

## CONTROLLABILITY AND POLE ASSIGNMENT FOR DISCRETE TIME LINEAR SYSTEMS DEFINED OVER ARBITRARY FIELDS\*

S. K. MITTER† AND R. FOULKES‡

**1. Introduction.** The theory of controllability and observability due to Kalman is certainly one of the most important conceptual contributions to linear systems theory. An account of the development of the ideas of controllability and observability as well as its implications on feedback control theory and realization theory may be found in the recent book of Kalman, Falb and Arbib [1].

It has been known for some time that for a linear continuous, finite-dimensional autonomous system with a scalar control variable, complete controllability is equivalent to being able to assign arbitrary poles to the closed loop transfer matrix by a suitable choice of state variable feedback gain matrix. This result was generalized to the vector control case by Wonham [2] and Simon and Mitter [3]. In [3] constructive recursive algorithms to achieve pole assignment were also presented. The objective of this note is to generalize this result to cover discrete-time, finite-dimensional, autonomous linear systems defined over arbitrary fields. The result can thus be applied to the feedback control of linear sequential machines [4]. By duality arguments the problem of state determination is also solved.

**2. Notation and system definition.** Let

- $T =$  time set  $= Z =$  (ordered Abelian group of) integers;
- $U =$  input values  $= F^m =$  vector space of  $m$ -tuples over the field  $F$ ;
- $X =$  state space  $= F^n$ ;
- $Y =$  output space  $= F^p$ ;
- $\Omega =$  input space of functions  $t \mapsto u(t)$ ; that is, arbitrary sequences  $u(-1), u(0), u(1), \dots$ , with  $u(t) \in U$ .

We shall be concerned with the discrete-time, autonomous, linear dynamical system  $\Sigma$  defined over a field  $F$ ,

$$(2.1) \quad \begin{aligned} x(t+1) &= Ax(t) + Bu(t), \\ y(t) &= Cx(t) \end{aligned}$$

with  $t \in Z$ ,  $x(t) \in F^n$ ,  $u(t) \in F^m$ ,  $y(t) \in F^p$  and where

$$(2.2) \quad \begin{aligned} A &: F^n \rightarrow F^n, \\ B &: F^m \rightarrow F^n, \\ C &: F^n \rightarrow F^p \end{aligned}$$

are  $F$ -homomorphisms.

---

\* Received by the editors July 10, 1969, and in revised form February 27, 1970. This work was supported in part by the National Science Foundation under Grant GK-3714 at Case Western Reserve University and the National Aeronautics and Space Administration under Grant NGL-22-009-124 at Electronic Systems Laboratory, Massachusetts Institute of Technology.

† Electronic Systems Laboratory, Massachusetts Institute of Technology, Cambridge, Massachusetts, and Systems Research Center, Case Western Reserve University, Cleveland, Ohio 44106.

‡ Systems Research Center, Case Western Reserve University, Cleveland, Ohio 44106.

We shall usually not make a distinction between  $(A, B)$  and  $(A, C)$  as a pair of  $F$ -homomorphisms or as a pair of matrices representing these homomorphisms with respect to a given basis.

With respect to the system (2.1) we make the assumption:

(i) the pair  $(A, B)$  is *completely reachable*, that is, the rank of the  $n \times nm$  matrix

$$(2.3) \quad H(A, B) = [B, AB, \dots, A^{n-1}B]$$

is  $n$ .

(ii) the pair  $(A, C)$  is *completely observable*, that is, the rank of the  $n \times np$  matrix

$$(2.4) \quad K(A, C) = [C^T, A^T C^T, \dots, (A^T)^{n-1} C^T]$$

is  $n$ .

**3. Statement of main theorem.**<sup>1</sup> The principal result of this paper is the following theorem.

**THEOREM 3.1.** *For the linear autonomous system (2.1),  $(A, B)$  is a completely reachable pair if and only if for every monic polynomial  $g$  of degree  $n$ , there exists an  $m \times n$  matrix  $K$  over  $F$  such that the characteristic polynomial of  $A + BK$  is precisely  $g$  (up to a factor of  $\pm 1$ ).*

The proof of the result proceeds via several propositions and is presented in the next section.

**4. Proof of main theorem.** The proof will be divided into three parts: necessity for the case when  $B$  is a column vector, necessity for a general  $B$  and sufficiency.

**PROPOSITION 4.1** (case  $m = 1$ ). *In (2.1) let  $B = b = n \times 1$  matrix. If  $(A, B)$  is a completely reachable pair, then there exists a  $1 \times n$  matrix  $k$  such that the characteristic polynomial of  $A + bk$  has an arbitrary preassigned form (of degree  $n$ ).*

The proof of this proposition essentially consists of transforming  $A$  to rational canonical form and is identical to the proof given for the field of real numbers (see, for example, [5, Theorems 7 and 9]).

We now consider the case where  $B$  is an  $n \times m$  matrix.

**PROPOSITION 4.2.** *If  $(A, B)$  is a completely reachable pair, then there exists a matrix  $K$  and a vector  $b$  such that  $(A + BK, b)$  is a completely reachable pair and  $b$  is in the column space of  $B$ .*

*Proof.* The proof presented is essentially the same as independently given by Heymann [6] and hence only an outline of the proof will be given.

Let  $b_j$  be the  $j$ th column of  $B$  and let  $E_j$  be the cyclic subspace of the coordinate space  $E = F^n$  generated by  $b_j$ . Since  $(A, B)$  is a completely reachable pair  $E = E_1 + \dots + E_m$ . In general, the  $E_i$  are not independent, that is,  $E_i \cap E_j \neq \emptyset$  for  $i \neq j$ . However, it is easy to see that there are subspaces  $S_i$  and a finite integer  $t$ ,  $0 < 1 \leq m$ , such that  $E = S_1 + \dots + S_t$  and  $S_i \cap S_j = \emptyset$  for  $i \neq j$ , that is,  $E$  is

---

<sup>1</sup> It was pointed out by the reviewer that a similar result has been obtained by R. E. Kalman in the unpublished notes: *Lectures on Controllability and Observability*, CIME Seminar, Italy, February 1969.

a direct sum of the subspaces  $S_i$ . A basis for  $E$  can now be obtained by combining the bases for the subspaces.

By rearranging the columns of  $B$  (hence the coordinates of the control) it can be assumed that the first  $t$  columns of  $B$  are used. Hence the basis is

$$b_1, \dots, A^{k_1-1}b_1, \dots, b_t, \dots, A^{k_t-1}b_t \quad \text{and} \quad \sum_{i=1}^t k_i = n.$$

Let  $R = [b_1, \dots, A^{k_1-1}b_1, \dots, b_t, \dots, A^{k_t-1}b_t]$  be the matrix whose columns are the above basis vectors. Clearly  $R$  is invertible.

Define an  $m \times n$  matrix  $S = [s_1 \cdots s_n]$ , where each column is an  $m$ -tuple defined as follows:

$$s_{r_j} = \varepsilon_{j+1}^{(m)} \quad \text{if} \quad r_j = \sum_{i=1}^j k_i \quad \text{and} \quad j = 1, \dots, t-1; \quad s_j = 0 \quad \text{otherwise,}$$

where  $\varepsilon_i^{(m)}$  is the  $i$ th standard basis vector of  $F^m$ .

Finally, let  $P = SR^{-1}$ . Clearly  $PA^{k_i-1}b_j = \varepsilon_{j+1}^{(m)}$ ,  $j = 1, \dots, t-1$ , and  $PA^i b_j = 0$  for all other powers of  $A$ .

Let  $\underline{A} = A + BP$ . Then the controllability matrix of the pair  $(\underline{A}, b_1)$  is  $\underline{H} = [b_1 \underline{A} b_1 \cdots \underline{A}^{n-1} b_1]$  and it has rank  $n$ . Clearly  $b_1$  is in the column space of  $B$ . The necessity part of the theorem now follows from Proposition 4.1.

We now prove sufficiency.

**PROPOSITION 4.3.** *Given an arbitrary monic polynomial  $g$  of degree  $n$ , if there exists an  $m \times n$  matrix  $K$  such that the characteristic polynomial of  $A + BK$  is precisely  $g$ , then  $(A, B)$  is a completely reachable pair.*

*Proof.* We first assume that the field  $F$  has a sufficient number of scalars  $a_1, \dots, a_n$  such that  $\det(A - a_i I) \neq 0$ ,  $i = 1, 2, \dots, n$ . From the above assumption and by hypothesis there is a  $K$  such that  $(A + BK)v_i = a_i v_i$  and  $v_i \neq 0$ . Since  $a_i I - A$  is invertible, we have

$$(4.1) \quad (a_i I - A)^{-1} B K v_i = v_i, \quad i = 1, 2, \dots, n.$$

Now for each  $a_i$  there are scalars  $b_j(a_i)$  such that

$$(4.2) \quad (a_i I - A)^{-1} = \sum_{j=1}^n b_j(a_i) A^{j-1}, \quad i = 1, 2, \dots, n.$$

Hence from (4.1) and (4.2), we obtain

$$(4.3) \quad \sum_{j=1}^n A^{j-1} B (b_j(a_i) K v_i) = v_i, \quad i = 1, 2, \dots, n.$$

Let  $H = [B, AB, \dots, A^{n-1}B]$  and  $y = (y_1 \cdots y_n)^T \in F^{nm}$ . Then

$$Hy = \sum_{j=1}^n A^{j-1} B y_j.$$

If we set  $y_j = b_j(a_i) K v_j$ , then (4.3) becomes

$$(4.4) \quad Hy_i^* = v_i, \quad i = 1, 2, \dots, n,$$

where  $y_i^* = (b_1(a_i) K v_1 \cdots b_n(a_i) K v_n)^T$ .

Since the eigenvalues of  $A + BK$  are distinct, the eigenvectors  $v_1, \dots, v_n$  are linearly independent and form a basis for  $F^n$ . Hence, by using (4.4), any  $v \in F^n$  can be written as  $v = H(\sum_i c_i y_i^*)$ . Therefore the range of  $H$  is  $F^n$  and hence  $(A, B)$  is a completely reachable pair.

Now if  $F$  does not contain enough distinct scalars, apply Proposition A.3 of the Appendix to  $f = \det(A - xI)$  and  $g = \det(A + BK - xI)$ . Then over some extension field  $F' \supset F$ ,  $g$  has  $n$  distinct roots none of which are roots of  $f$ . Now from the proof of Proposition 4.3,  $H$  considered as a linear transformation of  $(F')^m \rightarrow (F')^n$  has rank  $n$ . But  $H$  is a matrix over  $F \subset F'$ ; hence it has rank  $n$  over  $F$  also.

For finite fields (containing at least 2 elements) the following stronger result can be proved.

**THEOREM 4.4.** *The following statements are equivalent:*

- (i)  $(A, B)$  is a completely reachable pair;
- (ii) Given a monic polynomial  $g$  of degree  $n$ , there exists a matrix  $K$  such that the characteristic polynomial of  $A + BK$  is precisely  $g$ ;
- (iii)  $B \neq 0$ , and given an irreducible polynomial  $p$  of degree  $n$ , there exists a matrix  $K$  such that the characteristic polynomial of  $A + BK$  is  $p$ .

*Proof.* The theorem will be proved by showing that the statements (i) and (iii) are equivalent.

(i)  $\Rightarrow$  (iii) from Theorem 3.1.

We now prove the reverse implication. For  $n = 1$ , the result is obvious. For  $n > 1$ , by Proposition A.1 we can construct an irreducible polynomial of degree  $n$ .

Let  $\mathcal{R}$  denote the range of  $H(A, B)$ . Define the map

$$\bar{A}: F^n/\mathcal{R} \rightarrow F^n/\mathcal{R}$$

by

$$\bar{A}\bar{x}_i = \overline{(A + BK)x_i}, \quad i = 1, 2, \dots, n,$$

where  $F^n/\mathcal{R}$  is the quotient space,  $\{x_1, \dots, x_n\}$  is a basis for  $F^n$  and  $\bar{x}$  denotes the coset of  $x$  in the quotient space  $F^n/\mathcal{R}$  and  $K$  is an  $m \times n$  matrix. This is a well-defined map since  $\mathcal{R}$  is an  $A$ -invariant subspace of  $F^n$ .

Let  $p(x) = \sum_{i=0}^n p_i x^i$  be the characteristic polynomial of  $A + BK$ . Then by the Cayley–Hamilton theorem  $p(A + BK) = 0$ . It is easily verified that  $\bar{A}$  is an endomorphism of  $F^n/\mathcal{R} \rightarrow F^n/\mathcal{R}$ . Using an induction on  $k$  we can show

$$\bar{A}^k \bar{x} = \overline{(A + BK)^k x}$$

and we may verify that  $p(\bar{A}) = 0$  (that is, the zero map on  $F^n/\mathcal{R}$ ).

Let  $m$  be the minimal polynomial of  $\bar{A}$ . Then  $m$  divides  $p$  since  $p(\bar{A}) = 0$ . Since by hypothesis  $p$  is irreducible, either  $m = 1$  or  $m = \pm p$ .

Since  $B \neq 0$ ,  $\deg m < n = \deg p$ , so  $m = 1$ . But  $m(\bar{A}) = 0$ ; this means that the identity map on  $F^n/\mathcal{R}$  is equal to the zero map and hence  $F^n/\mathcal{R} = \mathcal{R}$ . Therefore  $\mathcal{R} = F^n$  and  $(A, B)$  is a completely reachable pair.

**5. An example.** As an example, consider the following three-state circuit over the field  $Z_3$ .



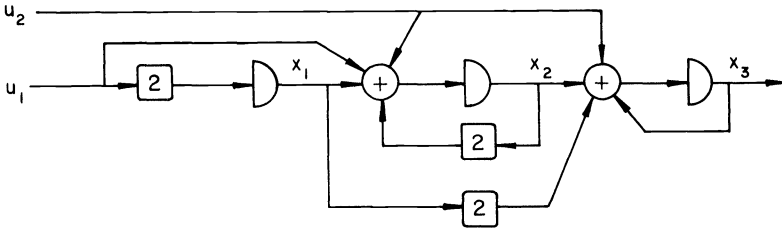


FIG. 1

By inspection of the diagram,

$$\begin{aligned} x_1(n+1) &= 2u_1(n), \\ x_2(n+1) &= x_1(n) + 2x_2(n) + u_1(n) + u_2(n), \\ x_3(n+1) &= 2x_1(n) + x_2(n) + x_3(n) + u_2(n). \end{aligned}$$

Letting  $x^T(n) = [x_1(n) \ x_2(n) \ x_3(n)]$  and  $u^T(n) = [u_1(n) \ u_2(n)]$ ,

$$x(n+1) = \begin{pmatrix} 0 & 0 & 0 \\ 1 & 2 & 0 \\ 2 & 1 & 1 \end{pmatrix} x(n) + \begin{pmatrix} 2 & 0 \\ 1 & 1 \\ 0 & 1 \end{pmatrix} u(n) = Ax(n) + Bu(n).$$

By direct calculation,

$$H = \begin{pmatrix} 2 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 2 & 2 & 1 \\ 0 & 1 & 2 & 2 & 0 & 1 \end{pmatrix},$$

which has rank 3. Following the construction in Proposition 4.2,

$$R = [b_1 \quad Ab_1 \quad A^2b_1] = \begin{pmatrix} 2 & 0 & 0 \\ 1 & 1 & 2 \\ 0 & 2 & 0 \end{pmatrix}.$$

Then  $R^{-1} = \begin{pmatrix} 2 & 0 & 0 \\ 0 & 0 & 2 \\ 2 & 2 & 2 \end{pmatrix}$ ; also,  $S = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{pmatrix}$ . Therefore  $P = SR^{-1} = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 2 \\ 2 & 2 & 2 \end{pmatrix}$ .

Again by direct calculation,  $\underline{A} = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 1 & 2 \\ 1 & 0 & 0 \end{pmatrix}$  and  $\underline{H} = \begin{pmatrix} 2 & 0 & 0 \\ 1 & 1 & 2 \\ 0 & 2 & 0 \end{pmatrix}$ , which is the

same as  $R$  in this case, so  $\underline{H}$  has rank 3. Letting  $p_1 = (p_{11}p_{12}p_{13})$ , the characteristic polynomial of  $\underline{A} + b_1p_1$  is  $-x^3 + (2p_{11} + p_{12} + 1)x^2 + (2p_{13} + p_{11})x + (p_{12} + p_{13})$ . To see that these coefficients may be chosen arbitrarily, it suffices to note that the following determinant is nonzero:

$$\begin{vmatrix} 2 & 1 & 0 \\ 1 & 0 & 2 \\ 0 & 1 & 1 \end{vmatrix} = 1 \neq 0.$$

**6. Observability and state reconstruction.** Since the pair  $(A, C)$  has been assumed to be a completely observable pair, it follows from Theorem 3.1 that for every monic polynomial  $g$  of degree  $n$ , there exists a  $p \times n$  matrix  $-D^T$  such that the characteristic polynomial of  $A^T - C^T D^T$  is precisely  $g$ . Hence the characteristic polynomial of  $A - DC$  can be made arbitrary.

Now consider an observer [7]:

$$(6.1) \quad \hat{x}(t+1) = X\hat{x}(t) + Dy(t) + Bu(t), \quad t = 0, 1, 2, \dots,$$

where  $X$  and  $D$  are  $n \times n$  and  $n \times p$  matrices respectively. Let  $D$  be chosen such that  $A^T - C^T D^T$  has arbitrary characteristic polynomial and let  $X = A - DC$ .

Then

$$(6.2) \quad \begin{aligned} \hat{x}(t+1) - x(t+1) &= (A - DC)[\hat{x}(t) - x(t)], \\ \hat{x}(0) - x(0) &= \text{given.} \end{aligned}$$

From (6.2) it follows that

$$\hat{x}(n) - x(n) = (A - DC)^n[\hat{x}(0) - x(0)].$$

Since the characteristic polynomial of  $A - DC$  can be made arbitrary, the matrix  $A - DC$  can in particular be made nilpotent and hence the observer reconstructs the initial state in at most  $n$  steps.

**Appendix.** In this Appendix some results on finite fields which are used in the proof of Proposition 4.3 are presented.

The following results are needed. The proof of the first two propositions are consequences of well-known results on finite fields (see Lang [8, Chap. VII, § 5]).

**PROPOSITION A.1** (see [9, p. 128]). *If  $F$  is a finite field consisting of at least two elements, then the polynomial ring  $P(F)$  contains irreducible polynomials of every degree  $\geq 2$ .*

**PROPOSITION A.2.** *For every irreducible polynomial over a finite field  $F$ , there is an extension field  $F'$  such that the given polynomial has  $n$  distinct roots in  $F'$ , where  $n$  is the degree of the polynomial.*

**PROPOSITION A.3.** *Let  $F$  be a finite field and  $f$  a given polynomial of degree  $n$  over  $F$ . Then there is a polynomial  $g$  of degree  $n$  over  $F$  and some extension field  $F' \supset F$  such that  $g$  has  $n$  distinct roots in  $F'$ , none of which are roots of  $f$ .*

*Proof.* First, consider the case when  $f$  has at least one root in  $F$ . Then  $f = f'f''$ , where  $f'$  is a product of linear factors and  $f''$  has no roots in  $F$ . Also,  $\deg f' \geq 1$ , so  $\deg f'' < n$ .

By Proposition A.1, there is an irreducible polynomial  $g$  over  $F$  of degree  $n$ . Then  $\gcd(f'', g) = 1$  since  $\deg f'' < \deg g$ . Clearly  $\gcd(f', g) = 1$ , so  $\gcd(f, g) = 1$ .

By Proposition A.2, there is an extension field  $F' \supset F$  such that  $g$  has  $n$  distinct roots in  $F'$ . But  $\gcd(f, g) = 1$  in  $F'$  also, so no root of  $g$  is a root of  $f$  in  $F'$ .

Next, consider the case when  $f$  has no roots in  $F$ . Now, either  $F$  has  $n$  distinct scalars or not. If it has, let  $g = (a_1 - x)(a_2 - x) \cdots (a_n - x)$ , where the  $a_i$  are distinct scalars in  $F$ . Then no root of  $g$  is a root of  $f$ , and  $F$  is the desired extension field.

If  $F$  does not have  $n$  distinct scalars, consider the prime factorization of  $f$ :  $f = p_1 p_2 \cdots p_k$ , where each  $p_i$  is an irreducible polynomial of degree at least 2, say  $\deg p_i = m_i$ , and  $\sum_{i=1}^k m_i = n$ . (Note that the  $p_i$  may not be distinct.)

If  $k > 1$ , then each  $m_i < n$ . Pick a polynomial  $g$  over  $F$  irreducible of degree  $n$ . Then  $\gcd(f, g) = 1$ . Again, let  $F' \supset F$  be an extension of  $F$  containing  $n$  distinct roots of  $g$ . Since  $\gcd(f, g) = 1$ , no root of  $g$  is a root of  $f$ .

If  $k = 1$ , then  $f = p_1$  and  $m_1 = n$ . Let  $F$  have  $p$  distinct scalars, and choose  $g' = (a_1 - x) \cdots (a_p - x)$ , where the  $a_i \in F$  are distinct. Also choose  $g''$  over  $F$  irreducible of degree  $n - p$ . Then  $\gcd(f, g'') = 1$  and  $\gcd(f, g') = 1$ ; hence,  $\gcd(f, g'g'') = 1$ .

Let  $F' \supset F$  be an extension field in which  $g''$  has  $n - p$  distinct roots. Therefore  $g = g'g''$  has  $n$  distinct roots in  $F'$ . Also, since  $\gcd(f, g) = 1$ , no root of  $g$  is a root of  $f$ .

**Acknowledgment.** The authors are indebted to the reviewer for constructive comments and suggestions. In particular the statement and proof of Theorem 4.4 was obtained as a result of comments by the reviewer.

*Note added in proof.* It was recently pointed out to me by R. W. Brockett that V. M. Popov proved the result on pole assignment earlier in his paper: *Hyperstability and optimality of automatic systems with several control functions*, Rev. Roumaine Sci. Techn.-Electrotechn. Energ., 9 (1964), pp. 629–690.

#### REFERENCES

- [1] R. E. KALMAN, P. L. FALB AND M. A. ARBIB, *Topics in Mathematical System Theory*, McGraw-Hill, New York, 1967.
- [2] W. M. WONHAM, *On pole-assignment in multi-input controllable linear systems*, IEEE Trans. Automatic Control, AC-12 (1967), pp. 660–665.
- [3] J. D. SIMON AND S. K. MITTER, *A theory of modal control*, Information and Control, 13 (1968), pp. 316–353.
- [4] A. GILL, *Linear Sequential Circuits: Analysis, Synthesis and Applications*, McGraw-Hill, New York, 1967.
- [5] E. B. LEE AND L. MARKUS, *Foundations of Optimal Control Theory*, John Wiley, New York, 1967.
- [6] M. HEYMANN, *Comments on "Pole Assignment in Multi-Input Controllable Linear Systems,"* IEEE Trans. Automatic Control, AC-13 (1969), pp. 748–749.
- [7] D. G. LUENBERGER, *Observers for multivariable systems*, Ibid., AC-11 (1966), pp. 190–197.
- [8] S. LANG, *Algebra*, Addison-Wesley, Reading, Mass., 1965.
- [9] I. T. ADAMSON, *Introduction to Field Theory*, Oliver and Boyd, Edinburgh and London, 1964.

## STABILITY OF A NONAUTONOMOUS DIFFERENTIAL EQUATION OF FOURTH ORDER\*

A. S. C. SINHA AND R. G. HOFT†

**1. Introduction.** In this paper, the object is to establish sufficient conditions for stability of (1) in the cases  $p(t) \equiv 0$ , and  $p(t) \not\equiv 0$ , respectively, which follows from some fundamental results due to LaSalle [1]. Further, it is worth noting that, in view of LaSalle's results, it is not necessary to show that the Lyapunov function is positive definite (and this would be a tedious computation) in order to conclude global asymptotic stability

$$(1) \quad \ddot{x} + f(\ddot{x})\ddot{x} + \phi(\dot{x}, \ddot{x})\dot{x} + \psi(\dot{x}) + \theta(x) = p(t).$$

It is easily seen that the results in this analysis become Routh-Hurwitz conditions when functions  $f$ ,  $\psi$ ,  $\phi$  and  $\theta$  are constants in the case  $p(t) \equiv 0$ . The conditions obtained are generalizations of the work of Ezeilo [2] and Harrow [3] for a class of fourth order systems. The Lyapunov function given here has been constructed by forming linear combinations of line integrals. However, the constants for the linear combinations are obtained by a trial method. The functions  $f$ ,  $\phi$ ,  $\psi$  and  $\theta$  are real-valued with continuous first partials. A smoothness property is also assumed to ensure the existence of integrals appearing in the analysis.

**2. Stability.** Consider the differential equation (1) and the associated system (Case 1.  $p(t) \equiv 0$ ):

$$(2) \quad \begin{aligned} \dot{x} &= y, & \dot{y} &= z, & \dot{z} &= w, \\ \dot{w} &= -f(z)w - \phi(y, z)z - \psi(y) - \theta(x), \end{aligned}$$

where  $\psi(0) = \theta(0) = 0$ , so that the origin is a critical point.

It is assumed that there exist positive constants such that

$$(3) \quad b \geq f(z) \geq f^o \geq a > 0,$$

$$(4) \quad \phi(y, z) \geq \phi^o > 0,$$

$$(5) \quad \psi(y)/y \geq \psi^o > 0,$$

$$(6) \quad (\theta^o \delta)^{1/2} \geq \theta(x)/x \geq \theta^o > 0,$$

$$(7) \quad \psi^o \leq \psi'(y) \leq \delta_1, \quad \theta^o \leq \theta'(x) \leq \theta'(0) \leq \delta,$$

---

\* Received by the editors June 10, 1969.

† College of Engineering, University of Missouri, Columbia, Missouri 65201. This work was supported in part by the National Science Foundation under Grant NSF GK805.

where the superscript  $o$  designates the evaluation of the given function at the origin. Let the following relations between the constants of (3) to (7) hold:

$$(8) \quad a\phi^o\psi^o\theta^o - \delta ab\theta^o - \delta_1^2\delta = \beta > 0,$$

$$(9) \quad \frac{\beta}{(\psi^o)^2(2a\psi^o\theta)^{1/2}} > \left[ \frac{\delta_1}{\psi^o} - \frac{\theta^o}{\delta} \right],$$

$$(10) \quad \frac{1}{z} \int_0^z f(s) ds - f(z) < \frac{\beta}{a^2\psi^o\theta^o} \quad \text{for all } z \neq 0.$$

The following  $V$ -functions are introduced:

$$(11) \quad V_1(y, z, w) = \frac{z}{F(z)} \left[ w + F(z) + \frac{\delta}{\psi^o} \frac{F(z)}{z} y \right]^2 + \left[ \frac{1}{a} - \frac{z}{F(z)} \right] w^2,$$

where  $\int_0^z f(s) ds = F(z)$ ;

$$(12) \quad V_2(y, z) = \frac{2}{a} \int_0^z s \{ \phi(y, s) - \phi^o \} ds + \frac{2\delta}{\psi^o} \int_0^y s \{ \phi(s, 0) - \phi^o \} ds,$$

$$(13) \quad V_3(x, y, z) = \frac{\theta^o\delta}{\psi^o} \left\{ x + \frac{\psi^o}{\theta^o\delta} \left[ \frac{\theta(x)}{x} \right] y + \frac{\psi^o}{a\theta^o\delta} \left[ \frac{\theta(x)}{x} \right] z \right\}^2 + \frac{2\delta}{\psi^o} \int_0^x s \left\{ \frac{\theta(s)}{s} - \theta^o \right\} ds,$$

$$(14) \quad V_4(x, y, z) = 2 \int_0^y \left\{ \frac{\psi(s)}{s} - \frac{\psi^o}{\theta^o\delta} \left[ \frac{\theta(x)}{x} \right]^2 \right\} s ds + \left\{ \frac{\phi^o\delta}{\psi^o} - \frac{\delta^2}{(\psi^o)^2} \left[ \frac{F(z)}{z} \right] - \frac{\theta'(0)}{a} \right\} y^2 + \left\{ \frac{\phi^o}{a} - \frac{\delta}{\psi^o} - \frac{\psi^o}{\delta\theta^o a^2} \left[ \frac{\theta(x)}{x} \right]^2 \right\} z^2 + 2 \int_0^z s f(s) ds - zF(z) + \frac{2}{a} \left\{ \frac{\psi(y)}{y} - \frac{\psi^o}{\theta^o\delta} \left[ \frac{\theta(x)}{x} \right]^2 \right\} yz.$$

Consider the following inequalities which are used to show that  $V \rightarrow \infty$  as  $|x| \rightarrow \infty$ . Now  $0 < F(z)/z < b$  and therefore from (7) it follows that

$$(15) \quad \alpha(x, y) = \frac{\psi(y)}{y} - \frac{\psi^\circ}{\theta^\circ \delta} \left[ \frac{\theta(x)}{x} \right]^2 \geq \frac{\psi(y)}{y} - \psi^\circ \geq 0,$$

$$(16) \quad \begin{aligned} \frac{\phi^\circ \delta}{\psi^\circ} - \frac{\delta^2}{(\psi^\circ)^2} \left[ \frac{F(z)}{z} \right] - \frac{\theta'(0)}{a} &\geq \frac{\delta}{a(\psi^\circ)^2 \theta^\circ} [a\phi^\circ \psi^\circ \theta^\circ - \delta ab\theta^\circ - (\psi^\circ)^2 \delta] \\ &\geq \frac{\delta}{a\psi^\circ \theta^\circ} [a\phi^\circ \psi^\circ \theta^\circ - \delta ab\theta^\circ - \delta_1^2 \delta] = \frac{\delta \beta}{a\psi^\circ \theta^\circ} > 0. \end{aligned}$$

From the mean value theorem,

$$\psi^\circ \leq \frac{\psi(y)}{y} = \int_0^1 \psi'(yt) dt \leq \delta_1, \quad \theta^\circ \leq \frac{\theta(x)}{x} = \int_0^1 \theta'(xt) dt \leq \delta$$

and (3), (6) and (7) give

$$(17) \quad \begin{aligned} \frac{\phi^\circ}{a} - \frac{\delta}{\psi^\circ} - \frac{\psi^\circ}{\delta \theta^\circ a^2} \left[ \frac{\theta(x)}{x} \right]^2 &\geq \frac{1}{a^2 \psi^\circ \theta^\circ} [a\phi^\circ \psi^\circ \theta^\circ - \delta a^2 \theta^\circ - (\psi^\circ)^2 \theta^\circ] \\ &\geq \frac{1}{a^2 \psi^\circ \theta^\circ} [a\phi^\circ \psi^\circ \theta^\circ - \delta ab\theta^\circ - \delta_1^2 \delta] = \frac{\beta}{a^2 \psi^\circ \theta^\circ} > 0, \end{aligned}$$

$$(18) \quad \begin{aligned} 2 \int_0^z sf(s) ds - zF(z) &= \int_0^z sf(s) ds - \int_0^z F(s) ds \\ &= - \int_0^z \left\{ \frac{F(s)}{s} - f(s) \right\} s ds > - \frac{\beta}{2a^2 \psi^\circ \theta^\circ} z^2. \end{aligned}$$

Therefore, using (15), (16), (17) and (18) in (14) gives

$$\begin{aligned} V_4(x, y, z) &\geq \int_0^y \alpha(x, s) s ds + \frac{\delta \beta}{a(\psi^\circ)^2 \theta^\circ} y^2 + \frac{\beta}{2a^2 \psi^\circ \theta^\circ} z^2 + \frac{2}{a} \alpha(x, y) yz \\ &\geq 2 \int_0^y \alpha(x, s) s ds + \frac{\beta}{a} \left[ \frac{y}{\psi^\circ} + \frac{\psi^\circ}{\beta} \alpha(x, y) z \right]^2 + Cz^2, \end{aligned}$$

where

$$C = \frac{\beta}{2a^2 \psi^\circ \theta^\circ} - \frac{(\psi^\circ)^2}{a\beta} \alpha^2(x, y) > 0.$$

Now  $\alpha(x, y) \leq \psi^\circ[\delta_1/\psi^\circ - \theta^\circ/\delta]$ ; therefore, from (9),  $\alpha^2(x, y) < \beta^2/[(\psi^\circ)^2 2a\psi^\circ\theta^\circ]$ .

Consider the Lyapunov function

$$\begin{aligned} 2V(x, y, z, w) &= V_1 + V_2 + V_3 + V_4 \\ &\geq \frac{z}{F(z)} \left[ w + F(z) + \frac{\delta}{\psi^\circ} \frac{F(z)}{z} y \right]^2 + \left[ \frac{1}{a} - \frac{z}{F(z)} \right] w^2 \\ &\quad + \frac{\theta^\circ \delta}{\psi^\circ} \left\{ x + \frac{\psi^\circ}{\theta^\circ \delta} \left[ \frac{\theta(x)}{x} \right] y + \frac{\psi^\circ}{a\theta^\circ \delta} \left[ \frac{\theta(x)}{x} \right] z \right\}^2 \\ &\quad + \frac{2\delta}{\psi^\circ} \int_0^x s \left\{ \frac{\theta(s)}{s} - \theta^\circ \right\} ds + 2 \int_0^y \alpha(x, s) s ds \\ &\quad + \frac{\beta}{a} \left[ \frac{y}{\psi^\circ} + \frac{\psi^\circ}{\beta} \alpha(x, y) z \right]^2 + Cz^2. \end{aligned}$$

Thus  $V \rightarrow \infty$  as  $|x| \rightarrow \infty$ . A little computation gives

$$\begin{aligned} 2V(x, y, z, w) &= \frac{2\delta}{\psi^\circ} \int_0^x \theta(s) ds + \left[ \frac{\phi^\circ \delta}{\psi^\circ} - \frac{\theta'(0)}{a} \right] y^2 \\ &\quad + \left( \frac{\phi^\circ}{a} - \frac{\delta}{\psi^\circ} \right) z^2 + \frac{1}{a} w^2 + 2wz + \frac{2\delta}{\psi^\circ} yw \\ (19) \quad &\quad + \frac{2\delta}{\psi^\circ} y \int_0^z f(s) ds + 2\theta(x)y + \frac{2}{a} \theta(x)z \\ &\quad + \frac{2}{a} \psi(y)z + 2 \int_0^y \psi(s) ds + \frac{2}{a} \int_0^z s \{ \phi(y, s) - \phi^\circ \} ds \\ &\quad + \frac{2\delta}{\psi^\circ} \int_0^y s \{ \phi(s, 0) - \phi^\circ \} ds + 2 \int_0^z s f(s) ds. \end{aligned}$$

Evaluating the time derivative of (19) along the trajectories of (2) yields

$$\begin{aligned} \dot{V}(x, y, z, w) &= \frac{1}{a} [f(z) - a] w^2 - \left[ \phi(y, z) - \frac{\psi'(y)}{a} - \frac{\delta}{\psi^\circ} \frac{F(z)}{z} \right] z^2 \\ &\quad - \left[ \frac{\delta}{\psi^\circ} \frac{\psi(y)}{y} - \theta'(x) \right] y^2 - \frac{1}{a} [\theta'(0) - \theta'(x)] yz \\ &\quad + \frac{1}{a} z \int_0^z s \phi_{,y}(y, s) ds - \frac{\delta}{\psi^\circ} [\phi(y, z) - \phi(y, 0)] yz \\ (20) \quad &\leq -\frac{1}{a} [f(z) - a] w^2 - \left\{ \phi^\circ - \frac{1}{a} \psi'(y) - \frac{b\delta}{\psi^\circ} - \frac{[\theta'(0) - \theta'(x)]}{4a^2} \right\} z^2 \\ &\quad - \delta_2 y^2 - [\theta'(0) - \theta'(x)] \left( y + \frac{z}{2a} \right)^2 \\ &\quad + \frac{1}{a} z \int_0^z s \phi_{,y}(y, s) ds - \frac{\delta}{\psi^\circ} \left[ y \int_0^1 \frac{\partial \phi}{\partial z}(y, zt) dt \right] \end{aligned}$$

and with  $\delta = \theta'(0) + \delta_2$ ,  $\delta_2 > 0$ ,

$$\phi(y, z) - \phi(y, 0) = z \int_0^1 \frac{\partial \phi}{\partial z}(y, zt) dt.$$

Then (20) is negative semidefinite if the following conditions hold :

$$(21) \quad z\phi_y(y, z) \leq 0,$$

$$(22) \quad y \int_0^1 \frac{\partial \phi}{\partial z}(y, zt) dt \geq 0,$$

$$(23) \quad \theta'(0) > 0, \quad \theta'(0) \geq \theta'(x) \geq \theta'(0) - \gamma_1, \quad \gamma_1 \geq 0,$$

where  $\gamma_1 = 4a^2(\beta/a\psi^\circ\theta^\circ - \gamma)$  for some  $\gamma \leq \beta/a\psi^\circ\theta^\circ$ . Using (7), (8) and (23) gives

$$\begin{aligned} & \left[ \phi^\circ - \frac{1}{a}\psi'(y) - \frac{b\delta}{\psi^\circ} \right] - \left[ \frac{\theta'(0) - \theta'(x)}{4a^2} \right] \\ &= \frac{1}{a\psi^\circ\theta^\circ} [a\phi^\circ\psi^\circ\theta^\circ - \psi^\circ\theta^\circ\psi'(y) - ab\delta\theta^\circ] - \left[ \frac{\theta'(0) - \theta'(x)}{4a^2} \right] \\ &\geq \frac{\beta}{a\psi^\circ\theta^\circ} - \frac{\gamma_1}{4a^2} = \gamma > 0. \end{aligned}$$

The conditions of Theorems 2 and 3 of LaSalle's paper [1] are satisfied; therefore global asymptotic stability has been proved. Thus these results are combined into the following theorem.

**THEOREM 1.** *If there exist positive constants as defined in (3)–(7) of the system (2) such that the following conditions hold:*

- (i)  $\psi(0) = \theta(0) = 0$ , and  $f, \phi, \psi, \theta$  have continuous first partials;
- (ii) conditions (8)–(10) and (21)–(23) hold;
- (iii)  $\dot{V} \neq 0$  along every nontrivial solution;

*then the system (2) is asymptotically stable.*

**3. Nonautonomous systems.** In this section, the system (1) is analyzed to give sufficient conditions for which every solution of (1) is bounded for all  $t \geq 0$ . Here it is assumed (Case 2) that  $p(t)$  is an integrable function.

**THEOREM 2.** *If the conditions (i) and (ii) of Theorem 1, along with the following conditions hold:*

- (i)  $f(z) \geq a + \gamma \geq a + \frac{1}{2}|p(t)|$ ;
- (ii)  $\delta_2 \geq (\delta/2\psi^\circ)|p(t)|$ ;
- (iii)  $\int_0^\infty |p(s)| ds \leq c < \infty$ ;

*then every solution of (1) is bounded in the future [1].*

*Proof.* To prove Theorem 2, consider the Lyapunov function (19) which was shown to be such that  $V \rightarrow \infty$  as  $|x| \rightarrow \infty$ . Evaluating the time derivative along



the trajectories of (1) results in the addition of  $p(t)$  terms to (20).

$$\begin{aligned}
 \dot{V}(x, y, z, w) &\leq -\frac{1}{a}[f(z) - a]w^2 - \left\{ \phi^o - \frac{1}{a}\psi'(y) - \frac{b\delta}{\psi^o} - \frac{[\theta'(0) - \theta'(x)]}{4a^2} \right\} z^2 \\
 &\quad - \delta_2 y^2 - [\theta'(0) - \theta'(x)] \left( y + \frac{z}{2a} \right)^2 \\
 &\quad + \frac{1}{a} z \int_0^s s \phi_y(y, s) ds - \frac{\delta}{\psi^o} \left[ y \int_0^z \frac{\partial \phi}{\partial z}(y, zt) dt \right] z^2 \\
 (24) \quad &\quad + \frac{1}{a} p(t)w + \frac{\delta}{\psi^o} p(t)y + p(t)z \\
 &\leq -\frac{1}{a} \left[ f(z) - a - \frac{1}{2}|p(t)| \right] w^2 - \left[ \gamma - \frac{1}{2}|p(t)| \right] z^2 \\
 &\quad - \left[ \delta_2 - \frac{\delta}{2\psi^o}|p(t)| \right] y^2 + \frac{1}{2} \left( \frac{1}{a} + \frac{\delta}{\psi^o} + 1 \right) |p(t)| \\
 &\leq \frac{1}{2} \left( \frac{1}{a} + \frac{\delta}{\psi^o} + 1 \right) |p(t)|
 \end{aligned}$$

if the conditions (i), (ii) and (iii) hold.

Now integrating (24), we have

$$\begin{aligned}
 V(x, y, z, w) &\leq V(0) + \frac{1}{2} \left( \frac{1}{a} + \frac{\delta}{\psi^o} + 1 \right) \int_0^t |p(s)| ds \\
 (25) \quad &\leq V(0) + \frac{1}{2} \left( \frac{1}{a} + \frac{\delta}{\psi^o} + 1 \right) c \leq k.
 \end{aligned}$$

Therefore, every solution of (1) is bounded in the future.

**4. Example.** A very simple case of differential equation (1) is when the functions  $f$  and  $\phi$  are constants. In this case the previous results become the Routh–Hurwitz conditions. Consider

$$(26) \quad \dot{x} + a_1 \ddot{x} + a_2 \ddot{x} + a_3 \dot{x} + a_4 x = p(t).$$

Let  $p(t) = 0$ ; then the conditions (8)–(10) result in only one nontrivial relation (27) which is clearly a Routh–Hurwitz condition for the system.

$$(27) \quad a_1 a_2 a_3 a_4 - a_4^2 a_1^2 - a_3^2 a_4 = \beta > 0.$$

Let  $p(t) \neq 0$  and define the constants  $a = a_1/2$ ,  $b = f^o = a_1$  and  $\delta = 2a_4$ . Then (26) has bounded solutions in the future if

- (i)  $\frac{1}{2} a_1 a_2 a_3 a_4 - a_4^2 a_1^2 - 2a_3^2 a_4 = \beta > 0$ ,
- (ii)  $a_1 \geq |p_1(t)|$ ,  $a_3 \geq |P(t)|$

and the conditions of Theorem 2 are satisfied.

#### REFERENCES

- [1] J. P. LASALLE, *Stability theory for ordinary differential equations*, Differential Equations, 4 (1968), pp. 57–65.

- [2] J. O. C. EZEILO, *Further results for the solutions of third order differential equations*, Proc. Cambridge Philos. Soc., 63 (1967), pp. 147–154.
- [3] M. HARROW, *Further results for the solution of third order differential equations*. Proc. London Math. Soc., 43 (1968), pp. 587–592.
- [4] J. LASALLE AND S. LEFCHETZ, *Stability by Liapunov's Direct Method with Applications*, Academic Press, New York, 1961.
- [5] J. O. C. EZEILO, *Further results for the solutions of a third-order differential equation*, Proc. Cambridge Philos. Soc., 59 (1963), pp. 111–116.
- [6] N. N. KRASOVSKII, *Stability of Motion*, Stanford University Press, Stanford, California, 1963.
- [7] W. LEIGHTON, *On the construction of Liapunov functions for certain autonomous nonlinear equations*, Contributions to Differential Equations, 4 (1963), pp. 367–383.
- [8] A. I. OGURTSOV, *The stability of solutions of two nonlinear differential equations of the third and fourth orders*, J. Appl. Math. Mech., 23 (1959), pp. 247–251.

## ERROR BOUNDS OF HIGH ORDER ACCURACY FOR THE STATE REGULATOR PROBLEM VIA PIECEWISE POLYNOMIAL APPROXIMATIONS\*

W. E. BOSARGE, JR. AND O. G. JOHNSON†

**Abstract.** Consider the linear quadratic cost control problem  $\dot{x} = A(t)x + B(t)u$ ,  $x(0) = x_0$ , with a cost functional  $J[u] = \frac{1}{2} \int_0^T [\langle x, Q(t)x \rangle + \langle u, R(t)u \rangle] dt$ . Let  $S$  be a suitable space of piecewise cubic polynomials on a mesh of norm  $h$  on the interval  $[0, T]$ . Then it is shown that a so-called Ritz–Treffitz method for minimizing  $J[\cdot]$  over  $S$  leads to an approximation to  $J[\cdot]$  of order  $O(h^7)$ . Further, a computable error bound can be exhibited. It is also shown that the computed pair  $(\bar{u}, \bar{x})$  converges to the optimal pair  $(u^*, x^*)$  with order  $O(h^3)$ . Similar statements are made for piecewise polynomial approximation of arbitrary positive order.<sup>1</sup>

**1. Introduction.** In the past considerable effort has been devoted to the use of direct methods in calculus of variations and control problems (see, for example, [4], [12], [6]). Specifically, the effort was directed towards finding the minimum of a functional  $J[u; x]$  (defined for some admissible controls and states) over a suitable finite-dimensional subspace of basis functions  $\{u_n(t), x_n(t); n = 0, 1, \dots, N\}$ . In this paper we consider the application of a modified Ritz direct method to the so-called state regulator control problem using finite-dimensional piecewise polynomial bases. We show that, for the “well-behaved” linear quadratic cost control problem, a numerically useful computational algorithm can be obtained for computing very high order approximations to the optimal control, the corresponding optimal state, and the associated cost functional. In particular, for the “smooth” problem, we show that a so-called Ritz–Treffitz method employing  $\alpha$ -order splines delivers approximations to the optimal control and optimal state with an error  $O(h^\alpha)$ , where  $h$  represents the mesh size for the spline basis. The corresponding approximate optimal cost estimates the true (optimal) cost to within  $O(h^{2\alpha})$ . We also exhibit an  $O(h^{2\alpha+1})$  error bound using the Hermite space  $H^\alpha(\Pi)$ . In addition, we represent an explicit derivation of the numerical Ritz–Treffitz algorithm. In a forthcoming paper we will discuss the practical utility of the method including results on its numerical stability. We will also present some comparisons of our method to the classical approach where numerical integration of a nonlinear Riccati differential equation is required. (Note  $O(h^\alpha)$  is a function less than  $Kh^\alpha$  in norm for some  $K > 0$ , where  $K$  is independent of  $h$ .)

**2. Problem description.** In this section we present a detailed description of the infinite-dimensional state regulator problem and recall a number of useful analytic results which we shall require in the sequel.

Consider the linear time-varying system

$$(2.1) \quad \begin{aligned} \dot{x}(t) &= A(t)x(t) + B(t)u(t), \\ x(0) &= x_0 \end{aligned}$$

\* Received by the editors October 9, 1969, and in revised form May 26, 1970.

† International Business Machines Scientific Center, Houston, Texas 77025 and Department of Mathematical Sciences, Rice University, Houston, Texas.

<sup>1</sup> A detailed numerical study of the algorithm resulting from this modified Ritz approach will appear in a forthcoming paper.

and the cost functional  $J[\cdot]$  defined by

$$(2.2) \quad J[u] = \frac{1}{2} \int_0^T [\langle x(t), Q(t)x(t) \rangle + \langle u(t), R(t)u(t) \rangle] dt.$$

Assume that  $x(\cdot)$  is an  $n$ -dimensional state vector,  $u(\cdot)$  an  $r$ -dimensional control vector,  $A(\cdot)$  an  $n \times n$  matrix,  $B(\cdot)$  an  $n \times r$  matrix,  $Q(\cdot)$  an  $n \times n$  positive definite matrix, and  $R(\cdot)$  an  $r \times r$  positive definite matrix. In addition, assume that  $A, B, Q, R$  are (elementwise) piecewise  $C^\alpha[0, T]$ ,  $\alpha \geq 0$ . We require that the terminal time  $T$  be fixed, that  $0 < r \leq n$ , and that  $u(\cdot)$  be unconstrained.

The problem we treat in this paper can be described in the following equivalent ways.

**DEFINITION 2.1 (Problem 1).** Given the linear system (2.1) and the cost functional (2.2) subject to the assumptions stated above, find the optimal control, i.e., the control which will drive the system (2.1) so as to minimize the cost functional (2.2). Thus we seek a  $u^*$  such that  $J[u^*] = \min_{u \in \mathcal{A}_u} J[u]$  with  $\dot{x}^* = A(t)x^*(t) + B(t)u^*(t)$  and  $x^*(0) = x_0$ . Here  $\mathcal{A}_u$  is some set of admissible vector-valued functions on  $[0, T]$ .

Since Problem 1 can be viewed as a standard calculus of variations minimization problem subject to nonholonomic constraints, we may introduce Lagrange multipliers  $\lambda(\cdot)$  (time-varying  $n$ -vector) and  $\gamma$  ( $n$ -vector) and define the Lagrangian  $L[u, x; \lambda, \gamma]$  by

$$(2.3) \quad L[u, x, \lambda, \gamma] \equiv J[u, x] + \int_0^T \langle \lambda(t), -\dot{x} + A(t)x + B(t)u \rangle dt + \langle \gamma, (x(0) - x_0) \rangle,$$

where  $J[u, x]$  is defined in (2.2) but is viewed as a function of two variables. Thus, an alternative way of describing Problem 1 is given by the following definition.

**DEFINITION 2.2 (Problem 1').** Given the linear system (2.1) and the cost functional (2.2) subject to the assumptions of Definition 2.1, find the  $u^*, x^*, \lambda^*$  and  $\gamma^*$  such that the Lagrangian is extremized; that is,

$$(2.4) \quad L[u^*, x^*; \lambda^*, \gamma^*] = \max_{\substack{\lambda \in \mathcal{A}_\lambda \\ \gamma \in R_n}} \min_{\substack{u \in \mathcal{A}_u \\ x \in \mathcal{A}_x}} L[u, x; \lambda, \gamma],$$

where  $R_n$  is real Euclidean  $n$ -space, and  $\mathcal{A}_x, \mathcal{A}_\lambda$  and  $\mathcal{A}_u$  are innerrelated.

One of the essential properties of the multipliers  $\lambda(\cdot)$  and  $\gamma$  is that, in the process of extremizing  $L$  over  $u, x, \lambda$  and  $\gamma$ , the Lagrangian is maximized over the multipliers  $\lambda$  and  $\gamma$ , and minimized with respect to  $u$  and  $x$  (see [10]). This fact is often misunderstood (see [12], [9] or [5]) and sometimes permits erroneous conclusions to be drawn in the convergence analysis of direct methods which have been applied to problems of the above type (see, for example, [12]).

We introduce still a third formulation of Problem 1 which is useful in the sequel.

**DEFINITION 2.3 (Problem 1'').** Under the assumptions of Definition 2.2, find the  $\lambda^*$  and  $\gamma^*$  such that

$$(2.5) \quad L[u^*, x^*; \lambda^*, \gamma^*] = \max_{\substack{\lambda \in \mathcal{A}_\lambda \\ \gamma \in R_n}} L[u, x; \lambda, \gamma]$$

subject to

$$(2.6) \quad \frac{\partial L}{\partial u}[u, x; \lambda, \gamma] = 0, \quad \frac{\partial L}{\partial x}[u, x; \lambda, \gamma] = 0,$$

where  $\partial L/\partial u$  and  $\partial L/\partial x$  are partial Fréchet derivatives of the scalar Langrangian.

The equivalence of the three formulations of Problem 1 follows from well-known "principles" of the calculus of variations (see, for example, [4]).

By smoothness of Problem 1 we shall mean the following.

DEFINITION 2.4. We say Problem 1  $\in PC^\alpha[0, T]$  if and only if  $A, B, Q, R$  have (elementwise) piecewise continuous  $\alpha$ th order derivatives. Equivalently we say Problem 1 is  $\alpha$ th order *smooth* if Problem 1  $\in PC^\alpha[0, T]$ ,  $\alpha \geq 0$ .

Basically, the solution of the state regulator problem leads to an optimal feedback system with the property that the state vector is "kept near zero" without excessive expenditure of control energy. Due to the linear nature of the system equations and the quadratic nature of the cost functional, the problem can be "handled analytically" in the sense that the optimal control is an explicit linear function of the state, i.e., is of the form

$$(2.7) \quad u(t) = G(t)x(t), \quad t \in [0, T],$$

with  $G(\cdot)$  an  $r \times n$  matrix-valued function defined by

$$(2.8) \quad G(t) = -R^{-1}(t)B^T(t)K(t)x(t),$$

where  $K(t)$  is the  $n \times n$  symmetric matrix solution of the well-known Riccati equation. Therefore, a complete solution is obtained (in general) only after  $K(\cdot)$  has been numerically approximated. We pursue this point further in a later section of the paper.

We now present a number of well-known results for Problem 1.

THEOREM 2.1. *Let us assume that an optimal control  $u^*$  for Problem 1 exists for the state  $x_0$ . Then, in order that  $u^*$  be optimal, it is necessary that there exist a vector function  $\lambda^*(\cdot)$  (Lagrange multiplier) such that*

(a)  $\lambda^*(\cdot)$  corresponds to  $u^*(\cdot)$  and  $x^*(\cdot)$  with  $\lambda^*(\cdot)$  and  $x^*(\cdot)$  solutions of the equation pair

$$(2.9) \quad \dot{x}^*(t) = A(t)x^*(t) + B(t)u^*(t)$$

and

$$(2.10) \quad \dot{\lambda}^*(t) = -Q(t)x^*(t) - A^T(t)\lambda^*(t),$$

subject to the boundary condition  $x^*(0) = x_0$  and the transversality condition  $\lambda^*(T) = 0$ ;

(b) along the optimal trajectory  $x^*(\cdot)$  we must have

$$(2.11) \quad R(t)u^*(t) + B^T(t)\lambda^*(t) = 0.$$

*Proof.* See Athans and Falb [2].

THEOREM 2.2 (Analytic solution of the state regulator problem 1). *Given the linear system (2.1) and the cost functional (2.2), where  $u(\cdot)$  is not constrained,  $T$  is specified and  $Q(\cdot)$  and  $R(\cdot)$  are positive definite matrices. Then an optimal control exists, is unique and is defined by (2.7) and (2.8), where the  $n \times n$  symmetric*

matrix  $K$  is the solution of the Riccati equation

$$(2.12) \quad \dot{K}(t) = -K(t)A(t) - A^T(t)K(t) + K(t)B(t)R^{-1}(t)B^T(t)K(t) - Q(t)$$

with  $K(T) = 0$ . The state  $x^*(\cdot)$  of the optimal system is then the solution of the linear system

$$(2.13) \quad \begin{aligned} \dot{x}(t) &= [A(t) - B(t)R^{-1}(t)B^T(t)K(t)]x(t), \\ x(0) &= x_0. \end{aligned}$$

*Proof.* See [2].

The notation  $x^*$ ,  $u^*$ ,  $\lambda^*$  and  $J^*$  will be used throughout the paper to denote optimal quantities. We now state a well-known theorem which we shall require in the sequel.

**THEOREM 2.3.** *Given the linear system (2.1) and the cost functional  $J$  of (2.2), let  $J^*$  denote the minimum value of  $J$ ; then  $J^*$  is given by*

$$(2.14) \quad J^*[x(t); t] = \frac{1}{2} \langle x^*(t), \lambda^*(t) \rangle, \quad 0 \leq t \leq T,$$

where  $\lambda^*(t) = K(t)x^*(t)$  with  $K(\cdot)$  the solution of (2.12).

*Proof.* See [2].

We now prove a theorem which tells us that the optimal quantities  $u^*$ ,  $x^*$  and  $\lambda^*$  for Problem 1 are in  $PC^\alpha[0, T]$  provided Problem 1  $\in PC^{\alpha-1}[0, T]$ .

**THEOREM 2.4.** *Assume that Problem 1 is  $(\alpha - 1)$ st order smooth. Then the optimal control  $u^*(\cdot)$ , the optimal state  $x^*(\cdot)$  and the optimal costate  $\lambda^*(\cdot)$  are each in  $PC^\alpha[0, T]$ .*

*Proof.* From (2.9), (2.10) and (2.11) we see that  $x^*(\cdot)$  and  $\lambda^*(\cdot)$  satisfy the  $2n$ -equation system

$$(2.15) \quad \begin{bmatrix} \dot{x}^*(t) \\ \dot{\lambda}^*(t) \end{bmatrix} = \begin{bmatrix} A(t) & -S(t) \\ -Q(t) & -A^T(t) \end{bmatrix} \begin{bmatrix} x^*(t) \\ \lambda^*(t) \end{bmatrix}$$

subject to  $x^*(0) = x_0$  and  $\lambda^*(T) = 0$ . Here  $S(t)$  is defined by

$$(2.16) \quad S(t) = B(t)R^{-1}(t)B^T(t).$$

Since  $A, B, Q, R \in PC^{\alpha-1}[0, T]$  and  $R$  is positive definite, it follows that  $R^{-1} \in PC^{\alpha-1}[0, T]$ . Thus,  $S \in PC^\alpha[0, T]$ . From the theory of ordinary differential equations we now conclude that  $x^*$  and  $\lambda^* \in PC^\alpha[0, T]$ . From (2.11) we deduce that  $u^* \in PC^\alpha[0, T]$  also and thus the theorem is proved.

**3. Finite-dimensional analogue of Problem 1.** In this section we present a special finite-dimensional analogue to Problem 1 using certain subspaces of piecewise polynomials. We first describe a number of useful piecewise polynomial spaces.

**3.1. Spaces of piecewise polynomials.** Let  $S_m^\alpha$  ( $\alpha \geq 1; m \geq 1$ , with  $\alpha, m$  integers) be an  $m$ -dimensional space of piecewise polynomials of fixed order  $\alpha - 1$ . Then we assume  $S_m^\alpha$  is characterized by the following properties:

P1. There exists a linear operator  $L_m: PC^\alpha[0, T] \rightarrow S_m^\alpha$ .

P2. For all  $f \in PC^\alpha[0, T]$ ,  $\|L_m f - f\|_2 = O(m^{-\alpha})$ .

P3. For all  $f \in PC^\alpha[0, T]$ ,  $\|(d/dt)(L_m f - f)\|_2 = O(m^{-\alpha+1})$ .

P4.  $(L_m f)^{(i)}(0) = f^{(i)}(0)$ ,  $(L_m f)^{(i)}(T) = f^{(i)}(T)$ ,  $i = 0, 1$ .

P5.  $S_m^\alpha \subset \mathcal{A}_\lambda$ .

It is well known that there exist spaces of piecewise polynomials possessing properties P1–P5 provided, of course,  $\mathcal{A}_\lambda$  is piecewise smooth. As an example, suppose we define  $\pi$  to be a mesh on  $[0, T]$  such that  $h = |\Pi|$ , i.e.,  $\Pi: 0 = t_0 < t_1 < \dots < t_v = T$  and

$$(3.1) \quad h = \max_{i=1, \dots, v} |t_i - t_{i-1}|.$$

Now suppose  $f \in PC^\alpha[0, T]$  and let  $\Pi$  contain all points of discontinuity of  $f^{(\alpha)}$  in  $[0, T]$ . Then the closure of the graph of  $f \in C^\alpha[t_{i-1}, t_i]$ ,  $i = 1, \dots, v$ . Thus,  $f^{(\alpha-1)}$  is absolutely continuous and  $f^{(\alpha)} \in L_\infty$ . Now, if  $m_0$  is an integer in  $[\alpha/2, \alpha + 1]$ , then there exists, for example, a Hermite interpolate to  $f$  (which we denote by  $L_m^H f$ ) of order  $2m_0 - 1$ . Further, it is well known in this case that (see, for example, [3])

$$(3.2) \quad \|L_m^H f - f\|_2 = O(h^{\alpha+1/2})$$

and

$$(3.3) \quad \left\| \frac{d}{dt} (L_m^H f - f) \right\| = O(h^{\alpha-1/2}).$$

We remark that the order of convergence is actually better (by  $\frac{1}{2}$ ) than we require.

Another extremely useful piecewise polynomial space possessing properties P1–P4 is the space of splines of order  $\alpha - 1$  (see, for example, [1]). Spline subspaces have been used frequently in recent years in the development of practical and efficient numerical algorithms for attacking wide classes of problems. In fact, for many practical problems, spline subspaces “deliver” the best results for an equivalent amount of computation, compared with an alternative finite dimension space of piecewise polynomials (see [8]). We shall therefore emphasize spline subspaces  $S_m^\alpha$  for our numerical comparisons which we present in a forthcoming paper.

**3.2. Application to state regulator Problem 1.** We now consider a special finite-dimensional approximation to Problem 1 (Definition 2.1) over a typical  $S_m^\alpha$ . The approximation is, in a sense, a Ritz method, but differs from the classical Ritz procedure in a fundamental way.

We begin by stating a special finite-dimensional analogue for Definition 2.1. Instead of requiring that the differential equation side condition (2.1) be satisfied identically, we relax the constraint as indicated in the following definition (see Trefftz [11]).

DEFINITION 3.1 (Problem 2). Given the linear system (2.1) and the cost functional (2.2) subject to the standard assumptions of Definition 2.1, find the optimal control  $\bar{u}$  such that

$$(3.4) \quad J[\bar{u}] = \min_{u \in \mathcal{A}_u} J[u]$$

subject to the side conditions

$$(3.5a) \quad \bar{x}(0) = x_0,$$

$$(3.5b) \quad \int_0^T \langle w_j(t), (-\dot{\bar{x}} + A(t)\bar{x} + B(t)\bar{u})_i \rangle dt = 0$$

for each basis function  $w_j \in S_m^\alpha$ ,  $i = 1, \dots, n$ ,  $j = 1, \dots, m$ .

The analogues of Problems 1' and 1'' are given in Definitions 3.2 and 3.3, respectively.

**DEFINITION 3.2 (Problem 2').** Under the same assumptions as in Definition 3.1, find the  $\bar{u}$ ,  $\bar{x}$ ,  $\bar{\lambda}$  and  $\bar{\gamma}$  such that the Lagrangian (2.3) is extremized when  $\lambda$  is restricted to  $S_m^\alpha$  ( $\lambda \in S_m^\alpha$  will henceforth mean that each  $\lambda_i \in S_m^\alpha$ ). Thus

$$(3.6) \quad L[\bar{u}, \bar{x}; \bar{\lambda}, \bar{\gamma}] = \max_{\substack{\lambda \in S_m^\alpha \\ \gamma \in R_n}} \min_{\substack{u \in \mathcal{A}_u \\ x \in \mathcal{A}_x}} L[u, x; \lambda, \gamma].$$

**DEFINITION 3.3 (Problem 2'').** Under the assumptions of Definition 3.1, find the  $\bar{\lambda}$  and  $\bar{\gamma}$  such that

$$(3.7) \quad L[\bar{u}, \bar{x}; \bar{\lambda}, \bar{\gamma}] = \max_{\substack{\lambda \in S_m^\alpha \\ \gamma \in R_n}} L[u, x; \lambda, \gamma]$$

subject to the constraints

$$(3.8) \quad \frac{\partial L}{\partial x}[u, x; \lambda, \gamma] = 0,$$

$$(3.9) \quad \frac{\partial L}{\partial u}[u, x; \lambda, \gamma] = 0.$$

The equivalence of Problems 2 and 2' follows immediately from the fundamental theorem of Lagrange multipliers. Further, the equivalence of Problems 2' and 2'' is guaranteed by fundamental principles of the calculus of variations (see [4]). The three formulations of the finite-dimensional problem will be required in subsequent sections. In particular, the formulation of Problem 2' will be required to generate certain key inequalities, essential in the convergence analysis of the numerical Ritz-Trefftz procedure. The third formulation (Problem 2'') is required in the development of a useful computational algorithm for the numerical solution of Problem 1.

**3.3. Development of numerical algorithm.** We consider now the development of a feasible computational algorithm using the problem formulation 2''. We begin by extremizing  $L[u, x; \lambda, \gamma]$  for  $\lambda_i \in S_m^\alpha$  and  $\gamma \in R_n$ , subject to the constraints (3.8) and (3.9). Thus, since  $\lambda_i \in S_m^\alpha$ , we can express  $\lambda$  as a linear combination of basis functions  $w_i$ ,  $i = 1, \dots, m$ . Consequently, we write

$$(3.10) \quad \lambda(t) = \sum_{j=1}^m c_j w_j(t), \quad 0 \leq t \leq T,$$

where each  $c_j$ ,  $j = 1, \dots, m$ , is an  $n$ -vector. Now, since  $x(T)$  is unconstrained, we conclude from (3.8) that

$$(3.11) \quad \begin{aligned} \dot{\lambda}(t) + A^T(t)\lambda(t) + Q(t)x(t) &= 0, & 0 \leq t \leq T, \\ \lambda(T) &= 0, & \gamma = -\lambda(0). \end{aligned}$$

Similarly, we deduce from condition (3.9) that

$$(3.12) \quad B^T(t)\lambda(t) + R(t)u(t) = 0, \quad 0 \leq t \leq T.$$



From (3.11) and (3.12) we observe that  $L$  can be written as a functional of the Lagrange multiplier  $\lambda(\cdot)$  only. Thus the problem is reduced to maximizing  $L$  over  $\lambda(\cdot)$ . We require, therefore, that

$$(3.13) \quad \frac{\partial L}{\partial \lambda} = 0,$$

where  $\lambda$  is defined in (3.10). Thus, we must first express  $L$  in terms of the  $c_j$ 's. We have

$$(3.14) \quad \begin{aligned} L[u, x; \lambda, \gamma] &= \frac{1}{2} \int_0^T \langle u, Ru \rangle dt + \frac{1}{2} \int_0^T \langle x, Qx \rangle dt \\ &\quad + \int_0^T \langle \lambda, -\dot{x} + Ax + Bu \rangle dt - \langle \lambda(0), x(0) - x_0 \rangle \\ &= J[u] - \langle \lambda(t), x(t) \rangle \Big|_0^T + \int_0^T \langle \dot{\lambda} + A^T \lambda, x \rangle dt \\ &\quad + \int_0^T \langle B^T \lambda, u \rangle dt - \langle \lambda(0), x(0) - x_0 \rangle. \end{aligned}$$

From conditions (3.11) and (3.12) we have

$$(3.15) \quad L[u, x; \lambda, \gamma] = -J[u] + \langle \lambda(0), x_0 \rangle.$$

From (3.10) we obtain

$$(3.16) \quad \int_0^T \langle u(t), R(t)u(t) \rangle dt = \sum_{i=1}^m \sum_{j=1}^m c_i^T \left[ \int_0^T E(t)w_i(t)w_j(t) dt \right] c_j,$$

where  $E$  is an  $n \times n$  symmetric matrix given by  $E(t) = B(t)R^{-1}(t)B^T(t)$ . Similarly

$$(3.17) \quad \int_0^T \langle x(t), Q(t)x(t) \rangle dt = \int_0^T [\langle \dot{\lambda}, Q^{-1}(t)\dot{\lambda} \rangle + 2\langle \lambda, F(t)\dot{\lambda} \rangle + \langle \lambda(t), G(t)\lambda \rangle] dt$$

or

$$(3.18) \quad \begin{aligned} \int_0^T \langle x(t), Q(t)x(t) \rangle dt &= \sum_{i=1}^m \sum_{j=1}^m c_i \left[ \int_0^T Q^{-1}(t)\dot{w}_i(t)\dot{w}_j(t) dt \right] c_j \\ &\quad + 2 \sum_{i=1}^m \sum_{j=1}^m c_i^T \left[ \int_0^T F(t)w_i(t)\dot{w}_j(t) dt \right] c_j \\ &\quad + \sum_{i=1}^m \sum_{j=1}^m c_i^T \left[ \int_0^T G(t)w_i(t)w_j(t) dt \right] c_j, \end{aligned}$$

where  $F(t) = A(t)Q^{-1}(t)$  and  $G(t) = A(t)Q^{-1}(t)A^T(t)$ . Finally, we have that

$$(3.18) \quad \langle \lambda(0), x_0 \rangle = \left\langle \sum_{i=1}^m c_i w_i(0), x_0 \right\rangle.$$

Maximizing  $\tilde{L}[c_1, \dots, c_m] = L[u(c_1, \dots, c_m), x(c_1, \dots, c_m); \lambda(c_1, \dots, c_m), \gamma(c_1, \dots, c_m)]$  over  $c_i$  ( $i = 1, \dots, m$ ), we obtain

$$(3.19) \quad \frac{\partial \tilde{L}}{\partial c_i} = - \sum_{j=1}^m \int_0^T [(E(t) + G(t))w_i(t)w_j(t) + F(t)w_i(t)\dot{w}_j(t) + F^T(t)\dot{w}_i(t)w_j(t) + Q^{-1}(t)\dot{w}_i(t)\dot{w}_j(t)] dt c_j + w_i(0)x_0 = 0$$

for  $i = 1, \dots, m$ . We define  $H_{ij}$  by

$$(3.20) \quad H_{ij} = \int_0^T [(E + G)w_i w_j + F w_i \dot{w}_j + F^T \dot{w}_i w_j + Q^{-1} \dot{w}_i \dot{w}_j] dt.$$

It follows that solving (3.19) for  $c_i$ ,  $i = 1, \dots, m$ , is equivalent to solving the linear system

$$(3.21) \quad Hy = b,$$

where the  $(m \times n) \times (m \times n)$  symmetric matrix  $H$  is defined by

$$(3.22) \quad H = \begin{bmatrix} H_{11} & H_{12} & \cdots & H_{1m} \\ H_{21} & H_{22} & \cdots & H_{2m} \\ \vdots & & & \vdots \\ H_{m1} & \cdots & & H_{mm} \end{bmatrix}, \quad H_{ij} = H_{ji}$$

and  $y$  and  $b$  are given by

$$(3.23) \quad y = \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_m \end{bmatrix}, \quad b = \begin{bmatrix} w_1(0)x_0 \\ w_2(0)x_0 \\ \vdots \\ w_m(0)x_0 \end{bmatrix}.$$

We note that  $H$  possesses certain desirable numerical properties (such as sparseness, definiteness, etc.), which insures that (3.21) is numerically "well-conditioned". We discuss these aspects of the numerical algorithm described by (3.21) in greater detail in a forthcoming paper. If now assume that  $\bar{y}$ , the solution of (3.21), is known exactly, then we write  $\bar{y} = [\bar{c}_1^T, \dots, \bar{c}_m^T]^T$  and set

$$(3.24) \quad \bar{\lambda}(t) = \sum_{i=1}^m \bar{c}_i w_i(t), \quad \bar{u} = u_{\bar{\lambda}}, \quad \bar{x} = x_{\bar{\lambda}},$$

where  $\bar{x}$  and  $\bar{u}$  are computed according to (3.11) and (3.12). In the remainder of the paper we shall be concerned with how "well"  $J[\bar{u}]$  approximates  $J[u^*]$ , and, likewise, how "well" the  $\bar{u}$  and  $\bar{x}$  approximate  $u^*$  and  $x^*$ , respectively.

**4. Cost functional convergence for pair  $(\bar{u}, \bar{x})$ .** In this section we state and prove the key convergence theorems for the cost functional  $J[\cdot]$ , interpreting the results in a number of useful piecewise polynomial subspaces.

**4.1. An order error bound for Ritz–Trefftz method.** From the equivalence of Problem 2 with Problem 2'' we can conclude that  $\bar{x}(0)$  is indeed  $x_0$ . Then it is clear that the final term  $\langle \gamma, \bar{x}(0) - x_0 \rangle$  may be omitted and we write  $L[\bar{u}, \bar{x}; \bar{\lambda}]$ , replacing  $L[\bar{u}, \bar{x}; \bar{\lambda}, \bar{\gamma}]$ .

**THEOREM 4.1.** *Assume Problem 1  $\in PC^{\alpha-1}[0, T]$ , and let  $\bar{u}$ ,  $\bar{x}$  and  $\bar{\lambda}$  be prescribed by the Ritz–Trefftz algorithm (defined in (3.21) to (3.24)) over an arbitrary  $S_m^\alpha$ . Suppose  $\lambda_s \in S_m^\alpha$ , where  $\lambda_s$  approximates  $\lambda^*$  to order  $h^\alpha$ ,  $h = 1/m$ . Further, suppose  $(u_s, x_s)$  is the pair generated by  $\lambda_s$  according to (3.11) and (3.12). Then*

$$(4.1) \quad \|\varepsilon_u\|_2 = O(h^\alpha), \quad \|\varepsilon_x\|_2 = O(h^{\alpha-1}),$$

$$(4.2) \quad J[u^*] \geq L[\bar{u}, \bar{x}; \bar{\lambda}] \geq J[u^*] + L[\varepsilon_u, \varepsilon_x; \varepsilon_\lambda],$$

where  $\varepsilon_u = u_s - u^*$ ,  $\varepsilon_x = x_s - x^*$  and  $\varepsilon_\lambda = \lambda_s - \lambda^*$ , with  $\|\cdot\|_2 = \left[ \int [\cdot]^2 \right]^{1/2}$ .

*Proof.* Since Problem 1  $\in PC^{\alpha-1}[0, T]$ , we have  $\lambda^* \in PC^\alpha[0, T]$ . Hence, from the hypothesis there exists a  $\lambda_s \in S_m^\alpha$  such that  $\|\lambda_s - \lambda^*\|_2 = O(h^\alpha)$ . Thus, from (3.11) and (3.12) we obtain

$$(4.3) \quad \|u_s - u^*\|_2 = \|R^{-1}B^T(\lambda_s - \lambda^*)\|_2 \leq O(h^\alpha),$$

$$(4.4) \quad \|x_s - x^*\|_2 = \|Q^{-1}[(\dot{\lambda}_s - \dot{\lambda}^*) + A^T(\lambda_s - \lambda^*)]\|_2 \leq O(h^{\alpha-1})$$

which proves (4.1).

For the proof of (4.2) we recall formulations 1'' and 2'', whence

$$(4.5) \quad J[u^*] = L[u^*, x^*; \lambda^*, \gamma^*] \geq L[\bar{u}, \bar{x}; \bar{\lambda}, \bar{\gamma}],$$

Since  $S_m^\alpha \subset \mathcal{A}_\lambda$  (in general,  $\mathcal{A}_\lambda$  will be a set of continuously differentiable functions on  $[0, T]$ ). And, from (3.7) we conclude that  $L[\bar{u}, \bar{x}; \bar{\lambda}] \geq L[u_s, x_s; \lambda_s]$ . But

$$(4.6) \quad \begin{aligned} L[u_s, x_s; \lambda_s] &= L[u^* + \varepsilon_u, x^* + \varepsilon_x; \lambda^* + \varepsilon_\lambda] \\ &= L[u^*, x^*; \lambda^*] + L[\varepsilon_u, \varepsilon_x; \varepsilon_\lambda] \\ &\quad + \int_0^T \langle u^*(t), R(t)\varepsilon_u(t) \rangle dt + \int_0^T \langle x^*(t), Q(t)\varepsilon_x(t) \rangle dt \\ &\quad + \int_0^T \langle \lambda^*(t), -\dot{\varepsilon}_x(t) + A(t)\varepsilon_x(t) + B(t)\varepsilon_u(t) \rangle dt. \end{aligned}$$

From (3.11), (3.12) and integration by parts we conclude that

$$(4.7) \quad \begin{aligned} &\int_0^T \langle \lambda^*(t), -\dot{\varepsilon}_x(t) + A(t)\varepsilon_x(t) + B(t)\varepsilon_u(t) \rangle dt \\ &= -\int_0^T \langle u^*(t), R(t)\varepsilon_u(t) \rangle dt - \int_0^T \langle x^*(t), Q(t)\varepsilon_x(t) \rangle dt \end{aligned}$$

(note that  $\varepsilon_x(0) = 0$  since  $\lambda_s(0) = \lambda^*(0)$ ,  $\lambda_s'(0) = \lambda^{*'}(0)$  by P4) which proves that

$$(4.8) \quad J[u^*] \geq L[\bar{u}, \bar{x}; \bar{\lambda}] \geq J[u^*] + L[\varepsilon_u, \varepsilon_x; \varepsilon_\lambda].$$

We consider now an important corollary to the previous theorem.

COROLLARY 4.2. *Suppose the conditions of Theorem 4.1 hold. Then*

$$(4.9) \quad 0 \leq J[u^*] - L[\bar{u}, \bar{x}; \bar{\lambda}] \leq O(h^{2(\alpha-1)}).$$

*Proof.* From (4.8) we observe that

$$(4.10) \quad 0 \leq J[u^*] - L[\bar{u}, \bar{x}; \bar{\lambda}] \leq -L[\varepsilon_u, \varepsilon_x; \varepsilon_\lambda].$$

But

$$(4.11) \quad \begin{aligned} L[\varepsilon_u, \varepsilon_x; \varepsilon_\lambda] &= \frac{1}{2} \left[ \int_0^T \langle \varepsilon_u, R(t)\varepsilon_u \rangle dt + \int_0^T \langle \varepsilon_x, Q(t)\varepsilon_x \rangle dt \right] \\ &\quad + \int_0^T \langle \varepsilon_\lambda, -\dot{\varepsilon}_x + A(t)\varepsilon_x + B(t)\varepsilon_u \rangle dt. \end{aligned}$$

Again integrating by parts and applying the order bounds to  $\varepsilon_\lambda, \varepsilon_x$  and  $\varepsilon_u$ , we obtain

$$\begin{aligned} 0 \leq J[u^*] - L[\bar{u}, \bar{x}; \bar{\lambda}] &\leq O(h^{2(\alpha-1)}) + O(h^{2\alpha}) \\ &= O(h^{2(\alpha-1)}), \end{aligned}$$

and the corollary is proved.

Two important piecewise polynomial subspaces which are useful computationally are  $S(\Pi)$  (the space of piecewise cubic splines on  $[0, T]$ ) and  $H(\Pi)$  (the piecewise cubic Hermite subspace on  $[0, T]$ ). If, for example, Problem 1  $\in PC^3[0, T]$ , then  $\alpha = 4$ , and we have

$$0 \leq J[u^*] - L[\bar{u}, \bar{x}; \bar{\lambda}] \leq O(h^6), \quad h = |\Pi|,$$

with  $\bar{u}, \bar{x}, \bar{\lambda}$  generated over  $S(\Pi)$ . Similarly, for  $\bar{u}, \bar{x}, \bar{\lambda}$  generated over  $H(\Pi)$ , we can show that

$$0 \leq J[u^*] - L[\bar{u}, \bar{x}; \bar{\lambda}] \leq O(h^6), \quad h = |\Pi|.$$

(Note that  $L[\bar{u}, \bar{x}; \bar{\lambda}]$  is just  $\frac{1}{2} \int_0^T \langle \bar{u}, R\bar{u} \rangle dt + \frac{1}{2} \int_0^T \langle \bar{x}, Q\bar{x} \rangle dt$ .)

**5. Norm convergence of the pair  $(\bar{u}, \bar{x})$ .** In this section we discuss the convergence of the computed pair  $\bar{u}$  and  $\bar{x}$  to the optimal pair  $u^*$  and  $x^*$ . We begin by proving a series of simple lemmas.

LEMMA 5.1. *Assume Problem 1  $\in PC^{\alpha-1}[0, T]$  and let  $\bar{u}, \bar{x}$  and  $\bar{\lambda}$  be prescribed by the Ritz-Trefftz procedure over the space  $S_m^\alpha$ . Then*

$$(5.1) \quad \begin{aligned} &\int_0^T \langle \lambda^*, -\dot{\bar{x}} + A(t)\bar{x} + B(t)\bar{u} \rangle dt \\ &= 2J[u^*] - \int_0^T \langle u^*, R(t)\bar{u} \rangle dt - \int_0^T \langle x^*, Q(t)\bar{x} \rangle dt. \end{aligned}$$

*Proof.* Integration of the left side of (5.1) by parts yields

$$\begin{aligned} & \int_0^T \langle \lambda^*(t), -\dot{\bar{x}} + A(t)\bar{x} + B(t)\bar{u} \rangle dt \\ &= -\langle \lambda^*, \bar{x} \rangle|_0^T + \int_0^T \langle \dot{\lambda}^* + A^T(t)\lambda^*, \bar{x} \rangle dt + \int_0^T B^T \lambda^*, \bar{u} \rangle dt. \end{aligned}$$

But, by transversality  $\lambda^*(T) = 0$ , and from (2.14) we conclude that  $2J[u^*] = \langle \lambda^*(0), \bar{x}(0) \rangle$ . Since  $\bar{x}(0) = x_0$ , we have

$$2J[u^*] = -\langle \lambda^*, \bar{x} \rangle|_0^T.$$

The desired result now follows from the standard necessary conditions for Problem 1 (see (2.9) and (2.10)).

LEMMA 5.2. *Assume the hypotheses of Lemma 5.1 hold. In addition suppose that  $\lambda_s$  is the  $S_m^a$  approximation to  $\lambda^*$ , and suppose  $(u_s, x_s)$  is the pair generated by  $\lambda_s$  according to (3.11) and (3.12). Then*

$$(5.2) \quad \int_0^T \langle \lambda^*, -\dot{\bar{x}} + A(t)\bar{x} + B(t)\bar{u} \rangle dt = O(h^{\alpha-1}).$$

*Proof.* From the formulation of Problem 2, we deduce that

$$\int_0^T \langle w_j, (-\dot{\bar{x}} + A(t)\bar{x} + B(t)\bar{u})_i \rangle dt = 0$$

for each  $w_j \in S_m^a$ ,  $i = 1, \dots, n$ ,  $j = 1, \dots, m$ . But  $\lambda_s(t) = \sum_{j=1}^m c_j^s w_j(t)$  for some set  $c_j^s$  of  $n$ -vectors,  $j = 1, \dots, m$ . Hence

$$(5.3) \quad \int_0^T \langle \lambda_s, -\dot{\bar{x}} + A(t)\bar{x} + B(t)\bar{u} \rangle dt = 0.$$

Hence

$$\begin{aligned} & \int_0^T \langle \lambda^*, -\dot{\bar{x}} + A(t)\bar{x} + B(t)\bar{u} \rangle dt \\ &= \int_0^T \langle \dot{\varepsilon}_\lambda, \bar{x} \rangle dt + \int_0^T \langle \varepsilon_\lambda, A\bar{x} \rangle dt + \int_0^T \langle \varepsilon_\lambda, B\bar{u} \rangle dt \\ &\leq \|\dot{\varepsilon}_\lambda\|_2 \|\bar{x}\|_2 + \|\varepsilon_\lambda\|_2 \|A\|_2 \|\bar{x}\|_2 + \|\varepsilon_\lambda\|_2 \|B\|_2 \|\bar{u}\|_2. \end{aligned}$$

But we know that  $\|\varepsilon_\lambda\|_2 = O(h^\alpha)$  and thus  $\|\dot{\varepsilon}_\lambda\|_2 = O(h^{\alpha-1})$ . Since  $\|\bar{x}\|_2$  and  $\|\bar{u}\|_2$  are each bounded uniformly in  $h$  (Theorem 4.1), (5.2) is established and the proof is complete.

LEMMA 5.3. *Assume the hypotheses of Lemma 5.1 hold. Then*

$$(5.4) \quad \|\bar{u} - u^*\|_R^2 \equiv \int_0^T \langle \bar{u} - u^*, R(t)(\bar{u} - u^*) \rangle dt = O(h^{\alpha-1}),$$

$$(5.5) \quad \|\bar{x} - x^*\|_Q^2 \equiv \int_0^T \langle \bar{x} - x^*, Q(t)(\bar{x} - x^*) \rangle dt = O(h^{\alpha-1}).$$

*Proof.* We can easily verify the identity

$$\begin{aligned} & \|u^* - \bar{u}\|_R^2 + \|x^* - \bar{x}\|_Q^2 \\ &= 2J[u^*] - 2 \left[ \int_0^T \langle \bar{u}, R(t)u^* \rangle dt + \int_0^T \langle \bar{x}, Q(t)x^* \rangle dt + 2J[\bar{u}, \bar{x}] \right]. \end{aligned}$$

Because of the equivalence of Problems 2, 2' and 2'' we recall that  $L[\bar{u}, \bar{x}; \bar{\lambda}, \bar{\gamma}] = J[\bar{u}, \bar{x}]$ . Hence

$$\begin{aligned} & \|u^* - \bar{u}\|_R^2 + \|x^* - \bar{x}\|_Q^2 \\ & \leq 4J[u^*] - 2 \left[ \int_0^T \langle \bar{u}, R(t)u^* \rangle dt + \int_0^T \langle \bar{x}, Q(t)x^* \rangle dt \right] + O(h^{2(\alpha-1)}). \end{aligned}$$

Now, applying the previous two lemmas, we obtain

$$0 \leq \|u^* - \bar{u}\|_R^2 + \|x^* - \bar{x}\|_Q^2 \leq 4J[u^*] - 4J[u^*] + O(h^{\alpha-1})$$

or

$$0 \leq \|u^* - \bar{u}\|_R^2 + \|x^* - \bar{x}\|_Q^2 \leq O(h^{\alpha-1})$$

from which (5.4) and (5.5) follow.

LEMMA 5.4. *Under the assumptions of Lemma 5.2 we assert that*

$$\begin{aligned} & \int_0^T \langle \lambda^*, -\dot{\bar{x}} + A(t)\bar{x} + B(t)\bar{u} \rangle dt \\ (5.6) \quad & \leq \|\dot{\varepsilon}_\lambda\|_2 \|\bar{x} - x^*\|_2 + \|\varepsilon_\lambda\|_2 \|A\|_2 \|\bar{x} - x^*\|_2 \\ & \quad + \|\varepsilon_\lambda\|_2 \|B\|_2 \|\bar{u} - u^*\|. \end{aligned}$$

*Proof.* The lemma follows immediately from (5.3) and integration by parts.

THEOREM 5.5. *Under the hypotheses of Lemma 5.2 we assert that*

$$(5.7) \quad \|\bar{u} - u^*\|_R = O(h^{\alpha-1}),$$

$$(5.8) \quad \|\bar{x} - x^*\|_Q = O(h^{\alpha-1}).$$

*Proof.* We observed in proving Lemma 5.3 that

$$\begin{aligned} & 0 \leq \|\bar{u} - u^*\|_R^2 + \|\bar{x} - x^*\|_Q^2 \\ & \leq 4J[u^*] - 2 \left[ \int_0^T \langle \bar{u}, R(t)u^* \rangle dt + \int_0^T \langle \bar{x}, Q(t)x^* \rangle dt \right] + O(h^{2(\alpha-1)}). \end{aligned}$$

Applying Lemmas 5.1, 5.2 and 5.4, we find that

$$\begin{aligned} & \|\bar{u} - u^*\|_R^2 + \|\bar{x} - x^*\|_Q^2 \\ & \leq 2[\|\dot{\varepsilon}_\lambda\|_2 \|\bar{x} - x^*\|_2 + \|\varepsilon_\lambda\|_2 \|A\|_2 \|\bar{x} - x^*\|_2 \\ & \quad + \|\varepsilon_\lambda\|_2 \|B\|_2 \|\bar{u} - u^*\|_2] + O(h^{2(\alpha-1)}). \end{aligned}$$

Now, since  $R(\cdot)$  and  $Q(\cdot)$  are positive definite, we can find positive constants  $v_1, v_2, q_1, q_2$  such that

$$(5.9) \quad v_1 \|\bar{u} - u^*\|_R^2 \leq \|\bar{u} - u^*\|_R^2 \leq v_2 \|\bar{u} - u^*\|_R^2$$

and

$$(5.10) \quad q_1 \|\bar{x} - x^*\|_2^2 \leq \|\bar{x} - x^*\|_Q^2 \leq q_2 \|\bar{x} - x^*\|_2^2.$$

Therefore, we obtain the inequality

$$(5.11) \quad \begin{aligned} & \|\bar{u} - u^*\|_R^2 + \|\bar{x} - x^*\|_Q^2 \\ & \leq O(h^{\alpha-1}) \|\bar{u} - u^*\|_R + O(h^{\alpha-1}) \|\bar{x} - x^*\|_Q + O(h^{2(\alpha-1)}). \end{aligned}$$

Set  $\delta_u = \|\bar{u} - u^*\|_R$  and  $\delta_x = \|\bar{x} - x^*\|_Q$ . Then (5.11) becomes

$$\delta_u^2 + O(h^{\alpha-1})\delta_u + \delta_x^2 + O(h^{\alpha-1})\delta_x \leq O(h^{2(\alpha-1)})$$

or, equivalently,

$$[\delta_u - O(h^{\alpha-1})]^2 + [\delta_x - O(h^{\alpha-1})]^2 \leq O(h^{2(\alpha-1)}).$$

Hence

$$\begin{aligned} \delta_u &= \|\bar{u} - u^*\|_R = O(h^{\alpha-1}), \\ \delta_x &= \|\bar{x} - x^*\|_Q = O(h^{\alpha-1}) \end{aligned}$$

and the theorem is established.

**COROLLARY 5.6.**  $\|\bar{u} - u^*\|_2 = O(h^{\alpha-1})$ ,  $\|\bar{x} - x^*\|_2 = O(h^{\alpha-1})$ .

*Proof.* The proof is trivial and is therefore omitted.

**6. Remarks on the numerical properties of the method.** The above theorems establish that the Ritz–Treffitz method has the unique theoretical property that it converges to the optimal value of the cost functional at a rate equal to the square of the order of the approximating subspaces. This advantage is not shared by any other approximate method, direct or indirect. In fact, the standard Ritz direct method is not known to converge for variational problems with nonholonomic constraints.

The numerical implications are obvious. The authors have a working program for the method and have compared it with the usual numerical approximation of the Riccati equation for numerous examples. The method is significantly faster, except on trivial problems where the two approaches yield approximately the same results. Since the resulting algebraic problem is definite, it can be solved by a Cholesky decomposition without pivoting, hence the band structure is preserved (providing a patch basis is used). Thus the numerical algorithm is extremely stable and fast. Also, reasonable bounds on the condition number can be computed.

In [13] we discuss the numerical utility of the algorithm in considerable detail. Additional convergence results for  $(\bar{u}, x_{\bar{u}})$  and  $J[\bar{u}; x_{\bar{u}}]$  were obtained and appear in [14].

#### REFERENCES

- [1] J. H. AHLBERG, E. N. NILSON AND J. L. WALSH, *The Theory of Splines and Their Applications*, Academic Press, New York, 1967.
- [2] M. ATHANS AND P. L. FALB, *Optimal Control: An Introduction to the Theory and Its Applications*, McGraw-Hill, New York, 1966.
- [3] G. BIRKHOFF, M. SCHULTZ AND R. VARGA, *Piecewise Hermite interpolation in one and two variables with applications to partial differential equations*, Numer. Math., 11 (1968), pp. 232–256.

- [4] R. COURANT AND D. HILBERT, *Methods of Mathematical Physics*, vol. I, Interscience, New York, 1953.
- [5] V. W. EVELEIGH, *Adaptive Control and Optimization Techniques*, McGraw-Hill, New York, 1967.
- [6] P. L. FALB, *Direct Methods in Optimal Control*, to appear.
- [7] G. HADLEY, *Nonlinear and Dynamic Programming*, Addison-Wesley, Reading, Mass., 1964.
- [8] O. G. JOHNSON, *Error bounds for Sturm–Liouville eigenvalues by several piecewise cubic Rayleigh–Ritz methods*, SIAM J. Numer. Anal., 6 (1969), pp. 317–333.
- [9] R. C. K. LEE, *Optimal Estimation, Identification and Control*, Research Monograph No. 28, M.I.T. Press, Cambridge, Mass., 1964.
- [10] D. G. LUENBINGER, *Optimization by Vector Space Methods*, John Wiley, New York, 1969.
- [11] E. TREFFTZ, *Konvergenz und Fehlerschätzung beim Ritzschen Verfahren*, Math. Ann., 100 (1928), pp. 503–521.
- [12] A. P. SAGE, *Optimum Systems Control*, Prentice-Hall, Englewood Cliffs, N.J., 1968.
- [13] W. E. BOSARGE, JR. AND O. G. JOHNSON, *Numerical properties of the Ritz–Treffitz algorithm for optimal control*, IBM DPD Tech. Rep. 320.2376, 1970; Math. Comp., to appear.
- [14] ———, *Direct method approximation to the state regulator control problem using a Ritz–Treffitz suboptimal control*, Proc. Joint Automatic Control Conference, 1970, pp. 51–56.



## BOUNDARY VALUE CONTROL OF THE HIGHER-DIMENSIONAL WAVE EQUATION\*

DAVID L. RUSSELL†

**1. Introduction.** Many important control processes can be described approximately by means of partial differential equations with control parameters appearing in the boundary conditions. For example, a triangular airplane wing may be equipped with ailerons on the trailing edge. An idealized model for such a plant would involve the partial differential equations which describe the motion of a plate with arbitrary control functions appearing in the boundary conditions along one of the sides of the triangle. Many other examples could be given.

Linear hyperbolic problems in one space dimension have been studied rather extensively; see, e.g., [1], [2], [3], [4], [5]. Here the theory is relatively uncomplicated. One can study questions of controllability using the geometric techniques based on characteristic curves or the more algebraic techniques based on the theory of nonharmonic Fourier series. One obtains not only theorems asserting the existence of controls transferring one state to another within a finite time period but also constructive proofs of these theorems which can be adapted to yield numerical techniques whereby the appropriate control functions can be calculated. The papers of Grainger [4] and Cirina [5] are noteworthy in this respect. Cirina's paper shows that such methods can even be used for quasilinear systems.

The theory is not nearly as complete for problems involving two or more space variables. The reasons why this should be so become apparent when one compares Chaps. V and VI of the treatise [6] of Courant–Hilbert. Some results have been obtained in this area by Fattorini [7] who considers, for the most part, boundary value control problems wherein the controls can be described by finitely many functions of the time  $t$ —physically the most realistic situation.

The purpose of the present paper is to study hyperbolic problems in several space dimensions using certain uniqueness theorems due to Holmgren [8] and John [9]. Using these results we can obtain very explicit estimates on the length of time required to transfer a given state into an arbitrarily small neighborhood of any other state using boundary value controls restricted to a subset of the boundary of the region in question. More specifically, we are able to show for the wave equations in 3 or fewer space variables that the system can be controlled in any time  $T$  which exceeds twice the wave propagation time from the boundary set where controls are applied to the rest of the physical medium. It should be noted that such a result is in agreement with known results [1], [2] for the case of a single space dimension.

**2. The control problem.** Let  $\Omega$  be a bounded, open connected domain in  $R^n$  whose boundary is an analytic surface  $\Gamma$  of dimension  $n - 1$ . We indicate

---

\* Received by the editors February 12, 1970.

† Departments of Mathematics and Computer Sciences, University of Wisconsin, Madison, Wisconsin 53706, and Research Consultant, Honeywell Inc., St. Paul, Minnesota 55113. This work was supported in part by the National Science Foundation under Grant GP-11495.

points in  $\Omega$  by

$$x = \begin{pmatrix} x^1 \\ x^2 \\ \vdots \\ x^n \end{pmatrix}.$$

The boundary surface  $\Gamma$  is parametrized by an  $(n - 1)$ -dimensional vector variable  $s$ . Integrals over  $\Omega$  will be denoted by  $\int_{\Omega} (\cdot) dx$  while integrals over  $\Gamma$  will be written  $\int_{\Gamma} (\cdot) ds$ . If we wish to indicate a point in  $R^n$  which lies on  $\Gamma$  we will write  $x(s)$ . The surface  $\Gamma$ , being analytic, has everywhere a unique unit outward normal vector which we will indicate by  $\eta(s)$ .

We consider a second order linear hyperbolic partial differential equation

$$(2.1) \quad L(w) = \rho(x)w_{tt} - \sum_{i,j=1}^n (\alpha_{ij}(x)w_{ij}) = 0.$$

The subscripts  $i, j$  indicate partial differentiation with respect to  $x^i, x^j$ , respectively. The coefficients  $\rho(x), \alpha_{ij}(x)$  are real analytic in some open subset of  $R^n$  which includes  $\Omega \cup \Gamma$ . Moreover, if  $A(x)$  is the  $n \times n$  matrix with entries  $\alpha_{ij}(x)$ ,  $A(x)$  is symmetric and there are positive numbers  $\rho_0$  and  $\alpha_0$  such that

$$\begin{aligned} \rho(x) &\geq \rho_0, & x \in \Omega, \\ v' A(x)v &\geq \alpha_0 \|v\|^2, & x \in \Omega, \quad v \in R^n. \end{aligned}$$

Let  $\Gamma_1$  be a relatively open subset of  $\Gamma$ . For  $T > 0$ , we denote by  $F$  the space of all  $C^\infty$  functions  $f: \Gamma \times [0, T] \rightarrow R^1$  with the property that  $f$  vanishes outside some compact subset of  $\Gamma_1 \times (0, T)$  (this set varies with  $f$ ). We pose for (2.1) the initial boundary value problem

$$(2.2) \quad \begin{aligned} w(x, 0) &\equiv w_t(x, 0) \equiv 0, & x \in \Omega, \\ w_x(x(s), t)A(x(s))\eta(s) &\equiv f(s, t), & (x(s), t) \in \Gamma \times [0, T]. \end{aligned}$$

The symbol  $w_x$  denotes the row vector of spatial partial derivatives of  $w$ :

$$w_x = (w_1, w_2, \dots, w_n).$$

With these assumptions it is known that the initial boundary value problem (2.1), (2.2) has a unique solution  $w^f(x, t)$  which lies in the class  $C^\infty((\Omega \cup \Gamma) \times [0, T])$ . The reader is referred to the papers of Friedrichs, Lax and Duff [10], [11], [12].

If  $f(s, t) \equiv 0$  for  $t_1 \leq t \leq t_2$ , then  $w^f(\cdot, t)$  can be considered as a vector-valued function with range in  $L^2(\Omega)$  which solves the evolution equation

$$(2.3) \quad d^2w/dt^2 + Bw = 0, \quad t_1 \leq t \leq t_2,$$

where  $B$  is the unbounded operator on  $L^2(\Omega)$  which is the unique self-adjoint extension of the operator  $\sum_{i,j=1}^n (\alpha_{ij}(x)w_{ij})$  defined on twice continuously

differentiable functions  $w(x)$  satisfying the boundary conditions

$$w_x(x(s))A(x(s))\eta(s) = 0, \quad x(s) \in \Gamma.$$

Let  $H_1$  denote the space of pairs of real-valued functions  $w(x)$ ,  $w_t(x)$  defined on  $\Omega$  with  $w_t(x)$  square integrable and  $w(x)$  having square integrable derivatives:

$$\int_{\Omega} \left( w_t(x)^2 + \sum_{i=1}^n w_i(x)^2 + w(x)^2 \right) dx < \infty.$$

Let  $H_E$  denote the space of equivalence classes of  $H_1$  modulo the zero energy states  $w_t(x) \equiv 0$ ,  $w(x) \equiv \text{const}$ . The “energy”

$$\mathcal{E}(w, w_t) = \int_{\Omega} (\rho(x)(w_t(x))^2 + w_x(x)A(x)w_x(x)) dx$$

is a constant on each such equivalence class.  $H_E$  is a Hilbert space with the inner product

$$\langle (w, w_t); (v, v_t) \rangle_E = \int_{\Omega} (\rho(x)w_t(x)v_t(x) + w_x(x)A(x)v_x(x)) dx$$

and resulting norm

$$\|(w, w_t)\|_E = \sqrt{\langle (w, w_t); (w, w_t) \rangle_E} = \sqrt{\mathcal{E}(w, w_t)}.$$

We shall not stress the distinction between  $H_1$  and  $H_E$  where unnecessary, and we shall say  $(w, w_t) \in H_E$  if the equivalence class of  $(w, w_t)$  is a member of  $H_E$ .

For each  $f \in F$  the corresponding solution  $w^f(x, t)$  of (2.1), (2.2) is such that  $(w^f(\cdot, T), w_t^f(\cdot, T)) \in H_E$ . In fact, if we put

$$R_T = \{(w^f(\cdot, T), w_t^f(\cdot, T)) | f \in F\},$$

then  $R_T$  is a subspace of  $H_E$  which we will call the *reachable space*. Following others we make the following definition.

**DEFINITION.** The control system  $\{(2.1), (2.2), f \in F\}$  is *approximately controllable* in time  $T > 0$  if  $R_T$  is dense in  $H_E$  relative to the topology induced by the norm  $\|\cdot\|_E$ .

We shall conclude this section with a theorem which relates approximate controllability to “observability.” (Cf. parallel results for ordinary differential equations [13].)

**THEOREM 1.** Let  $(\hat{v}, \hat{v}_t) \in H_E$  be such that both  $\hat{v}$  and  $\hat{v}_t$  lie in  $C^\infty(\Omega \cup \Gamma)$  and  $\hat{v}$  satisfies the consistency conditions

$$\hat{v}_x(x(s))A(x(s))\eta(x(s)) = 0, \quad \hat{v}_{t,x}(x(s))A(x(s))\eta(x(s)) = 0, \quad x(s) \in \Gamma.$$

Let  $v(x, t)$  be the unique  $C^\infty$  solution of  $L(v) = 0$  which satisfies the boundary conditions

$$(2.4) \quad v_x(x(s), t)A(x(s))\eta(x(s)) \equiv 0, \quad (x(s), t) \in \Gamma \times [0, T]$$

and the terminal conditions

$$(2.5) \quad v(x, T) \equiv \hat{v}(x), \quad v_t(x, T) \equiv \hat{v}_t(x).$$

Then  $(\hat{v}, \hat{v}_t) \in (R_T)^\perp$  in  $H_E$  if and only if  $v_t(x(s), t) \equiv 0$ ,  $(x(s), t) \in \Gamma_1 \times [0, T]$ .

*Proof.* For  $f \in F$  we have

$$(2.6) \quad \begin{aligned} 0 &= \int_{\Omega \times [0, T]} (v_t L(w^f) + w_t^f L(v)) dx dt \\ &= \int_{\Omega \times [0, T]} \operatorname{div}_{x,t} \begin{pmatrix} -v_t(w_x^f a_1) - w_t^f(v_x a_1) \\ \vdots \\ -v_t(w_x^f a_n) - w_t^f(v_x a_n) \\ \rho w_t^f v_t + w_x^f A v_x' \end{pmatrix} dx dt, \end{aligned}$$

where, for convenience, we have suppressed the arguments in the integrand and the column vectors  $a_i(x)$  are the columns of the symmetric matrix  $A(x)$ :

$$A(x) = (a_1(x), a_2(x), \dots, a_n(x)).$$

Applying the divergence theorem to the second member of (2.6) we obtain

$$(2.7) \quad \begin{aligned} 0 &= \int_{\Omega \times \{T\}} (\rho w_t^f v_t + w_x^f A v_x') dx - \int_{\Omega \times \{0\}} (\rho w_t^f v_t + w_x^f A v_x') dx \\ &\quad - \int_{\Gamma \times [0, T]} (v_t(w_x^f A \eta) + w_t^f(v_x A \eta)) ds dt. \end{aligned}$$

Using (2.5), (2.2), respectively, in the first two members of (2.7) and (2.2), (2.4) in the third member, we obtain

$$(2.8) \quad \int_{\Omega} (\rho(x) w_t^f(x, T) \hat{v}_t(x) + w_x^f(x, T) A(x) \hat{v}_x(x')) dx = \int_{\Gamma_1 \times [0, T]} (v_t f) ds dt.$$

From the definition of  $\langle \cdot, \cdot \rangle_E$ , we see that (2.8) becomes

$$(2.9) \quad \langle (w^f(\cdot, T), w_t^f(\cdot, T)); (\hat{v}, \hat{v}_t) \rangle_E = \int_{\Gamma_1 \times [0, T]} (v_t f) ds dt.$$

The right-hand side of (2.9) vanishes for all  $f \in F$  if and only if  $v_t(x(s), t) \equiv 0$ ,  $(x(s), t) \in \Gamma_1 \times [0, T]$ , and thus the proof is complete.

Theorem 1 is fundamental in the proofs of the controllability theorems of the subsequent sections.

**3. The time  $T_0$ .** In order to state and prove our theorems on approximate controllability of (2.1), (2.2) we must employ the concept of a *characteristic surface* for (2.1) in  $R^{n+1}$ . This concept is treated in detail in [6], for example, but we give a brief description to make our presentation somewhat self-contained.

Let  $S$  be a surface in  $R^{n+1}$  given by

$$S = \{(x, t) | \Phi(x, t) = 0\},$$

where  $\Phi(x, t)$  is a smooth real-valued function of  $n + 1$  variables. We define the *characteristic form*

$$\chi(\Phi, x, t) = \rho(x)(\Phi_t(x, t))^2 - \Phi_x(x, t)A(x)\Phi_x'(x, t).$$

The surface  $S$  is: characteristic if  $\chi(\Phi, x, t) \equiv 0$  for  $(x, t) \in S$ ; uniformly space-like if there exists a  $\delta > 0$  such that  $\chi(\Phi, x, t) \geq \delta$  for  $(x, t) \in S$ ; uniformly time-like if

there exists a  $\delta > 0$  such that  $\chi(\Phi, x, t) \leq -\delta$  for  $(x, t) \in S$ . For what is usually called the wave equation  $\rho(x) \equiv 1$  and  $A(x) \equiv I$ , the  $n \times n$  identity matrix, a surface is characteristic if and only if it everywhere makes an angle of  $45^\circ$  with any intersecting surface  $t = \text{const}$ . It is this special case that we shall use in our diagrams since it is less confusing than the general case.

Let  $(x_0, t_0) \in (\Omega \cup \Gamma) \times [0, T]$ . We define the *forward cone of influence* of  $(x_0, t_0)$  to be the subset  $K^+(x_0, t_0)$ , the largest closed subset of  $(\Omega \cup \Gamma) \times [t_0, T]$  which contains  $(x_0, t_0)$  and does not meet any uniformly space-like surface passing through  $(x_0, t_0)$ . Similarly we define  $K^-(x_0, t_0)$ , the *backward cone of influence* of  $(x_0, t_0)$ , by replacing  $[t_0, T]$  with  $[0, t_0]$ . It is easy to see that  $K^+(x_0, t_0)$  and  $K^-(x_0, t_0)$  have characteristic boundary surfaces. When  $\rho(x) \equiv 1$ ,  $A(x) \equiv I$ , we have  $K^+(x_0, 0) = \{(x, t) \in (\Omega \cup \Gamma) \times [0, T] | t^2 - \|x - x_0\|^2 \geq 0\}$ . If  $G$  is a subset of  $\Omega \cup \Gamma$ , we define forward and backward cones of influence of  $(G, t_0)$  by

$$K^+(G, t_0) = \overline{\bigcup_{x_0 \in G} K^+(x_0, t_0)}, \quad K^-(G, t_0) = \overline{\bigcup_{x_0 \in G} K^-(x_0, t_0)}.$$

Let  $t_0, t_1$  lie in  $[0, T]$  and let  $G \subseteq \Omega \cup \Gamma$ . We define

$$K(G, t_0, t_1) = K^+(G, t_0) \cap K^-(G, t_1).$$

Since the coefficients of the operator  $L$  do not depend upon  $t$ ,  $K(G, t_0, t_1)$  is symmetric about the plane  $t = \frac{1}{2}(t_0 + t_1)$ .

The fact that  $A(x)$  is uniformly positive definite can be used to prove that there is a least time  $T_0 > 0$  such that  $K^+(\Gamma_1, 0)$  includes the set  $\Omega \times \{T_0\}$ . Then  $K(\Gamma_1, 0, 2T_0)$  also includes  $\Omega \times \{T_0\}$ . If  $T > 2T_0$ , there is an  $\alpha > 0$  such that  $K(\Gamma_1, 0, T)$  includes  $\Omega \times \{t\}$  for  $|t - T/2| < \alpha$ . If  $T < 2T_0$ , the set

$$J(T/2) = \Omega \times \{T/2\} - K(\Gamma_1, 0, T)$$

is a nonempty set. For  $\rho(x) \equiv 1$ ,  $\Omega = \text{unit disc in } R^2$ , Figs. 1-4 illustrate the geometry of the situations described above both when  $\Gamma_1 = \Gamma$  and when  $\Gamma_1$  is a small subarc of  $\Gamma$ .

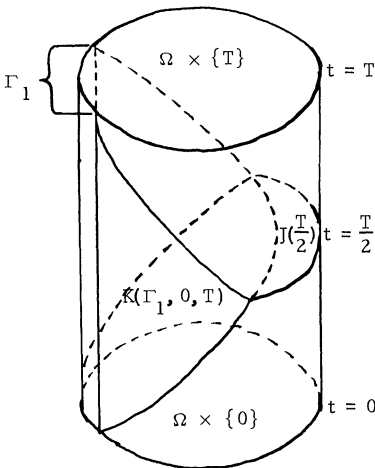


FIG. 1.  $\Gamma_1$  is a subarc of  $\Gamma$ ,  $T < 2T_0$

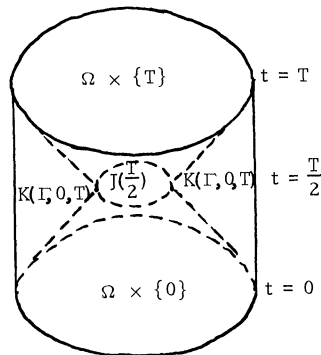


FIG. 2.  $\Gamma_1 = \Gamma$ ,  $T < 2T_0$

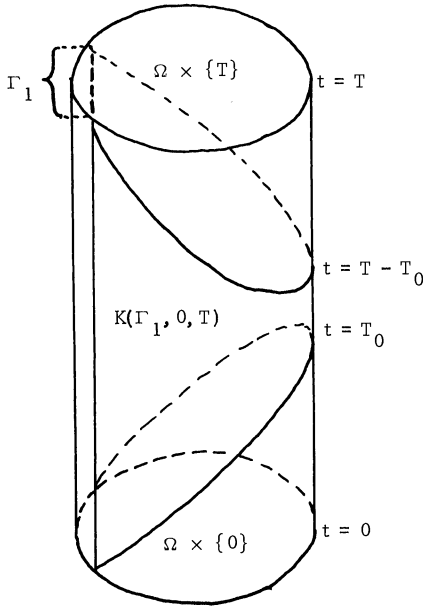


FIG. 3.  $\Gamma_1$  is a subarc of  $\Gamma$ ,  $T > 2T_0$

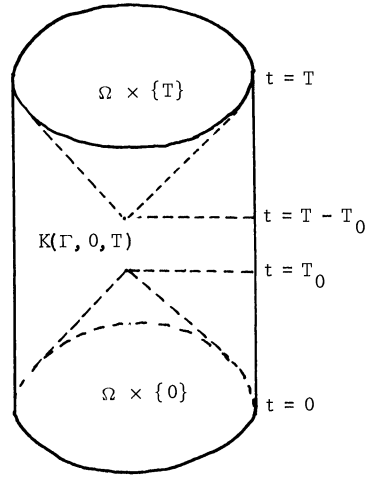


FIG. 4.  $\Gamma_1 = \Gamma$ ,  $T > 2T_0$

**4. Noncontrollability for  $T < 2T_0$ .** Since the reachable set  $R_T$  is a linear subspace of the Hilbert space  $H_E$ ,  $R_T$  fails to be dense in  $H_E$  just in case there is a nonzero element of  $H_E$  which is orthogonal to all elements of  $R_T$ . In view of our definition of  $H_E$  in terms of equivalence classes of states in  $H_1$  modulo the zero energy states, we see that  $R_T$  is dense in  $H_E$  if and only if the equations

$$(4.1) \quad \int_{\Omega} (\rho(x)w_t^f(x, T)\hat{v}_t(x) + w_x^f(x, T)A(x)\hat{v}_x(x)) dx = 0, \quad f \in F,$$

where  $\hat{v}, \hat{v}_t$  is a fixed element of  $H_1$ , imply that  $\hat{v}_t \equiv 0, \hat{v} = \text{const.}$

When  $T < 2T_0$  the subset  $J$  of  $\Omega$  given by

$$J = \{x | (x, T/2) \in J(T/2)\}$$

is a nonempty open set. Let  $\tilde{v}(x), \tilde{v}_t(x)$  be a state in  $H_1$  such that (i)  $\tilde{v}(x), \tilde{v}_t(x)$  has nonzero energy norm, (ii)  $\tilde{v}(x), \tilde{v}_t(x) \in C^\infty(\Omega)$  and vanish outside a compact subset of the interior of  $J$ .

We now permit the state  $\tilde{v}(x), \tilde{v}_t(x)$  to evolve, via the partial differential equation with

$$v\left(x, \frac{T}{2}\right) = \tilde{v}(x), \quad v_t\left(x, \frac{T}{2}\right) = \tilde{v}_t(x).$$

We put

$$(4.2) \quad \hat{v}(x) = v(x, T), \quad \hat{v}_t(x) = v_t(x, T)$$

and note that (2.1), (2.4) imply conservation of energy so that

$$\|(\hat{v}, \hat{v}_t)\|_E = \|(\tilde{v}, \tilde{v}_t)\|_E \neq 0.$$

Now our requirement (ii) on  $\tilde{v}$ ,  $\tilde{v}_i$  guarantees that  $v(x, t) \in C^\infty(\Omega \times [0, T])$ , as does  $w^f(x, t)$  for each  $f \in F$ . This enables us to apply Theorem 1 to show that (4.1) holds if and only if  $v_i(x(s), t) \equiv 0$  for  $(x(s), t) \in \Gamma_1 \times [0, T]$ .

It is easily seen that the fact that  $K(\Gamma_1, 0, T) \cap J = \emptyset$  implies that the interior of  $K(\Gamma_1, 0, T)$  does not meet  $K^+(J, T/2) \cup K^-(J, T/2)$ , the cone of influence of  $(J, T/2)$ . Well-known results for hyperbolic partial differential equations [6] then show that  $v(x, t) \equiv 0$ ,  $v_i(x, t) \equiv 0$  in  $K(\Gamma_1, 0, T)$ . But  $\Gamma_1 \times [0, T] \subseteq K(\Gamma_2, 0, T)$  so we conclude that  $v_i(x(s), t) \equiv 0$ ,  $(x(s), t) \in \Gamma_1 \times [0, T]$ . Therefore (4.1) must hold for the nonzero energy state  $\hat{v}$ ,  $\hat{v}_i$  and we have proved the following theorem.

**THEOREM 2.** *The system (2.1), (2.2) is not approximately controllable in time  $T$  if  $T < 2T_0$ .*

This theorem may be compared with comparable results [1], [2] for hyperbolic systems in one space dimension.

**5. The Holmgren–Fritz John uniqueness theorem.** The uniqueness theorems of Holmgren and Fritz John [8], [9], applied to the case we have in mind, reduce to the following theorem.

**THEOREM 3.** *Let  $u(x, t)$  be a twice continuously differentiable solution of  $L(u) = 0$  (cf. (2.1)) in  $K(\Gamma_1, t_0, t_1)$ ,  $[t_0, t_1] \subseteq [0, T]$ , with*

$$(5.1) \quad u_x(x(s), t)A(x(s))\eta(x(s)) \equiv 0, \quad u(x(s), t) \equiv 0, \\ (x(s), t) \in \Gamma_1 \times [t_0, t_1].$$

*Then*

$$(5.2) \quad u(x, t) \equiv 0 \quad \text{in } K(\Gamma_1, t_0, t_1).$$

The proof of this theorem is detailed in the works cited. However, we need to strengthen the theorem somewhat for our needs and this strengthening requires that we have some details of the proof. For this reason we give a short proof of Theorem 3. An important part of the proof is the lemma stated below, which we do not prove. See [9] for details in certain cases.

**LEMMA.** *If  $(\bar{x}, \bar{t})$  lies in the interior of  $K(\Gamma_1, t_0, t_1)$ , there is a uniformly time-like family of surfaces  $S(\lambda)$ ,  $0 \leq \lambda \leq 1$ , with the following properties:*

- (i)  $S(\lambda)$  is a compact subset of a relatively open analytic  $(n - 1)$ -dimensional surface;
- (ii)  $S(\lambda)$  varies analytically with respect to  $\lambda$ ,  $0 \leq \lambda \leq 1$ ;
- (iii)  $S(\lambda) \subseteq K(\Gamma_1, t_0, t_1)$ ,  $0 \leq \lambda \leq 1$ , and  $S(0)$  is a subset of the interior of  $\Gamma_1 \times [0, T]$ ;
- (iv) If  $0 \leq \lambda \leq 1$ , then  $S(0) \cup S(\lambda)$  is the boundary of an open subset  $D(\lambda) \subseteq K(\Gamma_1, t_0, t_1)$  and  $(\bar{x}, \bar{t}) \in D(1)$ .

Assuming this lemma, we proceed.

*Proof of Theorem 3.* The uniformly time-like character of the surfaces  $S(\lambda)$  together with the analyticity of these surfaces enables one to employ the Cauchy–Kowalewski theorem [6] to show that there are  $n$ -dimensional neighborhoods  $N(\lambda)$  of the surfaces  $S(\lambda)$  such that if analytic Cauchy data for  $z$  are prescribed on  $S(\lambda)$ , there will be a corresponding unique analytic solution  $z(x, t)$  of  $L(z) = 0$  in  $N(\lambda)$ . Since the equation  $L(z) = 0$  is linear,  $N(\lambda)$  depends only on  $S(\lambda)$ , not on the particular Cauchy data. Moreover,  $N(\lambda)$  varies continuously with  $\lambda$ ,  $0 \leq \lambda \leq 1$ .

Thus for sufficiently small  $\lambda > 0$ ,  $S(0) \subseteq N(\lambda)$  and the analytic solution  $z(x, t)$  is defined throughout the domain  $D(\lambda)$  which is bounded by  $S(0)$  and  $S(\lambda)$ . We consider the identity

$$\begin{aligned}
 (5.3) \quad 0 &= \int_{D(\lambda)} (uL(z) - zL(u)) \, dx \, dt \\
 &= \int_{D(\lambda)} \operatorname{div}_{x,t} \begin{pmatrix} -uz_x a_1 + zu_x a_1 \\ \vdots \\ -uz_x a_n + zu_x a_n \\ \rho uz_t - \rho zu_t \end{pmatrix} \, dx \, dt \\
 &= \int_{S(0) \cup S(\lambda)} (-uz_x A \eta + zu_x A \eta + \eta_0 \rho uz_t - \eta_0 \rho zu_t) \, d\sigma,
 \end{aligned}$$

where  $d\sigma$  denotes integration with respect to surface area on  $S(0) \cup S(\lambda)$  and  $\begin{pmatrix} \eta \\ \eta_0 \end{pmatrix}$  is the unit outward normal to  $S(0) \cup S(\lambda)$  in  $R^{n+1}$ , defined in the relative interiors of  $S(0)$  and  $S(\lambda)$ . We observe that  $\eta_0 \equiv 0$  on  $S(0)$ , and that (5.1) holds on  $S(0)$  since  $S(0) \subseteq \Gamma_1 \times [0, T]$ . Thus (5.3) reduces to

$$(5.4) \quad \int_{S(\lambda)} (u(\eta_0 \rho z_t - z_x A \eta) + z(\eta_0 \rho u_t - u_x A \eta)) \, d\sigma = 0.$$

We choose analytic Cauchy data for  $z$  on  $S(\lambda)$  as follows: we put  $z \equiv 0$  on  $S(\lambda)$  and we take the normal derivative of  $z$  across  $S(\lambda)$  to be an arbitrary real-valued analytic function  $\alpha$ , i.e.,

$$(5.5) \quad z(x(\sigma), t(\sigma)) = 0, \quad (x(\sigma), t(\sigma)) \in S(\lambda),$$

$$(5.6) \quad z_x(x(\sigma), t(\sigma))\eta(\sigma) + z_t(x(\sigma), t(\sigma))\eta_0(\sigma) = \alpha(\sigma), \quad (x(\sigma), t(\sigma)) \in S(\lambda).$$

The equations (5.5), (5.6) together imply that

$$(5.7) \quad (z_x(x(\sigma), t(\sigma)), z_t(x(\sigma), t(\sigma))) \equiv \alpha(\sigma)(\eta'(\sigma), \eta_0(\sigma)), \quad (x(\sigma), t(\sigma)) \in S(\lambda).$$

Substituting (5.7) and (5.5) in (5.4) we obtain

$$(5.8) \quad \int_{S(\lambda)} u\alpha(\rho(\eta_0)^2 - \eta' A \eta) \, d\sigma = 0.$$

Since  $S(\lambda)$  is uniformly time-like, we have

$$(5.9) \quad \rho(x(\sigma))(\eta_0(\sigma))^2 - \eta'(\sigma)A(x(\sigma))\eta(\sigma) \equiv \beta(\sigma) \leq -\beta_0 < 0$$

for all values of the vector  $\sigma$  parametrizing  $S(\lambda)$ . The equation (5.8) becomes

$$\int_{S(\lambda)} u\alpha\beta \, d\sigma = 0.$$

Since this equation holds for all real analytic functions  $\alpha$ , we conclude

$$u(x(\sigma), t(\sigma))\beta(\sigma) \equiv 0, \quad (x(\sigma), t(\sigma)) \in S(\lambda),$$



and, since (5.9) shows that  $\beta$  never vanishes, we have

$$u(x(\sigma), t(\sigma)) \equiv 0, \quad (x(\sigma), t(\sigma)) \in S(\lambda).$$

Repeating this argument on surfaces  $S(\mu)$ ,  $0 < \mu < \lambda$ , which sweep out the interior of  $D(\lambda)$ , we conclude

$$u(x, t) \equiv 0, \quad (x, t) \in D(\lambda).$$

We now let  $I$  denote the largest subinterval of  $[0, 1]$  which includes 0 and has the property that  $u \equiv 0$  on  $S(\lambda)$  if  $\lambda \in I$ . We have seen above that  $I$  is nonempty. Essentially the same technique can be used to show that  $I$  is open. But it is obvious that  $I$  is closed, and we conclude, from the connectedness of  $[0, 1]$ , that  $I = [0, 1]$ . Thus

$$u(x, t) \equiv 0, \quad (x, t) \in D(1),$$

and, since  $(\bar{x}, \bar{t}) \in D(1)$ , we have

$$u(\bar{x}, \bar{t}) \equiv 0.$$

Since  $(\bar{x}, \bar{t})$  is an arbitrary point in the interior of  $K(\Gamma_1, t_0, t_1)$  and since  $u$  is continuous in  $K(\Gamma_1, t_0, t_1)$ , we see that (5.2) follows and the proof is complete.

**6. Controllability for  $T > 2T_0$ ,  $n \leq 3$ .** Let the state  $(\hat{v}, \hat{v}_t)$  lie in the finite energy space  $H_E$ , and suppose that for all  $f \in F$  we have

$$(6.1) \quad \int_{\Omega} (\rho(x)w_t^f(x, T)\hat{v}_t(x) + W_x^f(x)A(x)\hat{v}_x'(x)) dx = 0.$$

If this implies  $\hat{v}_t(x) \equiv 0$ ,  $\hat{v}(x) \equiv \text{const.}$ , then  $R_T$  is dense in the Hilbert space  $H_E$  and we have approximate controllability.

Let  $v(x, t)$  be the generalized solution in  $\Omega \times [0, T]$  of the partial differential equation  $L(v) = 0$  corresponding to homogeneous boundary conditions (2.4). If  $v(x, t)$  were smooth, say  $v \in C^3((\Omega \cup \Gamma) \times [0, T])$ , the proof of our controllability result would not be difficult. Applying Theorem 1 we would get  $v_t(x(s), t) \equiv 0$  for  $(x(s), t) \in \Gamma_1 \times [0, T]$ . Putting  $u(x, t) \equiv v_t(x, t)$ , we would have a solution of  $L(u) = 0$  satisfying the hypotheses (5.1) of Theorem 3 and we could conclude  $v_t \equiv u \equiv 0$  in  $K(\Gamma_1, 0, T)$ . If  $T > 2T_0$ , the set  $K(\Gamma_1, 0, T)$  includes  $\Omega \times [T/2 - \varepsilon, T/2 + \varepsilon]$  for some  $\varepsilon > 0$ . If  $v_t$  vanishes in  $\Omega \times [T/2 - \varepsilon, T/2 + \varepsilon]$ , then  $v_{tt}(x, T/2) \equiv 0$  which would imply  $\sum_{i,j=1}^n (\alpha_{ij}(x)v_{it}(x, T/2))_j \equiv 0$ . Thus  $v(x, T/2)$  would be a solution of the elliptic boundary value problem

$$(6.2) \quad \sum_{i,j=1}^n \left( \alpha_{ij}(x)v_i \left( x, \frac{T}{2} \right) \right)_j = 0, \quad x \in \Omega,$$

$$(6.3) \quad v_x \left( x(s), \frac{T}{2} \right) A(x(s))\eta(s) = 0, \quad x(s) \in \Gamma.$$

It is clear that the only solutions of (6.2), (6.3) have the form  $v \equiv \text{const.}$  Thus we would have  $v_t(x, T/2) \equiv 0$ ,  $v_x(x, T/2) \equiv 0$  which would show that  $\mathcal{E}(v(\cdot, T/2), v_t(\cdot, T/2)) = 0$ . Since solutions of  $L(v) = 0$  with boundary conditions (2.4) conserve energy, we could then conclude  $\mathcal{E}(v(\cdot, T), v_t(\cdot, T)) = \mathcal{E}(\hat{v}, \hat{v}_t) = 0$ , so that  $\hat{v}_t \equiv 0$ ,

$\hat{v} \equiv \text{const.}$  and the proof would be complete. In fact, under these conditions we could obtain the result for  $T = 2T_0$  also and there is nothing special about  $n \leq 3$ .

Unfortunately we are not at all justified in assuming such smoothness for  $v(x, t)$ . A rigorous proof requires that we allow  $(\hat{v}, \hat{v}_i)$  to be an arbitrary finite energy state. All this gives us is that  $\hat{v}_i \in L^2(\Omega)$  and  $\hat{v}_i \in L^2(\Omega)$ ,  $i = 1, 2, \dots, n$ . The generalized solution  $v(x, t)$  is no smoother. For this reason it becomes a nontrivial task to justify the argument presented above. In the present paper we will give such justification only for  $n \leq 3$ ,  $T > 2T_0$ . The result undoubtedly remains true for larger values of  $n$  and for  $T = 2T_0$  but rather involved arguments seem to be required. Fortunately,  $n \leq 3$  includes most cases of physical interest.

**THEOREM 4.** *If  $(\hat{v}, \hat{v}_i) \in H_E$  is such that (6.1) holds for all  $f \in F$  then  $\hat{v}_i \equiv 0$ ,  $\hat{v} \equiv \text{const.}$  provided  $n \leq 3$ ,  $T > 2T_0$ . Thus the system (2.1), (2.2) is approximately controllable in time  $T > 2T_0$  when  $n \leq 3$ .*

*Proof.* Let  $v(x, t)$  be the generalized solution of  $L(v) = 0$  with homogeneous boundary conditions (2.4) and satisfying the terminal conditions  $v(x, t) \equiv \hat{v}(x)$ ,  $v_t(x, T) \equiv \hat{v}_t(x)$ . Let  $(\hat{v}^k, \hat{v}_t^k)$  be a sequence of states in  $H_E$  converging to  $(\hat{v}, \hat{v}_t)$  in the energy norm as  $k \rightarrow \infty$ . Moreover,  $(\hat{v}^k(x), \hat{v}_t^k(x)) \in C^\infty(\Omega \cup \Gamma)$  and satisfy the consistency conditions

$$(6.4) \quad \hat{v}_x^k(x(s))A(x(s))\eta(s) = 0, \quad \hat{v}_{tx}^k(x(s))A(x(s))\eta(s) = 0, \quad x(s) \in \Gamma.$$

One way in which this could be done is to expand  $\hat{v}(x)$ ,  $\hat{v}_t(x)$  in terms of the eigenfunctions  $\varphi_j(x)$  of the operator  $B$  introduced in (2.3):

$$\hat{v}(x) = \sum_{j=0}^{\infty} \alpha_j \varphi_j(x), \quad \hat{v}_t(x) = \sum_{j=0}^{\infty} \beta_j \varphi_j(x)$$

and take

$$\hat{v}^k(x) = \sum_{j=0}^k \alpha_j \varphi_j(x), \quad \hat{v}_t^k(x) = \sum_{j=0}^k \beta_j \varphi_j(x).$$

For  $k = 0, 1, 2, \dots$  let  $v^k(x, t)$  be the  $C^\infty$  solutions of  $L(v^k) = 0$  which satisfy the terminal and boundary conditions

$$(6.5) \quad \begin{aligned} v^k(x, T) &\equiv \hat{v}^k(x), & v_t^k(x, T) &\equiv \hat{v}_t^k(x), \\ v_x^k(x(s), t)A(x(s))\eta(s) &\equiv 0, & (x(s), t) &\in \Gamma \times [0, T]. \end{aligned}$$

It is known [6, Chap. VI] that for each fixed  $t_1$ ,  $0 \leq t_1 \leq T$ , the states  $(v^k(\cdot, t_1), v_t^k(\cdot, t_1))$  converge to  $(v(\cdot, t_1), v_t(\cdot, t_1))$  in the space  $H_E$ ; in fact, this is just a consequence of the energy conservation.

Since

$$(6.6) \quad \lim_{k \rightarrow \infty} \|(\hat{v}, \hat{v}_t) - (\hat{v}^k, \hat{v}_t^k)\|_E = 0$$

and (6.1) is assumed to hold, we have

$$\lim_{k \rightarrow \infty} \int_{\Omega} [\rho w_t^f \hat{v}^k + w_x^f A \hat{v}_x^k] dx = 0.$$

Performing a computation similar to that in the proof of Theorem 1, we conclude that for each fixed  $f \in F$ ,

$$(6.7) \quad \lim_{k \rightarrow \infty} \int_{\Gamma_1 \times [0, T]} v_t^k f \, ds \, dt = 0.$$

Let us now define functions  $D^{-j}v$ ,  $j = 1, 2, \dots$ , for  $x \in \Omega$  by

$$(D^{-1}v)(x, t) = \int_T^t v(x, \tau) \, d\tau, \quad (D^{-(j+1)}v)(x, t) = \int_T^t (D^{-j}v)(x, \tau) \, d\tau.$$

We define  $D^{-j}v^k$  similarly and verify without difficulty, since the  $v^k$  satisfy the evolution equation (2.3), that

$$(6.8) \quad \begin{aligned} \frac{d^2}{dt^2}(D^{-j}v^k) + B(D^{-j}v^k) &= \hat{v}_t^k \frac{(t-T)^{j-1}}{(j-1)!} + \hat{v}^k \frac{(t-T)^{j-2}}{(j-2)!}, \\ (D^{-j}v^k)(T) &= (D^{-j}v^k)_t(T) = 0, \end{aligned}$$

for  $j > 2$ .

A result proved in [14, Theorem 1.19, p. 486] shows that the inhomogeneous linear initial value problem

$$\frac{d^2 r}{dt^2} + Br = g(t), \quad r(T) = r_t(T) = 0,$$

where  $g: [0, T] \rightarrow L^2(\Omega)$  is continuously differentiable with respect to  $t$ , has a unique solution  $r$  such that  $r_{tt}(t_1) \in L^2(\Omega)$  and  $r(t_1)$  lies in the domain  $\Delta(B) \subseteq L^2(\Omega)$  of the unbounded self-adjoint operator  $B$  for  $0 \leq t_1 \leq T$ . Moreover it is also shown in the theorem cited that there are constants  $M_0, M_1$  such that

$$(6.9) \quad \|Br(t_1)\| < M_0 \sup_{0 \leq \tau \leq T} \|g(\tau)\| + M_1 \sup_{0 \leq \tau \leq T} \|g_t(\tau)\|$$

uniformly for all  $t_1 \in [0, T]$ , where  $\|\cdot\|$  denotes the usual norm in  $L^2(\Omega)$ , and it is also shown that  $Br(t)$  is continuous in  $t$  with respect to the norm  $\|\cdot\|$ .

Applying this theorem we see that, for  $j \geq 2$ ,  $D^{-j}v$  lies in  $\Delta(B)$  and solves

$$(6.10) \quad \begin{aligned} \frac{d^2}{dt^2}(D^{-j}v) + B(D^{-j}v) &= \hat{v}_t \frac{(t-T)^{j-1}}{(j-1)!} + \hat{v} \frac{(t-T)^{j-2}}{(j-2)!}, \\ (D^{-j}v)(T) &= (D^{-j}v)_t(T) = 0. \end{aligned}$$

Further, (6.9) together with (6.8), (6.10) shows that

$$(6.11) \quad \lim_{k \rightarrow \infty} \|B(D^{-j}v)(\cdot, t) - B(D^{-j}v^k)(\cdot, t)\| = 0$$

uniformly for  $0 < t < T$ .

Now the operator  $B$  is uniformly elliptic and we can apply known results from the theory of elliptic boundary value problems [15, Theorem 9.11, p. 132 and Remarks, p. 148] to show that the fact that  $D^{-j}v \in \Delta(B)$  together with (6.6) implies that  $D^{-j}v(\cdot, t)$  lies in the space  $H_2(\Omega)$  (for definition of  $H_m(\Omega)$ , see [15]) for  $0 \leq t \leq T$ ,  $\|(D^{-j}v)(\cdot, t)\|_{2, \Omega}$  is continuous and uniformly bounded for  $0 \leq t \leq T$ , and

$$(6.12) \quad \lim_{k \rightarrow \infty} \|(D^{-j}v)(\cdot, t) - (D^{-j}v^k)(\cdot, t)\|_{2, \Omega} = 0$$

uniformly for  $0 \leq t \leq T$ , where  $\|\cdot\|_{2,\Omega}$  is the norm in  $H_2(\Omega)$ , the sum of the integrals of the squares of the partial derivatives of order  $\leq 2$ .

The theorem of Sobolev [15, Theorem 3.9, p. 32] states that if  $r \in H_m(\Omega)$ , then  $r$  can be modified on a set of measure zero so that  $r \in C^1(\Omega \cup \Gamma)$ , provided  $l$  is an integer such that  $l < m - n/2$ . For  $m = 2$ , we have  $0 < 2 - n/2$  when  $n \leq 3$ , so, for such  $n$ ,  $r \in H_2(\Omega)$  implies  $r \in C^0(\Omega \cup \Gamma)$ . Moreover, if  $\|\cdot\|_s$  denotes the usual "sup" norm in  $C^0(\Omega \cup \Gamma)$ , we have

$$\|r\|_s \leq \alpha_0 \|r\|_{2,\Omega} + \alpha_1 \|r\|.$$

Applying these results with  $r = D^{-j}v$ ,  $j \geq 2$ , the uniform boundedness and continuity of  $\|(D^{-j}v)(\cdot, t)\|_{2,\Omega}$  and (6.11), we conclude that  $(D^{-j}v)(x, t)$  is continuous for  $(x, t) \in (\Omega \cup \Gamma) \times [0, T]$  and

$$(6.13) \quad \lim_{k \rightarrow \infty} (D^{-j}v^k)(x, t) = (D^{-j}v)(x, t)$$

uniformly for such  $(x, t)$ .

Having now obtained the continuity of  $(D^{-2}v)(x, t)$  we return to (6.7). Integrating by parts three times, we conclude that for all  $f \in F$ ,

$$\lim_{k \rightarrow \infty} \int_{\Gamma_1 \times [0, T]} (D^{-2}v^k) f_{iii} \, ds \, dt = 0$$

which with (6.13) implies

$$(6.14) \quad \int_{\Gamma_1 \times [0, T]} (D^{-2}v) f_{iii} \, ds \, dt = 0, \quad f \in F.$$

Taking account of the fact that  $f$  and all its derivatives vanish outside a compact subset of  $\Gamma_1 \times (0, T)$ , (6.14) implies that  $(D^{-2}v)(x(s), t)$  is a polynomial in  $t$  of degree at most 2 whose coefficients are continuous functions of  $x(s)$ , for all  $(x(s), t) \in \Gamma_1 \times [0, T]$ . Let  $\delta$  be a small positive number. We define the third order difference

$$\Delta^3 r(x, t) = r(x, t + 3\delta) - 3r(x, t + 2\delta) + 3r(x, t + \delta) - r(x, t)$$

for any function  $r$  defined on  $\Omega \times [0, T]$ . The function  $\Delta^3 r$  is defined on  $\Omega \times [0, T - 3\delta]$  and possesses all smoothness properties of  $r(x, t)$ . Applying this difference operator to  $D^{-2}v$  and  $D^{-2}v^k$  we obtain functions

$$(6.15) \quad \begin{aligned} \hat{u}(x, t) &= \Delta^3(D^{-2}v)(x, t), \\ u^k(x, t) &= \Delta^3(D^{-2}v^k)(x, t), \quad k = 0, 1, 2, \dots \end{aligned}$$

From (6.8), (6.10) we see that

$$L(\hat{u}) = \frac{d^2 \hat{u}}{dt^2} + B\hat{u} = 0, \quad L(u^k) = \frac{d^2 u^k}{dt^2} + Bu^k = 0,$$

the  $u^k$  satisfy the homogeneous boundary conditions

$$(6.16) \quad u_x^k(x(s), t)A(x(s))\eta(s) = 0, \quad (x(s), t) \in \Gamma \times [0, T - 3\delta],$$

while  $\hat{u}(x, t)$  is continuous,  $\hat{u}(\cdot, t)$  lies in  $\Delta(B)$  (which means  $\hat{u}(x, t)$  satisfies the boundary conditions  $\hat{u}_x(x(s), t)A(x(s))\eta(s) = 0$ ,  $(x(s), t) \in \Gamma \times [0, T - 3\delta]$  in some

sense which we need not specify) and, from the fact that  $D^{-2}v$  is a polynomial of degree at most 2 on  $\Gamma_1 \times [0, T]$ , we have

$$(6.17) \quad \hat{u}(x(s), t) \equiv 0, \quad (x(s), t) \in \Gamma_1 \times [0, T - 3\delta].$$

Now we refer back to Theorem 3, or, more precisely, its proof. We let  $t_0 = 0$ ,  $t_1 = T - 3\delta$  and define the surfaces  $S(\lambda)$  as we did there. We define  $z$  as we did there and put  $u^k$  in place of the function  $u$  of Theorem 3. Repeating the calculations following (5.3), we see that

$$\int_{S(0) \cup S(\lambda)} (-u^k z_x A \eta + z u_x^k A \eta + \eta_0 \rho u^k z_t - \eta_0 \rho z u_t^k) d\varphi = 0.$$

On  $S(0)$  we have  $\eta_0(s) \equiv 0$  and  $u_x^k A \eta \equiv 0$ . Defining  $\alpha$  and  $\beta$  as in the proof of Theorem 3 and recalling  $z \equiv 0$  on  $S(\lambda)$ , we have

$$\int_{S(\lambda)} u^k \alpha \beta d\varphi = \int_{S(0)} u^k z_x A \eta ds dt.$$

Now  $u^k$  converges uniformly to  $\hat{u}$  in  $(\Omega \cup \Gamma) \times [0, T - 3\delta]$ , so we have

$$\int_{S(\lambda)} \hat{u} \alpha \beta d\delta = \int_{S(0)} \hat{u} z_x A \eta ds dt = 0$$

since  $\hat{u}$  obeys (6.17). As in Theorem 3 we conclude that  $\hat{u} \equiv 0$  on  $S(\lambda)$ . A continuation process similar to that described in Theorem 3 can be used to show that  $\hat{u} \equiv 0$  on every surface  $S(\lambda)$ ,  $0 \leq \lambda \leq 1$ . Then, just as in Theorem 3, we conclude that

$$\hat{u}(x, t) \equiv 0, \quad (x, t) \in K(\Gamma_1, 0, T - 3\delta).$$

Now if  $T > 2T_0$ , we have

$$\Omega \times \left[ \frac{T}{2} - \varepsilon, \frac{T}{2} + \varepsilon \right] \subseteq K(\Gamma_1, 0, T - 3\delta)$$

if  $\varepsilon$  and  $\delta$  are both chosen sufficiently small. (We need  $\varepsilon + 3\delta < T/2 - T_0$ .) Thus we have, for small  $\varepsilon$  and  $\delta$ ,

$$(6.18) \quad \hat{u}(x, t) \equiv 0, \quad x \in \Omega, \quad \frac{T}{2} - \varepsilon \leq t \leq \frac{T}{2} + \varepsilon.$$

Returning to the definition (6.15) we see that if (6.18) holds for all small  $\delta$ , then it must be true that for  $T/2 - \varepsilon \leq t \leq T/2 + \varepsilon$ ,  $(D^{-2}v)(x, t)$  is a polynomial in  $t$  of degree not greater than 2 with coefficients which are functions of  $x$  lying in  $\Delta(B)$  (since  $D^{-2}v \in \Delta(B)$ ). Differentiating  $D^{-2}v$  twice with respect to  $t$  in  $T/2 - \varepsilon < t < T/2 + \varepsilon$ , we see that there is a function  $\tilde{v}(x)$  with  $\tilde{v} \in \Delta(B)$  such that

$$v(x, t) \equiv \tilde{v}(x), \quad x \in \Omega, \quad \frac{T}{2} - \varepsilon \leq t \leq \frac{T}{2} + \varepsilon.$$

But if  $v(x, t)$  is a generalized solution of  $d^2v/dt^2 + Bv = 0$  such that  $v \in \Delta(B)$  and  $v(x, t)$  is constant with respect to  $t$ , then we must have

$$Bv(x, t) \equiv B\tilde{v}(x) \equiv 0.$$

But the only elements  $\tilde{v} \in \Delta(B)$  for which  $B\tilde{v} = 0$  are of the form  $\tilde{v} \equiv \text{const.}$  Therefore, for  $T/2 - \varepsilon \leq t \leq T/2 + \varepsilon$ ,

$$v(x, t) \equiv \tilde{v}(x) \equiv \text{const.}, \quad v_t(x, t) \equiv 0.$$

Since the energy associated with the generalized solution  $v(x, t)$  is constant, we conclude that  $v(x, t) \equiv \hat{v}(x) \equiv \text{const.}$ ,  $v_t(x, T) \equiv \hat{v}_t(x) \equiv 0$  and the proof of Theorem 4 is complete.

**Acknowledgment.** I should like to express my appreciation to Professor J. L. Lions of the University of Paris whose suggestions in a 1966 letter provided the germinal idea for the proofs presented in this paper.

#### REFERENCES

- [1] D. L. RUSSELL, *On boundary-value controllability of linear symmetric hyperbolic systems*, Proc. Conference on Mathematical Theory of Control (Univ. of Southern California, Los Angeles, 1967), Academic Press, New York, 1967.
- [2] ———, *Nonharmonic Fourier series in the control theory of distributed parameter systems*, J. Math. Appl., 18 (1967), pp. 542–560.
- [3] M. A. CIRINA, *Boundary controllability of nonlinear hyperbolic systems*, this Journal, 7 (1969), pp. 198–212.
- [4] J. J. GRAINGER, *Boundary-value control of distributed systems characterized by hyperbolic differential equations*, Doctoral thesis, Electrical Engineering Dept., Univ. of Wisconsin, Madison, 1967.
- [5] R. H. THOMAS, *The identification and control of hyperbolic distributed processes*, Doctoral thesis, Electrical Engineering Dept., Univ. of Wisconsin, Madison, 1969.
- [6] R. COURANT AND D. HILBERT, *Partial Differential Equations*, Methods of Mathematical Physics, vol. II; Interscience, New York, 1962.
- [7] H. O. FATTORINI, *Boundary control systems*, this Journal, 6 (1968), pp. 349–388.
- [8] E. HOLMGREN, *Über Systeme von linearen partiellen Differentialgleichungen*, Ofversigt af kongl. Vetenskapsakad. Förel., 58 (1901), pp. 91–103.
- [9] F. JOHN, *On linear partial differential equations with analytic coefficients—Unique continuation of data*, Comm. Pure Appl. Math., 2 (1949), pp. 209–253.
- [10] K. O. FRIEDRICHS, *Symmetric hyperbolic linear differential equations*, Ibid., 7 (1954), pp. 345–392.
- [11] P. D. LAX, *On Cauchy's problem for hyperbolic equations and the differentiability of solutions of elliptic equations*, Ibid., 8 (1955), pp. 615–633.
- [12] G. F. D. DUFF, *Mixed problems for hyperbolic equations of general order*, Canad. J. Math., 2 (1959), pp. 195–221.
- [13] E. B. LEE AND L. MARKUS, *Foundations of Optimal Control Theory*, Interscience, New York, 1968.
- [14] T. KATO, *Perturbation Theory for Linear Operators*, Springer-Verlag, New York, 1966.
- [15] S. AGMON, *Lectures on Elliptic Boundary Value Problems*, Van Nostrand, Princeton, 1965.

## ON OPTIMAL CONTROLS FOR MEASURE DELAY- DIFFERENTIAL EQUATIONS\*

P. C. DAS† AND R. R. SHARMA‡

**1. Introduction.** In recent years it has been a matter of interest to study the systems described by differential equations containing impulses which cause discontinuous changes in the values of the state variables of the system. Such systems include pulse frequency modulation systems and models for biological neural nets. (See [10], [11] and the references given therein.) Schmaedeke [13] has considered the control system described by the equations

$$(1.1) \quad \frac{dx^i}{dt} = f^i(t, x^1, x^2, \dots, x^n, u^1, u^2, \dots, u^m) + \sum_{j=1}^m g_j^i(t) \frac{du^j}{dt},$$

$i = 1, 2, \dots, n,$

in which the control inputs  $u^j(t)$  may have discontinuities of the first kind giving rise to impulse control input  $du^j/dt$ . He developed the theory of the equations (1.1) when the control functions  $u^j(t)$  and the solutions  $x^i(t)$  are functions of bounded variation and the derivatives  $dx^i/dt$  and  $du^j/dt$  are taken in the sense of distribution derivatives, to be denoted by  $Dx^i$  and  $Du^j$  respectively, which can be identified with Stieltjes measures. He then considered the existence of optimal control for the system given by (1.1). But in this optimal control problem the coefficients  $g_j^i$  of the impulsive controllers may be expected in many cases to depend not only on  $t$  but also on the state variables  $x^i$ . Moreover, the system may have hereditary effect. Due to these considerations an attempt has been made in this paper to generalize the results in [13] by considering the control system described by delay-differential equations

$$(1.2) \quad \begin{aligned} Dx^i(t) = & f^i(t, x_t^1, x_t^2, \dots, x_t^n, u^1, u^2, \dots, u^m) \\ & + \sum_{j=1}^m g_j^i(t, x_t^1, x_t^1, \dots, x_t^n, u^1, u^2, \dots, u^m) Du^j(t), \quad i = 1, 2, \dots, n, \end{aligned}$$

where  $x_t^i$  represents the restriction of the function  $x^i(s)$  on the interval  $p(t) \leq s \leq q(t)$ ,  $p$  and  $q$  being real functions with the property  $p(t) \leq q(t) \leq t$  for each  $t$ ; for each fixed  $t$ ,  $f^i$  and  $g_j^i$  are functionals defined on the space  $BV([p(t), q(t)])$ ; and  $u^j(t)$  are right continuous functions of bounded variation. If the  $g_j^i$  depend only on  $t$ , and  $p(t) = q(t) = t$  so that  $x_t^i = x^i(t)$ , then (1.2) coincides with (1.1).

Throughout this paper, the (Stieltjes) integrals are taken to be Lebesgue (–Stieltjes) integrals. This requires in the proof of Theorem 1 a careful use of the integration by parts formula which does not always hold for Lebesgue–Stieltjes integrals. (In the corresponding theorem in [13] the integral used is actually Riemann–Stieltjes, though it has not been stated explicitly there.)

\* Received by the editors May 14, 1969, and in final revised form March 18, 1970.

† Department of Mathematics, Indian Institute of Technology, Kanpur, India.

‡ Department of Mathematics, Regional Institute of Technology, Jamshedpur, India.

**2. Preliminaries.** Let  $\Omega$  be a subset of  $n$ -dimensional Euclidean space  $E^n$ . We denote by  $C_c^\infty(\Omega)$  the class of infinitely partially differentiable complex functions, defined on  $\Omega$ , which have compact support.  $C_c^\infty(\Omega)$  is a normed linear space with addition, scalar multiplication and norm defined by

$$\begin{aligned}(\psi_1 + \psi_2)(x) &= \psi_1(x) + \psi_2(x), \\ (\alpha\psi)(x) &= \alpha\psi(x), \\ \|\psi\| &= \sup_{x \in \Omega} |\psi(x)|.\end{aligned}$$

A continuous linear functional defined on  $C_c^\infty(\Omega)$  is called a distribution on  $\Omega$ . It follows from the Riesz representation theorem that any distribution  $F$  on  $\Omega$  can be identified with a complex Borel measure  $\mu$  by the relation

$$(2.1) \quad F(\psi) = \int_{\Omega} \psi \, d\mu, \quad \psi \in C_c^\infty(\Omega).$$

If a distribution  $F$  is given by (2.1) and  $g$  is a  $\mu$ -integrable function, then we define the product  $gF$  by

$$(2.2) \quad (gF)(\psi) = \int_{\Omega} g\psi \, d\mu, \quad \psi \in C_c^\infty(\Omega).$$

It is easy to see that  $gF$  is a distribution on  $\Omega$ . A distribution  $F$  on an interval  $I$  is to be identified with the Lebesgue–Stieltjes measure  $dh(t)$  if for every closed finite interval  $J$  contained in  $I$ ,  $h(t)$  is of bounded variation on  $J$  and

$$(2.3) \quad F(\psi) = \int_J \psi(t) \, dh(t)$$

for all  $\psi \in C_c^\infty(J)$ . A distribution  $F$  on an interval  $I$  is to be identified with a point function  $f$  if for every closed finite interval  $J$  contained in  $I$ ,  $f$  is integrable on  $J$  and

$$(2.4) \quad F(\psi) = \int_J f(t)\psi(t) \, dt$$

for all  $\psi \in C_c^\infty(J)$ . The derivative  $DF$  of a distribution  $F$  on  $I$  is a distribution defined by

$$(2.5) \quad DF(\psi) = -F(\psi'), \quad ' \equiv d/dt.$$

For any vector  $x = (x^1, x^2, \dots, x^n) \in E^n$ , the norm will be defined by

$$(2.6) \quad |x| = \sum_{i=1}^n |x^i|.$$

The norm of an  $n \times m$  matrix  $G = (g_j^i)$  will be defined by

$$(2.7) \quad |G| = \sum_{i=1}^n \sum_{j=1}^m |g_j^i|.$$



The space  $BV(I)$  is defined for an interval  $I$  and consists of all scalar functions  $f$  on  $I$  which are of bounded variation. If  $a$  is the left endpoint of  $I$ , then the norm of  $f$  is

$$(2.8) \quad |f|_I = v(f, I) + |f(a+)|,$$

where  $v(f, I)$  denotes the total variation of  $f$  on  $I$ . With this norm the space  $BV(I)$  is a Banach space. The space  $BV(I)_n$  is defined for an interval  $I$  and consists of all vector functions  $f$  with values in  $E^n$  whose individual components belong to  $BV(I)$ . The norm of  $f$  is

$$\begin{aligned} \|f\|_I &= \sum_{i=1}^n |f^i|_I = \sum_{i=1}^n \{v(f^i, I) + |f^i(a+)\} \\ &= v(f, I) + |f(a+)|. \end{aligned}$$

With this norm  $BV(I)_n$  is a Banach space.

**3. Existence and uniqueness of solutions.** Let  $S$  be a domain (an open connected set) in  $E^n$ . The set of all functions in  $BV(I)_n$  with values in  $S$  will be denoted by  $BV(I, S)$ . Let  $\alpha, \beta$  and  $t_0$  be numbers such that

$$-\infty \leq \alpha < t_0 < \beta \leq \infty.$$

In what follows, the interval  $[\alpha, t_0]$  will be understood to be  $(-\infty, t_0]$  in case  $\alpha = -\infty$ ; similarly, if  $\beta = \infty$ , the interval  $[t_0, \beta]$  will mean  $[t_0, \infty)$ . Let  $p$  and  $q$  be two real functions defined on  $[t_0, \beta]$  and satisfying

$$\alpha \leq p(t) \leq q(t) \leq t$$

for each  $t$ . We define the interval

$$I_t = [p(t), q(t)].$$

Let  $x_t$  denote the restriction of the function  $x(\tau)$  on the interval  $I_t$ . For each  $t \in [t_0, \beta]$ , we define

$$(3.1) \quad R_t = \{(t, x_t) | x \in BV([\alpha, t], S)\}$$

and let  $R$  be defined by

$$(3.2) \quad R = \bigcup_{t \in [t_0, \beta]} R_t.$$

Let  $f(t, x_t)$  be an  $n$ -vector functional and  $G(t, x_t)$  an  $n \times m$  matrix functional defined on  $R$ . Let  $u(t)$  be a right continuous  $m$ -vector function of bounded variation defined on  $[t_0, \beta]$ . We assume that for each given  $x \in BV([\alpha, \beta], S)$ ,  $f(t, x_t)$  is Lebesgue measurable and  $G(t, x_t)$  is integrable with respect to the Lebesgue-Stieltjes measure  $du(t)$  on  $[t_0, \beta]$

Consider now the delay-differential equation

$$(3.3) \quad Dx = f(t, x_t) + G(t, x_t)Du, \quad t > t_0,$$

where the operations of differentiation are to be understood in the sense of distribution derivatives with respect to the real variable  $t$ . Since the distribution derivative  $Du$  of a function  $u$  of bounded variation can always be identified with

a Lebesgue–Stieltjes measure, we shall call (3.3) a measure delay-differential equation.

DEFINITION 1. A function  $x(t)$  is called a *solution of (3.3) on an interval  $I$* ,  $[\alpha, t_0] \subseteq I \subseteq [\alpha, \beta]$ , with the initial function  $\varphi \in BV([\alpha, t_0], S)$ , if

- (i)  $x \in BV(I, S)$ ,
- (ii)  $x(t) = \varphi(t)$  for  $t \in [\alpha, t_0]$ ,
- (iii)  $x(t)$  is continuous from the right on  $I \cap [t_0, \beta]$ ,
- (iv)  $x(t)$  satisfies (3.3) on  $I \cap (t_0, \beta]$ .

Consider the integral equation

$$(3.4) \quad x(t) = \begin{cases} \varphi(t) & \text{for } t \in [\alpha, t_0], \\ \varphi(t_0) + \int_{t_0}^t f(s, x_s) ds + \int_{t_0}^t G(s, x_s) du(s) & \text{for } t > t_0. \end{cases}$$

DEFINITION 2. A function  $x \in BV(I, S)$  is called a *solution of (3.4) on the interval  $I$* ,  $[\alpha, t_0] \subseteq I \subseteq [\alpha, \beta]$ , if it satisfies this equation for  $t \in I$ .

We shall now prove the following theorem.

THEOREM 1.  $x(t)$  is a solution of (3.4) if and only if it is a solution of (3.3) with initial function  $\varphi$ .

*Proof.* Let  $x(t)$  be a solution of (3.4) on the interval  $[\alpha, T]$ . Then conditions (i) and (ii) of Definition 1 are satisfied. The right continuity of  $u(t)$  implies that the integral  $\int_{t_0}^t G(s, x_s) du(s)$  is also a right continuous function of  $t$ . The integral

$\int_{t_0}^t f(s, x_s) ds$  is an absolutely continuous (and hence continuous) function of  $t$ .

Thus  $x(t)$  is right continuous. We shall show that  $x(t)$  satisfies (3.3) on  $(t_0, T]$ . Let  $F^i(\psi)$  be the distribution on  $[t_0, T]$  to be identified with the  $i$ th component  $x^i(t)$  of  $x(t)$ . Then for any closed interval  $J = [a, b]$  contained in  $[t_0, T]$ , we have

$$F^i(\psi) = \int_J \left\{ \varphi^i(t_0) + \int_{t_0}^t f^i(s, x_s) ds + \int_{t_0}^t [G(s, x_s) du(s)]^i \right\} \psi(t) dt$$

for all  $\psi \in C_c^\infty(J)$ . The derivative-distribution is

$$(3.5) \quad DF^i(\psi) = -F(\psi') = - \int_J \left\{ \varphi^i(t_0) + \int_{t_0}^t f^i(s, x_s) ds + \int_{t_0}^t \left( \sum_{j=1}^m g_j^i(s, x_s) du^j(s) \right) \right\} \psi'(t) dt,$$

where  $g_j^i(s, x_s)$  is the  $i, j$ th element of  $G(s, x_s)$  and  $u^j(s)$  is the  $j$ th component of  $u(s)$ .

Integration by parts yields

$$(3.6) \quad \begin{aligned} & \int_J \left\{ \varphi^i(t_0) + \int_{t_0}^t f^i(s, x_s) ds \right\} \psi'(t) dt \\ &= \left[ \left( \varphi^i(t_0) + \int_{t_0}^t f^i(s, x_s) ds \right) \psi(t) \right]_a^b - \int_J \psi(t) f^i(t, x_t) dt \\ &= - \int_J \psi(t) f^i(t, x_t) dt \quad (\text{since } \psi(b) = \psi(a) = 0). \end{aligned}$$

The function  $h(t) = \int_{t_0}^t g_j^i(s, x_s) du^j(s)$  is right continuous, and is of bounded variation on the interval  $J = [a, b]$ , by [6, § 52.18, p. 275]. We have

$$\int_J h(t)\psi'(t) dt = \int_J h(t) d\psi(t) = \int_{\{a\}} h(t) d\psi(t) + \int_{(a,b]} h(t) d\psi(t).$$

But

$$\begin{aligned} \int_{\{a\}} h(t) d\psi(t) &= h(a)(\psi(a) - \psi(a-)) \quad (\text{by [7, Example 5, p. 199]}) \\ &= 0 \quad (\text{since } \psi \text{ is continuous}); \end{aligned}$$

and

$$\begin{aligned} \int_{(a,b]} h(t) d\psi(t) &= h(b)\psi(b) - h(a)\psi(a) - \int_{(a,b]} \psi(t) dh(t) \quad (\text{by [7, Example n, p. 185]}) \\ &= - \int_{(a,b]} \psi(t) dh(t) \quad (\text{since } \psi(a) = \psi(b) = 0) \\ &= - \int_{[a,b]} \psi(t) dh(t) + \int_{\{a\}} \psi(t) dh(t) \\ &= - \int_{[a,b]} \psi(t) dh(t) + \psi(a)(h(a) - h(a-)) \\ &= - \int_{[a,b]} \psi(t) dh(t) \quad (\text{since } \psi(a) = 0). \end{aligned}$$

Therefore,

$$\int_J h(t)\psi'(t) dt = - \int_J \psi(t) dh(t);$$

i.e.,

$$\begin{aligned} &\int_J \left( \int_{t_0}^t g_j^i(s, x_s) du^j(s) \right) \psi'(t) dt \\ &= \int_J \psi(t) d \left\{ \int_{t_0}^t g_j^i(s, x_s) du^j(s) \right\} \\ &= \int_J \psi(t) g_j^i(t, x_t) du^j(t) \quad (\text{by [1, Corollary 6, p. 180]}); \end{aligned}$$

and, therefore,

$$(3.7) \quad \int_J \psi(t) d \left\{ \sum_{j=1}^m \int_{t_0}^t g_j^i(s, x_s) du^j(s) \right\} = \int_J \psi(t) \left( \sum_{j=1}^m g_j^i(t, x_t) du^j(t) \right).$$

From (3.5), (3.6) and (3.7), we obtain

$$DF'(\psi) = \int_J \psi(t) f^i(t, x_t) dt + \int_J \psi(t) [G(t, x_t) du(t)]^i.$$

Since  $G(t, x_i)$  is integrable with respect to  $du(t)$ , the last continuous linear functional in the above equation is identified with the measure  $[G(t, x_i) du(t)]^i$  (see (2.2)) while the first continuous linear functional in this equation is identified with the function  $f^i(t, x_i)$ . This holds for  $i = 1, 2, \dots, n$ , and therefore the derivative distribution  $DF(\psi)$  is identified with  $f(t, x_i) + G(t, x_i)Du$ . Hence,  $x(t)$  is also a solution of (3.3).

Conversely, let us suppose that  $x(t)$  is a solution of (3.3) with initial function  $\varphi(t)$ . Then for any closed interval  $J$  contained in  $[t_0, T]$  we have

$$(3.8) \quad \int_J \psi(t) Dx^i(t) = \int_J \psi(t) f^i(t, x_i) dt + \int_J \psi(t) [G(t, x_i) du(t)]^i$$

for  $\psi \in C_c^\infty(J)$ . By using [1, Corollary 6, p. 180] again we may write

$$\int_J \psi(t) [G(t, x_i) du(t)]^i = \int_J \psi(t) d \left( \int_{t_0}^t [G(s, x_s) du(s)]^i \right).$$

Integrating the three integrals in (3.8) in the way we have done above, we obtain

$$\int_J \psi'(t) (x^i(t) - \varphi^i(t_0)) dt = \int_J \psi'(t) \left\{ \int_{t_0}^t f^i(s, x_s) ds + \int_{t_0}^t [G(s, x_s) du(s)]^i \right\} dt.$$

Therefore,

$$x^i(t) = \varphi^i(t_0) + \int_{t_0}^t f^i(s, x_s) ds + \int_{t_0}^t [G(s, x_s) du(s)]^i$$

almost everywhere in  $J$ . But, since  $x^i(t)$  is continuous from the right, being a solution of (3.3), and since the right-hand side of the above equation is a right continuous function of  $t$ , equality holds everywhere in  $J$  in the above equation. Hence  $x(t)$  is a solution of (3.4). The proof of Theorem 1 is thus complete.

We now enumerate the hypotheses which will be used in this section.

H<sub>1</sub>.  $f(t, x_i)$  is locally Lipschitzian with respect to  $x$ ; i.e., for any  $t_1 > t_0$  and every  $k_1 < \infty$ , there exists  $k_2 = k_2(t_1, k_1)$  such that  $\|x\|_{[t_0, t_1]} \leq k_1$ ,  $\|y\|_{[t_0, t_1]} \leq k_1$  imply

$$|f(t, x_i) - f(t, y_i)| \leq k_2 \|x - y\|_{I_t}$$

for every  $t \in [t_0, t_1]$ .

H<sub>2</sub>.  $f(t, x_i)$  is continuous in  $t$  and  $x$ ; i.e.,  $\lim_{v \rightarrow \infty} t^{(v)} = t^*$  and  $\lim_{v \rightarrow \infty} x^{(v)}(t) = x(t)$  imply

$$\lim_{v \rightarrow \infty} f(t^{(v)}, x_t^{(v)}) = f(t^*, x_{t^*}).$$

H<sub>3</sub>. There exists a Lebesgue integrable function  $r(t)$  such that

$$|f(t, x_i)| \leq r(t)$$

uniformly with respect to  $x$ .

H<sub>4</sub>.  $G(t, x_i)$  is locally Lipschitzian with respect to  $x$ .

H<sub>5</sub>.  $G(t, x_i)$  is continuous in  $t$  and  $x$ .

H<sub>6</sub>. There exists a function  $w(t)$  integrable with respect to the Lebesgue–Stieltjes measure  $dV_u(t)$  ( $V_u(t) \equiv v(u, [t_0, t])$  denotes the total variation function

of  $u(t)$  such that

$$|G(t, x_t)| \leq w(t)$$

uniformly with respect to  $x$ .

**THEOREM 2 (Local existence and uniqueness).** *Let the hypotheses  $H_1, H_3, H_4, H_6$  be satisfied on the interval  $[t_0, T]$ . Then there exists a unique solution of (3.3) on an interval  $[\alpha, t_0 + a]$  with a given initial function  $\varphi \in BV([\alpha, t_0], S)$ .*

*Proof.* Denote by  $Q_t, t_0 \leq t \leq T$ , the space of all functions  $x(t)$  with the properties

- (i)  $x \in BV([\alpha, t])_n$ ,
- (ii)  $x(s) = \varphi(s)$  for  $s \in [\alpha, t_0]$ ,
- (iii)  $v(x, [t_0, t]) \leq b, b > 0$ .

Suppose that  $Q_t \subseteq BV([\alpha, t], S)$ . This is always possible if  $b$  is suitably chosen.

Choose  $a, 0 < a \leq T - t_0$ , such that

$$(3.9) \quad \int_{t_0}^{t_0+a} r(t) dt + \int_{t_0}^{t_0+a} w(t) dV_u(t) \leq b.$$

Since  $\int_{t_0}^t r(s) ds$  is a continuous function of  $t$ , and since the right continuity of  $u(t)$  and hence of  $dV_u(t)$  imply that  $\int_{t_0}^t w(s) dV_u(s)$  is also a right continuous function  $t$ , it is possible to choose such an  $a$ .

Now consider  $Q_{t_0+a}$  which forms a complete metric space. Let  $A$  be the mapping defined on  $Q_{t_0+a}$  through the relations

$$(3.10) \quad (Ax)(t) = \begin{cases} \varphi(t) & \text{for } t \in [\alpha, t_0], \\ \varphi(t_0) + \int_{t_0}^t f(s, x_s) ds + \int_{t_0}^t G(s, x_s) du(s) & \text{for } t \in (t_0, t_0 + a]. \end{cases}$$

We shall show that  $A$  maps  $Q_{t_0+a}$  into itself. Since both of the integrals on the right of (3.10) are functions of  $t$  which are of bounded variation on  $[t_0, t_0 + a]$ , it follows that  $Ax$  is a function in  $BV([\alpha, t_0 + a])_n$ . Furthermore,

$$(3.11) \quad \begin{aligned} v(Ax, [t_0, t_0 + a]) &\leq v\left(\int_{t_0}^t f(s, x_s) ds, [t_0, t_0 + a]\right) \\ &\quad + v\left(\int_{t_0}^t G(s, x_s) du(s), [t_0, t_0 + a]\right) \\ &\leq \int_{t_0}^{t_0+a} |f(s, x_s)| ds + \int_{t_0}^{t_0+a} |G(s, x_s)| dV_u(s) \\ &\leq \int_{t_0}^{t_0+a} r(s) ds + \int_{t_0}^{t_0+a} w(s) dV_u(s) && \text{(by } H_3 \text{ and } H_6) \\ &\leq b && \text{(by (3.9)).} \end{aligned}$$

Thus  $a$  maps  $\mathcal{Q}_{t_0+a}$  into itself.

We shall now show that  $A$  is a contraction. We have

$$\begin{aligned}
 \|Ax - Ay\|_{[\alpha, t_0+a]} &\leq \left\| \int_{t_0}^t [f(s, x_s) - f(s, y_s)] ds \right\|_{[t_0, t_0+a]} \\
 &\quad + \left\| \int_{t_0}^t [G(s, x_s) - G(s, y_s)] du(s) \right\|_{[t_0, t_0+a]} \\
 (3.12) \qquad &\leq \int_{t_0}^{t_0+a} |f(s, x_s) - f(s, y_s)| ds \\
 &\quad + \int_{t_0}^{t_0+a} |G(s, x_s) - G(s, y_s)| dV_u(s).
 \end{aligned}$$

We have, for every  $x \in \mathcal{Q}_{t_0+a}$ ,

$$\begin{aligned}
 \|x\|_{[\alpha, t_0+a]} &= \|\varphi\|_{[\alpha, t_0]} + v(x, [t_0, t_0+a]) \\
 &\leq \|\varphi\|_{[\alpha, t_0]} + b = k_1, \quad \text{say.}
 \end{aligned}$$

Therefore, by  $H_1$  and  $H_4$ , there exist constants  $k_2 = k_2(k_1)$ ,  $k_3 = k_3(k_1)$  such that

$$\begin{aligned}
 (3.13) \qquad |f(s, x_s) - f(s, y_s)| &\leq k_2 \|x - y\|_{I_s} \leq k_2 \|x - y\|_{[\alpha, s]}, \\
 |G(s, x_s) - G(s, y_s)| &\leq k_3 \|x - y\|_{I_s} \leq k_3 \|x - y\|_{[\alpha, s]}.
 \end{aligned}$$

From (3.12) and (3.13), we obtain

$$\begin{aligned}
 (3.14) \qquad \|Ax - Ay\|_{[\alpha, t_0+a]} &\leq \int_{t_0}^{t_0+a} k_2 \|x - y\|_{[\alpha, s]} ds + \int_{t_0}^{t_0+a} k_3 \|x - y\|_{[\alpha, s]} dV_u(s) \\
 &\leq \{ak_2 + k_3v(u, [t_0, t_0+a])\} \|x - y\|_{[\alpha, t_0+a]}.
 \end{aligned}$$

Since  $u(t)$  is right continuous,  $a$  can be chosen such that

$$ak_2 + k_3v(u, [t_0, t_0+a]) < 1,$$

and then  $A$  is, by (3.14), a contraction. Hence, by the principle of contraction mapping there is a unique fixed point. This completes the proof.

Our next theorem shows that the local existence of a solution can also be proved if the hypotheses  $H_1$  and  $H_4$  in the above theorem are replaced by the hypotheses  $H_2$  and  $H_5$ .

**THEOREM 3 (Local existence).** *Let the hypotheses  $H_2, H_3, H_5$  and  $H_6$  be satisfied. Then there exists a solution of (3.3) on an interval  $[\alpha, t_0+a]$  with a given initial function  $\varphi \in BV([\alpha, t_0], S)$ .*

*Proof.* Define  $\mathcal{Q}_t$  as in the proof of Theorem 2 and choose  $a > 0$  so as to satisfy (3.9).

Now define

$$(3.15) \qquad \bar{\varphi}(t) = \begin{cases} \varphi(t) & \text{for } t \in [\alpha, t_0], \\ \varphi(t_0) & \text{for } t \in (t_0, t_0+a] \end{cases}$$

and

$$(3.16) \quad x^{(v)}(t) = \begin{cases} \bar{\varphi}(t) & \text{for } t \in [\alpha, t_0 + a/v], \\ \varphi(t_0) + \int_{t_0+a/v}^t f\left(s - \frac{a}{v}, x_{s-a/v}^{(v)}\right) ds \\ \quad + \int_{t_0+a/v}^t G\left(s - \frac{a}{v}, x_{s-a/v}^{(v)}\right) du(s) & \text{for } t \in (t_0 + a/v, a]. \end{cases}$$

$v = 1, 2, \dots$

For any  $v \geq 2$ , the first expression in (3.16) defines  $x^{(v)}(t)$  on  $[\alpha, t_0 + a/v]$ , and then the second expression of (3.16) defines  $x^{(v)}(t)$  on  $(t_0 + a/v, t_0 + 2a/v]$ . Let us assume that  $x^{(v)}(t)$  is defined on  $[\alpha, t_0 + ja/v]$  for  $1 < j < v$ . Then the second expression of (3.16) defines  $x^{(v)}(t)$  on  $(t_0 + ja/v, t_0 + (j + 1)a/v]$ .  $x^{(v)}(t)$  is thus defined on  $[\alpha, t_0 + a]$ .

From  $H_3$ ,  $H_6$  and (3.9), we shall obtain, by calculations similar to (3.11),

$$(3.17) \quad v(x^{(v)}, [t_0, t_0 + a]) \leq b, \quad v = 1, 2, \dots$$

Thus,

$$x^{(v)} \in Q_{t_0+a}, \quad v = 1, 2, \dots,$$

and hence  $x^{(v)}(t)$  are of uniformly bounded variation and are also uniformly bounded. Hence by Helly's selection principle [2, Chap. XII, Theorem 33], there exists a subsequence  $x^{(v_j)}(t)$  and a function  $x^*(t)$  of bounded variation such that

$$(3.18) \quad \lim_{j \rightarrow \infty} x^{(v_j)}(t) = x^*(t),$$

and moreover,

$$v(x^*, [\alpha, t_0 + a]) \leq \liminf_{j \rightarrow \infty} v(x^{(v_j)}, [\alpha, t_0 + a]) \leq b$$

by (3.17). Obviously  $x^*(t) = \varphi(t)$  for  $t \in [\alpha, t]$ . Hence  $x^* \in Q_{t_0+a} \subset BV([\alpha, t_0], s)$ .

By  $H_2$ ,  $H_5$  and (3.18) we obtain

$$(3.19) \quad \lim_{j \rightarrow \infty} f(t - a/v_j, x_{t-a/v_j}^{(v_j)}) = f(t, x_t^*),$$

$$(3.20) \quad \lim_{j \rightarrow \infty} G(t - a/v_j, x_{t-a/v_j}^{(v_j)}) = G(t, x_t^*).$$

From  $H_3$  and (3.19), we obtain, by using Lebesgue's dominated convergence theorem,

$$(3.21) \quad \lim_{j \rightarrow \infty} \int_{t_0}^t f\left(s - \frac{a}{v_0}, x_{s-a/v_j}^{(v_j)}\right) ds = \int_{t_0}^t f(s, x_s^*) ds.$$

Let  $du^+(t)$  and  $du^-(t)$  be upper and lower variations (also called positive and negative variations) of the Lebesgue–Stieltjes measure  $du(t)$ . Then

$$\begin{aligned} \left| \int_{t_0}^t G(s, x_s) du^+(s) \right| &\leq \int_{t_0}^t |G(s, x_s)| du^+(s) \\ &\leq \int_{t_0}^t w(s) du^+(s) \quad (\text{by } H_6) \\ &\leq \int_{t_0}^t w(s) dV_d(s), \end{aligned}$$

since  $w(s) \geq 0$ ,  $dV_u(s) = du^+(s) + du^-(s)$ . Since  $w(s)$  is integrable with respect to  $dV_u(s)$ ,  $G(t, x_t)$  is integrable with respect to  $du^+(s)$ . Similarly  $G(t, x_t)$  is also integrable with respect to  $du^-(s)$ . Therefore,

$$(3.22) \quad \int_{t_0}^t G\left(s - \frac{a}{v_j}, x_{s-\frac{a}{v_j}}^{(v_j)}\right) du(s) = \int_{t_0}^t G\left(s - \frac{a}{v_j}, x_{s-\frac{a}{v_j}}^{(v_j)}\right) du^+(s) \\ - \int_{t_0}^t G\left(s - \frac{a}{v_j}, x_{s-\frac{a}{v_j}}^{(v_j)}\right) du^-(s)$$

(see [4, Example (7), § 29]; also see [7, p. 184]). Now take limits as  $v_j \rightarrow \infty$ . Due to  $H_6$  and (3.20), Lebesgue's dominated convergence theorem can be used on the integrals on the right-hand side of (3.22) to obtain

$$(3.23) \quad \lim_{j \rightarrow \infty} \int_{t_0}^t G\left(s - \frac{a}{v_j}, x_{s-\frac{a}{v_j}}^{(v_j)}\right) du(s) = \int_{t_0}^t G(s, x_s) du^+(s) - \int_{t_0}^t G(s, x_s) du^-(s) \\ = \int_{t_0}^t G(s, x_s) du(s).$$

Also

$$\left| \int_{t_0}^{t_0 + a/v_j} G\left(s - \frac{a}{v_j}, x_{s-\frac{a}{v_j}}^{(v_j)}\right) du(s) \right| \leq \int_{t_0}^{t_0 + a/v_j} w(s) dV_u(s) \quad (\text{by } H_6) \\ \rightarrow 0 \quad \text{as } j \rightarrow \infty,$$

since  $\int_{t_0}^t w(s) dV_u(s)$  is a right continuous function of  $t$ . So

$$(3.24) \quad \lim_{j \rightarrow \infty} \int_{t_0}^{t_0 + a/v_j} G\left(s - \frac{a}{v_j}, x_{s-\frac{a}{v_j}}^{(v_j)}\right) du(s) = 0.$$

Similarly,

$$(3.25) \quad \lim_{j \rightarrow \infty} \int_{t_0}^{t_0 + a/v_j} f\left(s - \frac{a}{v_j}, x_{s-\frac{a}{v_j}}^{(v_j)}\right) ds = 0.$$

Now, by (3.16), we have

$$x^{(v_j)}(t) = \begin{cases} \bar{\varphi}(t) & \text{for } t \in [\alpha, t_0 + a/v], \\ \varphi(t_0) + \int_{t_0}^t f\left(s - \frac{a}{v_j}, x_{s-\frac{a}{v_j}}^{(v_j)}\right) ds - \int_{t_0}^{t_0 + a/v_0} f\left(s - \frac{a}{v_j}, x_{s-\frac{a}{v_j}}^{(v_j)}\right) ds \\ \quad + \int_{t_0}^t G\left(s - \frac{a}{v_j}, x_{s-\frac{a}{v_j}}^{(v_j)}\right) du(s) - \int_{t_0}^{t_0 + a/v_j} G\left(s - \frac{a}{v_j}, x_{s-\frac{a}{v_j}}^{(v_j)}\right) du(s) & \text{for } t \in (t_0 + a/v, t_0 + a]. \end{cases}$$



Taking limits as  $j \rightarrow \infty$  and making use of (3.18), (3.19), (3.20), (3.24) and (3.25), we obtain

$$x^*(t) = \begin{cases} \varphi(t) & \text{for } t \in [\alpha, t_0], \\ \varphi(t_0) + \int_{t_0}^t f(s, x_s^*) ds + \int_{t_0}^t G(s, x_s^*) du(s) & \text{for } t \in (t_0, t_0 + a]. \end{cases}$$

$x^*(t)$  is thus the solution of (3.3) on the interval  $[\alpha, t_0 + a]$  with initial function  $\varphi(t)$ . This completes the proof.

**THEOREM 4 (Extended existence).** *Let the hypotheses  $H_2, H_3, H_5$  and  $H_6$  be satisfied on  $[t_0, \beta]$ , and let  $\varphi \in BV([\alpha, t_0], S)$  be given. If there exists a solution  $x(t, t_0, \varphi)$  of (3.3) on  $[\alpha, T)$ , where  $T$  is a point of continuity of  $u(t)$ , and if  $T < \beta$  and  $T$  cannot be increased, then for any compact set  $F \subset S$  there exists a sequence  $t_0 < t_1 < t_2 \cdots < t_k < \cdots \rightarrow T$  such that*

$$x(t_k) \in S - F \quad \text{for } k = 1, 2, \dots$$

*Proof.* Let  $x(t, t_0, \varphi)$  be a solution of (3.3) on  $[\alpha, T)$  with initial function  $\varphi$ , where  $T \in (t_0, \beta)$  is a point of continuity of  $u(t)$ . Suppose that there exists a compact set  $F \subset S$  such that  $x(t) \in F$  for  $t \in [\alpha, T]$ . We have

$$x(t) = \begin{cases} \varphi(t) & \text{for } t \in [\alpha, t_0], \\ \varphi(t_0) + \int_{t_0}^t f(s, x_s) ds + \int_{t_0}^t G(s, x_s) du(s) & \text{for } t \in (t_0, T). \end{cases}$$

If  $t_0 \leq t_1 < t_2 < T$ , then

$$\begin{aligned} (3.26) \quad |x(t_1) - x(t_2)| &\leq \int_{t_1}^{t_2} |f(s, x_s)| ds + \int_{t_1}^{t_2} |G(s, x_s)| dV_u(s) \\ &\leq \int_{t_1}^{t_2} r(s) ds + \int_{t_1}^{t_2} w(s) dV_u(s). \end{aligned}$$

Since  $u(t)$  is continuous at the point  $T$ , the function  $\int_{t_0}^t w(s) dV_u(s)$  is also continuous at  $T$ . So if  $\varepsilon > 0$ , we can choose  $t_1$  so close to  $T$  (but  $t_1 < T$ ) that

$$\int_{t_1}^T w(s) dV_u(s) < \frac{\varepsilon}{2},$$

and then we choose  $t_2$  so close to  $t_1$  that

$$\int_{t_1}^{t_2} r(s) ds < \frac{\varepsilon}{2}.$$

Thus, we obtain, from (3.26),

$$|x(t_1) - x(t_2)| < \varepsilon$$

for all  $t_0 \leq t_1 < t_2 < T$  such that  $t_1$  is sufficiently close to  $T$ . Therefore, by Cauchy's criterion,

$$x(T-) \equiv \lim_{t \rightarrow T-} x(t)$$

exists. Define

$$x(T) = x(T-).$$

Let

$$\varphi^*(t) = x(t) \quad \text{for } t \in [\alpha, T].$$

Obviously,  $\varphi^*(t) \in F \subset S$ . Replacing  $t_0$  by  $T$ ,  $\varphi$  by  $\varphi^*$ , we can prove, as in Theorem 3, that there exists a solution  $x^*(t, T, \varphi^*)$  of (3.3) on  $[\alpha, T + \delta]$ , where  $\delta > 0$ , which is represented by

$$x^*(t) = \begin{cases} \varphi^*(t) & \text{for } t \in [\alpha, T], \\ \varphi^*(T) + \int_T^t f(s, x_s) ds + \int_T^t G(s, x_s) du(s) & \text{for } t \in (T, T + \delta). \end{cases}$$

Obviously  $x^*(t)$  is also a solution of (3.3) with the initial condition  $x^*(t) = \varphi(t)$  for  $t \in [\alpha, t_0]$ ; i.e.,  $x^*(t)$  is the extension of  $x(t)$  over the interval  $[\alpha, T + \delta]$ . Thus, the value of  $T$  can be increased. This proves the theorem.

**4. Formulation of an optimal control problem.** Let  $[t_0, \beta]$  be a fixed interval and  $S$  a domain of  $E^n$ . Let  $B = \bigcup_{t \in [t_0, \beta]} BV([\alpha, t], S)$ . We shall denote by  $U$  the set of all right continuous functions  $u$  of bounded variation on  $[t_0, \beta]$  into a non-empty compact subset  $Q$  of  $E^m$ .

Consider a control process governed by the measure delay-differential equation

$$(4.1) \quad Dx = f(t, x_t, u(t)) + G(t, x_t, u(t)) Du, \quad t > t_0,$$

satisfying the following assumptions:

A<sub>1</sub>. The functional  $f$  with range in  $E^n$  is defined for all  $t \in [t_0, \beta]$ , for all  $x \in B$  and for all  $u \in U$ .

A<sub>2</sub>.  $f(t, x_t, u(t))$  is continuous in  $t, x$  and  $u$ .

A<sub>3</sub>. There exists a Lebesgue integrable real function  $r(t)$  for  $t \in [t_0, \beta]$  such that

$$|f(t, x_t, u(t))| \leq r(t)$$

uniformly with respect to  $x \in B, u \in U$ .

A<sub>4</sub>.  $G$  is an  $n \times m$  matrix defined for  $t \in [t_0, \beta]$ , for all  $x \in B$  and for all  $u \in U$ .

A<sub>5</sub>.  $G(t, x_t, u(t))$  is continuous in  $t, x$  and  $u$ .

A<sub>6</sub>. There exists a constant  $K$  such that

$$|G(t, x_t, u(t))| \leq K$$

uniformly with respect to  $x \in B, u \in U$ .

Let  $\varphi \in BV([\alpha, t_0], S)$  be given. Under the assumptions A<sub>1</sub> to A<sub>6</sub>, (4.1) has a bounded variation solution  $x(t)$ , for each choice of  $u \in U$ , satisfying the initial condition  $x(t) = \varphi(t)$  for  $t \in [\alpha, t_0]$ . This follows from Theorem 3. If  $u \in U$  defined on the interval  $[t_0, t_1]$  is such that the corresponding solution  $x(t, t_0, \varphi)$  of (4.1) is also defined on  $[t_0, t_1]$ , then  $u(t)$  is called an admissible control and  $x(t)$  is called

the corresponding trajectory. We define the cost of the control  $u(t)$  by

$$J(u) = \int_{t_0}^{t_1} f^0(t, x(t), u(t)) dt,$$

where  $f^0$  is a given continuous real function defined on  $[t_0, \beta] \times S \times Q$ . The functional  $J$  is called the cost functional of the system (4.1).

A target set is a family  $\mathcal{T}$  of nonempty compact sets  $T_t \subset E^n$  defined for  $t \in [t_0, \beta]$ . We say that an admissible control  $u(t)$  defined on  $[t_0, t_1]$  transfers the function  $\varphi \in BV([\alpha, t_0], S)$  to the target set  $\mathcal{T}$ , if the trajectory  $x(t) = x(t, t_0, \varphi, u)$ , corresponding to the initial function  $\varphi(t)$  and the control  $u(t)$ , satisfies the relation  $x(t_1) \in T_{t_1}$ . Let  $A$  be the set of all admissible controls which transfer  $\varphi$  to  $\mathcal{T}$ . A control  $u^* \in A$  is called an optimal (minimal) control if

$$J(u^*) \leq J(u) \quad \text{for each } u \in A.$$

The optimization problem consists in finding such an optimal control.

### 5. Existence of an optimal control.

**THEOREM 5.** *We are given the control problem with the following data:*

(a) 
$$Dx = f(t, x_t, u(t)) + G(t, x_t, u(t))Du, \quad t \in [t_0, \beta],$$

with  $f$  and  $G$  satisfying the assumptions  $A_1$  to  $A_6$  of § 4;

(b) the fixed initial function  $\varphi \in BV([\alpha, t_0], S)$ ;

(c) a target set  $\mathcal{T}$  of nonempty compact sets  $T_t \subset E^n$  defined on  $[t_0, \beta]$  and upper semicontinuous with respect to inclusion;

(d) the set  $A$  of admissible controls  $u(t)$  defined on subintervals  $[t_0, t_1]$  contained in  $[t_0, \beta]$  with the same left endpoint (and perhaps different right endpoint  $t_1 > t_0$ ) which transfer  $\varphi$  to  $\mathcal{T}$ , which is such that for all  $u \in A$ ,

$$|\Delta u| \leq \Delta h$$

on each subinterval of  $[t_0, t_1]$ , where  $h$  is a given nondecreasing right continuous function defined on  $[t_0, t_1]$ ; (the symbol  $\Delta h$  on the interval, say,  $[t_0, t_1]$  denotes  $h(t_1) - h(t_0)$ );

(e) the cost functional

$$J(u) = \int_{t_0}^{t_1} f^0(t, x(t), u(t)) dt,$$

where  $f^0$  is a continuous real function defined on  $[t_0, \beta] \times S \times Q$ .

Then if  $A$  is nonempty, there exists an optimal control in  $A$ .

*Proof.* Since for each  $u \in A$ ,  $|\Delta u| \leq \Delta h$  on each subinterval of  $[t_0, t_1]$ , therefore  $u \in A$  are of uniform bounded variation. We shall show that the corresponding trajectories  $x(t)$  are also of uniform bounded variation.

We have

$$(5.1) \quad x(t) = \begin{cases} \varphi(t) & \text{for } t \in [\alpha, t_0], \\ \varphi(t_0) + \int_{t_0}^t f(s, x_s, u(s)) ds + \int_{t_0}^t G(s, x_s, u(s)) du(s) & \text{for } t \in (t_0, t_1]. \end{cases}$$

As in (3.11),

$$\begin{aligned}
 v(x, [t_0, t_1]) &\leq \int_{t_0}^{t_1} |f(s, x_s)| ds + \int_{t_0}^{t_1} |G(s, x_s)| dV_u(s) \\
 (5.2) \qquad \qquad &\leq \int_{t_0}^{t_1} r(s) ds + Kv(u, [t_0, t_1]) \\
 &\leq \int_{t_0}^{t_1} r(s) ds + K(h(\beta) - h(t_0)),
 \end{aligned}$$

which shows that  $v(x, [t_0, t_1])$  are uniformly bounded for all the trajectories  $x(t)$  corresponding to controls in  $A$ . Hence the trajectories  $x(t)$  are also uniformly bounded.

Now, since  $A$  is nonempty, and the corresponding trajectories are uniformly bounded,  $\inf J(u) = \tilde{m} > -\infty$  for all  $u \in A$ . Either  $A$  is a finite set, in which case the theorem is trivially true, or we can select from  $A$  a sequence of controls  $u^{(k)}$  on the intervals  $[t_0, t_1^{(k)}]$  for which  $J(u^{(k)})$  decreases monotonically to  $\tilde{m}$ . Select a subsequence (which we shall still denote by  $u^{(k)}$ ) such that  $t_1^{(k)} \rightarrow t_1^*$  monotonically. Consider the case when  $\{t_1^{(k)}\}$  is monotonic decreasing (the case when it is increasing will be considered later). Next choose  $\hat{t}_1$  such that  $\hat{t}_1 = t_1$  if  $t_1^{(k)} = t_1$  for all  $k$ ; otherwise let  $t_1^* < t_1^{(k_0+1)} \leq \hat{t}_1 < t_1^{(k_0)} < \beta$  for some  $k_0$ . (From now on all reference to the index  $k$  tacitly assumes  $k > k_0$ .) We define the extended control  $\hat{u}^{(k)}$  on  $[t_0, \hat{t}_1]$  by

$$(5.3) \qquad \hat{u}^{(k)}(t) = \begin{cases} u^{(k)}(t) & \text{if } t \in [t_0, t_1^{(k)}], \\ u^{(k)}(t_1^{(k)}) & \text{if } t \in (t_1^{(k)}, \hat{t}_1]. \end{cases}$$

Since  $u^{(k)}(t) \in Q$  for all  $t \in [t_0, t_1^{(k)}]$ ,  $\hat{u}^{(k)}(t)$  also belongs to  $Q$  for  $t \in [t_0, \hat{t}_1]$ . Evidently  $\hat{u}^{(k)}$  is right continuous on  $[t_0, \hat{t}_1]$ , and  $|\Delta \hat{u}^{(k)}| \leq \Delta h$  on every subinterval of  $[t_0, \hat{t}_1]$ . Now, by Helly's principle of choice, there exists a subsequence (for which we shall use the same notation  $\hat{u}^{(k)}$ ) and a function of bounded variation  $u^*$  such that

$$(5.4) \qquad \lim_{k \rightarrow \infty} \hat{u}^{(k)}(t) = u^*(t)$$

everywhere on  $[t_0, \hat{t}_1]$ . Our aim is to show that  $u^* \in A$  and is an optimal control.

We shall first show that  $u^*$  is right continuous. Let  $\tau$  be any point in  $[t_0, \hat{t}_1]$ . Since  $|\Delta \hat{u}^{(k)}| \leq \Delta h$  on every subinterval of  $[t_0, \hat{t}_1]$ , we have

$$|\hat{u}^{(k)}(\tau + s) - \hat{u}^{(k)}(\tau)| \leq h(\tau + s) - h(\tau), \quad \tau \leq \tau + s \leq \hat{t}_1;$$

and since  $h$  is right continuous at  $\tau$ , it follows that given any  $\varepsilon > 0$ , there exists a  $\delta > 0$  such that for any  $k$ ,

$$|\hat{u}^{(k)}(\tau + s) - \hat{u}^{(k)}(\tau)| < \varepsilon \quad \text{for } 0 \leq s \leq \delta.$$

Therefore,

$$(5.5) \qquad \lim_{s \rightarrow 0^+} \hat{u}^{(k)}(\tau + s) = \hat{u}^{(k)}(\tau) \quad \text{uniformly for } k.$$

Also

$$\lim_{k \rightarrow \infty} \hat{u}^{(k)}(\tau + s) = u^*(\tau + s).$$

Hence, applying the Moore theorem on interchange of order of repeated limits, we obtain

$$\lim_{s \rightarrow 0+} \lim_{k \rightarrow \infty} \hat{u}^{(k)}(\tau + s) = \lim_{k \rightarrow \infty} \lim_{s \rightarrow 0+} \hat{u}^{(k)}(\tau + s);$$

i.e.,

$$u^*(\tau + 0) = u^*(\tau),$$

Thus,  $u^*$  is right continuous at  $\tau$ . But  $\tau$  was taken arbitrarily in  $[t_0, \hat{t}_1]$ . Hence  $u^*$  is right continuous at each  $t \in [t_0, \hat{t}_1]$ .

Since  $\hat{u}^{(k)}(t) \in Q$  for  $t \in [t_0, \hat{t}_1]$  and  $Q$  is compact, it follows from (5.4) that  $u^*(t)$  also belongs to  $Q$  for each  $t \in [t_0, \hat{t}_1]$ . Since  $\Delta u^* = \lim_{k \rightarrow \infty} \Delta \hat{u}^{(k)}$  on every subinterval of  $[t_0, \hat{t}_1]$ , it follows that, given any  $\varepsilon > 0$ ,

$$\begin{aligned} \Delta u^* &\leq |\Delta \hat{u}^{(k)}| + \varepsilon \\ &\leq \Delta h + \varepsilon \quad \text{for sufficiently large } k. \end{aligned}$$

Since  $\varepsilon$  is arbitrary, we get

$$(5.6) \quad \Delta u^* \leq \Delta h$$

on every subinterval of  $[t_0, \hat{t}_1]$ .

Let  $x^{(k)}$  be the trajectory defined on  $[t_0, t_1^{(k)}]$  corresponding to the control  $u^{(k)}$ . The extended control  $\hat{u}^{(k)}$  coincides with  $u^{(k)}$  on  $[t_0, t_1^{(k)}]$  and is constant and, therefore, continuous on  $[t_1^{(k)}, \hat{t}_1]$  (see (5.3)). Hence, by Theorem 4,  $x^{(k)}$  can be extended to  $[t_0, \hat{t}_1]$ . The extended trajectory  $\hat{x}^{(k)}$  is given by

$$(5.7) \quad \hat{x}^{(k)}(t) = \begin{cases} \varphi(t) & \text{for } t \in [\alpha, t_0], \\ \varphi(t_0) + \int_{t_0}^t f(s, \hat{x}_s^{(k)}, \hat{u}^{(k)}(s)) ds & \\ \quad + \int_{t_0}^t G(s, \hat{x}_s^{(k)}, \hat{u}^{(k)}(s)) d\hat{u}^{(k)}(s) & \text{for } t \in (t_0, \hat{t}_1]. \end{cases}$$

Since

$$v(\hat{u}^{(k)}, [t_0, \hat{t}_1]) = v(u^{(k)}, [t_0, t_1^{(k)}]),$$

total variations of  $\hat{x}^{(k)}$  can be seen, as in (5.2), to be uniformly bounded. The  $\hat{x}^{(k)}$  are also uniformly bounded. Hence there exists a subsequence (still labeled  $\hat{x}^{(k)}$ ) and a function  $x^*$  such that

$$(5.8) \quad \lim_{k \rightarrow \infty} \hat{x}^{(k)}(t) = x^*(t)$$

everywhere on  $[t_0, \hat{t}_1]$ . By selecting the corresponding subsequence from  $\hat{u}^{(k)}$  we do not change any of the preceding limiting operations satisfied by  $\hat{u}^{(k)}$ .

Since  $\hat{u}^{(k)}(t) \rightarrow u^*(t)$ ,  $\hat{x}^{(k)}(t) \rightarrow x^*(t)$  everywhere on  $[t_0, \hat{t}_1]$ , and the assumptions  $A_2$  and  $A_3$  of § 4 hold, we obtain, by Lebesgue's dominated convergence theorem,

$$(5.9) \quad \lim_{k \rightarrow \infty} \int_{t_0}^t f(s, \hat{x}_s^{(k)}, \hat{u}^{(k)}(s)) ds = \int_{t_0}^t f(s, x_s^*, u^*(s)) ds$$

for all  $t \in [t_0, \hat{t}_1]$

Now,  $|\Delta \hat{u}^{(k)}| \leq \Delta h, |\Delta u^*| \leq \Delta h$  on every subinterval of  $[t_0, \hat{t}_1]$ ; and  $\hat{u}^{(k)}(t) \rightarrow u^*(t)$  everywhere on  $[t_0, \hat{t}_1]$ . For each  $k$ ,  $G(t, \hat{x}_t^{(k)}, \hat{u}^{(k)}(t))$  is a continuous function of  $t$  (by  $A_5$  of § 4), and therefore each element  $g_j^i(t, \hat{x}_t^{(k)}, \hat{u}^{(k)}(t))$  of the matrix function  $G(t, \hat{x}_t^{(k)}, \hat{u}^{(k)}(t))$  is integrable with respect to  $dh(t)$ . Moreover,

$$\begin{aligned} \left| \int_{t_1}^{t_2} g_j^i(t, \hat{x}_t^{(k)}, \hat{u}^{(k)}(t)) dh(t) \right| &\leq \int_{t_1}^{t_2} |g_j^i(t, \hat{x}_t^{(k)}, \hat{u}^{(k)}(t))| dh(t) \\ &\leq \int_{t_1}^{t_2} |G(t, \hat{x}_t^{(k)}, \hat{u}^{(k)}(t))| dh(t) \leq \int_{t_1}^{t_2} K dh(t) \quad (\text{by } A_6) \\ &= K(h(t_2) - h(t_1)), \end{aligned}$$

and therefore the integrals  $\int g_j^i(t, \hat{x}_t^{(k)}, \hat{u}^{(k)}(t)) dh(t)$  are absolutely continuous with respect to  $h(t)$  uniformly in  $k$ , and bounded uniformly. Hence by [2, Theorem 27, p. 285] and since  $A_5$  holds, we obtain

$$(5.10) \quad \lim_{k \rightarrow \infty} \int_{t_0}^t G(s, \hat{x}_s^{(k)}, \hat{u}^{(k)}(s)) d\hat{u}^{(k)}(s) = \int_{t_0}^t G(s, x_s^*, u^*(s)) du^*(s)$$

for all  $t \in [t_0, \hat{t}_1]$ . By (5.7), (5.8), (5.9), (5.10), we obtain

$$(5.11) \quad x^*(t) = \begin{cases} \varphi(t) & \text{for } t \in [\alpha, t_0], \\ \varphi(t_0) + \int_{t_0}^t f(s, x_s^*, u(s)) ds \\ \quad + \int_{t_0}^t G(s, x_s^*, u^*(s)) du^*(s) & \text{for } t \in (t_0, \hat{t}_1]. \end{cases}$$

We shall now show that

$$(5.12) \quad \lim_{k \rightarrow \infty} \hat{x}^{(k)}(t_1^{(k)}) = x^*(t_1^*).$$

We have

$$(5.13) \quad |x^{(k)}(t_1^{(k)}) - x^*(t_1^*)| \leq |x^{(k)}(t_1^{(k)}) - x^*(t_1^{(k)})| + |x^*(t_1^{(k)}) - x^*(t_1^*)|.$$

Since  $x^*$  is continuous from the right (because  $u^*$  is),  $x^*(t_1^{(k)}) \rightarrow x^*(t_1^*)$ ; and hence the last term on the right in (5.13) is zero as  $k \rightarrow \infty$ . We have

$$\begin{aligned} x^{(k)}(t_1^{(k)}) - x^*(t_1^{(k)}) &= \int_{t_0}^{t_1^{(k)}} f(s, x_s^{(k)}, u^{(k)}(s)) ds + \int_{t_0}^{t_1^{(k)}} G(s, x_s^{(k)}, u^{(k)}(s)) du^{(k)}(s) \\ &\quad - \int_{t_0}^{t_1^{(k)}} f(s, x_s^*, u^*(s)) ds - \int_{t_0}^{t_1^{(k)}} G(s, x_s^*, u^*(s)) du^*(s) \end{aligned}$$

$$\begin{aligned}
&= \int_{t_0}^{t_1^*} [f(s, x_s^{(k)}, u^{(k)}(s)) - f(s, x_s^*, u^*(s))] ds \\
&\quad + \left[ \int_{t_0}^{t_1^*} G(s, x_s^{(k)}, u^{(k)}(s)) du^{(k)}(s) - \int_{t_0}^{t_1^*} G(s, x_s^*, u^*(s)) du^*(s) \right] \\
&\quad + \int_{t_1^*}^{t_1^{(k)}} [f(s, x_s^{(k)}, u^{(k)}(s)) - f(s, x_s^*, u^*(s))] ds \\
&\quad + \left[ \int_{t_1^*}^{t_1^{(k)}} G(s, x_s^{(k)}, u^{(k)}(s)) du^{(k)}(s) - \int_{t_1^*}^{t_1^{(k)}} G(s, x_s^*, u^*(s)) du^*(s) \right] \\
&= I_1 + I_2 + I_3 + I_4, \quad \text{say.}
\end{aligned}$$

By (5.9) and (5.10), we have

$$I_1 \rightarrow 0, \quad I_2 \rightarrow 0 \quad \text{as } k \rightarrow \infty.$$

We have, further,

$$\begin{aligned}
|I_3| &\leq \int_{t_1^*}^{t_1^{(k)}} |f(s, x_s^{(k)}, u^{(k)}(s))| ds + \int_{t_1^*}^{t_1^{(k)}} |f(s, x_s^*, u^*(s))| ds \\
&\leq 2 \int_{t_1^*}^{t_1^{(k)}} r(s) ds \rightarrow 0 \quad \text{as } k \rightarrow \infty, \\
|I_4| &\leq \int_{t_1^*}^{t_1^{(k)}} |G(s, x_s^{(k)}, u^{(k)}(s))| dh(s) + \int_{t_1^*}^{t_1^{(k)}} |G(s, x_s^*, u^*(s))| dh(s) \\
&\leq 2K[h(t_1^{(k)}) - h(t_1^*)] \rightarrow 0 \quad \text{as } k \rightarrow \infty,
\end{aligned}$$

since  $h$  is right continuous. Thus (5.12) is established.

Now,  $x^{(k)}(t_1^{(k)}) \in T_{t_1^{(k)}}$  for  $k = 1, 2, \dots$ , and

$$x^*(t_1^*) = \lim_{k \rightarrow \infty} x^{(k)}(t_1^{(k)}).$$

If  $x^*(t_1^*)$  were not in  $T_{t_1^*}$ , then there would exist a neighborhood  $N$  of the compact set  $T_{t_1^*}$  so that  $x^*(t_1^*)$  is not in the closure  $\bar{N}$  of  $N$ . But, since  $T_t$  is upper semicontinuous,  $T_t \subset N$  for  $t$  sufficiently near  $t_1^*$ . Thus  $x^{(k)}(t_1^{(k)}) \in N$  for large  $k$  and yet  $x^*(t_1^*) \notin \bar{N}$ . This contradiction shows that  $x^*(t_1^*) \in T_{t_1^*}$ . Hence the control  $u^*$  on  $[t_0, t_1^*]$  belongs to  $A$ .

We have

$$\begin{aligned}
J(u^{(k)}) &= \int_{t_0}^{t_1^{(k)}} f^0(s, x^{(k)}(s), u^{(k)}(s)) ds \\
&= \int_{t_0}^{t_1^*} f^0(s, x^{(k)}(s), u^{(k)}(s)) ds + \int_{t_1^*}^{t_1^{(k)}} f^0(s, x^{(k)}(s), u^{(k)}(s)) ds.
\end{aligned}$$

Since  $f^0(s, x^{(k)}(s), u^{(k)}(s))$  is uniformly bounded on  $[t_0, t_1^{(k)}]$  and so the last integral in the above equation approaches zero, and since

$$x^{(k)}(t) \rightarrow x^*(t), \quad u^{(k)}(t) \rightarrow u^*(t)$$

everywhere on  $[t_0, t_1^*]$ , therefore, by applying Legesgue's dominated convergence theorem to the first integral, we obtain

$$\lim_{k \rightarrow \infty} J(u^{(k)}) = \int_{t_0}^{t_1^*} f^0(s, x^*(s), u^*(s)) ds = J(u^*).$$

But  $\lim_{k \rightarrow \infty} J(u^{(k)}) = \tilde{m}$ . Therefore,  $J(u^*) = \tilde{m}$ . Hence  $u^*$  on  $[t_0, t_1^*]$  is an optimal control.

We shall now consider the case when  $\{t_1^{(k)}\}$  is monotonic increasing. We extend all the controls  $u^{(k)}$  to the interval  $[t_0, \hat{t}_1]$ , where  $\hat{t}_1 = t_1^* + \delta$  for appropriately small  $\delta > 0$ , by defining

$$(5.14) \quad \hat{u}^{(k)}(t) = \begin{cases} u^{(k)}(t) & \text{for } t \in [t_0, t_1^{(k)}], \\ u^{(k)}(t_1^{(k)}) & \text{for } t \in (t_1^{(k)}, \hat{t}_1]. \end{cases}$$

As before, there exists a subsequence  $\hat{u}^{(k)}$  (without changing the notation) and a right continuous function of bounded variation  $u^*$  such that everywhere on  $[t_0, \hat{t}_1]$ ,

$$(5.15) \quad \lim_{k \rightarrow \infty} \hat{u}^{(k)}(t) = u^*(t)$$

and

$$(5.16) \quad \lim_{k \rightarrow \infty} \hat{x}^{(k)}(t) = x^*(t),$$

where  $\hat{x}^{(k)}$  and  $x^*$  are the trajectories corresponding to  $\hat{u}^{(k)}$  and  $u^*$  respectively.

We shall show that

$$(5.17) \quad \lim_{k \rightarrow \infty} x^{(k)}(t_1^{(k)}) = x^*(t_1^*).$$

Consider

$$(5.18) \quad |x^{(k)}(t_1^{(k)}) - x^*(t_1^*)| \leq |x^{(k)}(t_1^{(k)}) - \hat{x}^{(k)}(t_1^*)| + |\hat{x}^{(k)}(t_1^*) - x^*(t_1^*)|.$$

Since  $\hat{x}^{(k)}(t_1^*) \rightarrow x^*(t_1^*)$ , the last term on the right in (5.18) approaches zero as  $k \rightarrow \infty$ . For the first term on the right, we write

$$\begin{aligned} x^{(k)}(t_1^{(k)}) - \hat{x}^{(k)}(t_1^*) &= \int_{t_0}^{t_1^{(k)}} f(s, x_s^{(k)}, \hat{u}^{(k)}(s)) ds + \int_{t_0}^{t_1^{(k)}} G(s, x_s^{(k)}, \hat{u}^{(k)}(s)) d\hat{u}^{(k)}(s) \\ &\quad - \int_{t_0}^{t_1^*} f(s, \hat{x}_s^{(k)}, \hat{u}^{(k)}(s)) ds - \int_{t_0}^{t_1^{(k)}} G(s, \hat{x}_s^{(k)}, \hat{u}^{(k)}(s)) d\hat{u}^{(k)}(s) \\ &= - \int_{t_1^{(k)}}^{t_1^*} f(s, \hat{x}_s^{(k)}, \hat{u}^{(k)}(s)) ds - \int_{t_1^{(k)}}^{t_1^*} G(s, \hat{x}_s^{(k)}, \hat{u}^{(k)}(s)) d\hat{u}^{(k)}(s) \\ &= I_1 + I_2, \quad \text{say.} \end{aligned}$$



Now,

$$\begin{aligned} |I_1| &\leq \int_{t_1^{(k)}}^{t_1^*} |f(s, \hat{x}_s^{(k)}, \hat{u}^{(k)}(s))| ds \\ &\leq \int_{t_1^{(k)}}^{t_1^*} r(s) ds \rightarrow 0 \quad \text{as } k \rightarrow \infty \end{aligned}$$

and

$$I_2 = - \int_{t_1^{(k)}}^{t_1^*} G(s, \hat{x}_s^{(k)}, \hat{u}^{(k)}(s)) d\hat{u}^{(k)}(s) = 0$$

for each  $k$ , since  $\hat{u}^{(k)}(t) = u^{(k)}(t_1^{(k)})$  for  $t \in [t_1^{(k)}, t_1^*]$  and is thus constant on  $[t_1^{(k)}, t_1^*]$ . Thus (5.17) is established. Exactly as before,  $x^*(t_1^*) \in T_{t_1^*}$ ; and therefore  $u^*$  on  $[t_0, t_1^*]$  belongs to  $A$ . It can be seen, as before, that  $J(u^*) = \bar{m}$ . Hence  $u^*$  is an optimal control. This completes the proof.

**Acknowledgment.** The authors are very thankful to the referee for his comments which helped them to improve this paper.

#### REFERENCES

- [1] N. DUNFORD AND J. SCHWARTZ, *Linear Operators. Part I*, Interscience, New York, 1964.
- [2] L. M. GRAVES, *The Theory of Functions of Real Variables*, McGraw-Hill, New York, 1956.
- [3] A. HALANAY, *Differential Equations: Stability, Oscillations, Time Lags*, Academic Press, New York, 1965.
- [4] P. R. HALMOS, *Measure Theory*, East-West ed., Van Nostrand, New Delhi, 1964.
- [5] I. HALPERIN, *Introduction to the Theory of Distributions*, University of Toronto Press, Toronto, 1952.
- [6] E. J. MCSHANE, *Integration*, Princeton University Press, Princeton, 1944.
- [7] M. E. MUNROE, *Introduction to Measure and Integration*, Addison-Wesley, Reading, Mass., 1959.
- [8] I. P. NATANSON, *Theory of Functions of a Real Variable*, vol. I, Frederick Ungar, New York, 1955.
- [9] M. N. OGUZTORELI, *Time-Lag Control Systems*, Academic Press, New York, 1966.
- [10] T. PAVLIDIS, *Optimal control of pulse frequency modulated systems*, IEEE Trans. Automatic Control, AC-11 (1966), pp. 676–684.
- [11] ———, *Stability of systems described by differential equations containing impulses*, Ibid., AC-12 (1967), pp. 43–45.
- [12] W. RUDIN, *Real and Complex Analysis*, McGraw-Hill, New York, 1966.
- [13] W. W. SCHMAEDEKE, *Optimal control theory for nonlinear vector differential equations containing measures*, this Journal, 3 (1965), pp. 331–380.

## NECESSARY AND SUFFICIENT OPTIMALITY CONDITIONS FOR A CLASS OF DISTRIBUTED PARAMETER CONTROL SYSTEMS\*

EARL R. BARNES†

**Abstract.** A class of optimal control problems arising in systems with distributed parameters is considered. Pontryagin's maximum principle is shown to be a necessary condition for optimality. Under certain convexity assumptions, it is shown that the maximum principle gives a sufficient condition for optimality.

The maximum principle of L. S. Pontryagin [4] gives a necessary condition for optimality for a large class of optimal control problems involving systems with lumped parameters. E. B. Lee has shown (cf. [1]) that the maximum principle is also a sufficient condition for optimality for a certain subclass of such problems. In this paper we shall derive analogous necessary and sufficient conditions for optimality in the form of a maximum principle for certain optimal control problems arising in systems with distributed parameters.

The paper is divided into two sections. The first section is concerned with the optimal control of systems (such as vibrating strings) that can be described by second order linear hyperbolic partial differential equations. In the second section, we discuss the problem of optimally controlling vibrating beams. The results obtained here are analogous to those obtained in Section 1 except now we restrict ourselves to problems in one space dimension.

In all cases we have tried to formulate problems that generalize certain control problems that have been studied for vibrating systems having finitely many degrees of freedom and lumped parameters. Other problems for vibrating strings and beams have been considered by Russell and Komkov in [2] and [3], respectively. These problems involve systems having quadratic cost functionals. A fairly extensive treatment of similar systems has been given by Lions in [11].

Necessary optimality conditions similar to ours have been obtained by Egorov in [5], [6] and [7] for a different class of distributed parameter control problems. However, our approach is different than Egorov's and we obtain certain advantages over his method of deriving necessary conditions. In the first place, our proofs are shorter and simpler. But more importantly, our technique of proof illustrates that Egorov's completeness assumption on the class of admissible controls (cf. [6]) can be dropped. This is desirable since there is apparently no technique available for determining when a given class of admissible controls is complete in the sense of Egorov.

---

\* Received by the editors August 12, 1969, and in revised form June 1, 1970.

† IBM Watson Research Center, Yorktown Heights, New York 10598. A portion of this paper was taken from the author's doctoral dissertation, written under the directorship of Professor G. S. Jones at the Institute for Fluid Dynamics and Applied Mathematics, University of Maryland, 1968. The author's graduate research was supported in part by an NDEA Fellowship and in part by a Minta Martin Fellowship.

## 1. The optimal control of a vibrating string.

**1.1. Introduction.** The problem of controlling vibrating systems is of considerable importance to engineers. Many vibrating systems can be approximated by lumped parameter systems which can in turn be described by a set of second order linear ordinary differential equations. Such systems have received considerable attention in the literature. See, for example, [8], [9] and [10, pp. 444–456].

The physical model for a vibrating system having lumped parameters is usually taken to be a system consisting of a finite number of springs and masses in combination. The control functions are then taken to represent forces that can be applied to the masses so as to bring the system to an equilibrium position in some optimal fashion. In this section we shall discuss the optimal control of the simplest vibrating systems having distributed parameters. Such systems can be described by one-dimensional wave equations. We assume that a physical model for our system can be taken to be an inhomogeneous string stretched between two fixed points and undergoing small planar vibrations. We also assume that we are able to apply a distributed force which acts in the plane of vibration and normal to the string, along the entire length of the string, in an effort to control the vibrations. Our aim is to determine a force which eliminates the vibrations in some optimal fashion.

Assume that the string is stretched between the points  $x = 0$  and  $x = l$  on the  $x$ -axis and is vibrating in a vertical plane. Denote the displacement at the point  $x$  and time  $t$  by  $u(x, t)$  and denote the externally applied force by  $f$ . Then, assuming that the frictional forces acting on the string are negligible,  $u$  and  $f$  are related by the differential equation

$$(1.1) \quad \mu(x) \frac{\partial^2 u}{\partial t^2} = \frac{\partial}{\partial x} \left( p(x) \frac{\partial u}{\partial x} \right) - q(x)u + f(x, t),$$

$$0 \leq x \leq l, \quad t > 0.$$

$\mu$  is the variable density of the string,  $p$  is the tension, and  $q$  the elastic restoring force.  $u$  satisfies boundary and initial conditions of the form

$$(1.2) \quad \begin{aligned} u(0, t) &= 0, & u(l, t) &= 0, \\ u(x, 0) &= \varphi(x), & u_t(x, 0) &= \psi(x). \end{aligned}$$

The functions  $\mu$ ,  $p$  and  $q$  in (1.1) are assumed to satisfy  $\mu, p \in C^3$ ,  $q \in C^1$ , and  $\mu(x), q(x), p(x) > 0$  for  $0 \leq x \leq l$ . The initial conditions  $\varphi$  and  $\psi$  are elements of  $L_2[0, l]$ . The function  $f$  denotes a control parameter.  $f$  is assumed to be admissible in the sense of the following definition.

**DEFINITION 1.1.** Let  $F \subset E_1$  be a given constraint set. A square integrable function  $f$  defined on  $[0, l] \times [0, T]$  and taking values in  $F$  is called an *admissible control*.

For a given admissible control  $f$ , let  $\{f_n\}$  be a sequence of continuous functions converging to  $f$  in the mean square. Similarly, let  $\{\varphi_n\}$  and  $\{\psi_n\}$  be sequences of  $C^2$  and  $C^3$  functions, respectively, converging in the mean square to  $\varphi$  and  $\psi$ . Then it is well known that there exist twice continuously differentiable functions  $u_n$  satisfying (1.1)–(1.2) with  $f$ ,  $\varphi$  and  $\psi$  replaced by  $f_n$ ,  $\varphi_n$  and  $\psi_n$ , respectively (cf. [12, p. 157]).

Moreover, each  $u_n$  satisfies the energy inequality

$$(1.3) \quad \int_0^l \left[ \mu(x) \left( \frac{\partial u_n}{\partial t}(x, t) \right)^2 + p(x) \left( \frac{\partial u_n}{\partial x}(x, t) \right)^2 + q(x) u_n^2(x, t) \right] dx \\ \leq c \int_0^t \int_0^l f_n^2(x, \tau) dx d\tau \\ + c \int_0^l \left[ \mu(x) \psi_n^2(x) + p(x) \left( \frac{\partial \varphi_n}{\partial x}(x) \right)^2 + q(x) \varphi_n^2(x) \right] dx.$$

$c$  is a constant, independent of the functions  $u_n$ .

If we apply the energy inequality to the differences  $u_k - u_n$ , we see that the limits

$$\lim_{k \rightarrow \infty} \frac{\partial u_k}{\partial t}(x, t), \quad \lim_{k \rightarrow \infty} \frac{\partial u_k}{\partial x}(x, t), \quad \lim_{k \rightarrow \infty} u_k(x, t)$$

exist in  $L_2[0, l]$  uniformly in  $t$  for  $0 \leq t \leq T$ . We shall denote these limits by

$$\frac{\partial u}{\partial t}(x, t), \quad \frac{\partial u}{\partial x}(x, t) \quad \text{and} \quad u(x, t),$$

respectively, and shall refer to this convergence as convergence in the energy norm.

The function  $u$  will be called a solution of (1.1)–(1.2) having finite energy.  $\partial u/\partial t$  and  $\partial u/\partial x$  will be called generalized derivatives of  $u$ .

In the sequel we shall obtain certain results by formally integrating by parts. The integrations will be valid for classical solutions, and hence, the results obtained will be valid for solutions with finite energy, by passing to the limit in the result for classical solutions.

**1.2. Formulation of the control problem.** The system (1.1)–(1.2) will be at rest at time  $T > 0$  if the condition  $u(x, T) \equiv 0$  and  $u_t(x, T) \equiv 0$  are satisfied. These conditions are equivalent to

$$(1.4) \quad \int_0^l u^2(x, T) dx = 0 \quad \text{and} \quad \int_0^l u_t^2(x, T) dx = 0.$$

The constraints (1.4) are special cases of more general terminal constraints assumed in the discussion to follow.

Let  $h_1(x, u)$ ,  $h_2(x, u)$ ,  $g_1(x, u)$ ,  $g_2(x, u)$ ,  $F_{-2}(x, t, u, f)$ ,  $F_{-1}(x, t, u, f)$ ,  $F_0(x, t, u, f)$ ,  $\dots$ ,  $F_n(x, t, u, f)$  be continuous functions, each having continuous first and second derivatives with respect to  $u$ . Consider the following optimal control problem.

Find an admissible control  $f^0$  which minimizes the functional

$$(1.5) \quad J_0(f) = \int_0^l [g_1(x, u(x, T)) + g_2(x, u_t(x, T))] dx \\ + \int_0^T \int_0^l F_0(x, t, u(x, t), f(x, t)) dx dt$$

subject to the differential equation (1.1)–(1.2) and the constraints

$$\begin{aligned}
 & \int_0^l h_2(x, u_i(x, T)) dx + \int_0^T \int_0^l F_{-2}(x, t, u(x, t), f(x, t)) dx dt = c_{-2}, \\
 & \int_0^l h_1(x, u(x, T)) dx + \int_0^T \int_0^l F_{-1}(x, t, u(x, t), f(x, t)) dx dt = c_{-1}, \\
 (1.6) \quad & \int_0^T \int_0^l F_i(x, t, u(x, t), f(x, t)) dx dt \leq c_i, \quad 1 \leq i \leq n', \\
 & \int_0^T \int_0^l F_i(x, t, u(x, t), f(x, t)) dx dt = c_i, \quad n' < i \leq n.
 \end{aligned}$$

The discussion to follow is still applicable if all these constraints are of one type, either equality or inequality.

DEFINITION 1.2. An admissible control which solves this control problem is said to be an *optimal control*.

**1.3. Necessary and sufficient conditions for optimality.** A set of necessary conditions for optimality is contained in the following theorem. We shall show that the conditions stated in the theorem are also sufficient for a certain class of problems.

THEOREM 1.1 (Maximum principle). *In order that an admissible control  $f^0$ , with the response  $u^0$ , be optimal, it is necessary that there exist constants  $\lambda_{-2}, \lambda_{-1}, \lambda_0, \dots, \lambda_n$  and a solution  $v(x, t)$  of*

$$\begin{aligned}
 \mu(x) \frac{\partial^2 v}{\partial t^2} &= \frac{\partial}{\partial x} \left( p(x) \frac{\partial v}{\partial x} \right) - q(x)v \\
 &+ \sum_{i=-2}^n \lambda_i \frac{\partial F_i}{\partial u}(x, t, u^0(x, t), f^0(x, t)), \\
 (1.7) \quad v(0, t) &= 0, \quad v(l, t) = 0, \\
 v(x, T) &= \frac{1}{\mu(x)} \lambda_0 \frac{\partial g_2}{\partial u}(x, u^0(x, T)) + \frac{1}{\mu(x)} \lambda_{-2} \frac{\partial h_2}{\partial u}(x, u^0(x, T)), \\
 \frac{\partial v}{\partial t}(x, T) &= -\frac{1}{\mu(x)} \lambda_{-1} \frac{\partial h_1}{\partial u}(x, u^0(x, T)) - \frac{1}{\mu(x)} \lambda_0 \frac{\partial g_1}{\partial u}(x, u^0(x, T))
 \end{aligned}$$

such that

$$\begin{aligned}
 \max_{f \in F} [v(x, t)f + \sum_{i=-2}^n \lambda_i F_i(x, t, u^0(x, t), f)] \\
 = v(x, t)f^0(x, t) + \sum_{i=-2}^n \lambda_i F_i(x, t, u^0(x, t), f^0(x, t))
 \end{aligned}$$

for almost all points  $(x, t) \in (0, l) \times (0, T)$ . The  $\lambda_i$  can be chosen such that  $\lambda_i \leq 0$  for  $0 \leq i \leq n'$  and some  $\lambda_i \neq 0$ .

*Remark.* The conclusion of this theorem will be deduced from three lemmas. Before stating and proving the lemmas, we introduce a special perturbation of an optimal control.

Suppose that  $f^0$  is an optimal control with the corresponding response  $u^0$ . Let  $(x_1, t_1), \dots, (x_N, t_N)$  be  $N$  arbitrary points in the open region  $(0, l) \times (0, T)$  and let  $f_1, \dots, f_N$  be  $N$  arbitrary points in  $F$ . We may assume that  $x_1 \leq x_2 \leq \dots \leq x_N$ . Choose  $\delta > 0$  such that  $x_i + N\sqrt{\delta} < x_j$  if  $x_i < x_j$ , and such that  $x_N + N\sqrt{\delta} < l$  and  $t_i + \sqrt{\delta} < T$  for each  $i$ . Let  $\varepsilon_1, \dots, \varepsilon_N$  be real parameters satisfying  $0 \leq \varepsilon_j \leq \delta$  ( $1 \leq j \leq N$ ). Let

$$X_1 = x_1 \quad \text{and} \quad X_j = x_j + \sqrt{\varepsilon_1} + \dots + \sqrt{\varepsilon_{j-1}}$$

for  $1 < j \leq N$ . Clearly, the intervals  $X_j \leq x \leq X_j + \sqrt{\varepsilon_j}$  are nonoverlapping. Therefore, the rectangles

$$R_j: [X_j, X_j + \sqrt{\varepsilon_j}] \times [t_j, t_j + \sqrt{\varepsilon_j}]$$

are nonoverlapping.

Define the admissible control  $f_\varepsilon$  on  $[0, l] \times [0, T]$  by

$$(1.8) \quad f_\varepsilon(x, t) = \begin{cases} f^0(x, t) & \text{if } (x, t) \notin \bigcup_{j=1}^N R_j, \\ f_j & \text{if } (x, t) \in R_j, \quad j = 1, \dots, N. \end{cases}$$

Here  $\varepsilon$  denotes the vector  $(\varepsilon_1, \dots, \varepsilon_N) \in E_N$ . The norm of  $\varepsilon$  is defined by

$$|\varepsilon| = \varepsilon_1 + \dots + \varepsilon_N.$$

LEMMA 1.1. Let  $u_\varepsilon$  denote the response of the system (1.1)–(1.2) due to  $f_\varepsilon$ . Let  $\Delta u = u_\varepsilon - u^0$  and let  $\Delta f = f_\varepsilon - f^0$ . Then

$$\int_0^l \Delta u^2(x, T) dx = o(\varepsilon), \quad \int_0^l \Delta u_t^2(x, T) dx = o(\varepsilon)$$

and

$$\int_0^T \int_0^l \Delta u^2(x, t) dx dt = o(\varepsilon).$$

$o(\varepsilon)$  is a quantity such that

$$\lim_{\varepsilon \rightarrow 0} (o(\varepsilon)/|\varepsilon|) = 0.$$

*Proof.* Let

$$E(t) = \int_0^l \left\{ \mu(x) \left( \frac{\partial \Delta u}{\partial t}(x, t) \right)^2 + p(x) \left( \frac{\partial \Delta u}{\partial x}(x, t) \right)^2 + q(x) \Delta u^2(x, t) \right\} dx.$$

Proceeding formally as in [14, p. 145], we obtain

$$(1.9) \quad E(t) = \int_0^t \frac{dE(\tau)}{d\tau} d\tau = 2 \int_0^t \int_0^l \frac{\partial \Delta u}{\partial t}(x, \tau) \Delta f(x, \tau) dx d\tau.$$

This equality is valid for weak solutions since the latter is the limit of classical solutions in the energy norm.

Applying the Cauchy–Schwarz inequality to the space integral, we obtain

$$\begin{aligned} E(t) &\leq 2 \int_0^t \left( \int_0^l \left( \frac{\partial \Delta u}{\partial t} \right)^2 dx \right)^{1/2} \left( \int_0^l (\Delta f)^2 dx \right)^{1/2} d\tau \\ &\leq \int_0^T 2E^{1/2}(\tau) \left( \int_0^l (\Delta f)^2 dx \right)^{1/2} d\tau. \end{aligned}$$

It follows that

$$\begin{aligned} (1.10) \quad \sup_{0 \leq t \leq T} E(t) &\leq 2 \sup_{0 \leq t \leq T} E^{1/2}(t) \int_0^T \left( \int_0^l (\Delta f)^2 dx \right)^{1/2} d\tau \\ &= 2 \sup_{0 \leq t \leq T} E^{1/2}(t) \sum_{i=1}^N O(\varepsilon_i^{3/4}). \end{aligned}$$

$O(r)$  is a quantity such that

$$\lim_{r \rightarrow 0^+} (O(r)/r) = \text{const.}$$

Inequality (1.10) shows that, for each  $t \in [0, T]$ , we have

$$0 \leq E^{1/2}(t) \leq 2 \sum_{i=1}^N O(\varepsilon_i^{3/4})$$

or

$$(1.11) \quad E(t) = o(\varepsilon).$$

The conclusion of the lemma now follows immediately from (1.11), since the coefficients  $q$  and  $\mu$  are bounded away from zero.

Let  $L$  denote the differential operator defined by

$$Lu \equiv \mu(x) \frac{\partial^2 u}{\partial t^2} - \frac{\partial}{\partial x} \left( p(x) \frac{\partial u}{\partial x} \right) + q(x)u.$$

The formal adjoint of  $L$  has the same definition. Thus we shall write

$$Mv = \mu(x) \frac{\partial^2 v}{\partial t^2} - \frac{\partial}{\partial x} \left( p(x) \frac{\partial v}{\partial x} \right) + q(x)v$$

to denote the adjoint of  $L$ . We use  $M$  instead of  $L$  to distinguish the roles of  $u$  and  $v$  in the discussion to follow.

LEMMA 1.2. *Let  $u$  and  $v$  be two functions for which the expressions  $Lu$  and  $Mv$  are square integrable. Thus  $u$  and  $v$  are weak solutions in the sense mentioned above. If  $u$  and  $v$  satisfy the boundary conditions (1.2), then*

$$\begin{aligned} (1.12) \quad \int_0^T \int_0^l [vLu - uMv] dx dt &= \int_0^l \left\{ \mu(x)v(x, T) \frac{\partial u}{\partial t}(x, T) - \mu(x)u(x, T) \frac{\partial v}{\partial t}(x, T) \right. \\ &\quad \left. + \mu(x)u(x, 0) \frac{\partial v}{\partial t}(x, 0) - \mu(x)v(x, 0) \frac{\partial u}{\partial t}(x, 0) \right\} dx. \end{aligned}$$

*Proof.* Formally we have

$$vLu - uMv = -\frac{\partial}{\partial x} \left( vp(x) \frac{\partial u}{\partial x} - up(x) \frac{\partial v}{\partial x} \right) + \frac{\partial}{\partial t} \left( -\mu(x)u \frac{\partial v}{\partial t} + \mu(x)v \frac{\partial u}{\partial t} \right).$$

Therefore, the conclusion of the lemma can be obtained for  $c^2$  functions by integrating both sides of this identity over the region  $[0, l] \times [0, T]$  and making use of the boundary conditions in (1.2). The lemma holds for weak solutions  $u$  and  $v$  by passing to the limit in (1.12) for  $c^2$  functions.

DEFINITION 1.3. Let  $f^0$  be any admissible control with the response  $u^0$ . Let  $v$  be the solution of system (1.9) for arbitrary constants  $\lambda_{-2}, \lambda_{-1}, \lambda_0, \dots, \lambda_n$ . Let  $g$  be any one of the following functions:

$$\begin{aligned} g(x, t) &= F_i(x, t, u^0(x, t), f), \quad f \in F \text{ is fixed,} \\ g(x, t) &= F_i(x, t, u^0(x, t), f^0(x, t)), \\ g(x, t) &= v(x, t)f, \quad f \in F \text{ is fixed,} \\ g(x, t) &= v(x, t)f^0(x, t). \end{aligned}$$

We shall call  $(\bar{x}, \bar{t})$  a regular point of  $f^0$  if

$$\int_{\bar{x}}^{\bar{x}+\sqrt{\varepsilon}} \int_{\bar{t}}^{\bar{t}+\sqrt{\varepsilon}} g(x, t) dx dt = \varepsilon g(\bar{x}, \bar{t}) + o(\varepsilon)$$

for any  $\varepsilon > 0$  and sufficiently small.

It follows from [15, Theorem 6.3, p. 118], that almost all points of  $[0, l] \times [0, T]$  are regular for each admissible control  $f$ .

LEMMA 1.3. Let the functions  $f^0$  and  $f_\varepsilon$  be defined as in Lemma 1.1 where we assume now that the points  $(x_i, t_i)$  are regular,  $i = 1, \dots, N$ . If  $N = 1$ , then there exist constants  $\lambda_{-2}, \lambda_{-1}, \lambda_0, \lambda_1, \dots, \lambda_n$  (not all zero) such that

$$(a) \quad \lambda_i \leq 0 \quad \text{for } 0 \leq i \leq n',$$

$$(b) \quad \lim_{\varepsilon \rightarrow 0^+} \frac{\bar{J}(f_\varepsilon) - \bar{J}(f^0)}{\varepsilon} \leq 0,$$

where  $\bar{J}(f)$ , for any admissible control  $f$ , is defined by

$$\begin{aligned} \bar{J}(f) &= \lambda_0 \int_0^l [g_1(x, u(x, T)) + g_2(x, u_t(x, T))] dx \\ &\quad + \int_0^l [\lambda_{-2} h_2(x, u_t(x, T)) + \lambda_{-1} h_1(x, u(x, T))] dx \\ &\quad + \int_0^T \int_0^l \sum_{i=-2}^n \lambda_i F_i(x, t, u(x, t), f(x, t)) dx dt. \end{aligned}$$



*Proof.* Define the functionals  $J_{-2}, J_{-1}, J_1, \dots, J_n$  on the class of admissible controls by

$$J_{-1}(f) = \int_0^l h_1(x, u(x, T)) dx + \int_0^T \int_0^l F_{-1}(x, t, u(x, t), f(x, t)) dx dt,$$

$$J_{-2}(f) = \int_0^l h_2(x, u_t(x, T)) dx + \int_0^T \int_0^l F_{-2}(x, t, u(x, t), f(x, t)) dx dt,$$

$$J_i(f) = \int_0^T \int_0^l F_i(x, t, u(x, t), f(x, t)) dx dt,$$

$i = 1, \dots, n$ .  $J_0$  is defined by (1.5). Let  $J$  denote the vector-valued functional  $(J_{-2}, J_{-1}, J_0, \dots, J_n)$ . Let  $Y$  denote the set

$$Y = \{J(f) : f \text{ admissible}\} \subset E_{n+3}.$$

We assume the reader is familiar with the terminology and the generalized multiplier rule introduced by Hestenes in [13, Chap. IV]. We shall construct a derived set  $K$  for the set  $Y$  at  $J(f^0)$ .

Define functions  $v_{-2}, v_{-1}, v_0, \dots, v_n$  by requiring the following conditions to be satisfied:

$$(1.13a) \quad \mu(x) \frac{\partial^2 v_j}{\partial t^2} = \frac{\partial}{\partial x} \left( p(x) \frac{\partial v_j}{\partial x} \right) - q(x)v_j + \frac{\partial F_j}{\partial u}(x, t, u^0(x, t), f^0(x, t))$$

$$0 \leq x \leq l, \quad 0 \leq t \leq T, \quad -2 \leq j \leq n,$$

$$(1.13b) \quad v_j(0, t) = 0, \quad v_j(l, t) = 0, \quad -2 \leq j \leq n,$$

$$(1.13c) \quad v_{-2}(x, T) = \frac{1}{\mu(x)} \frac{\partial h_2}{\partial u}(x, u_t^0(x, T)), \quad \frac{\partial v_{-2}}{\partial t}(x, T) = 0,$$

$$(1.13d) \quad v_{-1}(x, T) \equiv 0, \quad \frac{\partial v_{-1}}{\partial t}(x, T) = -\frac{1}{\mu(x)} \frac{\partial h_1}{\partial u}(x, u^0(x, T)),$$

$$(1.13e) \quad v_0(x, T) = \frac{1}{\mu(x)} \frac{\partial g_2}{\partial u}(x, u_t^0(x, T)),$$

$$\frac{\partial v_0}{\partial t}(x, T) = -\frac{1}{\mu(x)} \frac{\partial g_1}{\partial u}(x, u^0(x, T)),$$

$$(1.13f) \quad v_j(x, T) = 0, \quad \frac{\partial v_j}{\partial t}(x, T) = 0, \quad j = 1, \dots, n.$$

Thus, each  $v_j$  is defined to be a solution of a boundary value problem satisfying certain terminal conditions. The transformation  $x \rightarrow x, t \rightarrow T - t$  transforms each of these problems into an initial-boundary value problem which has a solution for the same reason that the system (1.1)–(1.2) has a solution for each admissible control.

For each point  $(x, t) \in (0, l) \times (0, T)$  and for each point  $\bar{f} \in F$ , define the scalars

$$(1.14) \quad \begin{aligned} k^i(x, t, \bar{f}) &= v_i(x, t)(\bar{f} - f^0(x, t)) \\ &+ F_i(x, t, u^0(x, t), \bar{f}) - F_i(x, t, u^0(x, t), f^0(x, t)) \end{aligned} \quad i = -2, -1, 0, \dots, n.$$

We shall show that the set

$$K = \{k | k = (k^{-2}(x, t, \bar{f}), \dots, k^n(x, t, \bar{f})), (x, t) \text{ a regular point of } f^0, \bar{f} \in F\}$$

is a derived set for  $Y$  at  $J(f^0)$ .

Let  $k_1, \dots, k_N$  be an arbitrary finite collection of vectors from  $K$ . We must show that there exist points  $J_\varepsilon \in Y$  depending continuously on the vector parameter  $\varepsilon = (\varepsilon_1, \dots, \varepsilon_N)$ , for all sufficiently small positive values of  $\varepsilon$ , such that

$$J_\varepsilon = J(f^0) + \sum_{j=1}^N k_j \varepsilon_j + o(\varepsilon).$$

Since  $k_j \in K$ ,  $j = 1, \dots, N$ , there exist points  $(x_1, t_1), \dots, (x_N, t_N)$  of regularity of  $f^0$  and points  $f_1, \dots, f_N \in F$  such that

$$k_j = (k^{-2}(x_j, t_j, f_j), \dots, k^n(x_j, t_j, f_j)),$$

$j = 1, \dots, N$ . We shall show that  $J_\varepsilon$  can be defined by  $J_\varepsilon = J(f_\varepsilon)$ , where  $f_\varepsilon$  is the admissible control defined in Lemma 1.1.

We have, for  $i = 1, \dots, n$ ,

$$(1.15) \quad \begin{aligned} J_i(f_\varepsilon) - J_i(f^0) &= \int_0^T \int_0^l [F_i(x, t, u_\varepsilon(x, t), f_\varepsilon(x, t)) \\ &\quad - F_i(x, t, u^0(x, t), f^0(x, t))] dx dt \\ &= \int_0^T \int_0^l [F_i(x, t, u_\varepsilon(x, t), f_\varepsilon(x, t)) \\ &\quad - F_i(x, t, u^0(x, t), f_\varepsilon(x, t))] dx dt \\ &\quad + \int_0^T \int_0^l [F_i(x, t, u^0(x, t), f_\varepsilon(x, t)) \\ &\quad - F_i(x, t, u^0(x, t), f^0(x, t))] dx dt \\ &= \int_0^T \int_0^l \frac{\partial F_i}{\partial u}(x, t, u^0(x, t), f^0(x, t)) \Delta u(x, t) dx dt \\ &\quad + \sum_{j=1}^N \varepsilon_j [F_i(x_j, t_j, u^0(x_j, t_j), f_j) \\ &\quad - F_i(x_j, t_j, u^0(x_j, t_j), f^0(x_j, t_j))] + \sum_{j=1}^N o(\varepsilon_j). \end{aligned}$$

In (1.15) we used the fact that  $f^0$  is regular at each of the points  $(x_j, t_j)$  and the fact that

$$\int_0^T \int_0^l \Delta u^2(x, t) dx dt = o(\varepsilon).$$

If now we use the fact that

$$Mv_i = \frac{\partial F_i}{\partial u}(x, t, u^0(x, t), f^0(x, t)), \quad i = 1, \dots, n,$$

and substitute in (1.12), we obtain

$$\begin{aligned} \int_0^T \int_0^l \Delta u(x, t) Mv_i dx dt &= \int_0^T \int_0^l v_i(x, t)(f_\varepsilon(x, t) - f^0(x, t)) \\ &= \sum_{j=1}^N \varepsilon_j v_i(x_j, t_j)(f_j - f^0(x_j, t_j)) + o(\varepsilon). \end{aligned}$$

It now follows from (1.14) and (1.15) that

$$(1.16) \quad J_i(f_\varepsilon) - J_i(f^0) = \sum_{j=1}^N \varepsilon_j k_j^i + o(\varepsilon), \quad i = 1, \dots, n,$$

where  $k_j^i$  denotes the  $i$ th component of  $k_j$ .

For  $i = 0$ , we have

$$\begin{aligned} J_0(f_\varepsilon) - J_0(f^0) &= \int_0^l \left[ \frac{\partial g_1}{\partial u}(x, u^0(x, T)) \Delta u(x, T) + \frac{\partial g_2}{\partial u}(x, u_t^0(x, T)) \Delta u_t(x, T) \right] dx \\ &\quad + \sum_{j=1}^N \varepsilon_j [F_0(x_j, t_j, u^0(x_j, t_j), f_j) \\ &\quad \quad \quad - F_0(x_j, t_j, u^0(x_j, t_j), f^0(x_j, f_j))] \\ &\quad + \int_0^T \int_0^l (Mv_0) \Delta u(x, t) dx dt + o(\varepsilon). \end{aligned}$$

If now we observe that by (1.12) and (1.13e)

$$\begin{aligned} \int_0^T \int_0^l \Delta u(x, t) Mv_0 dx dt &= \int_0^T \int_0^l v_0(x, t)(f_\varepsilon(x, t) - f^0(x, t)) dx dt \\ &\quad - \int_0^l \left[ \frac{\partial g_1}{\partial u}(x, u^0(x, T)) \Delta u(x, T) \right. \\ &\quad \quad \quad \left. + \frac{\partial g_2}{\partial u}(x, u_t(x, T)) \Delta u_t(x, T) \right] dx, \end{aligned}$$

we obtain (1.16) for  $i = 0$ .

Equations (1.12) and (1.13c)–(1.13d) can be used in a similar fashion to prove (1.16) for  $i = -1$  and  $-2$ . Thus we can write

$$J(f_\varepsilon) - J(f^0) = \sum_{j=1}^N \varepsilon_j k_j(x_j, t_j, f_j) + o(\varepsilon).$$

Define  $J_\varepsilon = J(f_\varepsilon)$ . The proof that  $K$  is a derived set for  $Y$  at  $J(f^0)$  is now complete.

It now follows from Theorem 3.1 in [13, p. 177] that there exist multipliers  $\lambda_{-2}, \lambda_{-1}, \lambda_0, \dots, \lambda_n$  with  $\lambda_i \leq 0$  for  $0 \leq i \leq n'$  and some  $\lambda_i \neq 0$ , such that

$$(1.17) \quad \sum_{i=-2}^n \lambda_i k^i \leq 0$$

for any vector  $k = (k^{-2}, k^{-1}, k^0, \dots, k^n)$  in  $K$ .

To obtain the conclusion of Lemma 1.3, we take  $N = 1$  in the foregoing discussion and put

$$\bar{J} = \sum_{i=-2}^n \lambda_i J_i.$$

We then have, by (1.16),

$$\bar{J}(f_\varepsilon) - \bar{J}(f^0) = \varepsilon \sum_{i=-2}^n \lambda_i k^i + o(\varepsilon)$$

for any vector  $k = (k^{-2}, k^{-1}, k^0, \dots, k^n)$  in  $K$ . The proof of Lemma 1.3 now follows from the fact that

$$\lim_{\varepsilon \rightarrow 0^+} \frac{\bar{J}(f_\varepsilon) - \bar{J}(f^0)}{\varepsilon} = \sum_{i=-2}^n \lambda_i k^i(x, t, f) \leq 0.$$

*Proof of Theorem 1.1.* Assume that the admissible control  $f^0$  is optimal. Let  $(x, t)$  be a regular point of  $f^0$ . Then, by Lemma 1.3, there exist constants  $\lambda_{-2}, \lambda_{-1}, \lambda_0, \dots, \lambda_n$  independent of  $(x, t)$ , with  $\lambda_i \leq 0$  for  $0 \leq i \leq n'$  and some  $\lambda_i \neq 0$ , such that

$$\sum_{i=-2}^n \lambda_i [v_i(x, t)(f - f^0(x, t)) + F_i(x, t, u^0(x, t), f) - F_i(x, t, u^0(x, t), f^0(x, t))] \leq 0$$

for any point  $f \in F$ . The  $v_i$ 's are defined in (1.13). This says that the function

$$\sum_{i=-2}^n \lambda_i [v_i(x, t)f + F_i(x, t, u^0(x, t), f)]$$

attains its maximum value on  $F$  at  $f = f^0(x, t)$ . If we put  $v(x, t) = \sum_{i=-2}^n \lambda_i v_i(x, t)$ , we obtain the conclusion of the maximum principle. This completes the proof of Theorem 1.1.

**THEOREM 1.2.** Consider the control system (1.1), (1.2), (1.5), (1.6). Assume

- (a)  $g_1, g_2, F_0, \dots, F_{n'}$ , are convex functions of  $u$ ;
- (b) the functions  $F_i$  are of the form

$$F_i(x, t, u, f) = F^i(x, t, u) + H^i(x, t, f), \quad i = -2, \dots, n.$$

Suppose there exist an admissible control  $f^0$  and constants  $\lambda_{-2}, \lambda_{-1}, \lambda_0, \dots, \lambda_n$  such that

$$\max_{f \in F} \left[ v(x, t)f + \sum_{i=-2}^n \lambda_i H^i(x, t, f) \right] = v(x, t)f^0(x, t) + \sum_{i=-2}^n \lambda_i H^i(x, t, f^0(x, t)),$$

where  $v$  is a nonzero solution of

$$Mv = \sum_{i=-2}^n \lambda_i \frac{\partial F^i}{\partial u}(x, t, u^0(x, t))$$

with the terminal and boundary conditions (1.9). Assume further that

- (c)  $\lambda_0 < 0, \lambda_i \leq 0$  for  $i = 1, \dots, n'$ ;
- (d) the constraints (1.6) are satisfied by  $f^0$ ;
- (e) if strict inequality holds in (1.6), the corresponding multiplier  $\lambda_i$  is zero;
- (f)  $-\lambda_i F^i, -\lambda_{-2} h_2, -\lambda_{-1} h_1$  are convex functions of  $u$  for  $n' < i \leq n$ .

Then  $f^0$  is optimal.

*Remarks.* Condition (e) is a necessary assumption since, as Hestenes shows in the proof of his multiplier rule [13], if an optimal control satisfies a strict inequality in (1.6), the corresponding multiplier must be zero. (f) is satisfied when the functions  $h_1, h_2, F_i, n' < i \leq n$ , are linear in  $u$ .

*Proof of Theorem 1.2.* If  $f$  and  $u$  satisfy (1.6), then by (e),

$$\begin{aligned} & \int_0^T \int_0^l \lambda_i [F^i(x, t, u(x, t)) - F^i(x, t, u^0(x, t))] dx dt \\ & + \int_0^T \int_0^l \lambda_i [H^i(x, t, f(x, t)) - H^i(x, t, f^0(x, t))] dx dt \geq 0 \end{aligned}$$

for  $i = 1, \dots, n$ . A similar result holds for  $i = -2$  and  $-1$ .

It follows that

$$\begin{aligned} (1.18) \quad -\lambda_0 [J_0(f) - J_0(f^0)] & \geq -\lambda_0 \int_0^l [g_2(x, u_i(x, T)) - g_2(x, u_i^0(x, T)) + g_1(x, u(x, T)) \\ & \qquad \qquad \qquad - g_1(x, u^0(x, T))] dx \\ & - \sum_{i=-2}^n \int_0^T \int_0^l \lambda_i [F^i(x, t, u(x, t)) - F^i(x, t, u^0(x, t))] dx dt \\ & - \sum_{i=-2}^n \int_0^T \int_0^l \lambda_i [H^i(x, t, f(x, t)) - H^i(x, t, f^0(x, t))] dx dt \\ & - \lambda_{-1} \int_0^l [h_1(x, u(x, T)) - h_1(x, u^0(x, T))] dx \\ & - \lambda_{-2} \int_0^l [h_2(x, u_i(x, T)) - h_2(x, u_i^0(x, T))] dx. \end{aligned}$$

Therefore, by the convexity assumption (f),

$$\begin{aligned}
-\lambda_0[J_0(f) - J_0(f^0)] &\geq -\lambda_0 \int_0^l \left[ \frac{\partial g_1}{\partial u}(x, u^0(x, T)) \Delta u(x, T) \right. \\
&\quad \left. + \frac{\partial g_2}{\partial u}(x, u_t(x, T)) \Delta u_t(x, T) \right] dx \\
&\quad - \int_0^T \int_0^l \sum_{i=-2}^n \lambda_i \frac{\partial F^i}{\partial u}(x, t, u^0(x, t)) \Delta u(x, t) dx dt \\
&\quad + \int_0^T \int_0^l \sum_{i=-2}^n \lambda_i [H^i(x, t, f^0(x, t)) - H^i(x, t, f(x, t))] dx dt \\
&\quad - \lambda_{-1} \int_0^l \frac{\partial h_1}{\partial u}(x, u^0(x, T)) \Delta u(x, T) dx \\
&\quad - \lambda_{-2} \int_0^l \frac{\partial h_2}{\partial u}(x, u_t^0(x, T)) \Delta u_t(x, T) dx \\
&= -\lambda_0 \int_0^l \left[ \frac{\partial g_1}{\partial u}(x, u^0(x, T)) \Delta u(x, T) \right. \\
&\quad \left. + \frac{\partial g_2}{\partial u}(x, u_t^0(x, T)) \Delta u_t(x, T) \right] dx \\
&\quad - \int_0^T \int_0^l (Mv) \Delta u(x, t) dx dt \\
&\quad + \int_0^T \int_0^l \sum_{i=-2}^n \lambda_i [H^i(x, t, f^0(x, t)) - H^i(x, t, f(x, t))] dx dt \\
&\quad - \int_0^l \left[ \lambda_{-1} \frac{\partial h_1}{\partial u}(x, u^0(x, T)) \Delta u(x, T) \right. \\
&\quad \left. + \lambda_{-2} \frac{\partial h_2}{\partial u}(x, u_t^0(x, T)) \Delta u_t(x, T) \right] dx.
\end{aligned}$$

If we now apply Lemma 1.2 and the conditions (1.7), we obtain

$$\begin{aligned}
(1.19) \quad -\lambda_0[J_0(f) - J_0(f^0)] &\geq \int_0^T \int_0^l \{v(x, t)[f^0(x, t) - f(x, t)] \\
&\quad + \sum_{i=-2}^n \lambda_i [H^i(x, t, f^0(x, t)) - H^i(x, t, f(x, t))]\} dx dt.
\end{aligned}$$

The right-hand side of (1.19) is nonnegative. Hence,  $J_0(f) - J_0(f^0) \geq 0$ .

**1.4. Generalizations.** We wish to point out briefly how the results of the previous section can be extended to more general systems. Thus, let us consider a

system governed by the general linear hyperbolic partial differential equation

$$(1.20) \quad \frac{\partial^2 u}{\partial t^2} = \sum_{i,j=1}^m \frac{\partial}{\partial x_i} \left( p_{ij}(x) \frac{\partial u}{\partial x_j} \right) + \sum_{i=1}^m r_i(x) \frac{\partial u}{\partial x_i} + q(x)u + f(x, t),$$

$(x, t) \in \Omega \times [0, T] = Q,$

together with boundary and initial conditions of the form

$$(1.21) \quad \begin{aligned} u(x, 0) &= \varphi(x), & \frac{\partial u}{\partial t}(x, 0) &= \psi(x) & \text{in } \Omega, \\ u(x, t) &= 0 & \text{on } \partial\Omega \times [0, T]. \end{aligned}$$

Here  $\Omega$  represents an open connected set in  $E_m$  and  $x = (x_1, \dots, x_m)$  is a point in  $\Omega$ .  $\partial\Omega$  denotes the boundary of  $\Omega$ .

The functions  $p_{ij}$ ,  $\partial p_{ij}/\partial x_k$ ,  $r_i$ ,  $\partial r_i/\partial x$  and  $q$  are assumed to be bounded and measurable in  $\Omega$ . Moreover, we assume that  $p_{ij} = p_{ji}$  and that there exist  $c$  and  $C$  such that

$$c \sum_{i=1}^m \xi_i^2 \leq \sum_{i,j=1}^m p_{ij}(x) \xi_i \xi_j \leq C \sum_{i=1}^m \xi_i^2$$

for all  $x \in \Omega$  and all  $\xi = (\xi_1, \dots, \xi_m) \in E_m$ .  $c$  and  $C$  are constants satisfying  $0 < c \leq C < \infty$ .

With these assumptions holding, the system (1.20)–(1.21) has a weak solution having finite energy for every choice of  $f \in L_2(Q)$  and  $\varphi, \psi \in L_2(\Omega)$ . This follows from Theorem 5.1 in [17].

Now let us imagine that in § 1.2 we have replaced the interval  $(0, l)$  by the set  $\Omega$  in  $E_m$  and that  $x$  is now an  $m$ -dimensional vector. The optimal control problem formulated in § 1.2 now makes sense for the system (1.20)–(1.21). Moreover, the statement of the necessary and sufficient optimality conditions remain the same except that now in place of the differential equations (1.1) and (1.7) we must substitute the differential equations

$$Lu = f$$

and

$$Mv = \sum_{i=1}^n \lambda_i \frac{\partial F_i}{\partial u}(x, t, u^0(x, t), f^0(x, t)),$$

where  $L$  and  $M$  are defined by (1.22) and (1.23) below. We must also take  $v(x, t) = 0$ ,  $x \in \partial\Omega$ , and  $\mu(x) \equiv 1$  in (1.7).

The reader who has understood the proofs in § 1.2 and who is familiar with the elementary theory of second order linear hyperbolic partial differential equations as discussed, for example, in [17] will have no trouble constructing the proofs of these facts for himself. Nevertheless, we shall state and indicate the proofs of the needed results. In the first place, we recall that Theorem 1.1 follows from Lemmas 1.1, 1.2 and 1.3. Lemma 1.3 is a consequence of Lemmas 1.1, 1.2 and Theorem 3.1 in [13, p. 177]. To obtain Lemma 1.2 for the system (1.20)–(1.21),

we must define

$$(1.22) \quad Lu \equiv \frac{\partial^2 u}{\partial t^2} - \sum_{i,j=1}^m \frac{\partial}{\partial x_i} \left( p_{ij}(x) \frac{\partial u}{\partial x_j} \right) - \sum_{i=1}^m r_i(x) \frac{\partial u}{\partial x_i} - q(x)u$$

and

$$(1.23) \quad Mv \equiv \frac{\partial^2 v}{\partial t^2} - \sum_{i,j=1}^m \frac{\partial}{\partial x_i} \left( p_{ij}(x) \frac{\partial v}{\partial x_j} \right) + \sum_{i=1}^m \frac{\partial}{\partial x_i} \left( r_i \frac{\partial v}{\partial x_i} \right) - q(x)v.$$

Lemma 1.2 now follows immediately from the identity (which is given by (4.10) in [17])

$$\begin{aligned} \int_0^T \int_{\Omega} vLu \, dx \, dt &= \int_{\Omega} \left\{ \frac{\partial u}{\partial t}(x, T)v(x, T) - \frac{\partial u}{\partial t}(x, 0)v(x, 0) \right\} dx \\ &\quad - \int_0^T \int_{\Omega} \left\{ \frac{\partial u}{\partial t} \frac{\partial v}{\partial t} - \sum_{i,j=1}^m p_{ij} \frac{\partial u}{\partial x_i} \frac{\partial v}{\partial x_j} + \sum_{i=1}^m r_i \frac{\partial u}{\partial x_i} + quv \right\} dx \, dt \end{aligned}$$

and its analogue for

$$\int_0^T \int_{\Omega} uMv \, dx \, dt.$$

We turn now to a discussion of the  $m$ -dimensional version of Lemma 1.1. To this end, let  $(x^1, t^1), \dots, (x^N, t^N)$  be  $N$  arbitrary points in  $\Omega \times (0, T)$  and let  $f_1, \dots, f_N$  be  $N$  arbitrary points in  $F$ . Assume that the  $(x^i, t^i)$  are labeled so that  $x_1^1 \leq x_1^2 \leq \dots \leq x_1^N$ , where  $x_i^j$  denotes the  $i$ th component of  $x^j$ . Choose  $\delta > 0$  such that  $x_1^i + N\delta < x_1^j$  if  $x_1^i < x_1^j$  and such that the points  $(y^j, \tau^j) \in E_{m+1}$  defined by

$$y_i^j = x_i^j + N\delta^{1/(m+1)}, \quad \tau^j = t^j + \delta^{1/(m+1)},$$

$i = 1, \dots, m, j = 1, \dots, N$ , belong to  $\Omega \times (0, T)$ . Let  $\varepsilon_1, \dots, \varepsilon_N$  be  $N$  real parameters satisfying  $0 \leq \varepsilon_j \leq \delta$  ( $1 \leq j \leq N$ ), and let  $\hat{m} = m + 1$ . Let

$$X_1 = x_1^1 \quad \text{and} \quad X_j = x_1^j + \sqrt[\hat{m}]{\varepsilon_1} + \dots + \sqrt[\hat{m}]{\varepsilon_{j-1}}$$

for  $1 < j \leq N$ . The intervals  $X_j \leq x \leq X_j + \sqrt[\hat{m}]{\varepsilon_j}$  are nonoverlapping, and therefore, the  $(m + 1)$ -dimensional cubes

$$R_j: [X_j, X_j + \sqrt[\hat{m}]{\varepsilon_j}] \times [x_2^j, x_2^j + \sqrt[\hat{m}]{\varepsilon_j}] \times \dots \times [x_m^j, x_m^j + \sqrt[\hat{m}]{\varepsilon_j}] \times [t^j, t^j + \sqrt[\hat{m}]{\varepsilon_j}],$$

$j = 1, \dots, N$ , are nonintersecting.

Let  $f_{\varepsilon}$  be defined on  $\Omega \times [0, T]$  according to definition (1.8). To obtain the analogue of Lemma 1.1 for the system (1.20)–(1.21), we assert the inequality

$$(1.24) \quad \begin{aligned} \int_{\Omega} \left\{ \left( \frac{\partial \Delta u}{\partial t}(x, t) \right)^2 + \sum_{i,j=1}^m p_{ij} \frac{\partial \Delta u}{\partial x_i}(x, t) \frac{\partial \Delta u}{\partial x_j}(x, t) + \Delta u^2(x, t) \right\} dx \\ \leq 2e^{Kt} \int_0^t \int_{\Omega} e^{-K\tau} \frac{\partial \Delta u}{\partial t}(x, \tau) \Delta f(x, \tau) \, dx \, d\tau \end{aligned}$$

which is obtained from the energy inequality (8.1) in [17].  $K$  is a constant depending only on  $m, c$  and the bounds for the coefficients of  $L$ .  $\Delta u$  and  $\Delta f$  are defined as in Lemma 1.1.



Denoting the left-hand side of (1.24) by  $E(t)$ , we obtain from the Cauchy–Schwarz inequality,

$$E(t) \leq 2e^{Kt} \int_0^t e^{-K\tau} E^{1/2}(\tau) \left( \int_{\Omega} (\Delta f)^2 dx \right)^{1/2} d\tau.$$

From this, it follows that

$$\begin{aligned} E^{1/2}(t) &\leq 2 \int_0^t \left( \int_{\Omega} (\Delta f)^2 dx \right)^{1/2} d\tau \\ &= \sum_{i=1}^N \varepsilon_i^{1/(m+1)} O(\varepsilon_i^{m/(2(m+1))}) = \sum_{i=1}^N O(\varepsilon_i^{(m+2)/(2(m+1))}). \end{aligned}$$

Now since  $(m+2)/(2(m+1)) > 1/2$ , it follows that  $E(t) = o(\varepsilon)$  uniformly in  $t$  for  $0 \leq t \leq T$ . This completes the development necessary for the extension of the results in § 1.2 to the system (1.20)–(1.21).

*Example.* Let the controlled process be described by the differential equation

$$(1.25) \quad \frac{\partial^2 u}{\partial t^2} = \frac{\partial^2 u}{\partial x^2} + f(x, t), \quad 0 \leq x \leq \pi, \quad 0 \leq t \leq \pi,$$

with the initial and boundary conditions

$$(1.26) \quad u(x, 0) = \sin x, \quad \partial u(x, 0)/\partial t = -\sin x, \quad u(0, t) = u(\pi, t) = 0.$$

Find an admissible control  $f^0$  which minimizes the functional

$$J_0(f) = \int_0^\pi [u^2(x, \pi) + u_t^2(x, \pi)] dx + \int_0^\pi \int_0^\pi f^2(x, t) dx dt,$$

subject to the differential equation (1.25)–(1.26). For this example, we shall assume that the set  $F$  appearing in the definition of an admissible control is the whole real line.

To find an admissible control  $f^0$  satisfying the condition of the maximum principle (Theorem 1.1), we look for a function which maximizes an expression of the form

$$(1.27) \quad v f + \lambda_0 f^2,$$

where  $\lambda_0 \leq 0$ . We assume  $\lambda_0 < 0$ . Then we may assume without loss of generality that  $\lambda_0 = -1$ . Since  $f^0$  maximizes (1.27), we must have

$$f^0 = -\frac{1}{2}v.$$

Putting this in (1.25) we are faced with the problem of solving the system

$$\frac{\partial^2 u}{\partial t^2} = \frac{\partial^2 u}{\partial x^2} - \frac{1}{2}v, \quad \frac{\partial^2 v}{\partial t^2} = \frac{\partial^2 v}{\partial x^2},$$

$$u(0, t) = u(\pi, t) = v(0, t) = v(\pi, t) = 0,$$

$$u(x, 0) = \sin x, \quad \partial u(x, 0)/\partial t = -\sin x,$$

$$v(x, \pi) = -2u(x, \pi), \quad v_t(x, \pi) = 2u(x, \pi).$$

The solution of this system is

$$v(x, t) = \frac{2}{1 + \pi/2} (\sin t + \cos t) \sin x,$$

$$u(x, t) = (\cos t - \sin t) \sin x$$

$$- \frac{\sin x}{1 + \pi/2} \int_0^t \sin(t - \tau) (\sin \tau + \cos \tau) d\tau.$$

Hence, an optimal control is

$$f^0(x, t) = -\frac{1}{2} v(x, t) = -\frac{1}{1 + \pi/2} (\sin t + \cos t) \sin x,$$

since the maximum principle is a sufficient condition for optimality for this example.

The minimum cost is  $4\pi/(2 + \pi)^2 + 2\pi^2/(2 + \pi)^2$ .

**2. The optimal control of a vibrating beam.** In this section we discuss the optimal control of the system studied by Komkov in [3]. Thus we let  $u$  denote the transverse deflection of an inhomogeneous beam subject to an external force  $f$ . Then  $u$  and  $f$  are related by the differential equation

$$(2.1) \quad \rho(x)A(x) \frac{\partial^2 u}{\partial t^2} + \frac{\partial^2}{\partial x^2} \left( E(x)I(x) \frac{\partial^2 u}{\partial x^2} \right) = f(x, t), \quad 0 \leq x \leq l, \quad 0 \leq t.$$

$l$  denotes the length of the beam.  $\rho(x)$  is the material density,  $A(x)$  the cross-sectional area,  $E(x)$  Young's modulus and  $I(x)$  is the moment of inertia of the cross-sectional area about the neutral axis.

We assume the boundary conditions for (2.1) to be given by

$$(i) \quad u(0, t) = u(l, t) = 0, \quad \frac{\partial u}{\partial x}(0, t) = \frac{\partial u}{\partial x}(l, t) = 0,$$

or by

$$(ii) \quad u(0, t) = u(l, t) = 0, \quad EI \frac{\partial^2 u}{\partial x^2}(0, t) = EI \frac{\partial^2 u}{\partial x^2}(l, t) = 0,$$

or by

$$(iii) \quad EI \frac{\partial^2 u}{\partial x^2}(0, t) = EI \frac{\partial^2 u}{\partial x^2}(l, t) = 0,$$

$$\frac{\partial}{\partial x} \left[ EI \frac{\partial^2 u}{\partial x^2}(0, t) \right] = \frac{\partial}{\partial x} \left[ EI \frac{\partial^2 u}{\partial x^2}(l, t) \right] = 0.$$

Conditions (i) correspond to a beam with built-in ends. Conditions (ii) correspond to a beam with freely supported ends. Conditions (iii) correspond to a beam with free ends.

The initial condition for our system is of the form

$$(2.2) \quad u(x, 0) = \varphi(x), \quad \frac{\partial u}{\partial t}(x, 0) = \psi(x),$$

where  $\varphi$  and  $\psi$  are continuously differentiable functions on  $[0, l]$ . As in [3], we assume the coefficients  $\rho$ ,  $A$ ,  $E$  and  $I$  to be twice continuously differentiable and positively bounded away from zero on  $[0, l]$ .

**2.1. Formulation of the control problem.** We shall formulate the control problem for the vibrating beam analogously to the formulation given above for the problem of controlling the vibrating string.

Let functions  $h_1, h_2, g_1, g_2, F_{-2}, \dots, F_n$  be given as in § 1.2.

We wish to find an admissible control  $f^0$  which minimizes the functional

$$(2.3) \quad J_0(f) = \int_0^l [g_1(x, u(x, T)) + g_2(x, u_t(x, T))] dx \\ + \int_0^T \int_0^l F_0(x, t, u(x, t), f(x, t)) dx dt$$

subject to the differential equation (2.1) with initial conditions (2.2) and subject to any set of boundary conditions (i), (ii) or (iii) and the side constraints (1.6).  $T > 0$  is fixed.

An admissible control for this problem is defined by Definition 1.1 in § 1.2.

For each admissible control  $f$  the system (2.1)–(2.2) has a weak solution  $u$  which is continuously differentiable on  $[0, l] \times [0, T]$ . Moreover, the energy function

$$(2.4) \quad E(t) = \frac{1}{2} \int_0^l \left\{ \rho(x)A(x) \left( \frac{\partial u(x, t)}{\partial t} \right)^2 + E(x)I(x) \left( \frac{\partial^2 u}{\partial x^2}(x, t) \right)^2 \right\} dx$$

is continuous in  $t$  and uniformly bounded for  $0 \leq t \leq T$ . It is easy to show that each weak solution is the limit of a sequence of classical solutions in the energy norm (2.4). This justifies certain results obtained below by formally integrating by parts.

**DEFINITION 2.1.** An admissible control which minimizes the functional (2.3) subject to the above constraints is said to be an *optimal control*.

**THEOREM 2.1 (Maximum principle).** *In order that an admissible control  $f^0$  with response  $u^0$  be optimal, it is necessary that there exist constants  $\lambda_{-2}, \lambda_{-1}, \lambda_0, \dots, \lambda_n$  and a solution  $v(x, t)$  of*

$$(2.5) \quad \rho(x)A(x) \frac{\partial^2 v}{\partial t^2} + \frac{\partial^2}{\partial x^2} \left( E(x)I(x) \frac{\partial^2 v}{\partial x^2} \right) = \sum_{i=-2}^n \lambda_i \frac{\partial F_i}{\partial u}(x, t, u^0(x, t), f^0(x, t)), \\ v(x, T) = \frac{\lambda_0}{\rho(x)A(x)} \frac{\partial g_2}{\partial u}(x, u_t^0(x, T)) + \frac{\lambda_{-2}}{\rho(x)A(x)} \frac{\partial h_2}{\partial u}(x, u_t^0(x, T)), \\ \frac{\partial v}{\partial t}(x, T) = -\frac{\lambda_0}{\rho(x)A(x)} \frac{\partial g_1}{\partial u}(x, u^0(x, T)) - \frac{\lambda_{-1}}{\rho(x)A(x)} \frac{\partial h_1}{\partial u}(x, u^0(x, T))$$

such that

$$\begin{aligned} \max_{f \in F} \left[ v(x, t)f + \sum_{i=-2}^n \lambda_i F_i(x, t, u^0(x, t), f) \right] \\ = v(x, t)f^0(x, t) + \sum_{i=-2}^n \lambda_i F_i(x, t, u^0(x, t), f^0(x, t)). \end{aligned}$$

$v$  satisfies the same boundary conditions as  $u$ . Not all of the  $\lambda_i$  are zero, and  $\lambda_i \leq 0$  for  $0 \leq i \leq n'$ .

The proof of this theorem closely parallels the proof of Theorem 1.1. We shall merely state and indicate the proofs of the necessary lemmas. The other details of the proof are left to the reader.

LEMMA 2.1. *Let the admissible control  $f^0$  be optimal for the system (2.1)–(2.2). Let  $u^0$  denote the optimal response corresponding to  $f^0$ . Let the functions  $f_\varepsilon$  and  $u_\varepsilon$  be defined as in Lemma 1.1. Then*

$$\int_0^l \Delta u^2(x, T) dx = o(\varepsilon), \quad \int_0^l \Delta u_t^2(x, T) dx = o(\varepsilon)$$

and

$$\int_0^T \int_0^l \Delta u^2(x, t) dx dt = o(\varepsilon),$$

where  $\Delta u = u_\varepsilon - u_0$ .

*Proof.* Equation (14) in [3] shows that

$$\begin{aligned} \frac{1}{2} \int_0^l \left\{ \rho(x)A(x) \left( \frac{\partial \Delta u}{\partial t}(x, t) \right)^2 + E(x)I(x) \left( \frac{\partial^2 \Delta u}{\partial x^2}(x, t) \right)^2 \right\} dx \\ = \int_0^t \int_0^l \frac{\partial \Delta u}{\partial t}(x, t) \Delta f(x, t) dx dt.^1 \end{aligned}$$

The portion of the lemma concerning  $\Delta u_t$  now follows as in the proof of Lemma 1.1. To obtain the remaining portion of the lemma we observe that

$$|\Delta u(x, t)| = \left| \int_0^t \Delta u_t(x, \tau) d\tau \right| \leq T^{1/2} \left( \int_0^T (\Delta u_t)^2 d\tau \right)^{1/2}$$

and so

$$\int_0^l \Delta u^2(x, t) dx = o(\varepsilon)$$

uniformly in  $t$  for  $0 \leq t \leq T$ . This completes the proof of the lemma.

Let  $L$  denote the differential operator defined by

$$Lu \equiv \rho(x)A(x) \frac{\partial^2 u}{\partial t^2} + \frac{\partial^2}{\partial x^2} \left( E(x)I(x) \frac{\partial^2 u}{\partial x^2} \right).$$

<sup>1</sup> This derivation in [3] is purely formal, but it can be justified by using the approximating series technique that by now has become standard in this paper.

The formal adjoint of  $L$  is defined by

$$Mv \equiv \rho(x)A(x) \frac{\partial^2 v}{\partial t^2} + \frac{\partial^2}{\partial x^2} \left( E(x)I(x) \frac{\partial^2 v}{\partial x^2} \right).$$

Again we are dealing with a self-adjoint system. As before, we shall use  $L$  and  $M$  to distinguish the roles of  $u$  and  $v$  in what follows.

LEMMA 2.2. *Let  $u$  and  $v$  be two functions for which  $Lu$  and  $Mv$  are defined and are square integrable. Thus  $u$  and  $v$  are weak solutions of certain differential equations. If  $u$  and  $v$  satisfy any set of the boundary conditions (i), (ii) or (iii), then*

$$(2.6) \quad \int_0^T \int_0^l [vLu - uMv] dx dt = \int_0^l \rho(x)A(x) \left[ v(x, T) \frac{\partial u}{\partial t}(x, T) - u(x, T) \frac{\partial v}{\partial t}(x, T) \right] dx \\ + \int_0^l \rho(x)A(x) \left[ u(x, 0) \frac{\partial v}{\partial t}(x, 0) - v(x, 0) \frac{\partial u}{\partial t}(x, 0) \right] dx.$$

*Proof.* For sufficiently smooth functions  $u$  and  $v$ , the lemma is proved by integrating the left-hand side of (2.6) by parts and making use of the boundary conditions (i), (ii) or (iii). For weak solutions  $u$  and  $v$ , the lemma is obtained by applying (2.6) to a sequence of smooth functions converging to  $u$  and  $v$  in the energy norm (2.4).

LEMMA 2.3. *Lemma 1.3 is valid for the problem of this section.*

There is no essential change in the proof given for Lemma 1.3. We shall, therefore, leave the proof to the reader. It is probably worthwhile to point out here that the functions  $v_j$  defined in (1.13a)–(1.13f) are now defined as follows:

$$(2.7a) \quad \rho(x)A(x) \frac{\partial^2 v_j}{\partial t^2} + \frac{\partial^2}{\partial x^2} \left( E(x)I(x) \frac{\partial^2 v_j}{\partial x^2} \right) = \frac{\partial F_j}{\partial u}(x, t, u^0(x, t), f^0(x, t)), \\ 0 \leq x \leq l, \quad 0 \leq t \leq T, \quad -2 \leq j \leq n.$$

These  $v_j$  satisfy the same set of boundary conditions that have been prescribed for the function  $u$  in (2.1). The terminal conditions on the  $v_j$  are given by

$$(2.7b) \quad v_{-2}(x, T) = \frac{1}{\rho(x)A(x)} \frac{\partial h_2}{\partial u}(x, u_t(x, T)), \quad \frac{\partial v_{-2}}{\partial t}(x, T) = 0,$$

$$(2.7c) \quad v_{-1}(x, T) = 0, \quad \frac{\partial v_{-1}}{\partial t}(x, T) = -\frac{1}{\rho(x)A(x)} \frac{\partial h_1}{\partial u}(x, u_t^0(x, T)),$$

$$(2.7d) \quad v_0(x, T) = \frac{1}{\rho(x)A(x)} \frac{\partial g_2}{\partial u}(x, u_t^0(x, T)),$$

$$\frac{\partial v_0}{\partial t}(x, T) = -\frac{1}{\rho(x)A(x)} \frac{\partial g_1}{\partial u}(x, u^0(x, T)),$$

$$(2.7e) \quad v_j(x, T) = 0, \quad \frac{\partial v_j}{\partial t}(x, T) = 0, \quad j = 1, \dots, n.$$

The proof of Theorem 2.1 can now be carried out following the pattern set forth in the proof of Theorem 1.1.

Similarly, one can prove the analogue of Theorem 1.2 for the system (2.1)–(2.2). If this is done, one should apply the remarks following Theorem 1.2 to the new theorem.

## REFERENCES

- [1] E. B. LEE, *A sufficient condition in the theory of optimal control*, this Journal, 1 (1963), pp. 241–245.
- [2] D. L. RUSSELL, *Optimal regulation of linear symmetric hyperbolic systems with finite dimensional controls*, this Journal, 4 (1966), pp. 276–295.
- [3] V. KOMKOV, *The optimal control of a transverse vibration of a beam*, this Journal, 6 (1968), pp. 401–421.
- [4] L. S. PONTRYAGIN, V. BOLTYANSKII, R. GAMKRELIDZE AND E. MISHCHENKO, *The Mathematical Theory of Optimal Processes*, L. W. Neustadt, ed., Interscience, New York, 1962.
- [5] A. I. EGOROV, *On optimal control of processes in distributed objects*, J. Appl. Math. Mech., 27 (1963), pp. 1045–1058.
- [6] ———, *Necessary optimality conditions for distributed parameter systems*, this Journal, 5 (1967), pp. 352–408.
- [7] ———, *Optimal processes in distributed parameter systems and certain problems in invariance theory*, this Journal, 4 (1966), pp. 601–661.
- [8] G. S. JONES AND A. STRAUSS, *An example of optimal control*, SIAM Rev., 10 (1968), pp. 25–55.
- [9] D. W. BUSHAW, *Optimal discontinuous forcing terms*, Contributions to the Theory of Nonlinear Oscillations, vol. IV, Annals of Mathematics Studies, No. 41, Princeton Univ. Press, Princeton, 1958, pp. 29–52.
- [10] E. B. LEE AND L. MARKUS, *Foundations of Optimal Control Theory*, John Wiley, New York, 1967.
- [11] J. L. LIONS, *Contrôle optimal de systèmes gouvernés par des équations aux dérivées partielles*, Dunod, Paris, 1968.
- [12] G. HELLWIG, *Partial Differential Equations*, Blaisdell, New York, 1964.
- [13] M. HESTENES, *Calculus of Variations and Optimal Control Theory*, John Wiley, New York, 1966.
- [14] N. S. KOSHLYAKOV, M. M. SMIRNOV AND E. B. GLIBER, *Differential Equations of Mathematical Physics*, North-Holland, Amsterdam, 1964.
- [15] S. SAKS, *Theory of the Integral*, Hafner, New York, 1937.
- [16] E. R. BARNES, *The optimal control of systems with distributed parameters*, Doctoral thesis, Mathematics Department, University of Maryland, College Park, 1968.
- [17] C. WILCOX, *Initial-boundary value problems for linear hyperbolic partial differential equations of the second order*, Arch. Rational Mech. Anal., 10 (1962), pp. 361–400.

## CONVERGENCE OF A DISCRETIZATION FOR CONSTRAINED SPLINE FUNCTION PROBLEMS\*

JAMES W. DANIEL†

**1. Introduction.** In [10], some generalizations of the basic ideas of spline functions were developed by considering certain minimization problems under a mixture of discrete and continuous inequality constraints, extending concepts in [1], [8], [9] and [13]. Sufficient and sometimes necessary conditions for a function to solve the minimization problem were presented via optimal control techniques, but no computational methods were discussed. In the present paper, we shall analyze the convergence of simple discretizations of the problem, such discretizations in many cases being finitely solvable by standard quadratic programming methods. Let us first define the problem.

Let  $m$  be a positive integer and let  $1 < p < \infty$ . For  $i = 1, 2, \dots, k$ , let  $M_i$  be not identically zero linear differential operators on  $[0, 1]$  of degree less than  $m$ , and similiary for  $N_i, i = 1, 2, \dots, n$ ; we write

$$M_i x = \sum_{j=0}^{m-1} b_{ij}(t)x^{(j)}(t), \quad N_i x = \sum_{j=0}^{m-1} c_{ij}(t)x^{(j)}(t).$$

We allow  $k = 0$  and  $n = 0$ . Let  $L$  be a linear differential operator on  $[0, 1]$  of exact degree  $m$ ,

$$Lx = \sum_{j=0}^m a_j(t)x^{(j)}(t), \quad a_m(t) \neq 0 \text{ in } [0, 1].$$

Let  $W^{m,p}$  be the (Sobolev) space of real-valued functions  $x$  on  $[0, 1]$  such that  $x^{(m-1)}$  is absolutely continuous and  $x^{(m)} \in L^p(0, 1)$ . Then our minimization problem is to

$$\begin{aligned} \text{minimize } f(x) &\equiv \int_0^1 |Lx(t)|^p dt \quad \text{over} \\ (1.1) \quad C &= \{x; x \in W^{m,p}, \alpha_i(t) \leq M_i x(t) \leq \beta_i(t) \text{ for } 0 \leq t \leq 1, \\ &\quad i = 1, \dots, k, \gamma_i \leq N_i x(\xi_i) \leq \delta_i \text{ for } i = 1, \dots, n\}, \end{aligned}$$

where  $\alpha_i$  and  $\beta_i$  are given functions, the  $\gamma_i$  and  $\delta_i$  are given scalars, and the  $\xi_i$  are points in  $[0, 1]$ . Some simple generalizations are possible by allowing one-sided constraints or by allowing the  $N_i$  to be difference operators, but we shall not consider these here. It is shown in [10] that, if  $C$  is nonempty, if  $a_j \in C^j[0, 1]$  for  $j = 0, \dots, m$ , if  $b_{ij} \in C[0, 1]$  for  $i = 1, \dots, k$  and  $j = 0, \dots, m-1$ , if  $c_{ij} \in C[0, 1]$  for  $i = 1, \dots, n$  and  $j = 0, \dots, m-1$ , and if  $\alpha_i$  and  $\beta_i$  lie in  $M_i W^{m,p}$  for  $i = 1, \dots, k$ , then there exists a solution  $x^*$  to the problem in (1.1). We shall assume the above hypotheses to hold throughout the ensuing discussion. Since the solution  $x^*$

\* Received by the editors March 20, 1970, and in revised form June 26, 1970.

† Computer Sciences Department, University of Wisconsin, Madison, Wisconsin 53706. This work was supported by the Office of Naval Research under Contract N00014-67-A-0128-0004.

may hit the boundary of  $C$  at unknown points, perhaps countably many times, the computation of  $x^*$  is difficult. One obvious way to handle this is as follows.

Let  $h > 0$  be some mesh size, say  $h = 1/Q$ , and let  $[0, 1]$  be partitioned by  $t_i = ih$ ,  $i = 0, \dots, 1/h = Q$ ; we suppose that all the points  $\xi_i$  lie on this mesh for all  $h$  to be used, that is,  $\xi_i/h$  is an integer (this assumption is not necessary but simplifies the notation). Our first discretization consists in merely replacing the continuous constraints by discrete ones, that is, we

$$\begin{aligned} & \text{minimize } f(x) = \int_0^1 |Lx(t)|^p dt \quad \text{over} \\ (1.2) \quad & C_1(h) \equiv \{x; x \in W^{m,p}, \alpha_i(t_j) \leq M_i x(t_j) \leq \beta_i(t_j) \text{ for} \\ & j = 0, \dots, Q, i = 1, \dots, k, \gamma_i \leq N_i x(\xi_i) \leq \delta_i \\ & \text{for } i = 1, \dots, n\}. \end{aligned}$$

As analyzed in [13], this problem can be solved in the common case of  $p = 2$  in finitely many steps by minimizing a quadratic function of  $2k(Q + 1) + m + 2n$  variables subject to  $2k(Q + 1) + m + 2n$  linear inequality constraints. We shall prove the following (§ 2, Theorem 2.1): *All weak limit points (in the  $W^{m,p}$  sense), at least one of which exists, of a sequence of solutions to the first discretization in (1.2) must solve the original problem of (1.1); if the solution to the original problem is unique, the approximating solutions converge to it  $W^{m,p}$ -weakly, and in particular the function and the first  $m - 1$  derivatives converge uniformly, that is, in  $C[0, 1]$ .*

If one must take  $h$  very small to obtain a reasonable approximation to  $x^*$ , one might well be satisfied to have only approximate values of  $x^*$  at the grid points  $t_i$  rather than throughout  $[0, 1]$ ; if  $x^*(t)$  were desired and  $x^*(t_i)$  was accurately known, unconstrained interpolation could be used to generate a reasonable approximation to  $x^*$ . Thus we are led to a second, more complete, discretization. If  $z$  is a function defined at least on the mesh points  $t_i = ih$ ,  $i = 0, \dots, 1/h = Q$ , let  $D = D_h$  be the operator such that  $Dz(t_i) = (z(t_{i+1}) - z(t_i))/h$ ,  $i = 0, 1, \dots, Q - 1$ . We then have

$$D^l z(t_i) = \frac{1}{h^l} \sum_{j=0}^l (-1)^{l-j} \binom{l}{j} z(t_{i+j})$$

as a natural analogue of the  $l$ th derivative of  $z$ . We therefore define

$$M_{i,h} z(t_j) = \sum_{l=0}^{m-1} b_{il}(t_j) D^l z(t_j), \quad 0 \leq j \leq Q - m + 1,$$

$$N_{i,h} z(t_j) = \sum_{l=0}^{m-1} c_{il}(t_j) D^l z(t_j), \quad 0 \leq j \leq Q - m + 1,$$

$$L_h z(t_j) = \sum_{l=0}^m a_{il}(t_j) D^l z(t_j), \quad 0 \leq j \leq Q - m.$$



Our second discretized problem is now to

$$\begin{aligned}
 \text{minimize } f_h(x_h) &= h \sum_{j=0}^{Q-m} |L_h x_h(t_j)|^p \quad \text{over} \\
 C_2(h) &\equiv \{x_h; x_h \equiv (x_h(0), x_h(t_1), \dots, x_h(1))^T \in R^{Q+1}, \\
 (1.3) \quad & -\varepsilon_h + \alpha_i(t_j) \leq M_{i,h} x_h(t_j) \leq \beta_i(t_j) + \varepsilon_h \\
 & \text{for } j = 0, \dots, Q - m + 1, i = 1, \dots, k, \\
 & -\varepsilon_h + \gamma_i \leq N_{i,h} x_h(\xi_i) \leq \delta_i + \varepsilon_h \text{ for } i = 1, \dots, n\}.
 \end{aligned}$$

Here  $\varepsilon_h$ , which tends to zero, gives a small expansion of the constraint set as  $h \rightarrow 0$ . If some such expansion is not allowed, the set  $C_2(h)$  can be empty [6], [7]; as will be clear from the use made of the expansion by  $\varepsilon_h$ , constraints of the form  $\alpha(t) \leq x(t) \leq \beta(t)$  or  $\alpha \leq x(\xi) \leq \beta$  need not be expanded. This problem in the common case of  $p = 2$  can be solved in finitely many steps since it involves a quadratic function of  $Q + 1$  variables subject to  $2k(Q - m + 2) + 2n$  linear inequalities. Since in general  $Q = 1/h$  is large, seeking only approximations to  $x^*(t_i)$  reduces the difficulty considerably. Under slight additional hypotheses on the problem, we shall prove the following (roughly stated) (§ 3, Theorem 3.1): *All  $W^{m,p}$  weak limit points (of certain "interpolations"), at least one of which exists, of a sequence of solutions to the discretization in (1.3) must solve the original problem of (1.1); if the solution to the original problem is unique, the ("interpolations" of the) approximate solutions converge to it  $W^{m,p}$ -weakly, and in particular the function values and the first  $m - 1$  divided differences converge uniformly, that is, in  $C[0, 1]$ , to  $x^*$  and its first  $m - 1$  derivatives.*

In our analysis of these discretizations, we shall follow the general approach of [6] and [7] for studying discretized optimization problems; in particular, we use the approach of [6] and [7] for studying discretized optimal control problems, since our problem can easily be stated in that formalism, as in our (2.6). Once we use the control theory formalism, however, other approaches for analyzing discretizations of control problems can be used; for example, the ideas of Cullum [2] and [3]. Our present problem has such special structure however (for example, the control variable is unconstrained but the state variable is subject to constraints) that the more general results give us too little information; we must make use of the problem's structure.

**2. Analysis of the first, simpler, discretization.** We have yet to define a norm on the space  $W^{m,p}$ ; two common norms are

$$\begin{aligned}
 \|x\|_0 &\equiv \left\{ \sum_{i=0}^m \int_0^1 |x^{(i)}(t)|^p dt \right\}^{1/p}, \\
 \|x\|_1 &\equiv \left\{ \sum_{i=1}^m |x(\theta_i)|^p + \int_0^1 |x^{(m)}(t)|^p dt \right\}^{1/p}
 \end{aligned}$$

for  $0 \leq \theta_1 < \theta_2 < \dots < \theta_m \leq 1$ . It is well known [14], [12] that these norms are equivalent; that is, there exist positive constants  $a, A$  such that  $a\|x\|_0 \leq \|x\|_1 \leq A\|x\|_0$  for all  $x$  in  $W^{m,p}$ . This equivalence is related to the famous Sobolev

inequalities which say that, for  $0 \leq i \leq m - 1$ , there exist constants  $A_i$  such that  $|x^{(i)}(t)| \leq A_i \|x\|_1$  for all  $t$  in  $[0, 1]$  and  $x$  in  $W^{m,p}$ . If a sequence  $\{x_n\}$  converges weakly in  $(W^{m,p}, \|\cdot\|_1)$  to  $x$ , then  $\|x_n\|_1$  is uniformly bounded and therefore, by the Sobolev inequalities, for  $0 \leq i \leq m - 1$ , the functions  $x_n^{(i)}$  are uniformly bounded in  $C[0, 1]$ . Hence the  $\{x_n^{(i)}\}$  are equicontinuous for  $0 \leq i \leq m - 2$  which implies that  $\{x_n^{(i)}\}$  converges uniformly, that is, in  $C[0, 1]$ , to  $x^{(i)}$  for  $0 \leq i \leq m - 2$ . A slightly more subtle argument, using

$$|x^{(m-1)}(t) - x^{(m-1)}(t')| \leq \left| \int_t^{t'} |x^{(m)}(\tau)| d\tau \right| \leq \left[ \int_0^1 |x^{(m)}(\tau)|^p d\tau \right]^{1/p} |t' - t|^{1-1/p},$$

shows as well that  $\{x_n^{(m-1)}\}$  converges uniformly to  $x^{(m-1)}$ . For  $1 < p < \infty$ ,  $W^{m,p}$  is reflexive and hence the unit sphere is weakly compact and weakly sequentially compact.

From the computational standpoint, serious difficulties arise if the original problem (1.1) admits solutions of arbitrarily large norm. For example, the functions  $x_n(t) = n$  form a minimizing sequence (in fact, they are all solutions) for  $\int_0^1 |x^{(2)}(t)|^2 dt$  over the set of  $x$  satisfying  $0 \leq x^{(1)}(t) \leq 1$  but have no convergent subsequence. In this situation, our analysis to follow could not guarantee that the approximate solutions have limit points; to avoid this we must eliminate problems admitting solutions of arbitrarily large norm. We pause to see what this means. For  $0 \leq l \leq k$ , let  $S_l \equiv \{x; x \in W^{m,p}, Lx \equiv 0, M_i x \equiv 0 \text{ for } 1 \leq i \leq l, \text{ and } N_i x(\xi_i) = 0 \text{ for } 1 \leq i \leq n\}$ . It is shown in [10] that, if  $d_{k+1}$  is the dimension of  $S_k$  and if  $d_{k+1-j}$  is the dimension of  $S_{k+1-j} - S_{k+2-j}$  for  $1 \leq j \leq k$ , then there exist points  $\theta_{l,i}$  with  $0 \leq \theta_{l,1} < \theta_{l,2} < \dots < \theta_{l,d_l} \leq 1$  for  $1 \leq l \leq k + 1$ , with the points  $\{\theta_{k+1,i}\}_{i=1}^{d_{k+1}}$  being completely arbitrary in  $[0, 1]$ , such that

$$\|x\| \equiv \left\{ \sum_{i=1}^k \sum_{j=1}^{d_i} |M_i x(\theta_{i,j})|^p + \sum_{j=1}^{d_{k+1}} |x(\theta_{k+1,j})|^p + \sum_{i=1}^n |N_i x(\xi_i)|^p + \int_0^1 |Lx(t)|^p dt \right\}^{1/p} \tag{2.1}$$

defines a norm on  $W^{m,p}$ . We remark that if for some  $l_0$  one has  $M_{l_0} x \equiv x$ , then one may take  $d_{l_0} = m$ , all other  $d_i = 0$ , eliminate the sum in  $N_i$  from (2.1), and take arbitrary distinct points for  $\theta_{l_0,i}$  to define the norm. By using the Sobolev inequalities, it is simple to show that this norm is in fact equivalent to  $\|\cdot\|_0$  and  $\|\cdot\|_1$ .

LEMMA 2.1.  $\|\cdot\|$  is equivalent to  $\|\cdot\|_0$  and  $\|\cdot\|_1$ .

*Proof.* The existence of an  $A'$  such that  $\|x\| \leq A' \|x\|_1$  for all  $x$  in  $W^{m,p}$  is clear from the Sobolev inequalities; we ask whether or not an  $a' > 0$  exists such that  $\|x\| \geq a' \|x\|_1$ . If not, we can find  $x_n \in W^{m,p}$  such that  $\|x_n\| \rightarrow 0$  but  $\|x_n\|_1 = 1$ ; by the weak sequential compactness of the sphere in  $(W^{m,p}, \|\cdot\|_1)$ , we may assume  $x_n$  converges  $(W^{m,p}, \|\cdot\|_1)$ -weakly to some  $x$  in  $W^{m,p}$ . Since  $L$  is bounded from  $(W^{m,p}, \|\cdot\|_1)$  into  $L^p(0, 1)$ , since  $M_i$  and  $N_i$  are bounded from  $(W^{m,p}, \|\cdot\|_1)$  into  $C[0, 1]$ , and since  $\int_0^1 |Lx_n(t)|^p dt$  tends to zero, we have  $\|x\| = 0$  and therefore  $x \equiv 0$ . Since  $x_n$  converges uniformly, that is, in  $C[0, 1]$ , to  $x$  and  $\|x_n\|_1 = 1$ , we have

$$\int_0^1 |x_n^{(m)}(t)|^p dt \text{ converging to 1. Then}$$

$$\left\{ \int_0^1 |Lx_n(t)|^p dt \right\}^{1/p} \geq \left\{ \int_0^1 |a_m(t)x_n^{(m)}(t)|^p dt \right\}^{1/p} - \left\{ \int_0^1 \left| \sum_{i=0}^{m-1} a_i(t)x_n^{(i)}(t) \right|^p dt \right\}^{1/p},$$

which then is bounded away from zero since  $|a_m(t)| \geq \varepsilon > 0$  for some  $\varepsilon$ , since  $\int_0^1 |x_n^{(m)}(t)|^p dt \rightarrow 1$ , and since  $x_n^{(i)}$  converges uniformly to  $x^{(i)} \equiv 0$  for  $0 \leq i \leq m - 1$ . This contradicts  $\|x_n\| \rightarrow 0$ . This completes the proof.

Now for any  $x$  in  $C$ , the values  $|M_l x(\theta_{l,i})|$  and  $|N_l x(\xi_i)|$  are bounded, uniformly in  $i, l$  and  $x$ . Given any  $\bar{x}$  in  $C$ , we would seek  $x^*$  only from among those  $x$  in  $C$  satisfying  $f(x) \leq f(\bar{x})$ , which in turn gives an upper bound for  $\int_0^1 |Lx(t)|^p dt$ .

From this then, we can deduce an a priori bound on  $\|x^*\|$  if and only if the  $\{|x(\theta_{k+1,j})|\}$  are uniformly bounded for  $1 \leq j \leq d_{k+1}$ . If  $d_{k+1} = 0$ , this is certainly true; if  $d_{k+1} \neq 0$ , there exists a nonzero function  $z$  in  $S_k$  and thus  $x^* + \alpha z \in C$  for all scalars  $\alpha$  and  $f(x^* + \alpha z) = f(x^*)$ . Therefore the original problem admits of an a priori bound on its solutions if and only if  $d_{k+1} = 0$ . We hereafter assume that  $d_{k+1} = 0$ . Computationally this may be accomplished simply by adding one continuous constraint  $|x(t)| \leq E$  for some large  $E$  or  $m$  discrete constraints  $|x(\theta_j)| \leq E$  and thus this is not a computationally significant restriction. Adding  $|x(t)| \leq E$  means that we may use as our norm the simple expression

$$(2.2) \quad \left\{ \sum_{i=1}^m |x(\theta_i)|^p + \int_0^1 |Lx(t)|^p dt \right\}^{1/p}.$$

We hereafter assume that the points  $\theta_{l,i}$  of (2.1) are mesh points for all  $h$  used. If some  $M_{l,0}x \equiv x$ , then we need only assume that the  $m$  points  $\theta_i$  of (2.2) are mesh points for all  $h$  used.

Now let, for  $Q = 1/h$ ,  $x_Q^*$  be a solution to the simple discretized problem of (1.2);  $C_1(h)$  is not empty since, in particular,  $x^* \in C_1(h)$ , where  $x^*$  solves the original problem. Since  $f(x_Q^*) \leq f(x^*)$ , since  $|M_l x_Q^*(\theta_{l,i})| \leq \max\{|\alpha_l(\theta_{l,i})|, |\beta_l(\theta_{l,i})|\}$  and since  $|N_l x_Q^*(\xi_i)| \leq \max\{|\delta_l|, |\gamma_l|\}$ , we conclude that  $\|x_Q^*\|$  is uniformly bounded. We note that  $x_Q^*$  is not necessarily in  $C$ ; it is however "near" to  $C$  as the following more general lemma demonstrates.

LEMMA 2.2. If  $x \in W^{m,p}$  and  $-\varepsilon + \alpha_i(t_j) \leq M_i x(t_j) \leq \beta_i(t_j) + \varepsilon$  for  $i = 1, \dots, k$  and  $j = 0, \dots, Q$ , then  $-(\varepsilon + \eta_i(h)) + \alpha_i(t) \leq M_i x(t) \leq \beta_i(t) + (\varepsilon + \eta_i(h))$  for  $0 \leq t \leq 1$ , where

$$\eta_i(h) = \max \left[ \left\{ \int_0^1 |\beta_i^{(1)}(t)|^p dt \right\}^{1/p}, \left\{ \int_0^1 |\alpha_i^{(1)}(t)|^p dt \right\}^{1/p} \right] h^{1/q} + B_i \|x\|_0 h^{1/q} + \|x\|_0 w_i(h),$$

where  $1/p + 1/q = 1$ ,  $|b_{ii}(t)| \leq B_i$ ,  $|b_{ii}(t_1) - b_{ii}(t_2)| \leq w_i(h)$  if  $|t_1 - t_2| \leq h$ ,  $l = 0, \dots, m - 1$ .

Proof. For the upper bound,

$$M_i x(t) - \beta_i(t) = M_i x(t) - M_i x(t_j) + M_i x(t_j) - \beta_i(t_j) + \beta_i(t_j) - \beta_i(t).$$

Since

$$\beta_i \in M_i W^{m,p} \subset W^{1,p},$$

$$|\beta_i(t) - \beta_i(t_j)| \leq \int_{t_j}^t |\beta_i^{(1)}(t)| dt \leq \left\{ \int_0^1 |\beta_i^{(1)}(t)|^p dt \right\}^{1/p} |t - t_j|^{1/q}.$$

Also

$$\begin{aligned} |M_i x(t) - M_i x(t_j)| &\leq \sum_{l=0}^{m-1} |b_{il}(t)x^{(l)}(t) - b_{il}(t_j)x^{(l)}(t_j)| \\ &\leq B_i \|x\|_0 |t - t_j|^{1/q} + \|x\|_0 w_i(|t - t_j|) \end{aligned}$$

arguing as for  $\beta_i$ . Thus, letting  $t_j$  be such that  $|t - t_j| \leq h$ , we have

$$M_i x(t) - \beta_i(t) \leq \varepsilon + \left[ \int_0^1 |\beta_i^{(1)}(t)|^p dt \right]^{1/p} + B_i \|x\|_0 h^{1/q} + \|x\|_0 w_i(h).$$

Similarly for the lower bound. This completes the proof.

We note that if the  $b_{il}$  are Hölder continuous with exponent greater than or equal to  $1/q$  (that is, if  $|b_{il}(t_1) - b_{il}(t_2)| \leq c|t_1 - t_2|^r$  for some constant  $c \geq 0$  and  $r \geq 1/q$ ) then  $|\eta_i(h)| \leq Fh^q$  for a constant  $F$  uniformly bounded whenever  $\|x\|$  is bounded.

**LEMMA 2.3.** *If functions  $x_Q \in W^{m,p}$  satisfy  $-\varepsilon_Q + \alpha_i(t) \leq M_i x_Q(t) \leq \beta_i(t) + \varepsilon_Q$  for  $i = 1, \dots, k$  and  $-\varepsilon_Q + \gamma_i \leq N_i x_Q(\xi_i) \leq \delta_i + \varepsilon_Q$  for  $i = 1, \dots, n$ , if  $\lim_{Q \rightarrow \infty} \varepsilon_Q = 0$  and if  $x_Q$  converges  $W^{m,p}$ -weakly to  $x$ , then  $x \in C$ .*

*Proof.* This is obvious since  $x_Q^{(j)}(t)$  converges to  $x^{(j)}(t)$  for each  $t$  in  $[0, 1]$  for  $0 \leq j \leq m - 1$ , and the constraints involve derivatives of order at most  $m - 1$ .

**LEMMA 2.4.** *As  $\varepsilon_Q \geq 0$  tends to zero, the minimum of  $f$  over the set  $C_{\varepsilon_Q} = \{x; x \in W^{m,p}, -\varepsilon_Q + \alpha_i(t) \leq M_i x(t) \leq \beta_i(t) + \varepsilon_Q$  for  $1 \leq i \leq k$  and  $0 \leq t \leq 1, -\varepsilon_Q + \gamma_i \leq N_i x(\xi_i) \leq \delta_i + \varepsilon_Q$  for  $1 \leq i \leq n\}$  converges to the minimum of  $f$  over  $C$ .*

*Proof.* Clearly each set  $C_{\varepsilon_Q}$  is weakly closed, and for the minimization problem over  $C_{\varepsilon_Q}$  we may restrict ourselves to those  $x$  satisfying  $f(x) \leq f(x^*)$ , since  $x^* \in C_{\varepsilon_Q}$ , where  $x^*$  minimizes  $f$  over  $C$ . Since, for all  $x$  in this set,  $\|x\|$  is uniformly bounded, the weakly lower semicontinuous functional  $f$  attains its minimum over the weakly compact set  $C_{\varepsilon_Q}$  at some point  $x_Q$ . Since  $\|x_Q\|$  is uniformly bounded, we may assume that  $x_Q$  converges weakly to some  $x$ , which must be in  $C$  by Lemma 2. Thus  $f(x) \leq \liminf_{Q \rightarrow \infty} f(x_Q) \leq \limsup f(x_Q) \leq f(x^*)$  since  $f(x_Q) \leq f(x^*)$  for all  $Q$ . Thus  $\lim_{Q \rightarrow \infty} f(x_Q) = f(x^*)$ .

We can now prove our discretization result for the simpler discretization.

**THEOREM 2.1.** *Let the general assumptions of § 1 hold and let the fixed points  $\theta_{i,i}$  defining the norm  $\|\cdot\|$  in (2.1) (or  $\theta_i$  in (2.2)) be mesh points in our discretization for all  $h$ . Let  $x_Q^* \in W^{m,p}$  solve the problem of (1.2), that is, minimize  $f(x)$  over  $C_1(h)$ . Then  $f(x_Q^*)$  converges to  $f(x^*)$  and all  $W^{m,p}$  weak limit points, at least one of which exists, of  $\{x_Q^*\}$  minimize  $f$  over  $C$ ; if  $x^*$  minimizing  $f$  over  $C$  is unique, then  $x_Q^*$  converges weakly to  $x^*$ , that is,  $x_Q^{*(i)}$  converges uniformly to  $x^{*(i)}$  for  $0 \leq i \leq m - 1$  and  $x_Q^{*(m)}$  converges  $L^p$ -weakly to  $x^{*(m)}$ .*

*Proof.* Arguing as in Lemma 2.4, we see that  $x_Q^*$  always exists,  $f(x_Q^*) \leq f(x^*)$  and there exists a constant  $E$  such that  $\|x_Q^*\| \leq E$ . Thus, by Lemma 2.2, there exist functions  $\eta_i(h)$  for  $1 \leq i \leq k$  tending to zero with  $h$ , and such that  $x_Q^* \in C_Q \equiv \{x; x \in W^{m,p}, -\eta_i(h) + \alpha_i(t) \leq M_i x(t) \leq \beta_i(t) + \eta_i(h) \text{ for } 1 \leq i \leq k, \gamma_i \leq N_i x(\xi_i) \leq \delta_i \text{ for } 1 \leq i \leq n\}$ . By Lemma 2.4,  $\zeta_Q = \min_C f - \min_{C_Q} f$  tends to zero. We write

$$(2.3) \quad f(x^*) = \min_C f = \min_{C_Q} f + \zeta_Q \leq f(x_Q^*) + \zeta_Q \leq f(x^*) + \zeta_Q.$$

This implies, since  $\zeta_Q \rightarrow 0$ , that  $f(x_Q^*) \rightarrow f(x^*)$ . Since  $\{x_Q^*\}$  is bounded, it has weak limit points. For any such weak limit point  $x'$  with  $x_Q^*$  weakly converging to  $x'$ , we have  $x' \in C$  by Lemma 2.3 and thus

$$f(x^*) \leq f(x') \leq \liminf_{Q' \rightarrow \infty} f(x_{Q'}^*) = \lim_{Q \rightarrow \infty} f(x_Q^*) = f(x^*).$$

This says that  $f(x') = f(x^*)$ , that is,  $x'$  minimizes  $f$  over  $C$ . The remainder follows from the definition of convergence in  $W^{m,p}$ .

We have not been able to estimate the rate of convergence as a function of  $h$ .

**3. Analysis of the more complete discretization.** We shall here use the norm  $\|\cdot\|$  defined in (2.1) via points  $\theta_{l,i}$  (or  $\theta_i$  in (2.2)) and we shall assume that the  $\theta_{l,i}$  are mesh points for all  $h$  used. We shall analyze our complete discretization, the relationship between the problems in (1.1) and (1.3), by the general discretization analysis of [5], [6] and [7]; for completeness our arguments are self-contained.

We wish to use roughly the arguments of Theorem 2.1 in this case also. If  $x_h^*$  minimizes  $f_h$  over  $C_2(h)$ , we unfortunately cannot talk about  $f(x_h^*)$  or  $f_h(x^*)$  as in (2.3) in the proof of Theorem 2.1, since these make no sense in our new situation. Instead, with  $x^*$  we shall associate a point  $y_h = r_h x^* \in C_2(h)$  by a "discretization" or "restriction" mapping  $r_h$  such that  $|f_h(r_h x^*) - f(x^*)|$  converges to zero with  $h$ . Similarly, with  $x_h^*$  we shall associate a  $z_h = p_h x_h^* \in W^{m,p}$  "converging into  $C$ " by an "interpolation" or "prolongation" mapping  $p_h$  and such that  $|f(p_h x_h^*) - f_h(x_h^*)|$  converges to zero with  $h$ . We can then imitate the proof of Theorem 2.1 by replacing (2.3) by roughly

$$\begin{aligned} f(x^*) &= \min_{C_Q} f + \zeta_Q \leq f(p_h x_h^*) + \zeta_Q \\ &= f_h(x_h^*) + \zeta_Q + [f(p_h x_h^*) - f_h(x_h^*)] \\ &\leq f_h(r_h x^*) + \zeta_Q + [f(p_h x_h^*) - f_h(x_h^*)] \\ &= f(x^*) + \zeta_Q + [f(p_h x_h^*) - f_h(x_h^*)] + [f_h(r_h x^*) - f(x^*)]. \end{aligned}$$

Having outlined our approach and the reasons for constructing certain mappings  $p_h$  and  $r_h$ , we now proceed with the technical details. We define the restriction  $r_h$  in the obvious manner. Let  $y_h = r_h x^*$  be the discrete mesh function (that is, defined at points  $t_i = ih$  only) defined by  $y_h(t_i) \equiv x^*(t_i)$ .

We need to develop some tools for using divided differences. For any  $l$  with  $0 \leq l \leq m$ , by Peano's theorem we can write

$$D^l x^*(t) = \frac{1}{(l-1)!} \int_0^1 D_t^l(t-\tau)_+^{l-1} x^{*(l)}(\tau) d\tau,$$

where the  $D_t$  indicates differences with respect to  $t$ , that is,  $D_t g(t, \tau) = (g(t + h, \tau) - g(t, \tau))/h$ , and where

$$(t - \tau)_+^{l-1} = \begin{cases} (t - \tau)^{l-1} & \text{for } t - \tau \geq 0, \\ 0 & \text{for } t - \tau \leq 0 \end{cases}$$

as usual.  $K_l(\tau - t) \equiv (1/(l-1)!)D_t^l(t - \tau)_+^{l-1}$  is a ‘‘basic spline’’ [4], that is,

$$K_l(s) \equiv \frac{1}{h^l(l-1)!} \sum_{i=0}^l (-1)^{l-i} \binom{l}{i} (ih - s)_+^{l-1}$$

vanishes identically for  $s \geq lh$  and  $s \leq 0$ , is strictly positive for  $s$  in  $(0, lh)$  and lies in  $C^{l-2}(-\infty, \infty)$ . Thus  $D^l t^l \equiv 1! = \int_0^{lh} K_l(s) ds$ , and we find

$$\int_0^{lh} K_l(s) ds = 1.$$

Also

$$K_l(s) = \frac{1}{h^l(l-1)!} \sum_{i=0}^l (-1)^{l-i} \binom{l}{i} \left(i - \frac{s}{h}\right)_+^{l-1} h^{l-1} \leq \frac{G}{h}$$

for some fixed  $G$ , since we have  $0 \leq s/h \leq l$ .

We now show that  $r_h x^* \in C_2(h)$  for large enough  $\varepsilon_h$ .

**LEMMA 3.1.**  $-\varepsilon_h + \alpha_i(t_j) \leq M_{i,h} y_h(t_j) \leq \beta_i(t_j) + \varepsilon_h$  for  $0 \leq j \leq Q - m + 1$  and  $1 \leq i \leq k$ , and  $-\varepsilon_h + \gamma_i \leq N_{i,h} y_h(\xi_i) \leq \delta_i + \varepsilon_h$  for  $1 \leq i \leq n$ , where  $\varepsilon_h = Gh^{1/q}$  and  $G \geq \|x^*\|_0 \sum_{i=0}^{m-1} \|b_{il}\|_\infty l^{1/q}$  for  $1 \leq i \leq k$ , and  $G \geq \|x^*\|_0 \sum_{l=0}^{m-1} \|c_l\|_\infty l^{1/q}$  for  $1 \leq i \leq n$ .

*Proof.* For  $0 \leq l \leq m - 1$ ,

$$\begin{aligned} |D^l y_h(t_j) - x^{*(l)}(t_j)| &= |D^l x^*(t_j) - x^{*(l)}(t_j)| \\ &= \left| \int_{t_j}^{t_j+lh} K_l(\tau - t_j) [x^{*(l)}(\tau) - x^{*(l)}(t_j)] dt \right| \\ &\leq \sup_{t_j \leq \tau \leq t_j+lh} |x^{*(l)}(\tau) - x^{*(l)}(t_j)|. \end{aligned}$$

Now

$$|x^{*(l)}(\tau) - x^{*(l)}(t_j)| \leq \int_{t_j}^{\tau} |x^{*(l+1)}(s)| ds \leq \|x^*\|_0 |\tau - t_j|^{1/q} \leq \|x^*\|_0 lh^{1/q}.$$

More generally,

$$\begin{aligned} |M_{i,h} y_h(t_j) - M_i x^*(t_j)| &= |M_{i,h} x^*(t_j) - M_i x^*(t_j)| \\ &\leq \sum_{l=0}^{m-1} |b_{il}(t_j)| |D^l x^*(t_j) - x^{*(l)}(t_j)| \\ &\leq \sum_{l=0}^{m-1} |b_{il}(t_j)| \|x^*\|_0 lh^{1/q} \end{aligned}$$

and similarly for  $N_i$  and  $N_{i,h}$ . The lemma follows since  $\alpha_i(t_j) \leq M_i x^*(t_j) \leq \beta_i(t_j)$  and similarly for  $N_i$ .

As our last step in treating  $r_h$ , we show that  $|f_h(r_h x^*) - f(x^*)|$  converges to zero.

LEMMA 3.2.  $\lim_{h \rightarrow 0} |f_h(r_h x^*) - f(x^*)| = 0$ .

*Proof.* Since the  $m$ -times continuously differentiable functions are dense in  $W^{m,p}$ , we can find such a function  $z$  arbitrarily near  $x^*$  and such that  $|f(x^*) - f(z)|$  is arbitrarily small. Since, for  $0 \leq l \leq m$ ,

$$\begin{aligned} |D^l x^*(t_j) - D^l z(t_j)| &\leq \int_{t_j}^{t_j+lh} K_l(\tau - t_j) |x^{*(l)}(\tau) - z^{(l)}(\tau)| d\tau \\ &\leq \left\{ \int_{t_j}^{t_j+lh} K_l(\tau - t_j) d\tau \right\}^{1/q} \\ &\quad \cdot \left\{ \int_{t_j}^{t_j+lh} K_l(\tau - t_j) |x^{*(l)}(\tau) - z^{(l)}(\tau)|^p d\tau \right\}^{1/p} \\ &\leq \left\{ Gl \int_0^1 |x^{*(l)}(\tau) - z^{(l)}(\tau)|^p d\tau \right\}^{1/p} \end{aligned}$$

which is arbitrarily small, and since the functions  $a_l$  are bounded, it is also clear that  $|f_h(r_h x^*) - f_h(r_h z)|$  can be made arbitrarily small independent of  $h$  by choosing  $z$  near  $x^*$ . Thus we are through if, after fixing  $z$ , we can show that  $|f(z) - f_h(r_h z)|$  tends to zero. By using the triangle inequality, we immediately find

$$\begin{aligned} |f_h(r_h z)^{1/p} - f(z)^{1/p}|^p &\leq \sum_{i=0}^{Q-m} \int_{ih}^{ih+h} \left| \sum_{l=0}^{m-1} [a_l(t_i) D^l z(t_i) - a_l(t) z^{(l)}(t)] \right|^p dt \\ &\quad + \int_{1-mh}^1 \left| \sum_{l=0}^{m-1} a_l(t) z^{(l)}(t) \right|^p dt, \end{aligned}$$

the latter term of which clearly tends to zero with  $h$  for fixed  $z$ . For the former term, since  $z^{(l)}$  is continuous, we have  $D^l z(t_i) = z^{(l)}(\lambda_i)$  for some  $\lambda_i$  in  $(t_i, t_i + lh)$ . Then the former term equals

$$\sum_{i=0}^{Q-m} \int_{ih}^{ih+h} \left| \sum_{l=0}^{m-1} [a_l(t_i) z^{(l)}(\lambda_i) - a_l(t) z^{(l)}(t)] \right|^p dt.$$

Since  $a_l$  and  $z^{(l)}$  are both continuous, the term under the integral sign for this fixed  $z$  is bounded by some function  $w(h)$  tending to zero with  $h$  and thus the whole expression is bounded by  $\sum_{i=0}^{Q-m} \int_{ih}^{ih+h} w(h) dt \leq w(h)$ .

*Remark.* The preceding lemma is of some independent interest. As a special case, it says that  $h \sum_{i=0}^{Q-m} |D^m x(t_i)|^p$  converges to  $\int_0^1 |x^{(m)}(t)|^p dt$  for all  $x$  in  $W^{m,p}$ ; since the sum looks something like a Riemann sum for the integral, it is interesting that convergence can be proved. This is vital for the work in this paper since [10] did not give broad necessary continuity conditions for  $x^*$ ; we know of examples in which  $x^{*(m)}$  has countably infinitely many finite jumps although the constraining functions are very smooth.

Next, we must consider a mapping  $p_h$  of  $x_h^*$  into  $p_h x_h^* = z_h \in W^{m,p}$  near  $C$  with  $|f(p_h x_h^*) - f_h(x_h^*)|$  converging to zero. Let  $v_h$  be an  $m$ -vector function on the mesh

points,  $v_h = (v_{h,0}, \dots, v_{h,m-1})^T$ , solving

$$\begin{aligned} v_{h,0}(t_{j+1}) &= v_{h,0}(t_j) + hv_{h,1}(t_j), \\ &\vdots \\ (3.1) \quad v_{h,m-2}(t_{j+1}) &= v_{h,m-2}(t_j) + hv_{h,m-1}(t_j), \end{aligned}$$

$$v_{h,m-1}(t_{j+1}) = v_{h,m-1}(t_j) + h \left[ - \sum_{i=0}^{m-1} (a_i(t_j)v_{h,i}(t_j)) + L_h x_h^*(t_j) \right]$$

for  $0 \leq j \leq Q - m$ , with  $v_{h,i}(0) = D^i x_h^*(0)$  for  $0 \leq i \leq m - 1$ . For convenience we have assumed  $a_m(t) \equiv 1$  without loss of generality. Clearly then, we have that

$$v_{h,i}(t_j) = D^i x_h^*(t_j) \quad \text{for } 0 \leq i \leq m - 1 \quad \text{and} \quad 0 \leq j \leq Q - m.$$

Consider the  $m$ -vector function  $V_h$  on  $[0, 1]$ ,  $V_h = (V_{h,0}, \dots, V_{h,m-1})^T$ , solving the system of differential equations

$$\begin{aligned} (3.2) \quad V_{h,0}^{(1)} &= V_{h,1}, \\ &\vdots \\ V_{h,m-2}^{(1)} &= V_{h,m-1}, \\ V_{h,m-1}^{(1)} &= - \sum_{i=0}^{m-1} a_i V_{h,i} + u_h \end{aligned}$$

with  $V_{h,i}(0) = v_{h,i}(0)$  for  $i = 0, \dots, m - 1$ , and  $u_h(t) = L_h x_h^*(t_i)$  for  $t_i \leq t < t_{i+1}$  and  $0 \leq i \leq Q - m$ ,  $u_h(t) = 0$  for  $t \geq 1 - (m - 1)h$ . We see immediately that the  $v_h$  is obtained by applying Euler's method to solve the system in (3.2) which, for convenience, we write as

$$(3.3) \quad V_h^{(1)} = AV_h + eu_h, \quad V_h(0) = v_h(0),$$

where  $A$  is the obvious matrix and  $e = (0, 0, \dots, 0, 1)^T$ .

We now define  $z_h = p_h x_h^* \equiv V_{h,0}$ . We notice that  $z_h = V_{h,0}$  solves the equation  $Lz_h = u_h$  and thus

$$(3.4) \quad \int_0^1 |Lz_h(t)|^p dt = \int_0^1 |u_h(t)|^p dt = h \sum_{i=0}^{Q-m} |L_h x_h^*(t_i)|^p,$$

that is,  $f(p_h x_h^*) = f_h(x_h^*)$ . Thus we have accomplished the goal of making  $|f(p_h x_h^*) - f_h(x_h^*)|$  tend to zero; we now check to see if  $z_h = p_h x_h^*$  is "near"  $C$  by relating  $V_h$  to  $v_h$ .

Now we write

$$V_h(t_{j+1}) = V_h(t_j) + \int_{t_j}^{t_{j+1}} [A(t)V(t) + eu_h(t)] dt$$

and

$$v_h(t_{j+1}) = v_h(t_j) + \int_{t_j}^{t_{j+1}} [A(t)v_h(t) + eu_h(t)] dt.$$



Letting  $e_h(t_j) = V_h(t_j) - v_h(t_j)$  and arguing in the usual way, we find, writing

$$\|e_h(t_j)\|_\infty = \max_{0 \leq i \leq m-1} |e_{h,i}(t_j)| \quad \text{and} \quad F = \max_{0 \leq t \leq 1} \|A(t)\|_\infty,$$

$$\|e_h(t_j)\|_\infty \leq \frac{\exp(F) - 1}{(1 + 1/q)F} h^{1/q} \int_0^1 |(A(t)V_h(t))^{(1)}|^p dt.$$

Thus,  $v_h$  and  $V_h$  will be uniformly close if  $AV_h \in W^{1,p}$  and is uniformly bounded in  $W^{1,p}$ . If  $A^{(1)} \in C[0, 1]$ , that is  $A \in C^1[0, 1]$ , then since  $(AV_h)^{(1)} = A^{(1)}V_h + A^2V_h + Ae u_h$  and  $\|A^{(1)}\|_\infty, \|A\|_\infty$  and  $\int_0^1 |u_h(t)|^p dt$  are uniformly bounded,  $AV_h$  will be uniformly bounded in  $W^{1,p}$  if  $V_h$  is uniformly bounded in  $L^p(0, 1)$ ; this finally is clearly true if  $V_h(0)$  is uniformly bounded in  $\mathcal{R}^m$ .

LEMMA 3.3. *If the coefficients  $a_i$  defining the operator  $L$  are in  $C^1[0, 1]$ , then there exists a constant  $K$  such that  $\|v_h(t_j) - V_h(t_j)\|_\infty \leq Kh^{1/q}$  for  $0 \leq j \leq Q - m$ ,  $1/p + 1/q = 1$ . Also  $\|V_{h,0}\| = \|p_h x_h^*\|$  is uniformly bounded.*

*Proof.* Because of the preceding arguments, we need only show that  $V_h(0) = v_h(0)$  is uniformly bounded in  $\mathcal{R}^m$ . Because of (3.1) we can write  $v_h(t_j)$  via

$$(3.5) \quad v_h(t_j) = h \sum_{i=0}^{j-1} [I + hA(t_{j-1})] \cdots [I + hA(t_{i+1})] u_h(t_i) \\ + [I + hA(t_{j-1})] \cdots [I + hA(t_0)] v_h(0).$$

Consider the term in this expression involving the sum; this, call it  $w_h$ , solves (3.1) with  $w_h(0) = 0$  and is therefore within  $O(h^{1/q})$  of the solution  $W_h$  to (3.3) with  $W_h(0) = 0$  by our preceding arguments. Note that, by applying the operators  $M_{l,h}$  at  $\theta_{l,r}$ , and  $N_{r,h}$  at  $\xi_r$  to the first components of the vectors on both sides of (3.5), we immediately see that  $v_h(0)$  solves a certain system of linear equations. Applying one of these operators to the first component of the left-hand side yields merely that operator applied to  $x_h^*$ , and these values are uniformly bounded at the  $\theta_{l,r}$  and  $\xi_r$ . Applying an operator  $D^l$  for  $0 \leq l \leq m - 1$  to  $w_{h,0}$  merely gives  $w_{h,l}$ , which is uniformly close to  $W_{h,l} = W_{h,0}^{(l)}$ , which is uniformly bounded; it then follows that application of one of the operators  $M_{l,h}$  or  $N_{r,h}$  to  $w_{h,0}$  gives uniformly bounded values.

Thus we have found that  $v_h(0)$  solves a linear system with right-hand side uniformly bounded in  $\mathcal{R}^m$ . A typical row in the matrix  $B_h$  of this system consists of, say,  $M_{l,h}$  applied at  $\theta_{l,r}$  to the components of the first row of the matrix function whose value at  $t_j$  is

$$[I + hA(t_{j-1})] \cdots [I + hA(t_0)].$$

Arguing as we have done above, it is easy to show that such an expression converges uniformly to the row (the collection of which forms a matrix  $B$ ) consisting of the application of  $M_l$  at  $\theta_{l,r}$  to the components of the first row of the matrix function whose value at  $t$  is

$$\exp \left[ \int_0^t A(\tau) d\tau \right].$$

A matrix  $B$  of such rows, however, must be of full rank, since by assumption there are no nonzero functions  $x \in W^{m,p}$  such that  $\|x\| = 0$ . If we only apply those operators at those points which in the limit give an  $m \times m$  nonsingular matrix, as we can always do since  $\text{rank}(B) = m$ , then for small  $h$ , the matrices multiplying  $v_h(0)$  are uniformly nonsingular, and therefore the  $v_h(0) = V_h(0)$  are uniformly bounded in  $\mathcal{R}^m$ . Since  $V_h(0)$  is uniformly bounded, it follows that  $\|V_{h,0}\|$  is also.

We can now prove convergence for the more complex discretization.

**THEOREM 3.1.** *Let the general assumptions of § 1 hold and let the fixed points  $\theta_{i,i}$  defining the norm  $\|\cdot\|$  in (2.1) (or the  $\theta_i$  in (2.2)) be mesh points in our discretization for all  $h$ . Let the problem in (1.1) not admit solutions of arbitrarily large  $W^{m,p}$ -norm, for example, some  $M_{i,x} \equiv x$ . Let  $x_h^* \in \mathcal{R}^{Q+1}$  solve the problem in (1.3), that is, minimize  $f_h(x_h)$  over  $C_2(h)$ , where  $\varepsilon_h \geq Gh^{1/q}$  and  $G$  is defined in Lemma 3.1; one may thus take  $h^{1/q} = o(\varepsilon_h)$  for small  $h$ . Suppose the functions  $a_i$  defining  $L$  lie in  $C^1[0, 1]$ . Let  $p_h x_h^* = z_h$  solve*

$$Lz_h = \begin{cases} L_h x_h^*(t_i) & \text{for } t_i \leq t < t_{i+1}, 0 \leq i \leq Q - m, \\ 0 & \text{for } t \geq 1 - (m - 1)h. \end{cases}$$

Then  $f_h(x_h^*)$  converges to  $f(x^*)$  and all  $W^{m,p}$  weak limit points, at least one of which exists, of  $\{p_h x_h^*\}$ , minimize  $f$  over  $C$ ; if  $x^*$  minimizing  $f$  over  $C$  is unique, then  $p_h x_h^*$  converges weakly to  $x^*$ . If (some subsequence of)  $p_h x_h^*$  converges weakly to a point  $x$ , then  $x_h^*$  and its first  $m - 1$  difference approximations  $D^l x_h^*$  evaluated at the points  $t_i = ih$ ,  $0 \leq i \leq Q - l$ , converge uniformly, that is, in  $C[0, 1]$ , to  $x$  and its first  $m - 1$  derivatives at the points  $t_i$ .

*Proof.* By Lemma 3.1 and the hypothesis on  $\varepsilon_h$ ,  $r_h x^* \in C_2(h)$ , so  $C_2(h)$  is not empty. Since  $C_2(h)$  is not empty,  $x_h^*$  exists. By Lemmas 3.3 and 2.2 and the facts that  $v_{h,0} = x_h^* \in C_2(h)$  and  $\|p_h x_h^*\|$  is uniformly bounded, there exist functions  $\eta_i(h)$  tending to zero with  $h$  and such that  $p_h x_h^* \in C_h \equiv \{x; x \in W^{m,p}, -\eta_i(h) - \varepsilon_h - Kh^{1/q} + \alpha_i(t) \leq M_i x(t) \leq \beta_i(t) + \eta_i(h) + \varepsilon_h + Kh^{1/q} \text{ for } 1 \leq i \leq k, -\varepsilon_h - Kh^{1/q} + \gamma_i \leq N_i x(\xi_i) \leq \delta_i + \varepsilon_h + Kh^{1/q} \text{ for } 1 \leq i \leq n\}$ . By Lemma 2.4,  $\zeta_h \equiv \min_C f - \min_{C_h} f$  tends to zero. We write

$$f(x^*) = \min_C f = \min_{C_h} f + \zeta_h \leq f(p_h x_h^*) + \zeta_h = f_h(x_h^*) + \zeta_h,$$

the last equality following from the construction of  $p_h$ . Thus, we have

$$\begin{aligned} f(x^*) &\leq f(p_h x_h^*) + \zeta_h = f_h(x_h^*) + \zeta_h \leq f_h(r_h x^*) + \zeta_h \\ &\leq f(x^*) + \zeta_h + [f_h(r_h x^*) - f(x^*)]. \end{aligned}$$

From Lemma 3.2 and this inequality, we conclude that  $f(x^*) = \lim_{h \rightarrow 0} f(p_h x_h^*) = \lim_{h \rightarrow 0} f_h(x_h^*)$ . Since, from Lemma 3.3,  $p_h x_h^*$  is bounded, it has weak limit points; for any such weak limit point  $x'$  we have  $x' \in C$  by Lemma 2.3 and thus

$$f(x^*) \leq f(x') \leq \liminf_{h \rightarrow 0} f(p_h x_h^*) = f(x^*),$$

which says that  $x'$  minimizes  $f$  over  $C$ . If  $p_h x_h^*$  converges to some  $x$  weakly in  $W^{m,p}$ , then  $p_h x_h^*$  and its first  $m - 1$  derivatives converge uniformly to  $x$  and its first  $m - 1$  derivatives. By Lemma 3.3, the numbers  $v_{h,i}(t_j) = D^l v_{h,0}(t_j) = D^l x_h^*(t_j)$  are uniformly close to  $V_{h,i}(t_j) = V_{h,0}^{(l)}(t_j) = p_h x_h^{*(l)}(t_j)$  for  $0 \leq l \leq m - 1$ .

We have not been able to estimate the rate of convergence as a function of  $h$ .

**4. An elementary example.** Consider the example in [10] with  $m = 1$ ,  $Lx \equiv x^{(1)}$ ,  $k = 1$ ,  $M_1x \equiv x$ ,  $\alpha_1(t) = t - t^2$ ,  $\beta_1(t) = t$ ,  $n = 1$ ,  $N_1x(\xi_1) = x(1)$ ,  $\delta_1 = \gamma_1 = c \in [0, 1]$ , that is,

$$\text{minimize } \int_0^1 |x^{(1)}(t)|^2 dt \quad \text{over}$$

$$C = \{x; x \in W^{1,2}, t - t^2 \leq x(t) \leq t, x(1) = c\}.$$

The unique solution to this problem is

$$x^*(t) = \begin{cases} t - t^2 & \text{for } 0 \leq t \leq 1 - \sqrt{c}, \\ (2\sqrt{c} - 1)t + (1 + c - 2\sqrt{c}) & \text{for } 1 - \sqrt{c} \leq t \leq 1, \end{cases}$$

as pictured in Fig. 4.1.

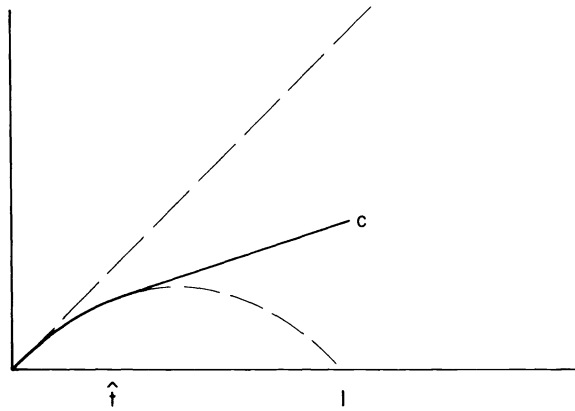


FIG. 4.1

Let  $\hat{t} = 1 - \sqrt{c}$ , the point at which  $x^*$  leaves the lower curve.

If we use the simple discretization and merely discretize the constraints at  $t_i = ih$ ,  $0 \leq i \leq Q = 1/h$ , then if  $\hat{t}_i = \max \{t_i; t_i \leq \hat{t}\}$ , then  $x_Q^*$  is just the piecewise linear interpolant of  $t - t^2$  at  $t_i$  for  $t \leq \hat{t}_i$  and is the linear interpolant between  $\hat{t}_i - \hat{t}_i^2$  and  $c$  for  $\hat{t}_i \leq t \leq 1$ . For small enough  $h$ , the solution for the complete discretization is also unique; such discrete variational splines are studied in [11].

If we define, for  $c > 1/2$ , the numbers

$$\alpha_h = \min \{t_i; t_i \geq \sqrt{2\varepsilon_h}\}, \quad \beta_h = \max \{t_i; t_i \leq 1 - \sqrt{c}\},$$

the unique solution  $x_h^*$  for the complete discretization with  $\varepsilon_h = h^{1/2-\delta}$ ,

$\delta \in (0, 1/2)$ , is

$$x_h^*(t_i) = \begin{cases} \varepsilon_h + t_i \left[ \frac{\alpha_h - \alpha_h^2 - 2\varepsilon_h}{\alpha_h} \right] & \text{if } 0 \leq t_i \leq \alpha_h, \\ t_i - t_i^2 - \varepsilon_h & \text{if } \alpha_h \leq t_i \leq \beta_h, \\ \beta_h - \beta_h^2 - \varepsilon_h + (t_i - \beta_h) \left[ \frac{c - \beta_h + \beta_h^2}{1 - \beta_h} \right] & \text{if } \beta_h \leq t_i \leq 1. \end{cases}$$

**Acknowledgment.** The author thanks O. L. Mangasarian and L. L. Schumaker for their suggestion of the problem and for their continuing interest.

#### REFERENCES

- [1] M. ATTEIA, *Fonctions "spline" définies sur un ensemble convexe*, Numer. Math., 12 (1968), pp. 192–210.
- [2] J. CULLUM, *Discrete approximations to continuous optimal control problems*, this Journal, 7 (1969), pp. 32–49.
- [3] ———, *An explicit procedure for discretizing continuous optimal control problems*, J. Optimization Theory and Applications, to appear.
- [4] H. B. CURRY AND I. J. SCHOENBERG, *On Polya frequency functions. IV: The fundamental spline functions and their limits*, J. Analyse Math., 17 (1966), pp. 71–107.
- [5] J. W. DANIEL, *On the approximate minimization of functionals*, Math. Comp., 23 (1969), pp. 573–582.
- [6] ———, *On the convergence of a numerical method for optimal control problems*, J. Optimization Theory and Applications, 4 (1969), pp. 330–342.
- [7] ———, *The Approximate Minimization of Functionals*, Prentice-Hall, Englewood Cliffs, N.J., to appear.
- [8] M. GOLUMB AND J. W. JEROME, *Ordinary differential equations with boundary conditions on arbitrary sets*, to appear.
- [9] J. W. JEROME AND L. L. SCHUMAKER, *On Lg-splines*, J. Approx. Theor., 2 (1969), pp. 29–49.
- [10] O. L. MANGASARIAN AND L. L. SCHUMAKER, *Splines via optimal control*, Approximations with Special Emphasis on Spline Functions, I. J. Schoenberg, ed., Academic Press, New York, 1969.
- [11] ———, *Discrete splines via mathematical programming*, this Journal, to appear.
- [12] C. B. MORREY, JR., *Multiple Integrals in the Calculus of Variations*, Springer-Verlag, New York, 1966.
- [13] K. RITTER, *Generalized spline interpolation and nonlinear programming*, Approximations with Special Emphasis on Spline Functions, I. J. Schoenberg, ed., Academic Press, New York, 1969.
- [14] S. L. SOBOLEV, *Applications of Functional Analysis in Mathematical Physics*, vol. 7, Trans. Math. Monographs, American Mathematical Society, Providence, R.I., 1963.

## INVARIANCE PROPERTIES IN THE THEORY OF ORDINARY DIFFERENTIAL EQUATIONS WITH APPLICATIONS TO STABILITY PROBLEMS\*

NELSON ONUCHIC†

**1. Introduction.** The main purpose of this paper is to find conditions under which every solution of a second order scalar differential equation tends to zero with its derivative as  $t \rightarrow \infty$ . The basic tool to carry out this objective is provided by Theorems 1, 2 and 3 given below.

Theorem 1 is a modified but closely related version of Yoshizawa's theorem 5 in [6]. Theorem 2 is essentially a special case of Miller's theorem 1 in [3]. Some remarks relating Theorems 1 and 2 and those ones by Yoshizawa and Miller will be done following each theorem. Theorem 3 is a particular case of La Salle's Theorem 1 in [1].

Theorem 4 is a result concerning the above posed problem obtained as an application of Theorems 1 and 2. This result is more general than the one by Yoshizawa in [6, Example 2]. See also [2] and [7, Example 14.1, p. 63]. Theorem 5, previously obtained by us [4], is here stated because it gives additional information on the problem under consideration which is intimately related to Theorem 4. Theorem 6 also gives a contribution to our problem and it is strongly dependent on Theorem 3.

**2. Invariance properties.** Consider a system of differential equations defined on an open set  $Q \subset R^n$ :

$$(1) \quad \dot{x} = H(x),$$

where  $H(x)$  is continuous on  $Q$ .

If  $M$  is a subset of  $Q$ , then  $M$  is called quasi-invariant with respect to (1) if and only if for each  $x_0 \in M$  there is a solution  $x(t)$  of (1) with  $x(0) = x_0$ , such that  $x(t)$  exists and remains in  $M$  for all real  $t$ . Let  $x(t)$  be a continuous function defined in the future, that is, for all  $t \geq$  some real  $t_0$ . A point  $p$  of  $R^n$  is said to be a limit point of  $x(t)$  if there exists a sequence  $\{t_m\}$ ,  $t_m \rightarrow \infty$  as  $m \rightarrow \infty$ , such that  $x(t_m) \rightarrow p$  as  $m \rightarrow \infty$ . The set of all limit points of  $x(t)$  is denoted by  $\Omega$  and is called the  $\omega$ -limit set of  $x(t)$ . If  $x(t)$  is bounded in the future, that is,  $x(t)$  is bounded on some interval  $[a, \infty)$ ,  $a > -\infty$ , it is easily seen that  $\Omega$  is a nonempty, connected and compact set with  $x(t) \rightarrow \Omega$  as  $t \rightarrow \infty$ , that is,  $\text{dist}(x(t), \Omega) \rightarrow 0$  as  $t \rightarrow \infty$ .

Consider the differential system defined on  $[0, \infty) \times Q$ ,  $Q$  being an open set of  $R^n$ ,

$$(2) \quad \dot{x} = f(t, x),$$

where  $f$  is assumed to be continuous for  $t \geq 0$ ,  $x \in Q$ .

Let  $V(t, x)$  be a real-valued  $C^1$  function defined for  $t \geq 0$ ,  $x \in Q$ . Define

$$\dot{V}_{(2)}(t, x) = \sum_{j=1}^n \frac{\partial V(t, x)}{\partial x_j} f_j(t, x) + \frac{\partial V(t, x)}{\partial t}.$$

---

\* Received by the editors June 17, 1970.

† Escola de Engenharia de São Carlos, São Carlos, SP Brazil. This research was supported in part by the Conselho Nacional de Pesquisas, Brazil.

Consider also the differential system defined on  $[0, \infty) \times Q$ :

$$(3) \quad \dot{x} = F(t, x) + G(t, x),$$

where  $F(t, x)$  and  $G(t, x)$  are continuous for  $t \geq 0, x \in Q$ .

THEOREM 1. *Suppose that the following hypotheses hold with respect to system (3):*

- (i)  $F(t, x)$  is bounded for all  $t \geq 0$  when  $x$  belongs to an arbitrary compact subset of  $Q$ .
- (ii) For every compact set  $B \subset Q$  and every continuous function  $z(t) \in B$ , defined on  $[0, \infty)$ , it follows that

$$(4) \quad \int_s^{s+t} G(v, z(v)) dv \rightarrow 0 \quad \text{as } s \rightarrow \infty,$$

uniformly on  $t \in [0, 1]$ .

(iii) *There are a real-valued nonnegative  $C^1$  function  $V(t, x)$  and a real-valued nonnegative continuous function  $W(x)$  such that*

$$\dot{V}_{(3)}(t, x) \leq -W(x), \quad t \geq 0 \quad \text{and} \quad x \in Q.$$

Let  $x(t)$  be a solution of (3), defined in the future, with  $x(t) \in K$  for  $t \geq$  some  $t_0$ , where  $K$  is a compact subset of  $Q$ .

Then  $\Omega \subset E = \{x \in Q | W(x) = 0\}$ , where  $\Omega$  is the  $\omega$ -limit set of  $x(t)$ .

Note. The main difference between the above theorem and Theorem 5 in [6] is that Yoshizawa assumes the condition that  $\int^\infty \|G(v, z(v))\| dv < \infty$ , which is stronger than the one given by (4).

A sufficient condition for (4) is given as follows:

(a) For every compact set  $B \subset Q$  there corresponds a scalar function  $\sigma_B(t)$  defined for  $t \geq 0$  such that

$$\int_t^{t+1} \sigma_B(s) ds \rightarrow 0 \quad \text{as } t \rightarrow \infty,$$

and  $\|G(t, x)\| \leq \sigma_B(t)$  for all  $t \geq 0, x \in B$ .

An example considered in [5] shows that condition (4) is not implied by (a).

The proof of Theorem 1 can be done by following essentially the same ideas contained in the proof of Theorem 5 in [6].

Consider the differential system defined on  $[0, \infty) \times Q$ :

$$(5) \quad \dot{x} = H(x) + S(t, x) + G(t, x),$$

where  $H(x), S(t, x)$  and  $G(t, x)$  are assumed to be continuous on  $[0, \infty) \times Q$ .

Let  $A$  be a fixed subset of  $Q$ . Assume that  $S(t, x)$  satisfies the following property with respect to the set  $A$ :

(b) For each  $\varepsilon > 0$  and each compact subset  $K$  of  $Q$  there corresponds  $\delta = \delta(\varepsilon, K) > 0$  and  $T_0 = T_0(\varepsilon, K) > 0$  such that  $t \geq T_0, x \in K$  and  $\text{dist}(x, A) < \delta$  imply  $\|S(t, x)\| < \varepsilon$ .

**THEOREM 2.** *Let hypotheses (4) and (b) hold. Let  $x(t)$  be a solution of (5) such that  $x(t) \in K$ ,  $T \leq t < \infty$ , where  $T \geq 0$  and  $K$  is a compact subset of  $Q$  and  $x(t) \rightarrow A$  as  $t \rightarrow \infty$ . Then the  $\omega$ -limit set  $\Omega$  of  $x(t)$  is a nonempty, connected, compact and quasi-invariant set with respect to (1).*

*Note.* As observed, this theorem is essentially a particular case of Miller's Theorem 1 in [3]. Miller deals with delay equations where the unperturbed system is uniformly almost periodic. However, condition (4) considered in Theorem 2 is weaker than the one in Miller's theorem. But Miller's proof also works well with condition (4).

Consider the system

$$(6) \quad \dot{x} = F(t, x),$$

where  $F$  is continuous for  $t \geq 0$ ,  $x \in R^n$  and satisfies any one of the conditions guaranteeing uniqueness of solutions.

**THEOREM 3.** *Let  $V(t, x)$  and  $W(x)$  be real-valued nonnegative  $C^1$  functions on  $[0, \infty) \times R^n$  such that  $\dot{V}_{(8)}(t, x) \leq -W(x)$  for all  $t > 0$  and  $x \in G$ , where  $G$  is a compact set of  $R^n$ . Let  $x(t)$  be a solution of (6) such that  $x(t) \in G$  for all  $t > t_0$  and let  $\dot{W}$  be bounded from above or below along the solution  $x(t)$ . Then  $\Omega \subset E = \{x \in G \mid W(x) = 0\}$ , where  $\Omega$  is the  $\omega$ -limit set of  $x(t)$ .*

**3. Applications.** The objective of this section is to study the problem posed in § 1 by using the results discussed in § 2. Specifically, our purpose is to find conditions under which, for all solutions  $x(t)$  of the second order scalar differential equation

$$(7) \quad \ddot{x} + h(t, x, \dot{x})\dot{x} + f(x) + g(t, x, \dot{x}) + p(t, x, \dot{x}) = 0,$$

we can guarantee that  $(x(t), \dot{x}(t)) \rightarrow (0, 0)$  as  $t \rightarrow \infty$ . To this end, consider the equivalent equation

$$(7') \quad \begin{aligned} \dot{x} &= y, \\ \dot{y} + h(t, x, y)y + f(x) + g(t, x, y) + p(t, x, y) &= 0 \end{aligned}$$

and the following set of assumptions with respect to the functions in (7').

( $H_1$ )  $h(t, x, y)$  is continuous in  $[0, \infty) \times R^2$  and  $h(t, x, y)$  is bounded when  $x^2 + y^2$  is bounded, and, moreover,  $h(t, x, y) \geq k(x, y) \geq 0$ , where  $k(x, y)$  is a continuous function in  $R^2$ .

( $H_2$ ) The sets  $R^+ = \{(x, 0) \mid x > 0\}$  and  $R^- = \{(x, 0) \mid x < 0\}$  are connected components with respect to the topological space  $\Gamma - \{(0, 0)\}$ , where  $\Gamma = \{(x, y) \in R^2 \mid yk(x, y) = 0\}$ .

( $H_3$ )  $f(x)$  is continuous,  $xf(x) > 0$  for all  $x \neq 0$ , and

$$\int_0^x f(s) ds \rightarrow \infty \quad \text{as } |x| \rightarrow \infty.$$

( $H_3$ )  $f(x)$  is continuous in  $R$  and there is a positive  $\rho$  such that  $\int_0^x f(s) ds > 0$  for  $0 < |x| \leq \rho$ .

$(H_3'')$   $f(x)$  is continuous in  $R$  and  $\int_0^x f(s) ds \rightarrow \infty$  as  $|x| \rightarrow \infty$ .

We see that  $(H_3)$  implies  $(H_3')$  and  $(H_3)$  implies  $(H_3'')$ .

$(H_4)$   $p(t, x, y)$  is continuous and  $|p(t, x, y)| \leq \beta(t)$  for all  $t \geq 0, x, y$  in  $R$ , where  $\beta(t)$  is continuous with  $\int_0^\infty \beta(t) dt < \infty$ .

$(H_5)$   $g(t, x, y)$  is continuous and  $yg(t, x, y) \geq 0$  for all  $t \geq 0, x, y$  in  $R$ .

$(H_6)$  For every compact subset  $B$  of  $R$  and for all continuous functions  $x(t)$  and  $y(t)$ , defined on  $[0, \infty)$  with values in  $B$ , it follows that

$$\int_s^{s+t} g(v, x(v), y(v)) dv \rightarrow 0 \quad \text{as } s \rightarrow \infty,$$

uniformly on  $t \in [0, 1]$ .

A sufficient condition for  $(H_6)$  is given as follows:

$(\tilde{H}_6)$  For every compact subset  $B$  of  $R$  there corresponds a real function  $\sigma_B(t)$ , defined for  $t \geq 0$ , such that

$$\int_t^{t+1} \sigma_B(s) ds \rightarrow 0 \quad \text{as } t \rightarrow \infty,$$

and  $|g(t, x, y)| \leq \sigma_B(t)$  for all  $t \geq 0, x, y$  in  $B$ .

LEMMA 1. Let hypotheses  $(H_3')$ ,  $(H_4)$  and  $(H_5)$  hold. Let  $h(t, x, y)$  be non-negative and continuous for  $t \geq 0, x, y$  in  $R$ . Then for each  $\varepsilon > 0$  there are positive  $T(\varepsilon)$  and  $\delta(\varepsilon)$  such that  $t \geq t_0 \geq T(\varepsilon)$  and  $(x_0, y_0) \in R^2$ , with  $|x_0| + |y_0| < \delta$ , imply  $|x(t)| + |y(t)| < \varepsilon$  for every solution  $(x(t), y(t))$  of (7') satisfying  $x(t_0) = x_0$  and  $y(t_0) = y_0$ .

*Proof.* Consider the function

$$V(t, x, y) = \left[ y^2 + 2 \int_0^x f(s) ds \right]^{1/2} + \int_t^\infty \beta(s) ds$$

defined for  $t \geq 0$  and  $|x| + |y| < \rho$ .

Given  $0 < \varepsilon < \rho$ , define

$$2m = \min_{|x| + |y| = \varepsilon} \left[ y^2 + 2 \int_0^x f(s) ds \right]^{1/2} > 0.$$

Then condition  $(H_3)$  implies

$$\inf_{\substack{t \geq 0 \\ |x| + |y| = \varepsilon}} V(t, x, y) \geq 2m > 0.$$

Let  $T(\varepsilon) > 0$  and  $0 < \delta(\varepsilon) < \varepsilon$  be chosen such that

$$\int_{T(\varepsilon)}^\infty \beta(t) dt > m \quad \text{and} \quad \left[ y^2 + 2 \int_0^x f(s) ds \right]^{1/2} < m$$

for  $|x| + |y| \leq \delta(\varepsilon)$ . Then

$$\max_{\substack{t \geq T(\varepsilon) \\ |x| + |y| = \delta(\varepsilon)}} V(t, x, y) < 2m.$$



An easy computation shows that

$$\dot{V}_{(7')}(t, x, y) \leq \frac{-y^2 h(t, x, y) - yg(t, x, y)}{\left[ y^2 + 2 \int_0^x f(s) ds \right]^{1/2}} + |p(t, x, y)| - \beta(t) \leq 0$$

for all  $t \geq 0$ ,  $0 < |x| + |y| < \rho$ .

Suppose that there are  $t_0 \geq T(\varepsilon)$  and  $(x_0, y_0) \in R^2$  with  $|x_0| + |y_0| < \delta(\varepsilon)$ , such that for some solution  $(x(t), y(t))$  of (7') satisfying  $(x(t_0), y(t_0)) = (x_0, y_0)$  and for some  $\tilde{t} > t_0$  we have  $|x(\tilde{t})| + |y(\tilde{t})| \geq \varepsilon$ . Then there are real numbers  $t_1$  and  $t_2$ ,  $t_0 < t_1 < t_2$ , such that  $|x(t_1)| + |y(t_1)| = \delta(\varepsilon)$ ,  $|x(t_2)| + |y(t_2)| = \varepsilon$  and  $\delta(\varepsilon) \leq |x(t)| + |y(t)| \leq \varepsilon$  for  $t_1 \leq t \leq t_2$ . But this implies  $V(t_1, x(t_1), y(t_1)) < V(t_2, x(t_2), y(t_2))$  and, since  $(x(t), y(t)) \neq (0, 0)$  on  $[t_1, t_2]$ , we have that  $V(t, x(t), y(t))$  is differentiable in  $[t_1, t_2]$ . Hence there is a real number  $s$ ,  $t_1 < s < t_2$ , satisfying  $\dot{V}(s, x(s), y(s)) > 0$ , leading to a contradiction. The proof is complete.

**COROLLARY 1.** *Suppose that some uniqueness condition with respect to the initial value problem holds for (7'). Suppose that all hypotheses of Lemma 1 are satisfied with  $p(t, 0, 0) = 0$  for  $t \geq 0$ . Then the null solution of (7') is uniformly stable.*

*Proof.* Corollary 1 follows from Lemma 1, taking into account that the origin is an equilibrium point of (7') and considering the usual continuous dependence argument.

We observe that  $g(t, 0, 0) = 0$  is a consequence of hypothesis  $(H_5)$ .

**LEMMA 2.** *Let hypotheses  $(H_3')$ ,  $(H_4)$  and  $(H_5)$  hold. Let  $h(t, x, y)$  be non-negative and continuous for  $t \geq 0$ ,  $x, y$  in  $R$ . Then every solution of (7') is bounded in the future.*

*Proof.* Let  $V(t, x, y) = \left[ y^2 + 2 \int_0^x f(s) ds + M \right]^{1/2} + \int_t^\infty \beta(s) ds$ , where  $M$  is chosen so that  $2 \int_0^x f(s) ds + M > 0$  for all  $x$  in  $R$ . It is easy to see that  $\dot{V}_{(7')}(t, x, y) \leq 0$  for all  $t \geq 0$  and  $x, y$  in  $R$ . Then as  $W(x, y) = \left[ y^2 + 2 \int_0^x f(s) ds + M \right]^{1/2} \leq V(t, x, y)$  and  $W(x, y) \rightarrow \infty$  as  $x^2 + y^2 \rightarrow \infty$  it follows that every solution of (7') is bounded in the future.

The proof is complete.

**THEOREM 4.** *Let hypotheses  $(H_1)$  through  $(H_6)$  hold. Then for every solution  $x(t)$  of (7) we have that  $x(t) \rightarrow 0$  and  $\dot{x}(t) \rightarrow 0$  as  $t \rightarrow \infty$ .*

*Proof.* Let  $(x(t), y(t))$  be any solution of (7'). Lemma 2 implies that this solution is bounded in the future. Then we know that the  $\omega$ -limit set  $\Omega$  of  $(x(t), y(t))$  is a nonempty, connected and compact set with  $(x(t), y(t)) \rightarrow \Omega$  at  $t \rightarrow \infty$ . We must have  $\Omega \cap R_x \neq \emptyset$ , where  $R_x$  is the  $x$ -axis, because, if this were not true, it would follow that  $|x(t)| \rightarrow \infty$  as  $t \rightarrow \infty$ , a contradiction.

Defining  $V(t, x, y) = \left[ y^2 + 2 \int_0^x f(s) ds + 1 \right]^{1/2} + \int_t^\infty \beta(s) ds$  and

$$W(x, y) = \frac{y^2 k(x, y)}{\left[ y^2 + 2 \int_0^x f(s) ds + 1 \right]^{1/2}},$$

a computation shows, by taking into account hypothesis  $(H_1)$ , that

$$\dot{V}_{(7)}(t, x, y) \leq -W(x, y) \leq 0$$

for  $t \geq 0, (x, y) \in Q = R^2$ . Then it follows from Theorem 1 that

$$\Omega \subset \Gamma = \{(x, y) \in R^2 | yk(x, y) = 0\}.$$

We claim that  $(0, 0) \in \Omega$ . Indeed, otherwise, since  $\Omega \cap R_x \neq \emptyset$ , it would follow that  $\Omega \subset \Gamma - \{(0, 0)\}$  and  $\Omega \cap (R^+ \cup R^-) \neq \emptyset$ . But, since  $\Omega$  is a connected set and  $(H_2)$  holds, it would follow that  $\Omega \subset R^+ \cup R^-$ . Then, by applying Theorem 2 with  $A = R_x, H = (y, -f(x)), S = (0, -h(t, x, y)y), G = (0, -g(t, x, y) - p(t, x, y))$  and  $Q = R^2$ , it would follow that  $\Omega$  is quasi-invariant with respect to the system

$$(8) \quad \begin{aligned} \dot{x} &= y, \\ \dot{y} + f(x) &= 0. \end{aligned}$$

But condition  $(H_3)$  implies that the unique quasi-invariant set for (8) contained in  $R_x$  is  $(0, 0)$ , leading to a contradiction. Thus  $(0, 0) \in \Omega$ , and hence there exists a sequence  $\{t_m\}, t_m \rightarrow \infty$  as  $m \rightarrow \infty$ , such that  $(x(t_m), y(t_m)) \rightarrow (0, 0)$  as  $m \rightarrow \infty$ . Then it follows from Lemma 1 that  $(x(t), y(t)) \rightarrow (0, 0)$  as  $t \rightarrow \infty$ , completing the proof.

**COROLLARY 2.** *Suppose that some uniqueness condition with respect to the initial value problem holds for (7'). Let hypotheses  $(H_1)$  through  $(H_6)$  hold and  $p(t, 0, 0) = 0$ . Then the origin is globally asymptotically stable for (7').*

*Proof.* This corollary follows from Corollary 1 and Theorem 4.

The next theorem, which is a result previously obtained by us [4, Theorem], gives sufficient conditions to guarantee that for every solution  $x(t)$  of the scalar equation

$$(9) \quad \ddot{x} + h(x, \dot{x}) + f(x) + g(t, x, \dot{x}) + p(t, x, \dot{x}) = 0,$$

we have  $x(t) \rightarrow 0$  and  $\dot{x}(t) \rightarrow 0$  as  $t \rightarrow \infty$ . Most of the ideas used in the proof of Theorem 5 are closely related to the ones in the proof of Theorem 4, but Theorem 5 depends also on the well-known Poincaré-Bendixon theorem for two-dimensional autonomous systems.

**THEOREM 5.** *Suppose that the following hypotheses hold with respect to (9):*

- (i)  $h(x, y)$  is continuous and  $yh(x, y) \geq 0$  for all real numbers  $x$  and  $y$ .
- (ii) For every Jordan curve  $\gamma$  in  $R^2$  containing the origin in its interior, there exists at least one point  $(x, y) \in \gamma$  such that  $y \neq 0$  and  $h(x, y) \neq 0$ .
- (iii)  $g(t, x, y)$  is continuous with  $y[h(x, y) + g(t, x, y)] \geq 0$  for all  $t \geq 0$  and all real numbers  $x$  and  $y$ .
- (iv) Hypotheses  $(H_3), (H_4)$  and  $(H_6)$  hold.

*Then for every solution  $x(t)$  of (9) we have that  $x(t) \rightarrow 0$  and  $\dot{x}(t) \rightarrow 0$  as  $t \rightarrow \infty$ .*

**COROLLARY 3.** *Suppose that some uniqueness condition with respect to the initial value problem holds for*

$$(9') \quad \begin{aligned} \dot{x} &= y, \\ \dot{y} + h(x, y) + f(x) + g(t, x, y) + p(t, x, y) &= 0. \end{aligned}$$

Suppose that all hypotheses of Theorem 5 are satisfied, with  $p(t, 0, 0) = 0$  for all  $t \geq 0$ . Then the origin is globally asymptotically stable for (9').

**THEOREM 6.** Let  $h(t, x, y)$  be continuous for  $t \geq 0, x, y$  in  $R$  and moreover,  $h(t, x, y) \geq k(x, y) \geq 0$ , where  $k(x, y)$  is a  $C^1$  function with  $y k_y(x, y) \geq 0$  for all  $x, y$  in  $R$ . (Here  $k_y$  denotes the partial derivative with respect to  $y$ .) Let hypotheses  $(H_2)$ ,  $(H'_3)$ ,  $(H''_3)$  and  $(H_5)$  hold. Let hypothesis  $(H_4)$  be satisfied with  $\beta(t)$  bounded on  $(0, \infty)$ . Suppose that any one of the conditions guaranteeing uniqueness of solutions is satisfied with respect to (7').

Then every solution  $x(t)$  of (7) is bounded in the future and  $\dot{x}(t) \rightarrow 0$  as  $t \rightarrow \infty$ . Furthermore if  $x(t)$  is a solution of (7) such that  $\liminf_{t \rightarrow \infty} |x(t)| = 0$ , then also  $x(t) \rightarrow 0$  as  $t \rightarrow \infty$ .

*Proof.* Let  $x(t)$  be any solution of (7). Then it follows from Lemma 2 that  $(x(t), y(t))$  is bounded in the future. Hence, there are  $t_0 \geq 0$  and a compact set  $G$  of  $R^2$  such that  $(x(t), y(t)) \in G$  for all  $t > t_0$ .

Define  $V(t, x, y)$  as in Lemma 2, that is,

$$V(t, x, y) = \left[ y^2 + 2 \int_0^x f(s) ds + M \right]^{1/2} + \int_t^\infty \beta(s) ds.$$

It is easily seen that

$$\dot{V}_{(7')} (t, x, y) \leq - \frac{y^2 k(x, y)}{\left[ y^2 + 2 \int_0^x f(s) ds + M \right]^{1/2}}.$$

By taking

$$\gamma = \inf_{(x,y) \in G} \frac{1}{\left[ y^2 + 2 \int_0^x f(s) ds + M \right]^{1/2}},$$

it follows that

$$\frac{y^2 k(x, y)}{\left[ y^2 + 2 \int_0^x f(s) ds + M \right]^{1/2}} \geq \gamma y^2 k(x, y)$$

for all  $(x, y) \in G$ . Therefore, defining  $W(x, y) = \gamma y^2 k(x, y)$ , we have that  $\dot{V}_{(7')} (t, x, y) \leq -W(x, y)$  for all  $t > 0, (x, y) \in G$ .

To apply Theorem 3, let us show that  $\dot{W}$  is bounded from above along each solution which remains in  $G$  for all  $t > t_0 \geq 0$ .

Let  $z(t)$  be a solution of (7) such that  $(z(t), \dot{z}(t)) \in G$  for all  $t > t_0$ . Then

$$\begin{aligned} \dot{W}(z(t), \dot{z}(t)) &= \gamma \frac{d}{dt} \{ [\dot{z}(t)]^2 k(z(t), \dot{z}(t)) \} \\ &= 2\gamma \dot{z}(t) \ddot{z}(t) k(z(t), \dot{z}(t)) + \gamma [\dot{z}(t)]^3 k_x(z(t), \dot{z}(t)) + \gamma [\dot{z}(t)]^2 k_y(z(t), \dot{z}(t)) \dot{z}(t) \\ &= \gamma [\dot{z}(t)]^3 k_x(z(t), \dot{z}(t)) + 2\gamma \dot{z}(t) k(z(t), \dot{z}(t)) [-h(t, z(t), \dot{z}(t)) \dot{z}(t) - f(z(t)) \\ &\quad - g(t, z(t), \dot{z}(t)) - p(t, z(t), \dot{z}(t))] \end{aligned}$$

$$\begin{aligned}
 & + \gamma[\dot{z}(t)]^2 k_y(z(t), \dot{z}(t))[-h(t, z(t)), \dot{z}(t)]\dot{z}(t) \\
 & \qquad - f(z(t)) - g(t, z(t), \dot{z}(t)) - p(t, z(t), \dot{z}(t))] \\
 \leq & \gamma[\dot{z}(t)]^3 k_x(z(t), \dot{z}(t)) - 2\gamma\dot{z}(t)k(z(t), \dot{z}(t))f(z(t)) \\
 & + 2\gamma|\dot{z}(t)|k(z(t), \dot{z}(t))\beta_0 - \gamma[\dot{z}(t)]^2 k_y(z(t), \dot{z}(t))f(z(t)) \\
 & + \gamma[\dot{z}(t)]^2 |k_y(z(t), \dot{z}(t))|\beta_0,
 \end{aligned}$$

where  $\beta_0 = \sup_{0 \leq t < \infty} \beta(t)$ .

Then as  $G$  is compact and  $(z(t), \dot{z}(t)) \in G$ , there is a positive number  $C = C(G)$  such that  $\dot{W}(z(t), \dot{z}(t)) \leq C$  for  $t > t_0$ .

Therefore Theorem 3 implies  $\Omega \subset \Gamma = \{(x, y) \in R^2 | yk(x, y) = 0\}$ , where  $\Omega$  is the  $\omega$ -limit set of  $(x(t), y(t))$ . We must have  $\Omega \cap R_x \neq \emptyset$ , because otherwise it would follow that  $|x(t)| \rightarrow \infty$  as  $t \rightarrow \infty$ , a contradiction.

Consider, then, the two possibilities:

(a)  $(0, 0) \in \Omega$ , and

(b)  $\Omega \subset \Gamma - \{(0, 0)\}$  and hence  $\Omega \cap (R^+ \cup R^-) \neq \emptyset$ .

Case (a) implies, taking into account Lemma 1, that  $(x(t), \dot{x}(t)) \rightarrow (0, 0)$  as  $t \rightarrow \infty$ .

Case (b) implies, by using  $(H_2)$ , that  $\Omega \subset R^+ \cup R^-$  and consequently  $\dot{x}(t) \rightarrow 0$  as  $t \rightarrow \infty$ .

In case  $\liminf_{t \rightarrow \infty} |x(t)| = 0$ , we have that  $(0, 0) \in \Omega$  and, consequently,  $(x(t), \dot{x}(t)) \rightarrow (0, 0)$  as  $t \rightarrow \infty$ . The proof is complete.

**Acknowledgment.** The author is indebted to Professor Taro Yoshizawa for his invaluable comments on the original manuscript.

REFERENCES

[1] J. P. LA SALLE, *An invariance principle in the theory of stability*, Differential Equations and Dynamical Systems, Academic Press, New York, 1967, pp. 277-286.

[2] J. J. LEVIN, *On the global asymptotic behavior of nonlinear systems of differential equations*, Arch. Rational Mech. Anal., 6 (1965), pp. 65-74.

[3] R. K. MILLER, *Asymptotic behavior of nonlinear delay-differential equations*, J. Differential Equations, 1 (1965), pp. 293-305.

[4] N. ONUCHIC, *Stability properties of a second order differential equation*, Acta Mexicana de Ciencia y Tecnologia, 3 (1969), pp. 6-11.

[5] A. STRAUSS AND J. A. YORKE, *On asymptotically autonomous differential equations*, Math. Systems Theor., 1, 2 (1967), pp. 175-182.

[6] T. YOSHIZAWA, *Asymptotic behavior of solutions of a system of differential equations*, Contributions to Differential Equations, 1 (1963), pp. 371-387.

[7] ———, *Stability theory by Lyapunov's second method*, J. Math. Soc. Japan, 1966.

## THE GENERATION OF LYAPUNOV FUNCTIONS FOR INPUT-OUTPUT STABLE SYSTEMS\*

JAN C. WILLEMS †

**Abstract.** This paper discusses the relationship between properties of input-output descriptions and state space models for dynamical systems. It is shown that a state space realization of an input-output stable dynamical system is globally asymptotically stable in the sense of Lyapunov if it is uniformly observable and if every state is reachable. This result is proved in the context of abstract dynamical systems and leads to the equivalence of input-output stability and asymptotic stability for uniformly controllable and uniformly observable linear finite-dimensional systems. The generation of Lyapunov functions is subsequently considered, and variational techniques for the construction of Lyapunov functions are presented. Passivity and related energy concepts are particularly exploited in this context. These results yield the Lyapunov functions used in the proofs of the circle criterion and the Popov criterion as particular cases. The generality of the approach, however, makes these ideas applicable to much more general situations. Examples illustrating the results and the unifying point of view are included.

**1. Introduction.** “Dynamical systems” as they are studied and defined<sup>1</sup> in modern system theory distinguish themselves from arbitrary operators in mathematics by one basic property: they are causal, i.e., nonanticipatory: future values of the input do not influence past values of the output. This basic realizability property of physical systems may be incorporated in the mathematical model in two ways: either by appropriately restricting the operator defining the input-output relationship, or by working with a state space description which will then automatically ensure this causality. This last approach has proved particularly useful in optimal control theory due to the fact that any deterministic optimal controller can always be implemented with a memoryless function of the state in the feedback. It is therefore very advantageous to work with a state space model from the very start.

This duality in the possible description of systems has reflected itself in other areas of system theory and is particularly prevalent in stability theory. The input-output approach leads to the concept of input-output stability and has been developed mainly in the last decade, especially following the work of Sandberg [4] and Zames [5]. The state space description leads to concepts such as global asymptotic stability in the sense of Lyapunov and poses the stability problem in a setting which does not involve inputs, thus making use of the theory of classical dynamical systems. Which of the two approaches is to be preferred depends on

---

\* Received by the editors January 20, 1970.

† Electronic Systems Laboratory, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139. Now at Department of Applied Mathematics and Theoretical Physics, University of Cambridge, Cambridge, England. This research was supported by the National Aeronautics and Space Administration under Contract NGL-22-009-124 and by the National Science Foundation under Grant GK-14152.

<sup>1</sup> There appears to be no agreement as to the use of the term “dynamical system.” For the purpose of this paper any causal input-output relation will be termed a dynamical system. It will be shown that this is equivalent to the existence of a state. Zadeh [1] and Balakrishnan [2] appear to reserve the term for systems in which the state evolution is governed by a differential equation. The dynamical systems studied in classical mechanics [3] correspond to the state evolution equations in the absence of inputs.

the particular application. (For a discussion of these issues, see the survey paper by the author [6].) It has become clear, however, that the input-output approach leads to more powerful and general results. There are, in fact, a number of interesting stability criteria available which have been obtained in an input-output stability setting, but for which no proofs using Lyapunov methods exist.

Notwithstanding this success of input-output stability, there are certain aspects of Lyapunov stability theory which make the study and development of this "internal" approach to stability theory both useful and important. Not the least of these advantages is the possibility of obtaining estimates on the domain of attraction of an equilibrium in the case of nonglobal stability, a concept which has not even been satisfactorily formulated, let alone developed in the context of input-output stability. The main stumbling block in applying Lyapunov methods to the stability analysis particularly of nonlinear systems remains the absence of general methods for the construction of Lyapunov functions. This paper in part addresses itself to this problem.

The paper is concerned with the implications of input-output stability to global stability of dynamical systems, and with the construction of Lyapunov functions for input-output stable systems. The converse question, i.e., the implications of global stability to input-output stability, will not be considered, in view of space limitations and in view of the fact that such implications are much easier to obtain.

The first part of the paper introduces the concept of a dynamical system, which is defined as a causal operator between signal spaces, and the concept of a realization in which state space concepts become relevant. Some important properties of dynamical systems and realizations are then introduced: they are those of stability, controllability, observability, reachability, connectedness, and irreducibility. These notions play an important role in the sequel.

The second part of the paper discusses the generation of Lyapunov functions for input-output stable systems. Particular emphasis is placed on passive systems and on concepts such as available energy, required energy, and cyclic energy, the latter of which is very reminiscent of certain notions in thermodynamics.

The third part of the paper is concerned with feedback systems. Feedback systems are very important in control, and their stability is, of course, the main qualification on the performance of a feedback structure as a controller. Moreover, for design purposes, it is extremely desirable that properties of feedback systems be concluded from considerations of the open-loop elements. This aspect makes the results of the previous sections not easily applicable to the analysis of feedback systems, and a somewhat different approach is thus required. The ensuing Lyapunov functions are defined in terms of variational problems.

The paper ends with a list of examples. They illustrate the viewpoint adopted here and lead to the Lyapunov functions used to prove the circle criterion and the Popov criterion.

The work reported here has been directly inspired by a very interesting paper by Baker and Bergen [7] which appeared recently. They indeed posed the problem of constructing Lyapunov functions as a variational problem, an approach which has been fully exploited in the context presented here. Some of these ideas already appeared in the work of Popov [8], Kalman [9], and Anderson [10], [11]. The

author obtained a great deal of insight from the work of Brockett on passivity and stability [12]. It is interesting to note that the ingenious independence of path argument as exploited by the latter author in his construction of Lyapunov functions follows here as a rather logical consequence of the variational problems which lead to the desired Lyapunov functions.

The paper also indicates what may be the basic reason why stability conditions appear to be easier to obtain using input-output methods than through the construction of suitable Lyapunov functions: both input-output stability and Lyapunov stability can be posed as variational minimization problems, and whereas Lyapunov methods need the explicit solution of these variational problems (thus the boundedness *and* the value of an infimum), input-output stability only requires the boundedness of this infimum. This observation is due to Zames (private communication).

**2. Dynamical systems.** A dynamical system is usually defined on a subset of the real line as a mapping between function spaces satisfying an appropriate set of axioms. This paper will be concerned with continuous time systems only. Moreover, it will be assumed that the inputs and the outputs take their values in appropriate inner product spaces and that their norm is a locally square integrable function of time. This restriction precludes a certain amount of generality and is made mainly for expository purposes since the results of the paper generalize to much more general situations. In particular, the assumption that the input and output spaces are inner product spaces is of no consequence to many of the results in the paper. One of the reasons for treating systems in this setting is the possibility of introducing and exploiting concepts related to energy and passivity of systems. Indeed, these have far-reaching implications in stability theory.

There are two main avenues for obtaining mathematical models of systems: the first one starts with an internal model in which physical laws and interconnections are used to describe the dynamics and which then yield the relation between the influence variables (the inputs) and the variables of interest (the outputs). The second approach starts with an input-output relation as the basic mathematical model to be used. Such a model is usually the logical consequence from identification experiments at the input-output terminals.

Besides inputs and outputs there is an additional set of variables which is of fundamental importance in the description of dynamical systems. These are the so-called states which summarize the effect of past inputs. The internal modeling approach, in fact, usually displays a state explicitly. More often than not the state has no immediate physical significance and there is never any uniqueness as to its choice. Although the basic mechanism of interest in system theory is the generation of outputs from inputs, it is very often advantageous, however, to view this process as taking place through this intermediate variable, the state. This point of view has been particularly useful in such fields as dynamic optimization theory and the study of Markov processes.

These concepts are formally introduced in the present section, and it is shown that input-output descriptions and state space descriptions of dynamical systems are essentially equivalent.

The first notion is that of signal spaces which will be the input and output function spaces.

Let  $V$  be an inner product space and let  $R$  denote the real line. Let  $f$  be a  $V$ -valued function defined on  $R$ . Then the *causal truncation* of  $f$  at  $T$  is defined to be the result of the projection operator  $P_T$  defined by

$$(P_T f)(t) \triangleq \begin{cases} f(t) & \text{for } t \leq T, \\ 0 & \text{otherwise.} \end{cases}$$

The *anticausal truncation* of  $f$  at  $T$  is defined as  $Q_T f = f - P_T f$ . Consider now the vector space of  $V$ -valued functions on  $R$  with  $\int_{-\infty}^{+\infty} \|f(t)\|_V^2 dt < \infty$ . This vector space is itself an inner product space with

$$\langle f_1, f_2 \rangle = \int_{-\infty}^{+\infty} \langle f_1(t), f_2(t) \rangle_V dt$$

as the inner product. It will be denoted by  $L_2(V)$  and is complete if  $V$  is. As usual, no attention is paid to the fact that a function in  $L_2(V)$  actually represents the equivalence class of functions which are equal to it almost everywhere with respect to Lebesgue measure.

As Wiener remarked when extending Fourier transforms,  $L_2(V)$  is not a very interesting class of functions since it consists of functions which were small in the remote past and are destined to become small in the remote future. This last aspect, in particular, makes this function space of very limited use in stability studies which precisely refer to this remote future, and any a priori limitations on the future would therefore be very inappropriate.

A useful extension of  $L_2(V)$  is its so-called *causal extension* denoted by  $L_{2e}(V)$ , which consists of all  $V$ -valued functions on  $R$  whose causal truncations belong to  $L_2(V)$ , i.e.,

$$L_{2e}(V) \triangleq \{f : R \rightarrow V \mid P_T f \in L_2(V), \text{ all } T \in R\}.$$

The *anticausal extension* of  $L_{2e}(V)$  is similarly defined as

$$\{f : R \rightarrow V \mid Q_T f \in L_2(V), \text{ all } T \in R\}.$$

Since all time functions considered in this paper will be assumed to start at some finite time, very little use of this anticausal extension will be made.

DEFINITION 1. Let  $L_{2e}(V)$  denote the causal extension of  $L_2(V)$ . Then the subspace of  $L_{2e}(V)$  defined by

$$S(V) \triangleq \{f \in L_{2e}(V) \mid Q_T f = 0 \text{ for some } T \in R\}$$

will be called a *signal space*. Elements of  $S(V)$  will be called *signals*, and elements of  $L_2(V) \cap S(V)$  will be called *small signals*.

Thus, signal spaces consist of functions which vanish in the remote past and which have, in a sense, no finite escape, but are otherwise quite arbitrary. As is customary in the related literature, it will be assumed that inputs are applied to systems starting at some finite time in the past. This time need not be a priori fixed and will, in general, be different for each experiment. Note that signal spaces are closed under concatenation and that any "reasonable" physical signal belongs to a signal space. For the purposes of this paper, signal spaces represent a very convenient abstraction of reality. The fact that signals are required to have their



support on a half-line is very important and results in mild conditions for the well-posedness of mathematical models. In other words, given a mathematical model for a system (e.g., an integral equation or a differential equation), it will, in general, be relatively easy to establish that inputs in a given class generate well-defined outputs, provided, however, that these inputs have their support on a half-line. For inputs defined on the whole real line  $(-\infty, +\infty)$ , establishing this existence and uniqueness of outputs usually leads to stringent conditions and requires typically input-output continuity of the system. This case is moreover of dubious physical significance. For details see [13], [14, § 4.6], [15].

Let  $U = S(V_u)$  and  $Y = S(V_y)$  be signal spaces.  $U$  will be called the *input space*, and  $Y$  will be called the *output space*. Elements of  $U$  and  $Y$  will be called respectively *input signals* and *output signals*. A mapping  $F$  from  $U$  into  $Y$  is said to be *causal* (or *nonanticipatory*) if for all  $T \in \mathbb{R}$  and all  $u_1, u_2 \in U$  with  $P_T u_1 = P_T u_2$ , the equality  $P_T F u_1 = P_T F u_2$  holds.<sup>2</sup> This condition is equivalent to requiring that  $P_T F P_T = P_T F$  on  $U$ .

Note that the signal spaces as introduced above could have been called somewhat more consistently causal signal spaces. The analogous concepts of anticausal signal spaces and anticausal operators thus become straightforward. No use will be made of these concepts, however. An additional notion which is of some importance is that of a *memoryless* operator. This would most logically be defined as an operator which is *both* causal and anticausal but is easiest (although equivalently) defined as an operator  $F$ , defined by an element  $r \in Y$  and an instantaneous map,  $f$ , from  $V_u \times \mathbb{R}$  into  $V_y$  with  $f(0, t) = 0$  for all  $t \in T$  and  $Fu \triangleq r + Nu$ , where  $(Nu)(t) \triangleq f(u(t), t)$  is such that any function  $u \in L_2(V_u)$  with compact support yields  $Nu \in L_2(V_y)$  (consequently also with compact support).

**DEFINITION 2.** A *dynamical system* is defined as a causal mapping from the input signal space  $U$  into the output signal space  $Y$ . If this mapping is memoryless, then the dynamical system will similarly be called memoryless.

The above setting for the study of input-output relations is similar to the one employed by Balakrishnan in [16]. The definition eliminates the possibility of studying differentiators, for instance, but for the purposes of this paper (stability) such a restriction is not very disturbing. In the study of networks, however, one clearly wants a more general definition which admits singularity functions in the impulse response. Zemanian [17] and Balakrishnan [18] have studied systems in which the inputs are assumed to be infinitely smooth functions and the outputs are distributions. Extended spaces appeared first in the context of stability theory as a result of the work of Sandberg [4] and Zames [5].

For many purposes, it is convenient to impose some smoothness conditions on the operators in question. Note that  $U$  and  $Y$  have, as signal spaces, no topology since they are, although derived from normed spaces, not normed themselves. However, causality enables one nevertheless to make a suitable definition of local continuity. Although simple continuity is the most logical smoothness condition to impose, it is very often advantageous to require somewhat stronger conditions, more specifically Lipschitz continuity. Recall that a (in general non-linear) map  $F$  between normed spaces is said to be *Lipschitz continuous* if there

<sup>2</sup> Note the abuse of notation in the fact that the symbol  $P_T$  is used to denote an operator on  $U$  and an operator on  $Y$ . This ambiguity, however, causes no difficulty.

exists a real constant  $K < \infty$  such that for all  $x_1, x_2$  in the domain of  $F$ ,  $\|Fx_1 - Fx_2\| \leq K\|x_1 - x_2\|$ . Let  $F$  be a causal map from the input space  $U$  into the output space  $Y$ . Then  $F$  is said to be *locally Lipschitz continuous* if for all  $t_0, t_1 \in \mathbb{R}$ ,  $P_{t_1}FQ_{t_0}$  is Lipschitz continuous as a map from  $L_2(V_u)$  into  $L_2(V_y)$ .

DEFINITION 3. A dynamical system is said to be *smooth* if the defining map  $G$  is locally Lipschitz continuous. It is said to be *uniformly smooth* if for any given  $T > 0$ ,  $P_{t+T}GQ_t$  is Lipschitz continuous uniformly in  $t$ .

*Convention.* For convenience it will be assumed that all dynamical systems under consideration are *unbiased*, i.e., that they map the zero input into the zero output; hence  $G0 = 0$ . This absence of a bias term can always be obtained by a trivial redefinition of  $G$  and assumes, for instance, that for memoryless operators the element  $r \in Y$  appearing in the definition is the zero element.

Now that the definition of input-output models of dynamical systems has been introduced, attention is focused on the formalism for the state space description of dynamical systems. Let  $\mathbb{R}_2^+$  denote the *causal sector* of  $\mathbb{R}_2$  defined as  $\mathbb{R}_2^+ \triangleq \{(t_2, t_1) | t_2, t_1 \in \mathbb{R}, t_2 \geq t_1\}$ .

DEFINITION 4. A (mathematical model of a) dynamical system is said to be in *state space form* if it is determined by an abstract set  $X$  (the *state space*) and two maps,  $\phi$ , the *state transition map*, and  $y$ , the *output reading map*, satisfying the following axioms:

- (i)  $\phi$  maps  $\mathbb{R}_2^+ \times X \times U$  into  $X$ ;
- (ii) (*Causality*):  $\phi(t, t_0, x_0, u) = \phi(t, t_0, x_0, P_t Q_{t_0} u)$  for all  $(t, t_0) \in \mathbb{R}_2^+$ ,  $x_0 \in X$ , and  $u \in U$ ;
- (iii) (*Consistency*):  $\phi(t_0, t_0, x_0, u) = x_0$  for all  $t_0 \in \mathbb{R}$ ,  $x_0 \in X$ , and  $u \in U$ ;
- (iv) (*Composition law or semi-group property*):  $\phi(t_2, t_0, x_0, u) = \phi(t_2, t_1, \phi(t_1, t_0, x_0, u), u)$  for all  $(t_1, t_0), (t_2, t_1) \in \mathbb{R}_2^+$ ,  $x_0 \in X$ , and  $u \in U$ ;
- (v)  $y$  maps  $\mathbb{R} \times X \times V_u$  into  $V_y$  and the value of the output at time  $t$  is given by  $y(t, x(t), u(t))$ ;
- (vi)  $X$  is a subset of an inner product space  $V_x$ ;
- (vii) (*Unbiasedness*):  $\phi(t, t_0, 0, 0) = 0$  for all  $(t, t_0) \in \mathbb{R}_2^+$ , and  $y(t, 0, 0) = 0$  for all  $t \in \mathbb{R}$ ;
- (viii) Let  $\tilde{X}$  denote the signal space induced by  $V_x$  (i.e.,  $\tilde{X} = S(V_x)$ ); it is then assumed that the functions

$$x(t) = Q_{t_0} \phi(t, t_0, x, u) = \begin{cases} \phi(t, t_0, x, u) & \text{for } t \geq t_0, \\ 0 & \text{otherwise,} \end{cases}$$

$$Q_{t_0} y(t, x(t), u(t)) = \begin{cases} y(t, x(t), u(t)) & \text{for } t \geq t_0, \\ 0 & \text{otherwise,} \end{cases}$$

belong to  $\tilde{X}$  and  $Y$  respectively for all  $t_0 \in \mathbb{R}$ ,  $x_0 \in X$ , and  $u \in U$ .

Axioms (i)–(v) are the usual axioms involved in describing dynamical systems from a state space point of view. Axiom (vi) induces a topology on the state space and will be needed in the definitions of Lyapunov stability, for instance.<sup>3</sup> Axiom

<sup>3</sup> The assumption that  $X$  is an inner product space is restrictive and inconvenient for many applications, more so than would appear at first sight. For a study of dynamical systems whose state space is a group manifold see [26].

(vii) is in keeping with the unbiasedness convention introduced above, and axiom (viii) guarantees the absence of a finite escape and results in the fact that locally square integrable inputs produce locally square integrable outputs.

A dynamical system in state space form thus views the generation of outputs from inputs as occurring through the composition of two maps,  $G_x$  and  $G_y$ , with  $G_x: U \rightarrow X$  and  $G_y: X \times U \rightarrow Y$ . The map  $G_x$  is a dynamical system in its own right (with  $X$  viewed as the output space) but satisfies a richer set of axioms than merely those implied by the dynamical system axioms: in addition, it is required that this map have the Markov property which results in a decoupling of the past from the future in the sense that the present value of the state has sufficient information in it so as to summarize the effect of past inputs. The state space thus represents an adequate memory-bank. The map  $G_y$  is memoryless (with input space  $X \times U$ ) and the dependence of  $y(t)$  on past values of  $u$  is obtained through the dependence on  $x(t)$ . It is a simple matter to verify that the composite system  $y = G_y(G_x u, u)$  is indeed a dynamical system in the input-output sense. Note also that  $G_x 0$  satisfies the axioms for dynamical systems without inputs as studied in classical mechanics and its extensions.

The next definitions refer to the smoothness of the state transition map and the output reading map. These smoothness conditions are generally quite important aspects of a particular dynamical system in state space form. For instance, it can be shown that otherwise any finite-dimensional dynamical system can be realized by a one-dimensional dynamical system if this latter is not required to have any smoothness. It suffices therefore to consider a one-to-one map from  $R^n$  into  $R$  and appropriately modify the state space and the maps defining the dynamical system.

**DEFINITION 5.** A dynamical system in state space form is said to be *smooth* if for any  $(t_1, t_0) \in R_2^+$  there exist  $K_1, K_2, K_3, K_4 < \infty$  such that

$$\begin{aligned} & \|P_{t_1} Q_{t_0}(\phi(t, t_0, x_1, u_1) - \phi(t, t_0, x_2, u_2))\| \\ & \leq K_1 \|x_1 - x_2\| + K_2 \|P_{t_1} Q_{t_0}(u_1 - u_2)\| \end{aligned}$$

for all  $x_1, x_2 \in X$  and  $u_1, u_2 \in U$ , and

$$\begin{aligned} & \|P_{t_1} Q_{t_0}(y(t, x_1(t), u_1(t)) - y(t, x_2(t), u_2(t)))\| \\ & \leq K_3 \|P_{t_1} Q_{t_0}(x_1 - x_2)\| + K_4 \|P_{t_1} Q_{t_0}(y_1 - y_2)\| \\ & \text{for all } x_1, x_2 \in \tilde{X} \text{ and } u_1, u_2 \in U. \end{aligned}$$

It is said to be *uniformly smooth* if for any  $T \geq 0$  and  $t_1 = t_0 + T$  in the above inequalities, the constants  $K_1, K_2, K_3, K_4$  may be taken independent of  $t_0$ .

These definitions are entirely analogous to those imposed for input-output systems. It is a simple matter to verify that (uniformly) smooth dynamical systems in state space form define (uniformly) smooth input-output dynamical systems. It is also clear that uniform smoothness and smoothness are equivalent for time-invariant systems.

The final discussion of this section involves the relationship between the above definitions of dynamical systems. As might be expected they are indeed equivalent.

DEFINITION 6. Consider the dynamical system  $G$  and a dynamical system in state space form with defining maps  $\phi$  and  $y$ . Then the dynamical system in state space form is said to be a *realization* of  $G$  if any  $u \in U$  with  $t_0 \in R$  such that  $P_{t_0}u = 0$ , yields  $(Gu)(t) = y(t, \phi(t, t_0, 0, u), u(t))$  for all  $t \in R$ . The dynamical system in state space form thus defines the same input-output relation as  $G$ .

THEOREM 1. *Every (smooth, uniformly smooth) dynamical system has a (smooth, uniformly smooth) realization in state space form.*

*Proof.* The proof proceeds by construction. The state space  $X = V_x$  will be taken to be the collection of all functions in  $L_2(V_u)$  with compact support, and the state at time  $t$  will be taken to be  $P_t u$ , where  $u \in U$  is the input to the system. Thus, for instance,

$$\phi(t, x_0, t_0, u) = S_t P_{t_0} u \stackrel{\Delta}{=} x_0 + S_t P_{t_0} Q_{t_0} u,$$

where  $S_T$  denotes the shift operator  $(S_T z)(t) = z(t - T)$ , and

$$y(t, x, u) = (Gu)(t) = (GP_t u)(t) = (GS_{-t} P_t u)(t) = \stackrel{\Delta}{=} (GS_{-t} x(t))(t).$$

It is left to the reader to verify that these maps indeed satisfy the axioms of dynamical systems in state space form. The smoothness claims are also easily verified directly in view of the simplicity of the state transition map and the fact that the output map and the original dynamical system are essentially identical.

The above theorem, although too trivial and general to be of significance in specific instances, yields a rather interesting canonical decomposition of nonlinear dynamical systems into a linear, time-invariant, reachable dynamical part followed by a memoryless nonlinear part.<sup>4</sup> Note also that the dynamical part in this decomposition may be described by the partial differential equation

$$\partial x(z, t) / \partial t = \partial x(z, t) / \partial z, \quad z \leq 0,$$

with the boundary control  $x(0, t) = u(t)$  and with solutions defined in an appropriate sense. The function  $x(z, t)$  for  $z \geq 0$  then plays the role of the state at time  $t$ , and the partial differential equation describes the evolution of the initial state  $x(z, t_0)$  resulting from the input  $u(t)$ . Notice also that in the above realization the map  $y$  inherits linearity and time-invariance of  $G$ .

It should be noted that the equivalence of a dynamical system and a state space realization of a dynamical system might nevertheless lead the latter model to produce an output which could not be the result of any input to the former model. Such outputs result from initial states which are somewhat artificial in the sense that they cannot be produced by past inputs. The equivalence of a dynamical system and one of its realizations is thus really a zero initial state equivalence.

*Notation.* Let  $G$  denote a dynamical system in state space form,  $x_0 \in X$ ,  $t_0 \in R$ , and  $u \in U$ . Then the function defined by

$$\begin{cases} y(t, \phi(t, t_0, x_0, u), u(t)) & \text{for } t \geq T \geq t_0, \\ 0 & \text{otherwise,} \end{cases}$$

---

<sup>4</sup> It was pointed out to the author that similar decomposition due to Wiener [19] and Balakrishnan [16] have appeared in the literature.

will be denoted by  $Q_T G(t_0, x_0, u)$  (or  $Q_T y(t_0, x_0, u)$  when there is no danger of confusion).

**3. Fundamental properties of dynamical systems.** A number of fundamental concepts related to dynamical systems and their state space realizations are introduced and discussed in this section: they relate to the influence of the control on the state (reachability, controllability, and connectedness), of the state on the output (observability and irreducibility), and of the input on the output (stability and continuity).

**DEFINITION 7.** The state space of a dynamical system in state space form is said to be *reachable* if given any  $x \in X$  and  $t \in R$ , there exists a  $t_0 \in R$ ,  $t_0 \leq t$ , and a  $u \in U$  such that  $\phi(t, t_0, 0, u) = x$ . A dynamical system in state space form is said to be *controllable* if given any  $x_0 \in X$  and  $t_0 \in R$ , there exists a  $t \in T$ ,  $t \geq t_0$ , and a  $u \in U$  such that  $\phi(t, t_0, x_0, u) = 0$ . The state space of a dynamical system in state space form is said to be *connected* if given any  $x_0, x_1 \in X$  there exists an element  $(t_1, t_0) \in R_2^+$  and a  $u \in U$  such that  $\phi(t_1, t_0, x_0, u) = x_1$ .

Reachability thus requires the map  $\phi(t, \cdot, 0, \cdot)$  to be onto  $X$ , whereas controllability requires that 0 be in the range space of  $\phi(\cdot, t_0, x_0, \cdot)$ . Note that reachability, controllability, and time-invariance imply connectedness.

**DEFINITION 8.** A dynamical system in state space form is said to be *observable* if for any  $t_0 \in R$ , knowledge of  $Q_{t_0} y(t_0, x_0, 0)$  (uniquely) determines  $x_0 \in X$ . The state space of a dynamical system in state space form is said to be *irreducible* if for any given  $t_0 \in R$  and  $x_0 \in X$  there exists a  $Q_{t_0} u \in U$ , such that knowledge of  $Q_{t_0} y(t_0, x_0, u)$  (uniquely) determines  $x_0 \in X$ .

Observability thus requires the map  $y(t_0, \cdot, 0)$  to be one-to-one on  $X$ , whereas irreducibility requires the map  $y(t_0, \cdot, u)$  to be one-to-one on  $X$  by choosing  $(Q_{t_0} u)(x)$ . It is clear that observability implies irreducibility and that the nomenclature “irreducible” is quite appropriate since if the state space is not irreducible, then there exist at least two initial states which will be completely indistinguishable under experimentation: these two states are thus entirely equivalent, and nothing will be lost by eliminating one of them from the state space.

The above nomenclature is common (although far from standard) in the related literature with the possible exception of irreducibility which is often taken to indicate the set-theoretic minimality of the state space. Observability and irreducibility are equivalent for linear systems with a finite-dimensional smooth state space realization. The simplest example of systems in which these concepts are different are systems with multiplicative control described, e.g., by  $\dot{x} = uAx$ . It should also be remarked that the above definitions, although natural, are not the most convenient ones for certain applications. Although it can be shown that every dynamical system has a realization<sup>5</sup> with a reachable state space, it is sometimes very difficult to discover exactly what states are reachable (and to define  $X$  then appropriately). For instance, in systems described by partial differential equations these reachable states have certain smoothness properties which are not a priori known; therefore, in certain applications it is much more convenient to adopt an “almost” reachability requirement. The same remark holds for the following

---

<sup>5</sup> See, for example, the proof of Theorem 1.

definitions which, in addition, require an appropriate choice of the topology on the state space.

DEFINITION 9. The state space of a dynamical system in state space form is said to be *uniformly reachable* if there exist a continuous function  $\alpha: R^+ \rightarrow R^+$  ( $R^+$  denotes the nonnegative real numbers) with  $\alpha(0) = 0$  and a constant  $T \geq 0$ , such that for any  $x \in X$  and  $t \in R$ , there exists a  $u \in U$  with  $\|P_t u\|^2 \leq \alpha(\|x\|)$  such that  $\phi(t, t - T, 0, u) = x$ . *Uniform controllability* and *uniform connectedness* are similarly defined. A dynamical system is said to be *uniformly observable* if there exist a strictly monotone increasing continuous function  $\beta: R^+ \rightarrow R$  with  $\beta(0) = 0$  and  $\lim_{\sigma \rightarrow +\infty} \beta(\sigma) = +\infty$  and a constant  $T \in R$ ,  $T \geq 0$ , such that for any  $x \in X$  and  $t_0 \in R$ ,

$$\|P_{t_0+T} Q_{t_0} y(t_0, x, 0)\|^2 \geq \beta(\|x\|).$$

The state space of a dynamical system in state space form is said to be *uniformly irreducible* if with  $\beta$  and  $T$  as before, the inequality

$$\|P_{t_0+T} Q_{t_0} (y(t_0, x_1, u) - y(t_0, x_2, u))\|^2 \geq \beta(\|x_1 - x_2\|)$$

holds for all  $x_1, x_2 \in X$ ,  $t_0 \in R$ , and some  $u \in U$ .

The above definitions differ somewhat from those in the literature. Most of the papers concerned with uniform controllability for linear systems follow Kalman's [20] original definition, which imposes many more restrictions than the definitions used here. In particular, it requires any control which makes the transfer from state 0 at time  $t - T$  to state  $x$  at time  $t$  to be such that  $\|P_t Q_{t-T} u\|^2 \geq \alpha(\|x\|) > 0$ .

The most efficient realization of a dynamical system is one in which the state space is reachable *and* irreducible. This indeed guarantees that every output which can be observed as a result of initial conditions and inputs could have been observed by properly choosing the past input and that two different initial conditions will lead to different outputs by properly choosing the input. Two realizations which are both reachable *and* irreducible are thus isomorphic. They differ in the sense that their state spaces are labeled differently. The one-to-one onto map between these state spaces may, in general, be a function of time, however. A realization of a dynamical system in which the state space is reachable and irreducible can thus properly be called *minimal*, a notion which has many more substantive implications for linear systems. In looking for reachable *and* irreducible realizations it is natural to consider as the candidate for the state space the equivalence classes of those inputs up to time  $t$  which yield the same output after time  $t$ , regardless of the input after time  $t$ ; more precisely, by considering the equivalence class  $\{P_t u \mid Q_t y \text{ is fixed for all } Q_t u\}$  as a typical element of the state space. The difficulty with this representation is that, in general, the state space itself then becomes a function of time. There are two methods of getting around this difficulty: one is to modify the original axioms and definitions so as to allow for a state space which is itself a function of time; the other is to define a dynamical system as a causal *and* a noncausal map depending on whether one considers time moving forward or backward from the initial time. The state is thus alternatively required to summarize past and future, and the state space thus has many more invariant properties with respect to time. This device has been used successfully by Kalman

[21] and others in their study of systems described by the Volterra integral equation

$$y(t) = \int_{t_0}^t w(t, \tau)u(\tau) d\tau$$

with separable kernel  $w$ . This principle rests on dubious physical grounds, however, and leads to technical difficulties for infinite-dimensional systems. The above problems do not occur in stationary systems.

Recall that a mapping between normed spaces is said to be *bounded* if it maps bounded sets into bounded sets. It is said to have a *finite gain* if there exists a  $K < \infty$  such that for any  $\rho \geq 0$  the ball with radius  $\rho$  gets mapped into the ball with radius  $K\rho$ . The infimum of all real numbers  $K$  achieving the above inequality is usually called the *gain* of the operator.

DEFINITION 10. A dynamical system,  $G$ , is said to be *input-output stable* if it maps bounded sets of small signals in  $U$  into bounded sets of small signals in  $Y$ . It is said to be *finite-gain input-output stable* if it is stable and if there exists a  $K < \infty$  such that for any small signal  $u \in U$ ,  $\|Gu\| \leq K\|u\|$ . The infimum of all such real numbers  $K$  will be denoted by  $\|G\|$ . A dynamical system  $G$  is said to be *input-output continuous* if it is stable and if the map  $G$  is continuous (in the topology induced by  $L_2(V_u)$  and  $L_2(V_y)$ ) as a map from  $U \cap L_2(V_u)$  into  $Y \cap L_2(V_y)$ . It is said to be *input-output Lipschitz continuous* if  $G$  is actually Lipschitz continuous.

It can be shown [14, § 2.4] that a dynamical system is finite-gain stable if and only if the gain of  $P_{t_1}GQ_{t_0}$  is bounded for all  $t_0, t_1 \in \mathbb{R}$ , uniformly in  $t_0$  and  $t_1$ . In fact,

$$\|G\| = \lim_{-t_0, t_1 \rightarrow \infty} \|P_{t_1}GQ_{t_0}\|,$$

and this limit is approached monotonically. A similar relationship holds for Lipschitz continuity.

Related, but not identical, are the following more familiar Lyapunov stability concepts for dynamical systems in state space form.

DEFINITION 11. The equilibrium state of a dynamical system in state space form is said to be *globally attractive* if for any  $x_0 \in X$  and  $t_0 \in \mathbb{R}$ ,

$$\lim_{T \rightarrow \infty} \phi(t_0 + T, t_0, x_0, 0) = 0.$$

It is said to be *uniformly globally attractive* if this limit is uniform in  $t_0$ . It is said to be *stable* if for any  $\varepsilon > 0$  and  $t_0 \in \mathbb{R}$  there exists a  $\delta(\varepsilon, t_0)$  such that  $\|\phi(t_0 + T, t_0, x_0, 0)\| \leq \varepsilon$  for all  $T \geq 0$  whenever  $\|x_0\| \leq \delta$ . It is said to be *uniformly stable* if  $\delta(\varepsilon, t_0)$  may be chosen independent of  $t_0$ . A dynamical system in state space form is said to be *bounded* if for any  $x_0 \in X$  and  $t_0 \in \mathbb{R}$ ,  $\phi(t_0 + T, t_0, x_0, 0)$  is bounded on the half-line  $T \geq 0$ . It is said to be *uniformly bounded* if this bound may be chosen independent of  $t_0$ . A dynamical system in state space form is said to be *globally asymptotically stable* if the equilibrium state is globally attractive and stable. It is said to be *uniformly globally asymptotically stable* if the equilibrium state is uniformly globally attractive, uniformly stable, and uniformly bounded.

The usual method of proving stability of systems in state space form is to consider an appropriate Lyapunov function. The following definition of a Lyapunov function is a convenient one for the present discussion.

DEFINITION 12. Let  $V$  be a mapping from  $X \times R$  into  $R^+$ , with  $V(x, t) = 0$  if and only if  $x = 0$ . Then  $V$  is said to be a *Lyapunov function* for a dynamical system in state space form if for any  $x_0 \in X$  and  $t_0 \in R$ ,

- (i)  $V(\phi(t, t_0, x_0, 0), t)$  is a monotone nonincreasing function of  $t$  for  $t \geq t_0$ ;
- (ii)  $\lim_{T \rightarrow \infty} V(\phi(t_0 + T, t_0, x_0, t), t_0 + T) = 0$ .

The function  $V$  will be called a *uniform Lyapunov function* if, in addition, the limit in (ii) is uniform in  $t_0$  and if  $V(x, t)$  is bounded in  $t$  for all  $x \in X$ . The function  $V$  is said to be *decreasing* if there exists a continuous function  $\alpha: R^+ \rightarrow R^+$  with  $\alpha(0) = 0$  such that  $V(x, t) \leq \alpha(\|x\|)$  for all  $x \in X$  and  $t \in R$ . It is said to be *positive definite* if there exists a monotone increasing continuous function  $\beta: R^+ \rightarrow R^+$  with  $\beta(0) = 0$  such that  $V(x, t) \geq \beta(\|x\|)$  for all  $x \in X$  and  $t \in R$ . It is said to be *radially unbounded* if there exists a continuous function  $\gamma: R^+ \rightarrow R^+$  with  $\lim_{\sigma \rightarrow +\infty} \gamma(\sigma) = +\infty$  such that  $V(x, t) \geq \gamma(\|x\|)$  for all  $x \in X$  and  $t \in R$ . If  $V$  is a Lyapunov function for a dynamical system in state space form, then the equilibrium state is globally asymptotically stable if  $V$  is positive definite, and the dynamical system is bounded if  $V$  is radially unbounded. If  $V$  is a uniform Lyapunov function, then the equilibrium state is uniformly globally attractive if  $V$  is positive definite, and uniformly stable if  $V$  is positive definite and decreasing; the system is uniformly bounded if  $V$  is radially unbounded, and uniformly globally asymptotically stable if  $V$  is radially unbounded, positive definite, and decreasing. Notice also that decreasency implies the last condition in the definition of a uniform Lyapunov function, and thus a decreasing Lyapunov function for a uniformly globally asymptotically stable system is a uniform Lyapunov function.

The main purpose of this paper is to study the relations between input-output stability and global stability. It seems reasonable to expect that an input-output stable system will be globally stable if inputs sufficiently influence states and if states sufficiently influence outputs. Then internal instability should reflect into external instability. That this can be made precise is shown in the next section.

**4. Input-output stability and global stability.** This section establishes the fundamental relationship between input-output stability and global stability. In trying to obtain these internal stability implications from external data, one defines certain functions which depend on the external variables only. In order for these functions to be well-defined and to qualify as suitable Lyapunov functions, a number of additional assumptions have to be made, and it is at this point that reachability, controllability, observability, and input-output stability become relevant. There are two natural functions to consider for this purpose:

- (i)  $V_r(x, t) \triangleq \inf \|P_t Q_{t_0} u\|^2$ , where the infimum is to be taken over all  $t_0 \leq t$  and  $u \in U$  with  $\phi(t, t_0, 0, u) = x$  (the infimum (supremum) over the void set is by assumption  $+\infty$  ( $-\infty$ )), and
- (ii)  $V_o(x, t) \triangleq \|Q_t y(t, x, 0)\|^2$ .

The symbolism is clear: the first function is inspired by reachability, and the second by observability.  $V_r$  is well-defined if the state space is reachable, and  $V_o$  is well-defined if the state space is reachable and if the dynamical system is input-output stable. Indeed, let  $u \in U$  be such that  $\phi(t, t_0, 0, u) = x$ . Then

$$\|Q_t y(t, x, 0)\|^2 \leq \|G P_t u\|^2 \leq \|G\|^2 \|P_t u\|^2.$$



If  $V_0(x, t)$  is well-defined, then it is clearly monotone nonincreasing along undriven solutions and approaches zero when  $t \rightarrow \infty$  since

$$V_0(\phi(t, t_0, x_0, 0), t) = \int_t^\infty \|y(\tau, \phi(\tau, t_0, x_0, 0), 0)\|_{V_y}^2 d\tau.$$

The function  $V_r$  is also monotone nonincreasing along undriven solutions, although this is not as immediate. The argument used in the demonstration of this fact will repeatedly be used in the sequel and will therefore only here be done explicitly. Thus, consider  $V_r(\phi(t_1, x, t, 0), t_1)$  with  $t_1 \geq t$ , and denote  $\phi(t_1, x, t, 0)$  by  $x_1$ . Then  $V_r(x_1, t_1) = \inf \|P_{t_1} Q_{t_0} u\|^2$ . The state of the dynamical system can be driven to  $x_1$  at time  $t_1$  by first driving it from 0 at time  $t_0$  to  $x$  at time  $t$  and then applying zero control until time  $t_1$ . This is, in general, a suboptimal control for reaching  $x_1$  at time  $t_1$ , even when it is driven to  $x$  at time  $t$  in an optimal fashion. This suboptimal strategy thus shows that  $V_r(x_1, t_1) \leq V_r(x, t)$ .

The basic relationships between input-output stability and global stability are stated in the following theorems.

**THEOREM 2.** *A uniformly observable realization of an input-output stable dynamical system with a reachable state space is globally asymptotically stable and bounded, and  $V_0$  is a positive definite radially unbounded Lyapunov function for it.*

*Proof.* It suffices to show that  $V_0$  is a positive definite radially unbounded Lyapunov function. By reachability and input-output stability,  $V_0$  is well-defined. By observability,  $V_0 = 0$  if and only if  $x = 0$ ; it is monotone nonincreasing and approaches zero along undriven solutions since

$$\int_t^\infty \|y(\tau, \phi(\tau, t, x, 0), 0)\|_{V_y}^2 d\tau < \infty.$$

It remains to be shown that  $V$  is positive definite and radially unbounded. This, however, is an immediate consequence of uniform observability.

Note that in the above theorem the uniform observability condition cannot simply be relaxed to only observability, even if at the same time one assumes uniform reachability rather than merely reachability.

**THEOREM 3.** *A uniformly observable realization of a finite-gain input-output stable dynamical system with a uniformly reachable state space is uniformly globally asymptotically stable, and  $V_r$  and  $V_0$  are positive definite radially unbounded decrescent uniform Lyapunov functions for it.*

*Proof.* From finite-gain input-output stability it follows that

$$\|P_t y(t, 0, u)\|^2 \leq \|G\|^2 \|P_t u\|^2$$

and thus that

$$V_0(x, t) \leq \|G\|^2 V_r(x, t).$$

$V_r$  is well-defined and decrescent by uniform reachability, and  $V_0$  is well-defined, positive definite, and radially unbounded by uniform observability. Thus both  $V_0$  and  $V_r$  are positive definite, radially unbounded, and decrescent, and the theorem follows if it can be shown that  $V_0(\phi(t_0 + T, t_0, x_0, 0), t_0 + T)$  approaches zero as  $T \rightarrow \infty$ , uniformly in  $t_0$ . Notice that with  $T$  as in the definition of uniform

observability it follows that

$$V_0(\phi(t_0 + T, t_0, x_0, 0), t + T) V_0(x_0, t_0) \leq 1 - \frac{\beta(\|x_0\|)}{\|G\|^2 \alpha(\|x_0\|)},$$

where  $\alpha$  and  $\beta$  are as in the definitions of uniform reachability and uniform observability. This implies uniform convergence of the Lyapunov function  $V_0$  since for any  $\varepsilon > 0$  and  $M < \infty$  there exists a  $\delta > 0$  such that  $\beta(\sigma)\alpha(\sigma) \geq \delta > 0$  for all  $0 < \varepsilon \leq \sigma \leq M < \infty$ . Hence the system is uniformly globally asymptotically stable, which in turn implies that  $V_r$  is a uniform Lyapunov function.

Note that finite-gain stability (or uniform smoothness) and uniform observability yield that every control  $u$  transferring state 0 at time  $t_0$  to state  $x$  at time  $t$  requires

$$\|P_t Q_{t_0} u\|^2 \geq \|G\|^{-2} \beta(\|x\|),$$

a condition which is usually part of the definition of uniform controllability [20]. As a final remark in this section, note that the fact that the inputs and outputs take their values in inner product spaces is inessential and that the results hold, mutatis mutandis, if these spaces are merely normed spaces. The inner product structure becomes very important in the next section, which is concerned with passivity.

**5. Lyapunov functions for passive systems.** The notions which will be introduced in this section are those of passivity and certain concepts related to energy. It will be assumed in this section that the inner product spaces under consideration are real.<sup>6</sup>

DEFINITION 13. Let  $U = Y$  and let  $G$  be a dynamical system from  $U$  into  $Y$ . Then  $G$  is said to be *passive* if for all  $u \in U$  and  $t \in R$ ,  $\langle P_t u, P_t G u \rangle \geq 0$ . It is said to be *strictly passive* if  $G - \varepsilon I$  is passive for some  $\varepsilon > 0$ .

This terminology is to be interpreted as follows:  $\langle u(t), y(t) \rangle_{V_u}$  represents the instantaneous power delivered to the system from the outside. Thus  $\langle P_t u, P_t y \rangle$  represents the total energy at time  $t$  delivered to the system from the outside. If regardless of the termination and in the absence of initial excitations this energy is nonnegative, then the system is passive viewed from its input-output terminals.

It can be shown (see [14, § 2.17]) that if  $G$  is input-output stable, then it is passive if and only if  $\langle u, G u \rangle \geq 0$  for all small input signals  $u$ .

DEFINITION 14. Let  $G$  be a dynamical system in state space form. Then the *required energy*,  $E_r$ , is defined on  $X \times R$  as

$$E_r(x, t) \triangleq \inf \langle P_t u, P_t G u \rangle,$$

where the infimum is to be taken over all  $t_0 \leq t$  and  $u \in U$  with  $P_{t_0} u = 0$  which yield  $\phi(t, t_0, 0, u) = x$ . The *available energy*,  $E_a$ , is defined on  $X \times R$  as

$$E_a(x, t) \triangleq \sup_{\substack{u \in \bar{U} \\ t_1 \geq t}} - \langle P_{t_1} Q_{t_1} u, P_{t_1} Q_{t_1} y(t, x, u) \rangle.$$

The *cycle energy*,  $E_c$ , is defined on  $X \times R$  as  $E_r - E_a$ . Thus

$$E_c(x, t) = \inf \langle P_{t_1} u, P_{t_1} G u \rangle,$$

---

<sup>6</sup> Complex inner product spaces can be treated equally well by considering the real part of the inner product in the definitions of passivity and energy.

where the infimum is to be taken over all  $t_0$  and  $t_1$  with  $t_0 \leqq t \leqq t_1$  and  $u \in U$  with  $P_{t_0}u = 0$ , which yields  $\phi(t, t_0, 0, u) = x$ .

The available energy is thus the maximum energy which can be extracted from a system, whereas the required energy is the energy needed to excite a system to a given set of initial conditions. The cycle energy is the minimum energy it takes to cycle a system between the equilibrium and a given state. Note that all of the above energies are defined in terms of input-output relations.

LEMMA 1. Consider a realization of a passive dynamical system and assume that the state space is reachable. Then  $E_a, E_r$  and  $E_c$  exist (i.e.,  $E_a, E_r, E_c < \infty$ ) and are nonnegative. Moreover,  $0 \leqq E_a, E_c \leqq E_r$ .

Proof. That  $E_r$  and  $E_c$  are finite and nonnegative follows immediately from passivity and reachability. Hence, since  $E_a + E_c = E_r, E_a \leqq E_r$ .

It remains to be shown that  $E_a$  is nonnegative. This follows by considering  $Q_t u = 0$ , which shows that the supremum in the definition of  $E_c$  is taken over a set which contains zero. This completes the proof of the lemma.

The inequality  $E_a \leqq E_r$  formalizes the intuitive notion that passive systems cannot supply more energy to the outside than has previously been supplied to them from the outside. Note that none of the above notions satisfactorily defines the stored energy,  $E_s(x, t)$ , which is an internal property of a dynamical system and thus usually a function of the realization. The passivity definition employed here is purely input-output. Similar definitions of internal passivity can be made, and the theory for linear time-invariant dissipative systems [22] is available. One can then pose the question of whether or not every input-output passive system has a passive realization. These ramifications fall beyond the scope of the present paper. It would be interesting to verify that the stored energy in a passive realization of a passive system satisfies the inequality  $E_a \leqq E_s \leqq E_r$ , as it should.

The cycle energy  $E_c$  is a measure of the degree of irreversibility of a system. This is the intuitive basis for the following definitions.

DEFINITION 15. A passive dynamical system in state space form is said to be *irreversible* if  $E_c(x, t) = 0$  only if  $x = 0$ . It is said to be *uniformly irreversible* if there exists a monotone increasing function  $\gamma: R^+ \rightarrow R^+$  with  $\gamma(0) = 0$  and  $\lim_{\sigma \rightarrow +\infty} \gamma(\sigma) = +\infty$  such that for all  $x \in X$  and  $t \in R, E_c(x, t) \geqq \gamma(\|x\|)$ . It is said to be *reversible* if  $E_c = 0$ , i.e., if  $E_r = E_a$ .

THEOREM 4. The available energy,  $E_a$ , and the required energy,  $E_r$ , are decrescent uniform Lyapunov functions for a uniformly observable realization of a passive finite-gain input-output stable dynamical system with a uniformly reachable state space.

Proof. It will first be shown that  $E_r$  is decrescent. By the Schwarz inequality,

$$|\langle P_t u, P_t G u \rangle| \leqq \|P_t u\| \|P_t G u\|.$$

It thus follows from finite-gain stability that there exists a constant  $K < \infty$  such that  $E_r(x, t) \leqq K \inf \|P_t u\|^2$ , where the infimum is to be taken over all  $t_0 \leqq t$  and  $u \in U$  with  $P_{t_0}u = 0, \phi(t, t_0, 0, u) = x$ . By uniform reachability,  $E_r(x, t)$  is thus decrescent. It will now be shown that  $E_r$  is a Lyapunov function. That  $E_r(x, t)$  is nonincreasing along undriven solutions follows from its definition and by letting  $u = 0$  from  $t_0$  until  $t_1$ , using an analogous argument to the one used in § 4 in

showing that  $V_r$  is monotone nonincreasing. By Theorem 3,

$$\phi(t_0 + T, t_0, x_0, 0) \rightarrow 0 \quad \text{as } T \rightarrow \infty,$$

uniformly in  $t_0$ , which by decreasence indeed implies that  $E_r$  is a uniform Lyapunov function. Since  $E_a \leq E_r$ , it remains to be shown that  $E_a$  is monotone nonincreasing along undriven solutions. This follows from an analogous argument to the one used to show that  $E_r$  is nonincreasing. This completes the proof.

Theorem 4 is not as convincing as one might like it to be since it does not make any claims about the positive definiteness of the Lyapunov functions. This positive definiteness can be obtained using somewhat stronger hypotheses. The available energy  $E_a$  will be positive definite if the feedback system with the dynamical system  $G$  in the forward and some constant gain  $k > 0$  in the feedback loop remains a well-defined dynamical system. This is the case under weak additional assumptions on  $G$ . The resulting control to be used to show definiteness of  $E_a$  is the solution  $e$  of the feedback equation  $Q_t(e + kG(t, x, e)) = 0$ . This corresponds to the input which results from a termination of the system with a positive resistor. The required energy  $E_r$  will be positive definite if  $E_a$  is or if the system is strictly passive (rather than merely passive). A third possibility is to require uniform irreversibility, since  $E_c \leq E_r$ .

**6. Feedback systems.** One of the main reasons for being interested in stability stems from its importance in feedback control. The canonical form of the feedback system considered in this paper is shown in Fig. 1, and the closed loop system is thus described by the implicit equations

$$(FE) \quad (I + G)e = u, \quad y = Ge.$$

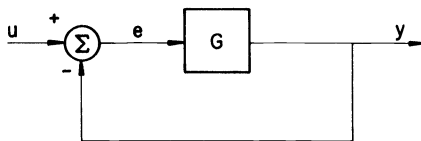


FIG. 1. *The feedback system under consideration*

It will be assumed that the input signal space  $U$  and the output signal space  $Y$  are the same and that  $G$  is a dynamical system from  $U$  into  $Y$ . Two questions related to the well-posedness of this feedback system arise: first, whether or not the closed loop feedback system still represents a well-defined dynamical system in its own right and, second, whether or not the state space induced by a realization of  $G$  will also qualify as the state space for the closed loop system. These issues fall beyond the scope of this paper, and it will be explicitly assumed rather than these well-posedness conditions are satisfied. It is thus assumed that:

(i)  $(I + G)^{-1}$  exists (as a map from  $U$  into itself) and is causal. This implies that the closed loop system  $G(I + G)^{-1} = I - (I + G)^{-1}$  is itself a dynamical system from the input space  $U$  into the output space  $Y$ .

(ii) If  $G$  is described in state space form with state space  $X$ , then  $X$  also qualifies as the state space for the closed loop dynamical system  $G(I + G)^{-1}$ , and a unique solution to the feedback equations exists for any initial condition  $x_0 \in X$ , initial time  $t_0 \in \mathbb{R}$ , and input  $u \in U$ .

These well-posedness questions have been investigated in the literature [14, § 4.2], [15] and the simplest sufficient conditions essentially impose a restriction on the feedthrough in  $G$ , in addition to some smoothness conditions on the open loop system. They are satisfied in most models and, in particular, whenever  $G$  contains a pure or generalized delay.

In the study of feedback systems it is important to establish conditions on the open loop operator in order to draw conclusions about the closed loop system. The first questions thus answered are those related to controllability and irreducibility.

**THEOREM 5.** *Consider the feedback system described by equations (FE). Then reachability (controllability, connectedness, irreducibility) of the state space realization of the open loop system implies reachability (controllability, connectedness, irreducibility) of the associated state space realization of the closed loop system. Uniform reachability (controllability, connectedness, irreducibility) of the state space realization of the open loop system implies uniform reachability (controllability, connectedness, irreducibility) of the associated state space realization of the closed loop system provided the open loop system is in addition uniformly smooth.*

*Proof.* Preservation of reachability, controllability, or connectedness is essentially obvious. Indeed, let  $u_1$  be a control which results in the desired transfer for the open loop system. Then the control  $u = u_1 + Gu_1$  will clearly result in the same transfer for the closed loop system. Irreducibility of the closed loop system will be established by contradiction. Assume therefore that there exist  $x_1, x_2 \in X, x_1 \neq x_2$ , and  $t_0 \in \mathbb{R}$  such that

$$Q_{t_0}G(I + G)^{-1}(t_0, x_1, u) = Q_{t_0}G(I + G)^{-1}(t_0, x_2, u) \quad \text{for all } u \in U.$$

This implies that

$$Q_{t_0}G(t_0, x_1, u_1) = Q_{t_0}G(t_0, x_2, u_1)$$

for all  $u_1$  which can be written as

$$u_1 = u - Q_{t_0}G(I + G)^{-1}(t_0, x_1, u), \quad u \in U.$$

Since  $u_1$  can thus be taken to be any element of  $U$  by choosing  $u = u_1 + Q_{t_0}G(t_0, x_1, u_1)$ , this shows that

$$Q_{t_0}G(t_0, x_1, u_1) = Q_{t_0}G(t_0, x_2, u_1) \quad \text{for all } u_1 \in U.$$

Hence, the open loop system is not irreducible whenever the closed loop system is not irreducible. To show uniform reachability, let  $u$  be a control such that  $\phi(t, t_0, 0, u) = x$  with  $t_0 = t - T$  and  $T$  as in the definition of uniform reachability. The control  $u + Gu$  then transfers the closed loop system from state 0 at  $t_0 = t - T$  to state  $x$  at  $t$ . Since

$$\|P_t Q_{t-T} u + P_t G Q_{t-T} u\| \leq (1 + \|P_t G Q_{t-T}\|) \|P_t Q_{t-T} u\|,$$

uniform reachability of the closed loop system thus becomes a consequence of uniform reachability of the open loop system if  $\|P_t G Q_{t-T}\|$  is uniformly bounded, which in turn is a consequence of the uniform smoothness assumption. Preservation of uniform irreducibility is shown in a similar way. This completes the proof.

The above theorem contains no surprises with the possible exception that it does not state the preservation of observability under feedback. This is in fact untrue, and it is necessary to consider nonlinear systems to obtain satisfactory counterexamples. The system  $\dot{x} = uAx$ ,  $y = Cx$  will lead to a counterexample for the contention that closed loop observability follows from open loop observability. This phenomenon is by and large a consequence of the definition of observability which is really observability under zero input. If one were to modify this definition and require that for all given inputs  $u \in U$  the response to different initial states should be different, then this observability under arbitrary inputs would indeed be preserved under feedback. These two types of observability are equivalent for linear systems. Note that the above theorem states the equivalence of reachability (controllability, connectedness, irreducibility) of the open loop and the closed loop system since putting positive unit feedback around the closed loop system gives back the open loop system. The fact that only feedback systems with unit feedback are being considered is also inessential to the basic result.

**7. Lyapunov functions for feedback systems.** In trying to define Lyapunov functions for input-output stable feedback systems, one can of course apply the techniques developed by § 4 and § 5. Such an approach is not very promising since the computation of some of the Lyapunov functions defined there requires detailed knowledge of the closed loop system, whereas it is desirable to pose the calculations and variational questions entirely in terms of the open loop dynamical system. This holds, in particular, for the function  $V_0$  defined in § 4.

It can be shown [13] that a feedback system is finite-gain input-output stable if and only if there exists a constant  $\varepsilon > 0$  such that the inequality  $\|P_t(I + G)u\| \geq \varepsilon\|P_t u\|$  holds for all  $u \in U$  and  $t \in R$ . In fact,  $\varepsilon^{-1}$  may be taken as any real number larger than the gain of  $(I + G)^{-1}$ . It should also be remarked that for linear feedback systems stability, continuity and finite-gain stability are equivalent.

Now consider the following two functions which are defined as variational problems and will lead to Lyapunov functions for feedback systems:

- (i)  $V_r(x, t) = \inf \|P_t Q_{t_0}(u + y(t_0, 0, u))\|^2$  when the infimum is to be taken over all  $u \in U$  and  $t_0 \in R$  such that  $\phi(t, t_0, 0, u) = x$  ( $\phi$  and  $y$  denote the state transition and output reading map of the open loop dynamical system).
- (ii)  $V_\varepsilon(x, t) = -\inf (\|P_{t_1} Q_{t_1}(u + y(t, x, u))\|^2 - \varepsilon^2 \|P_{t_1} Q_{t_1} u\|^2)$ , where the infimum is to be taken over all  $u \in U$  and  $t_1 \geq t$ .

**THEOREM 6.** *Assume that the feedback system described by equation (FE) is uniformly observable, finite-gain input-output stable, and that the state space is uniformly reachable. Let  $K$  denote the gain of  $(I + G)^{-1}$  and let  $0 < \varepsilon \leq K^{-1}$ . Then the feedback system is uniformly globally asymptotically stable and  $V_r$  and  $V_\varepsilon$  are positive definite radially unbounded decrescent uniform Lyapunov functions for it.*

*Proof.* Note that uniform global asymptotic stability and the claims about  $V_r$  follow from Theorem 3. It will now be shown that  $V_\varepsilon$  is finite.

Let  $u_1 \in U$  and  $t_0 \leqq t$  be such that  $P_{t_0}u_1 = 0$  and  $\phi(t, t_0, 0, u_1) = x$ . Then

$$P_{t_1}Q_{t_1}(u + y(t, x, u)) = P_{t_1}(Q_{t_1}u + P_{t_1}u_1 + y(t_0, 0, Q_{t_1}u + P_{t_1}u_1)) - P_{t_1}(u_1 + y(t_0, 0, P_{t_1}u_1)).$$

Thus,

$$\begin{aligned} \|P_{t_1}Q_{t_1}(u + y(t, x, u))\|^2 &\geqq \varepsilon^2\|P_{t_1}(Q_{t_1}u + P_{t_1}u_1)\|^2 - \|P_{t_1}(u_1 + y(t_0, 0, P_{t_1}u_1))\|^2 \\ &= \varepsilon^2\|P_{t_1}Q_{t_1}u\|^2 + \varepsilon^2\|P_{t_1}u_1\|^2 - \|P_{t_1}(u_1 + y(t_0, 0, P_{t_1}u_1))\|^2 \end{aligned}$$

and

$$-V_\varepsilon(x, t) \geqq \varepsilon^2\|P_{t_1}u_1\|^2 - \|P_{t_1}(u_1 + y(t_0, 0, P_{t_1}u_1))\|^2.$$

Since the right-hand side of this inequality depends on  $P_{t_1}u_1$  only, the result follows. The fact that  $V_\varepsilon$  is nonnegative follows from taking  $u = G(I + G)^{-1}(t, x, 0)$ . That  $V_\varepsilon$  is monotone nonincreasing along solutions follows from the usual argument explained earlier. To show that

$$V_\varepsilon(\phi(t_0 + T, t_0, x_0, 0), t_0 + T) \rightarrow 0 \quad \text{as } T \rightarrow \infty$$

uniformly in  $t_0$ , it suffices to show that  $V_\varepsilon$  is decrescent since  $\phi(t_0 + T, t_0, x_0, 0) \rightarrow 0$  as  $T \rightarrow \infty$  uniformly in  $t_0$ . It follows from the above inequality that

$$V_\varepsilon(x, t) \leqq \inf(\|P_{t_1}(u_1 + y(t_0, 0, P_{t_1}u_1))\|^2 - \varepsilon^2\|P_{t_1}u_1\|^2),$$

where the infimum is to be taken over all  $t_0 \leqq t_1$  and  $u_1 \in U$  with  $P_{t_0}u_1 = 0$  and  $\phi(t, t_0, 0, u_1) = x$ . Decrescence thus follows from uniform reachability and finite-gain stability. Positive definiteness and radial unboundedness follows by considering  $u = e = -G(I + G)^{-1}(t, x, 0)$  which yields  $V_\varepsilon(x, t) \geqq \beta(\|x\|)$  by uniform observability. This completes the proof.

The standard methods for proving stability of feedback systems is to show that the open loop gain is less than unity (small loop gain theorem) or to show that the open loop dynamical system may be viewed as the cascade of two passive systems (positive operator conditions). These cases admit special consideration and are treated in the remainder of this section.

Consider therefore the following two functions:

(i)  $V_1(x, t) = -\inf(\|Q_{t_1}u\|^2 - \|Q_{t_1}y(t, x, u)\|^2)$ , where the infimum is to be taken over all  $u \in U$  with  $\|u\| < \infty$ , and

(ii)  $V_2(x, t) = \inf(\|P_{t_1}u\|^2 - \|P_{t_1}Gu\|^2)$ , where the infimum is to be taken over all  $t_0 \in R$  and  $u \in U$  with  $P_{t_0}u = 0$  and  $\phi(t, t_0, 0, u) = x$ .

**THEOREM 7.** *Assume that the feedback system described by equations (FE) is uniformly observable, that the open loop dynamical system has gain less than unity and that the state space is uniformly reachable. Then the feedback system is finite-gain input-output stable and uniformly globally asymptotically stable and  $V_1$  and  $V_2$  are positive definite radially unbounded decrescent uniform Lyapunov functions for it.*

The case in which the open loop consists of the composition of a dynamical system  $G_1$  followed by a memoryless dynamical system  $G_2$  whose gain product is less than unity has received a great deal of attention and leads to a Lyapunov

function which only depends on  $G_1$ . Consider therefore the following functions:

(i)  $V_1^*(x, t) = -\inf(\|Q_t u\|^2 - \|G_1\|^{-2}\|Q_t G_1(t, x, u)\|^2)$ , where the infimum is to be taken over all  $u \in U$  with  $\|u\| < \infty$ , and

(ii)  $V_2^*(x, t) = \inf(\|P_t u\|^2 - \|G_1\|^{-2}\|P_t G_1 u\|^2)$ , where the infimum is to be taken over all  $t_0 \in R$  and  $u \in U$  with  $P_{t_0} u = 0$  and  $\phi(t, t_0, 0, u) = x$ .

**THEOREM 8.** *Assume that the feedback system described by equations (FE) is uniformly observable and that the state space is uniformly reachable. Let  $G = G_2 G_1$ , where  $G_1$  is a uniformly controllable, uniformly observable dynamical system and  $G_2$  is a memoryless dynamical system. Assume that the product of the gains of  $G_1$  and  $G_2$ ,  $\|G_1\| \|G_2\|$ , is less than unity. Then the feedback system is finite-gain input-output stable and uniformly globally asymptotically stable and  $V_1^*$  and  $V_2^*$  are positive definite radially unbounded decrescent uniform Lyapunov functions for it.*

*Proof.* The proofs of Theorems 7 and 8 offer no surprises considering the previous theorems, and the details will be omitted. The stability claims follow from the so-called small-gain theorem [5] for input-output stability and Theorem 2. Existence of  $V_1$ ,  $V_2$ ,  $V_1^*$  and  $V_2^*$  follows from the small gain condition, decrescence from uniform reachability, and positive definiteness by taking  $u = 0$ . Monotonicity along undriven solutions requires a minor modification of the usual argument. Consider, for instance, the function  $V_1^*$  of Theorem 8. Let  $t_1 \geq t_0$  and  $x_1 = \psi(t_1, t_0, x_0, 0)$ , with  $\psi$  the state transition map of the closed loop feedback system. Choose  $u$  on  $(t_0, t_1)$  to equal  $e = P_{t_1} Q_{t_0} (I + G)^{-1}(t_0, x_0, 0)$ . Thus,

$$V_1^*(x_0, t_0) \geq -\|e\|^2 + \|G_1\|^{-2}\|G_1(t_0, x_0, e)\|^2 + V_1^*(x_1, t_1).$$

Since, however,

$$\begin{aligned} -\|e\|^2 + \|G_1\|^{-2}\|P_{t_1} Q_{t_0} G_1(t_0, x_0, e)\|^2 &= -\|G_2 P_{t_1} Q_{t_0} G_1(t_0, x_0, e)\|^2 \\ &\quad + \|G_1\|^{-2}\|P_{t_1} Q_{t_0} G_1(t_0, x_0, e)\|^2 \\ &\geq (-\|G_2\|^2 + \|G_1\|^{-2})\|P_{t_1} Q_{t_0} G_1(t_0, x_0, e)\|^2 \end{aligned}$$

and  $\|G_1\|^{-2} - \|G_2\|^2 > 0$ ,  $V_1^*(x_0, t_0) \geq V_1^*(x_1, t_1)$  as desired. This completes the proof.

Theorem 8 is particularly useful, for instance, when  $G_1$  is linear and  $G_2$  is nonlinear or when  $G_1$  is linear and time invariant and  $G_2$  is time-varying. The variational problems which then result are indeed much simpler if one applies Theorem 8 than those needed in Theorem 7.

The two theorems which follow are the counterparts of the preceding ones, but with passivity conditions replacing the small gain condition. The stability theorem which lies at the basis of these results states that a feedback system is finite-gain input-output stable if the open loop dynamical system  $G$  is the composition of a passive system,  $G_1$ , and a strictly passive finite-gain input-output stable system,  $G_2$ . This decomposition is usually not the result of physical considerations, but rather a mathematical device which allows one to prove stability of the closed loop system.

Assume thus that  $G$ ,  $G_1$  and  $G_2$  with  $G = G_2 G_1$  are dynamical systems in state space form with state spaces  $X$ ,  $X_1$  and  $X_2$  respectively. The space  $X_1 \times X_2$  certainly qualifies as another state space for  $G$  but will, in general, be much larger than  $X$ , particularly if the latter is minimal (i.e., reachable and irreducible). In



general, this is true in stability applications since the factors  $G_1$  and  $G_2$  are usually not natural decompositions of  $G$  but are constructed with the aid of so-called “multipliers,” which usually results in this inflation of the state space. Assume now that the dynamical system  $G$  in state space form has a reachable and irreducible state space  $X$ . Consider the dynamical system  $G = G_2G_1$  with state space  $X_1 \times X_2$  and assume that the state  $(x_1, x_2) \in X_1 \times X_2$  is reachable at  $t \in R$ , i.e., that there exist a  $t_0 \leq t$  and  $u \in U$  with  $P_{t_0}u = 0$  such that  $\phi_1(t, t_0, 0, u) = x_1$  and  $\phi_2(t, t_0, 0, G_1u) = x_2$  ( $\phi_1$  and  $\phi_2$  denote the state transition maps of  $G_1$  and  $G_2$  respectively). Consider now on this subset of reachable states at time  $t$  the equivalence classes of those which yield the same output after time  $t$  for all  $u \in U$ , i.e., the reachable states  $(x'_1, x'_2)$  and  $(x''_1, x''_2)$  will be considered equivalent if

$$Q_t G_2(t, x'_2, G_1(t, x'_1, u)) = Q_t G_2(t, x''_2, G_1(t, x''_1, u)) \quad \text{for all } u \in U.$$

There is (by minimality) a one-to-one and onto correspondence between these equivalence classes and the space  $X$ . Denote by  $X_t(x_1, x_2)$  the element of  $X$  corresponding in this sense to the equivalence class derived from the reachable state  $(x_1, x_2)$ . The map  $X_t$  may in general depend explicitly on  $t$ . Assume furthermore that there exist constants  $k$  and  $K$  such that

$$k(\|x_1\|^2 + \|x_2\|^2) \leq \|X_t(x_1, x_2)\|^2 \leq K(\|x_1\|^2 + \|x_2\|^2) \quad \text{for all } t \in R$$

and reachable states  $(x_1, x_2) \in X_1 \times X_2$ . The decomposition of  $G$  into  $G = G_2G_1$  will then be called a *compatible factorization* of the dynamical system  $G$ .

The statement of the theorem which follows involves Lyapunov functions defined on  $X_1 \times X_2$ , but these can, by the above remarks, also be considered as Lyapunov functions on the state space  $X$  provided one only considers pairs  $(x_1, x_2)$  which are reachable. The following theorem statement then becomes clear.

**THEOREM 9.** *Assume that the feedback system described by equations (FE) is uniformly observable and that the state space is uniformly reachable. Assume also that the open loop dynamical system  $G$  has a uniformly reachable and uniformly irreducible state space and that it admits a compatible factorization  $G = G_2G_1$  into the uniformly observable dynamical systems  $G_1$  and  $G_2$  with uniformly reachable state spaces  $X_1$  and  $X_2$  respectively. Assume that one of these factors is passive and uniformly smooth and that the other is strictly passive and finite-gain input-output stable. Then the closed loop feedback system is finite-gain input-output stable and uniformly globally asymptotically stable, and the total available energy,  $E_a = E_a^{(1)} + E_a^{(2)}$ , and the total required energy,  $E_r = E_r^{(1)} + E_r^{(2)}$ , are decrescent uniform Lyapunov functions for it. (The superscripts refer to the dynamical systems composing  $G$ .)*

*Proof.* Decrescence of  $E_r$  on  $X_1 \times X_2$  follows from Theorem 4 with an appropriate modification in the proof in order to replace finite-gain stability by the uniform smoothness condition. Decrescence on  $X_1 \times X_2$  then implies decrescence on  $X$  by the inequality in the definition of a compatible factorization. Since  $E_a \leq E_r$ ,  $E_a$  is also decrescent. The stability claims about the feedback system are well known [5], and it remains to be shown that the energy functions are monotone nonincreasing along undriven solutions. This will only be shown for the required energy. The proof for the available energy is similar.

Let  $(x'_1, x'_2) \in X_1 \times X_2$  and  $t_0 \in \mathbb{R}$  be given, and let  $(x''_1, x''_2) \in X_1 \times X_2$  denote the state of the dynamical systems  $G_2G_1$  at time  $t_1 \geq t_0$  resulting from the transfer along solutions of the undriven feedback system. Then

$$\begin{aligned} & \inf_{\rightarrow x'_1, t_1} \langle P_{t_1} u_1, P_{t_1} G_1 u_1 \rangle + \inf_{\rightarrow x''_2, t_2} \langle P_{t_1} u_2, P_{t_1} G_2 u_2 \rangle \\ & \leq \inf_{\rightarrow x'_1, t_0} \langle P_{t_0} u_1, P_{t_0} G_1 u_1 \rangle + \inf_{\rightarrow x''_2, t_0} \langle P_{t_0} u_2, P_{t_0} G_2 u_2 \rangle \\ & \quad + \langle P_{t_1} Q_{t_0} e_1, P_{t_1} Q_{t_0} G_1(t_0, x'_1, e_1) \rangle + \langle P_{t_1} Q_{t_0} e_2, P_{t_1} Q_{t_0} G_2(t_0, x'_2, e_2) \rangle, \end{aligned}$$

where the notation  $\inf_{\rightarrow x'_1, t_1}$ , for instance, denotes the infimum over all  $t \leq t_1$  and  $u_1 \in U$  with  $P_t u_1 = 0$  and  $\phi_1(t_1, t, 0, u_1) = x'_1$ . The other symbolism is to be interpreted in an analogous way. The inputs  $e_1$  and  $e_2$  denote respectively  $(I + G)^{-1}(t_0, (x'_1, x'_2), 0)$  and  $G_1(I + G)^{-1}(t_0, (x'_1, x'_2), 0)$ . The desired result then follows if one notices that  $Q_{t_0} e_2 = Q_{t_0} G_1(t_0, x'_1, e_1)$  and that  $Q_{t_0} e_1 = -Q_{t_0} G_2(t_0, x'_2, G_1(t_0, x'_1, e_1))$  since this shows that the contributions of the last two terms in the above inequality cancel. This completes the proof.

The reader is referred to the remark following Theorem 4 for conditions to ensure positive definiteness and radial unboundedness. Notice again that positive definiteness on  $X_1 \times X_2$  suffices for positive definiteness on  $X$  by the definition of a compatible factorization. The case in which the operator  $G_2$  is memoryless leads, as in the small gain case, to a simplification. This is stated in the following final theorem.

**THEOREM 10.** *Assume that the feedback system described by equations (FE) is uniformly observable and that the state space is uniformly reachable. Assume that the open loop dynamical system,  $G = G_2G_1$ , consists of the composition of a uniformly observable, finite-gain input-output stable, strictly passive dynamical system,  $G_1$ , with a uniformly reachable state space, followed by a memoryless passive dynamical system,  $G_2$ . Then the closed loop feedback system is finite-gain input-output stable and uniformly globally asymptotically stable, and the available energy,  $E_a^{(1)}$ , and the required energy,  $E_r^{(1)}$ , are decrescent uniform Lyapunov functions for it.*

*Proof.* The proof combines the ideas in the proofs of Theorems 8 and 9 and is left to the reader.

The theorems developed here treat the small gain stability conditions and the passive operator stability conditions. The methods can, however, easily be extended to treat conic operators as well.

## 8. Examples.<sup>7</sup>

*Example 1.* Let  $G(s)$  be a  $p \times m$  matrix of rational functions of  $s$  with  $\lim_{s \rightarrow \infty} G(s) = 0$ , and assume that  $\{A, B, C\}$  is a minimal<sup>8</sup> realization of  $G(s)$ . Assume that

<sup>7</sup> The norms and inner products involved in these examples are the usual norms and inner products of Euclidean spaces. Prime denotes transposition. For the calculations involved in the solution of the variational problems in this section, see [23, §§ 21, 22, 23 and 25]. Although some of the problems are not treated explicitly there, the modifications merely require algebraic manipulation and no new methodology.

<sup>8</sup> Algebraically this means that  $A$  is an  $n \times n$  matrix, that the  $n \times nm$  and  $n \times np$  matrices  $(B, AB, \dots, A^{n-1}B)$  and  $(C, A'C, \dots, (A')^{n-1}C)$  are of full rank  $n$ , and that  $G(s) = C(Is - A)^{-1}B$ . The full

the poles of  $G(s)$ , which by minimality equal the eigenvalues of  $A$ , are in  $\text{Re } s < 0$ . The system  $\dot{x} = Ax + Bu; y = Cx$  is thus finite-gain input-output stable and globally asymptotically stable, and Theorem 3 yields as positive definite decrescent radially unbounded Lyapunov functions the quadratic forms  $x'K_1x$  and  $x'K_2x$ , where  $K_1$  is the (unique, positive definite) solution of the linear matrix equation  $A'X + XA = -C'C$ , and  $K_2^{-1}$  is the (unique, positive definite) solution of the linear matrix equation  $AX + XA' = -BB'$ . Clearly, these are only two of many possible Lyapunov functions for this asymptotically stable dynamical system.

*Example 2.* Let  $G(s)$  be an  $m \times m$  matrix of rational functions of  $s$  with  $\lim_{s \rightarrow \infty} G(s) < \infty$ , and assume that  $\{A, B, C, D\}$  is a minimal realization of  $G(s)$ . Assume that the poles of  $G(s)$  are in  $\text{Re } s < 0$ , that  $G(j\omega) + G'(-j\omega)$  is Hermitian positive definite for all  $\omega \in R$ , and that  $D + D'$  is positive definite. The  $n$ -dimensional system  $\dot{x} = Ax + Bu; y = Cx + Du$  and thus strictly passive, finite-gain input-output stable, and globally asymptotically stable. The available energy,  $E_a(x_0, t_0)$ , is given by

$$= \inf_{u \in L_2(0, \infty)} \eta, \quad \text{where } \eta = \int_0^\infty u'(t)y(t) dt,$$

subject to the constraint  $\dot{x} = Ax + Bu; y = Cx + Du, x(0) = x_0$ , and is independent of  $t_0$ . This variational problem is a least squares problem and, by Lemma 1, an infimum exists. This infimum is, in fact, attained by the feedback control

$$u = -(D + D')^{-1}(C + B'K)x$$

and

$$\min_{u \in L_2(0, \infty)} \eta = \eta^* = -x_0'Kx_0/2,$$

where  $K = K'$  is the (unique) negative definite solution of the algebraic Riccati equation

$$0 = -A'X - XA + (C + B'X)(D + D')^{-1}(C + B'X).$$

Note [8], [9], [10] that this implies the existence of an  $n \times n$  positive definite matrix  $P = P'$  ( $P = -K$ ), and  $n \times m$  matrix  $L$ , and an  $m \times m$  matrix  $W_0$  such that (Kalman-Yakubovich-Popov):

$$A'P + PA = -LL',$$

$$PB = C' - LW_0,$$

$$W_0'W_0 = D + D'.$$

---

rank condition on  $(B, AB, \dots, A^{n-1}B)$  is equivalent to controllability, reachability, and connectedness, and the full rank condition on  $(C, A'C, \dots, (A')^{n-1}C)$  is equivalent to observability and irreducibility, where these notions refer to the linear time-invariant finite-dimensional system  $\dot{x} = Ax + Bu; y = Cx$ . For these systems, the observability considered here is equivalent to observability under arbitrary inputs, and all of these properties hold uniformly whenever they hold. Global asymptotic stability requires all the eigenvalues of  $A$  to be in  $\text{Re } \lambda < 0$  and is equivalent to input-output continuity if the system is minimal.

Since, moreover, in this (linear) case the passivity is necessary and sufficient for the existence of the required infimum, and since passivity is equivalent to positive realness of  $G(s)$ , the above conditions (existence of  $K$  or  $P$ ) are also necessary for positive realness. Thus the available energy  $E_a(x, t) = x'Px/2$  is a positive definite radially unbounded decrescent Lyapunov function. The required energy,  $E_r(x_0, t_0)$ , is somewhat more involved to calculate and is defined by

$$\inf_{T \geq 0} \inf_{u \in L_2(-T, 0)} \eta,$$

where

$$\eta = \int_{-T}^0 u'(t)y(t) dt,$$

subject to the constraint<sup>9</sup>  $\dot{x} = Ax + Bu; y = Cx + Du, x(-T) = 0, x(0) = x_0$ , and is independent of  $t_0$ . The above variational problem is again a least squares problem, and by Lemma 1, an infimum exists. This infimum can be characterized as follows:

$$\eta^* = \inf_{T \geq 0} \min_{u \in L_2(-T, 0)} \eta = -x_0' \Sigma x_0 / 2,$$

where  $\Sigma = P + W^{-1}$  and  $P = P'$  is the (unique) positive definite solution of the algebraic Riccati equation

$$0 = A'X + XA + (C - B'X)(D + D')^{-1}(C - B'X).$$

In fact, this matrix is the same as the one appearing in the calculation of the available energy and is such that  $A_1 = A - B(D + D')^{-1}(C - B'P)$  is an asymptotically stable matrix.  $W$  is the (unique) solution of the linear matrix equation  $A_1X + XA_1' = -B(D + D')^{-1}B'$ , and is symmetric positive definite. The required energy  $E_r(x, t) = \frac{1}{2}x'Px + \frac{1}{2}x'W^{-1}x$  is also a positive definite radially unbounded decrescent Lyapunov function. The cycle energy  $E_c = E_r - E_a$  is given by  $E_c(x, t) = \frac{1}{2}x'W^{-1}x$ . The system is thus lossy.

*Example 3.* Let  $g(s)$  be a rational function of  $s$  with  $\lim_{s \rightarrow \infty} g(s) = 0$  and assume that  $\{A, b, c'\}$  is a minimal realization of  $g(s)$ . Let  $k$  be a scalar. Assume that the Nyquist locus of  $g(s)$  does not intersect but encircles the  $-1/k$  point in the complex plane  $-\rho$  times in the clockwise direction, where  $\rho$  is the number of poles of  $g$  in  $\text{Re } s \geq 0$ . Then the closed loop system  $\dot{x} = (A - kbc')x$  is globally asymptotically stable, and Theorem 6 yields as a positive definite radially unbounded decrescent Lyapunov function  $-x'Rx$ , where  $R = R'$  is the (unique) negative definite solution of the algebraic matrix Riccati equation

$$0 = -A'X - XA + \frac{(kc + Xb)(kc + Xb)'}{1 - \varepsilon^2} - cc'$$

with  $\varepsilon > 0$  such that  $|1 + kg(j\omega)| \geq \varepsilon_1 > \varepsilon$  for all  $\omega \in R$ .

---

<sup>9</sup> It is important to realize that this variational problem is *not* equivalent to the simpler one which asks to evaluate  $\inf_{u \in L_2(-\infty, 0)} \int_{-\infty}^0 u'(t)y(t) dt$  subject to  $\dot{x} = Ax + Bu; y = Cx + Du, x(0) = x_0$ . (This latter variational problem leads again to the available energy.)

*Example 4.* Let  $G(s)$  be a  $p \times m$  matrix of rational functions of  $s$  with  $\lim_{s \rightarrow \infty} G(s) = 0$  and assume that  $\{A, B, C\}$  is a minimal realization of  $G(s)$ . Assume that the poles of  $G(s)$  are in  $\text{Re } s < 0$ , and that for all  $\omega \in R$  the eigenvalues of the matrix  $G'(-j\omega)G(j\omega)$  are inside the open ball with radius  $\rho^{-2}$  in the complex plane. Let  $f(\sigma, t)$  be a  $R^m$ -valued function defined on  $R^p \times R$ , Lipschitz continuous on  $R^p$ , uniformly in  $t$ , and satisfying, for some  $\alpha < \rho$ , the inequality  $\|f(\sigma, t)\| \leq \alpha \|\sigma\|$  for all  $(\sigma, t) \in R^p \times R$ . Consider now the nonlinear differential equation

$$\dot{x}(t) = Ax(t) - Bf(Cx(t), t).$$

This differential equation may be viewed as the mathematical model of the undriven feedback system studied in § 6 with the open loop dynamical system described by the equations

$$\dot{x}(t) = Ax(t) + Bu(t), \quad y(t) = f(Cx(t), t),$$

and the closed loop dynamical system determined by the equations

$$\dot{x}(t) = Ax(t) - Bf(Cx(t), t) + Bu(t), \quad y(t) = f(Cx(t), t).$$

This system satisfies all the assumptions for Theorem 8 to be applicable, and the nonlinear differential equation is thus uniformly globally asymptotically stable by the small gain theorem. Consider now

$$\inf_{u \in L_2(0, \infty)} \int_0^\infty [u'(t)u(t) - \rho^2 y'(t)y(t)] dt,$$

subject to the constraint  $\dot{x} = Ax + Bu$ ;  $y = Cx$ ,  $x(0) = x_0$ . This infimum exists and is given by  $x_0'Kx_0$  when  $K$  is the (unique) negative definite solution of the matrix Riccati equation

$$0 = -A'X - XA + XBB'X + \rho^2 C'C.$$

Theorem 8 thus states that  $-x'Kx$  is a positive definite radially unbounded decrescent uniform Lyapunov function for this nonlinear differential equation. Theorem 8 yields as another positive definite radially unbounded decrescent uniform Lyapunov function  $x'(K + W^{-1})x$ , where  $K$  is as defined above and  $W$  is the (unique, positive definite) solution of the linear matrix equation

$$(A - BB'K)X + X(A - BB'K)' = -BB'.$$

*Example 5.* Let  $g(s)$  be a rational function of  $s$  with  $\lim_{s \rightarrow \infty} g(s) = 0$ , and assume that  $\{A, b, c'\}$  is a minimal realization of  $g(s)$ . Assume that the poles of  $g(s)$  are in  $\text{Re } s < 0$  and that there exists a real number  $\alpha \geq 0$  such that for some constant  $\varepsilon > 0$ ,  $\text{Re}(\alpha + j\omega)g(j\omega) \geq \varepsilon > 0$  for all  $\omega \geq 0$ . Let  $f(\sigma)$  be a real-valued function defined on  $R$ , Lipschitz continuous on  $R$ , and satisfying for some  $\delta > 0$  the inequality  $f(\sigma)/\sigma > \delta > 0$  for all  $\sigma \in R$ ,  $\sigma \neq 0$ . Consider now the nonlinear differential equation  $\dot{x}(t) = Ax(t) - bf(c'x(t))$ . This differential equation may be viewed as the mathematical model of the undriven feedback system studied in § 6 with the open loop dynamical system described by the equation

$$\dot{x} = Ax + bu; \quad y = f(c'x),$$

and the closed loop dynamical system described by the equations

$$\dot{x} = Ax - bf(c'x) + bu, \quad y = f(c'x).$$

The open loop dynamical system can be viewed as the cascade of the two systems: a linear time-invariant system with transfer function  $(s + \alpha)g(s)$  followed by a nonlinear system which is the cascade of a linear time-invariant system with transfer function  $1/(s + \alpha)$  followed by the memoryless nonlinearity  $f(\cdot)$ . The state equations of the first system are

$$\dot{x} = Ax + bu, \quad y = \alpha c'x + c'Ax + c'bu.$$

This system satisfies the assumptions of the system of Example 3 which thus yields expressions for the available energy, the required energy, and the cycle energy. The state equations for the second system are

$$\dot{z} = -\alpha z + u, \quad y = f(z).$$

The available energy  $E_a(z_0, t)$  for this system is independent of  $t_0$  and is defined by

$$E_a(z_0) = - \inf_{T \geq 0} \inf_{u \in L_2(0, T)} \eta,$$

where

$$\eta = \int_0^T u(t)y(t) dt,$$

subject to  $\dot{z} = -\alpha z + u$ ,  $y = f(z)$ ,  $z(0) = z_0$ . Thus

$$\eta = F(z(T)) + \alpha \int_0^T z(t)f(z(t)) dt - F(z_0)$$

with

$$F(z) = \int_0^z \sigma f(\sigma) d\sigma.$$

Since  $F(\sigma) \geq 0$  and  $\sigma f(\sigma) \geq 0$  for all  $\sigma \in R$ , and since the value of

$$\alpha \int_0^T z(t)f(z(t)) dt + F(z(T))$$

can be made arbitrarily small by proper choice of  $u$ , it follows that

$$E_a(z_0) = F(z_0) = \int_0^{z_0} \sigma f(\sigma) d\sigma.$$

Similarly, the required energy

$$E_r(z_0) = F(z_0) = \int_0^{z_0} \sigma f(\sigma) d\sigma.$$

The cycle energy  $E_c$  for this first order nonlinear system is thus zero, and the system is reversible. It is a simple matter to show that the above factorization of the system  $\dot{x} = Ax + bu$ ,  $y = f(c'x)$  is a compatible factorization as defined in § 7. Notice that reachable states satisfy the condition  $z_0 = c'x_0$ , which defines a hyperplane in the space  $R^n \times R$ . Theorem 9 thus yields as positive definite radially unbounded

decreasing Lyapunov functions for the nonlinear differential equation  $\dot{x} = Ax - bf(c'x)$  satisfying the conditions enumerated earlier (i.e., the conditions of the Popov criterion):

(i)  $V(x) = \frac{1}{2}x'Px + \int_0^{c'x} f(\sigma) d\sigma$ , where  $P = P'$  is the (unique) positive definite solution of the algebraic Riccati equation<sup>10</sup>

$$0 = A'X + XA + \frac{1}{2c'b}(ac' + c'A - b'X)(ac' + c'A - b'X).$$

(ii)  $V(x) + \frac{1}{2}x'W^{-1}x$ , where  $W$  is the (unique) solution of the linear matrix equation

$$A_1X + XA_1' = -\frac{bb'}{2c'b} \quad \text{with} \quad A_1 = A - \frac{b}{2c'b}(ac' + c'A - b'P),$$

and is symmetric positive definite.

*Example 6. Path integrals* [23, § 26], [25].

LEMMA 2 [23, p. 170]. Assume that  $x(t)$  is an  $n$  times differentiable function of  $t$  and that  $\alpha_{ij}, i, j = 0, 1, \dots, n$ , are constants. Then the

$$\eta = \int_{t_0}^{t_1} \sum_{i,j=0}^n \alpha_{ij} \frac{d^i x(t)}{dt^i} \frac{d^j x(t)}{dt^j} dt$$

is independent of path (i.e., it depends only on the values of  $x(t)$  and its derivatives at  $t = t_0$  and  $t = t_1$ ) if and only if the polynomial

$$h(s) = \sum_{i,j=0}^n \alpha_{ij} s^i (-s)^j + (-s)^i s^j$$

vanishes identically.

This lemma leads to rather specific formulas for the Lyapunov functions described in this paper. For instance, Theorem 3 thus yields as a Lyapunov function for the differential equation  $p(D)x(t) = 0$  with  $D = d/dt$  and  $p(s)$  a polynomial with all its roots in  $\text{Re } s < 0$ ,

$$V(x, x^{(1)}, \dots, x^{(n-1)}) = \inf_{T \geq 0} \inf_{\substack{x(t)|_{x(-T)=\dots=x^{(n-1)}(-T)=0} \\ x(0)=x, \dots, x^{(n-1)}(0)=x^{(n-1)}}} \eta$$

where  $\eta = \int_{-T}^0 (p(D)x(t))^2 dt$ . Let  $r(s)$  be a solution of the polynomial equation  $p(s)p(-s) = r(s)r(-s)$ , and let  $(p\bar{p})^+(s) = p(-s)$  and  $(p\bar{p})^-(s) = p(s)$  denote the solutions with poles respectively in  $\text{Re } s > 0$  and  $\text{Re } s < 0$ . Now rewriting  $\eta$  as

$$\eta = \int_{-T}^0 [(p(D)x(t))^2 - (r(D)x(t))^2] dt + \int_{-T}^0 (r(D)x(t))^2 dt,$$

one observes that by Lemma 2 the first integral is independent of path and thus depends on the values of  $x, x^{(1)}, \dots, x^{(n-1)}$  only. The integrand in the second integral is nonnegative and should hence be made as small as possible. By choosing

<sup>10</sup> Compare with the results of [24].

$r(s) = (p\bar{p})^+(s) = p(-s)$  and letting  $T \rightarrow +\infty$ , the contribution of this second integral can indeed be made arbitrarily small and yields a positive definite radially unbounded decrescent Lyapunov function for  $p(D)x(t) = 0$ , the quadratic form

$$V(x, \dot{x}, \dots, x^{(n-1)}) = \int_{-\infty}^0 [(p(D)x(t))^2 - (p(-D)x(t))^2] dt,$$

with  $x(t)$  any  $n$  times differentiable function such that

$$x(0) = x, \dots, x^{(n-1)}(0) = x^{(n-1)}$$

and

$$\lim_{t \rightarrow -\infty} x(t) = \dots = x^{(n-1)}(t) = 0.$$

If  $g(s)$  is chosen such that

$$\lim_{s \rightarrow \infty} \frac{q(s)}{p(s)} < \infty \quad \text{and} \quad \operatorname{Re} \frac{q(j\omega)}{p(j\omega)} \geq 0 \quad \text{for all } \omega,$$

then Theorem 4 yields as Lyapunov functions

$$E_r(x, \dot{x}, \dots, x^{(n-1)}) = \int_{-\infty}^0 [p(D)x(t)q(D)x(t) - \frac{1}{2}((p\bar{q} + \bar{p}q)^+(D)x(t))^2] dt$$

and

$$E_a(x, \dot{x}, \dots, x^{(n-1)}) = \int_0^{\infty} [p(D)x(t)q(D)x(t) - \frac{1}{2}((p\bar{q} + \bar{p}q)^-(D)x(t))^2] dt$$

with  $x(t)$  any  $n$  times differentiable function such that

$$x(0) = x, \dots, x^{(n-1)}(0) = x^{(n-1)}$$

and

$$\lim_{t \rightarrow \pm \infty} x(t) = \dots = x^{(n-1)}(t) = 0.$$

**9. Conclusions.** The development of the results and the techniques described in this paper evolves in three stages: the first one introduces and compares the input-output description with the state space description of dynamical systems and shows their equivalence. The second part in the development leads to the equivalence of input-output stability and global stability under appropriate controllability and observability conditions; the third issue is the construction of Lyapunov functions.

The methods for constructing Lyapunov functions involve, for the most part, variational problems and are posed in the framework of systems with inputs and outputs; this notwithstanding the fact that the system for which global asymptotic stability (in the sense of Lyapunov) is to be shown is an autonomous (undriven) system. The results thus obtained serve as a further relationship between the areas of dynamic optimization and stability theory and focus interest on a class of optimization problems, some of which will, in fact, lead to singular controls.



It is felt that the importance of this paper lies in its theoretical contribution in demonstrating the equivalence between global asymptotic stability and input-output stability, which, as expected, merely requires appropriate controllability and observability (more precisely: reachability and uniform observability). It also serves to unify the two main approaches to stability theory: input-output stability and Lyapunov stability. In this latter class it unifies and generalizes the various available results by posing the construction of these Lyapunov functions as variational problems.

The results of the paper could also serve as a starting point to develop techniques which will lead to suitable Lyapunov functions for estimating the domain of attraction for nonglobal stable systems. This is a problem of great practical importance, and the methods of the paper lead to tractable variational problems which could be used in such an analysis.

**Acknowledgment.** The author acknowledges Professor R. W. Brockett of Harvard University and Dr. J. L. Willems of the University of Ghent, Belgium (on leave at Harvard University) for some helpful discussions.

#### REFERENCES

- [1] L. ZADEH and C. A. DESOER, *Linear System Theory*, McGraw-Hill, New York, 1963.
- [2] A. V. BALAKRISHNAN, *On the state space theory of nonlinear systems*, Functional Analysis and Optimization, E. R. Caianiello, ed., Academic Press, New York, 1966, pp. 15–36.
- [3] V. V. NEMYTSKII AND V. V. STEPANOV, *Qualitative Theory of Differential Equations*, Princeton University Press, Princeton, N.J., 1960.
- [4] I. W. SANDBERG, *Some results on the theory of physical systems governed by nonlinear functional equations*, Bell System Tech. J., 44 (1965), pp. 871–898.
- [5] G. ZAMES, *On the input-output stability of time-varying nonlinear feedback systems. Part I: Conditions derived using concepts of loop gain, conicity, and positivity; Part II: Conditions involving circles in the frequency plane and sector nonlinearities*, IEEE Trans. Automatic Control, AC-11 (1966), pp. 228–238, 465–476.
- [6] J. C. WILLEMS, *A survey of stability of distributed parameter systems*, Control of Distributed Parameter Systems, 1969 Joint Automatic Control Conference, Boulder, Colo., ASME Publ., 1969, pp. 63–102.
- [7] R. A. BAKER AND A. R. BERGEN, *Lyapunov stability and Lyapunov functions of infinite dimensional systems*, IEEE Trans. Automatic Control, AC-14 (1969), pp. 325–334.
- [8] V. M. POPOV, *Hyperstability and optimality of automatic systems with several control functions*, Rev. Roumaine Sci. Tech. Electrotechn. et Energ., 9 (1964), no. 4, pp. 629–690.
- [9] R. E. KALMAN, *Lyapunov functions for the problem of Lur'e in automatic control*, Proc. Nat. Acad. Sci. U.S.A., 49 (1963), pp. 201–205.
- [10] B. D. O. ANDERSON, *A system theory criterion for positive real matrices*, this Journal, 5 (1967), pp. 171–182.
- [11] B. D. O. ANDERSON AND J. B. MOORE, *New results in linear system stability*, this Journal, 7 (1969), pp. 398–414.
- [12] R. W. BROCKETT, *Path integrals, Liapunov functions, and quadratic minimization*, Proc. 4th Annual Allerton Conf. on Circuit and System Theory, Monticello, Ill., 1966, pp. 685–698.
- [13] J. C. WILLEMS, *Stability, instability, invertibility and causality*, this Journal, 7 (1969), pp. 645–671.
- [14] ———, *The Analysis of Feedback Systems*, MIT Press, Cambridge, Mass., 1970.
- [15] G. ZAMES, *Realizability conditions for feedback systems*, IEEE Trans. Circuit Theory, CT-11 (1964), pp. 186–194.
- [16] A. V. BALAKRISHNAN, *On the controllability of a nonlinear system*, Proc. Nat. Acad. Sci. U.S.A., 55 (1966), pp. 465–468.
- [17] A. H. ZEMANIAN, *The Hilbert port*, SIAM J. Appl. Math., 18 (1970), pp. 98–138.

- [18] A. V. BALAKRISHNAN, *Foundations of the state-space theory of continuous systems. I*, J. Computer and System Sci., 1 (1967), pp. 91–116.
- [19] N. WIENER, *Nonlinear Problems in Random Theory*, MIT Press, Cambridge, Mass., 1958.
- [20] R. E. KALMAN, *Contributions to the theory of optimal control*, Bol. Soc. Mat. Mexicana, 5 (1960), pp. 102–119.
- [21] ———, *Canonical structure of linear dynamical systems*, Proc. Nat. Acad. Sci. U.S.A., 48 (1962), pp. 596–600.
- [22] G. LUMER AND R. S. PHILLIPS, *Dissipative operators in a Banach space*, Pacific J. Math., 11 (1961), pp. 679–698.
- [23] R. W. BROCKETT, *Finite Dimensional Linear Systems*, John Wiley, New York, 1970.
- [24] K. R. MEYER, *On the existence of Lyapunov functions for the problem of Lur'e*, this Journal, 3 (1966), pp. 373–383.
- [25] R. W. BROCKETT AND J. L. WILLEMS, *Frequency domain stability criteria. Parts I and II*, IEEE Trans. Automatic Control, AC-10 (1965), pp. 255–261, 401–413.
- [26] R. W. BROCKETT, *System theory on group manifolds and coset spaces*, submitted to this Journal.

## THE DISTRIBUTION OF $2 \times n$ GAME VALUES AND PROGRAM OPTIMA\*

DONALD J. SOULTS† AND H. T. DAVID‡

**1. Introduction.** In this paper we derive an expression (2.14) for the distribution of the value  $V$  of a  $2 \times n$  game when the  $n$  columns of the game matrix are distributed independently, each according to a bivariate distribution satisfying certain continuity conditions.

The sort of expression obtained is best illustrated under the further assumption that the columns are identically distributed. Suppose  $F$  is that common distribution, and let  $\alpha(p, v, F)$  be the probability under  $F$  of the half-plane  $px_1 + (1 - p)x_2 > v$ . Then, if  $F$  satisfies certain continuity conditions,  $P\{V > v\}$  is representable as an integral of  $\alpha^{n-1}(p, v, F)$ :

$$(1.1) \quad P\{V > v\} = \int_0^1 \alpha^{n-1}(p, v, F) dW(p).$$

Either of two functions  $W(p)$  may be used in (1.1). Of these, one typically has a discontinuity at 0 and the other typically has a discontinuity at 1.

The expectation of the distribution (2.14) is the value of a two-stage game in which Nature, at the first stage, randomly selects the game to be played at the second stage. The random selection is according to  $\prod_j F_j$ , where  $F_j$  is the distribution of the  $j$ th column. Hence the expectation of (2.14) is a possible valuation of a statistical population of  $2 \times n$  games of which a randomly selected one eventually is to be played.

By a familiar correspondence between games and programs, we obtain as well the distribution (3.2) of the optimum of the linear program :

$$(1.2) \quad \begin{aligned} & \min \{b_1 y_1 + b_2 y_2\}, \\ & a_{1j} y_1 + a_{2j} y_2 \geq c_j, \quad j = 1, 2, \dots, n, \quad y_1, y_2 \geq 0, \end{aligned}$$

where the  $b_i$  and the  $c_j$  are nonrandom and positive, and where the  $n$  vectors  $(a_{1j}, a_{2j})$  are random and independent. Because of the possibility that (1.2) is infeasible, this distribution typically is improper, with the probability deficit equal to the probability of infeasibility (cf.(3.3)).

The distribution (3.2) and probability (3.3) provide in effect closed expressions for the propagation of variability from the technological matrix to the optimum; hence, they seem of possible interest for sensitivity analysis.

As indicated below, the early parts of our argument apply to the general  $m \times n$  case. However, we have not managed to arrive at representations similar to (1.1), (2.14) or (3.2) for  $m \geq 3$ .

There is a vast literature in stochastic programming, indicated, for example, in [5]. We dwell solely on the distributional problem, to the exclusion of decision-

\* Received by the editors February 17, 1970, and in revised form June 12, 1970.

† Boeing-Aircraft Company, Renton, Washington 98055.

‡ Department of Statistics, Iowa State University, Ames, Iowa 50010.

making and economic considerations. References [1], [2], [4], [6], [7], [8] and the references cited therein have similar emphasis.

**2. The value distribution.** Let  $\sigma_m$  be the  $m$ -simplex, i.e., the set of points  $p: (p^{(1)} \dots, p^{(m)})$  with  $p^{(i)} \geq 0$  and  $\sum p^{(i)} = 1$ . Consider an  $m \times n$  game matrix  $\|X_{ij}\|, i = 1, \dots, m; j = 1, 2, \dots, n$ . The value of the corresponding game is

$$V = \max_{p \in \sigma_m} f(p),$$

where  $f(p)$ , continuous and concave on  $\sigma_m$ , is given by

$$f(p) = \min_{j=1,2,\dots,n} \sum_{i=1}^m p^{(i)} X_{ij}.$$

Hence the event  $V > v$  is equivalent to the event :

$$E: \text{for at least one } p \text{ in } \sigma_m, f(p) > v.$$

Let  $p_1, p_2, \dots, p_v$  be the points of a suitable grid  $G_v$  in  $\sigma_m$ . Define the approximating events

$$E^v = \bigcup_{k=0,1,\dots,v} E_k^v,$$

where  $E_k^v = f(p_k) > v$ .

It is readily verified using the continuity of  $f(p)$  that  $E$  is the set-theoretic limit of  $E^v$  if  $\|G_v\| \rightarrow 0$ , in which case [3, p. 40]

$$(2.1) \quad \lim_{v \rightarrow \infty} P\{E^v\} = P\{E\},$$

and inclusion-exclusion gives

$$(2.2) \quad \begin{aligned} P\{E^v\} = & \sum_{k=0}^v \langle k \rangle - \sum_{k_1 < k_2} \langle k_1, k_2 \rangle \\ & + \sum_{k_1 < k_2 < k_3} \langle k_1, k_2, k_3 \rangle \\ & \dots \dots \dots \\ & + (-1)^v \langle 0, 1, 2, \dots, v \rangle, \end{aligned}$$

where, for example,

$$\langle k_1, k_2, k_3 \rangle = P\{E_{k_1}^v \cap E_{k_2}^v \cap E_{k_3}^v\}.$$

We now discuss a reduction of (2.2), leading to (2.5) when  $m = 2$ . Consider any index set  $(k_1, \dots, k_h)$ , to be denoted by  $\mathcal{K}$ . Denote the corresponding probability  $\langle k_1, \dots, k_h \rangle$  by  $\langle \mathcal{K} \rangle$  and the corresponding subset  $(p_{k_1}, \dots, p_{k_h})$  of  $G_v$  by  $\mathcal{P}(\mathcal{K})$ . Let  $\mathcal{C}(\mathcal{K})$  be the convex hull of  $\mathcal{P}(\mathcal{K})$  and let  $\mathcal{S}(\mathcal{K})$  be the set of extreme points of  $\mathcal{C}(\mathcal{K})$ .

It is useful to single out sets  $\mathcal{K}$  with the property that

$$(2.3) \quad \mathcal{C}(\mathcal{K}) \cap [G_v - \mathcal{S}(\mathcal{K})] \neq \emptyset.$$

Probabilities  $\langle \mathcal{K} \rangle$  whose corresponding index sets  $\mathcal{K}$  satisfy (2.3) cancel in (2.2);

hence they can be excluded from the summations in (2.2). In other words, the summations in (2.2) need be extended only over index sets  $\mathcal{K}$  satisfying

$$(2.4) \quad \mathcal{C}(\mathcal{K}) \cap [G_v - \mathcal{S}(\mathcal{K})] = \emptyset.$$

The cancellation of terms  $\langle \mathcal{K} \rangle$  with  $\mathcal{K}$  satisfying (2.3) is due to two facts :

(i) In view of the concavity of  $f$  over  $\sigma_m$ ,  $[\mathcal{S}(\mathcal{K}_1) = \mathcal{S}(\mathcal{K}_2)]$  implies  $[\langle \mathcal{K}_1 \rangle = \langle \mathcal{K}_2 \rangle]$ . In other words, given  $\mathcal{S}$ ,  $\langle \mathcal{K} \rangle$  is the same for all index sets  $\mathcal{K}$  for which  $\mathcal{S}(\mathcal{K}) = \mathcal{S}$ .

(ii) For given  $\mathcal{S}$ , consider all index sets  $\mathcal{K}$  for which  $\mathcal{S}(\mathcal{K}) = \mathcal{S}$  and for which (2.3) is satisfied. Of the probabilities  $\langle \mathcal{K} \rangle$  in (2.2) corresponding to these index sets, half appear in subtracted sums, and half in added sums. For example, if  $G_v - \mathcal{S} = (p, p')$ , then the index sets  $\mathcal{K}$  with  $\mathcal{S}(\mathcal{K}) = \mathcal{S}$  and  $\mathcal{S}, \mathcal{S} \cup (p), \mathcal{S} \cup (p')$ , and  $\mathcal{S} \cup (p, p')$ , and it is clear that, if the term  $\langle \dots \rangle$  of (2.2) corresponding to  $\mathcal{S}$  appears in a subtracted sum, then so must the term  $\langle \dots \rangle$  corresponding to  $\mathcal{S} \cup (p, p')$ , while the terms  $\langle \dots \rangle$  corresponding to  $\mathcal{S} \cup (p)$  and  $\mathcal{S} \cup (p')$  must appear in added sums.

We next consider expression (2.2), and its reduction by cancellation, for the special case  $m = 2$ , with  $G_v$  the grid composed of the evenly spaced points  $p_k = (k/v, 1 - k/v)$ . For this case, the preceding considerations specialize as follows. Assuming first that  $\mathcal{K}$  is not a singleton, we have that  $\mathcal{C}(\mathcal{K})$  is the closed line segment  $\mathcal{I}(\mathcal{K})$  just containing all the points of  $\mathcal{P}(\mathcal{K})$ ,  $\mathcal{S}(\mathcal{K})$  consists of the endpoints of  $\mathcal{I}(\mathcal{K})$  and condition (2.4) specifies that, excepting its endpoints,  $\mathcal{I}(\mathcal{K})$  contains no points of  $G_v$ . Again, if  $\mathcal{K}$  contains only  $k$ , then  $\mathcal{C}(\mathcal{K}) = \mathcal{S}(\mathcal{K}) = (p_k)$ , and (2.4) is met. In other words, (2.4) is satisfied either when  $\mathcal{K}$  is a singleton or when  $\mathcal{P}(\mathcal{K})$  consists of two adjacent points of the grid, and the cancellation of terms  $\langle \dots \rangle$  corresponding to all other index sets  $\mathcal{K}$  reduces (2.2) to the more manageable expression

$$(2.5) \quad \sum_{k=0}^v \langle k \rangle - \sum_{k=0}^{v-1} \langle k, k + 1 \rangle.$$

Further, the summands of the second sum in (2.5) may be written

$$\begin{aligned} \langle k, k + 1 \rangle = P \left\{ \left( \frac{k}{v} \right) X_{11} + \left( 1 - \frac{k}{v} \right) X_{21} > v, \left( \frac{k + 1}{v} \right) X_{11} \right. \\ + \left( 1 - \left( \frac{k + 1}{v} \right) \right) X_{21} > v; \left( \frac{k}{v} \right) X_{12} + \left( 1 - \frac{k}{v} \right) X_{22} > v, \\ \left. \left( \frac{k + 1}{v} \right) X_{12} + \left( 1 - \left( \frac{k + 1}{v} \right) \right) X_{22} > v; \dots ; \left( \frac{k}{v} \right) X_{1n} \right. \\ + \left. \left( 1 - \frac{k}{v} \right) X_{2n} > v, \left( \frac{k + 1}{v} \right) X_{1n} + \left( 1 - \left( \frac{k + 1}{v} \right) \right) X_{2n} > v \right\}, \end{aligned}$$

and similarly for the summands  $\langle k \rangle$ .

Now define :

$F$  : A probability distribution on the plane,

$A(p, v)$  : the region  $px_1 + (1 - p)x_2 > v$ ,

$B(p, v)$  : the region  $px_1 + (1 - p)x_2 \leq v, x_2 \geq v$ ,

$C(p_1, p_2, v)$  : the region  $p_1x_1 + (1 - p_1)x_2 > v, p_2x_1 + (1 - p_2)x_2 > v$ ,

$$\begin{aligned} \alpha(p, v, F) &: P_F\{A(p, v)\}, \\ \beta(p, v, F) &: P_F\{B(p, v)\}, \\ \gamma(p_1, p_2, v, F) &: P_F\{C(p_1, p_2, v)\}. \end{aligned}$$

Then, assuming independence of columns, i.e., among pairs  $(X_{1j}, X_{2j})$ , and supposing that  $(X_{1j}, X_{2j})$  is distributed according to  $F_j$ , expression (2.5), and hence  $P\{E^v\}$ , may be written

$$\sum_{k=0}^{v-1} \left[ \prod_{j=1}^n \alpha\left(\frac{k}{v}, v, F_j\right) - \prod_{j=1}^n \gamma\left(\frac{k}{v}, \frac{k+1}{v}, v, F_j\right) \right] + \prod_{j=1}^n \alpha(1, v, F_j).$$

In accordance with (2.1), we must now take this to the limit with  $v$ . To this end we need the following lemma.

LEMMA 1. Suppose that  $\{\Phi_v(p)\}$  is a sequence of continuous functions on  $[0, 1]$ , converging uniformly to  $\Phi(p)$ . Let  $\{F_j; j = 1, \dots, n\}$  be any set of bivariate distributions such that the  $n$  functions  $\alpha(p, v, F_j)$  are continuous in  $p$  for  $p \in [0, 1]$  and the  $n$  functions  $\gamma(p_1, p_2, v, F)$  are uniformly (with respect to  $p_1 \in [0, 1]$ ) continuous in  $p_2$  at  $p_2 = p_1$ . Then

$$\begin{aligned} (2.6) \quad & \lim_{v \rightarrow \infty} \sum_{k=0}^{v-1} \Phi_v\left(\frac{k}{v}\right) \left[ \prod_{j=1}^n \alpha\left(\frac{k}{v}, v, F_j\right) - \prod_{j=1}^n \gamma\left(\frac{k}{v}, \frac{k+1}{v}, v, F_j\right) \right] \\ & = \sum_{j=1}^n \int_0^1 \Phi(p) \prod_{i \neq j} \alpha(p, v, F_i) d\beta(p, v, F_j). \end{aligned}$$

Before proceeding to the induction on  $n$  that verifies the lemma note that the introduction of the sequence  $\{\Phi_v\}$  is not an idle complication beyond what we actually require; rather, the  $\Phi_v$  seem required in the induction argument. Also, to abbreviate, set

$$(2.7) \quad \prod_{i \neq j} \alpha(p, v, F_i) \equiv \delta(p, j, n).$$

*Proof.* To begin with, the right-hand side of (2.6) is well-defined, since  $\Phi(z) \cdot \delta(z, j, n)$  is continuous and  $\beta(z, v, F_j)$  is nondecreasing.

For  $n = 1$ , we must show that the limit with  $v$  of

$$(2.8) \quad \left| \sum_{k=0}^{v-1} \Phi_v\left(\frac{k}{v}\right) \left[ \alpha\left(\frac{k}{v}, v, F\right) - \gamma\left(\frac{k}{v}, \frac{k+1}{v}, v, F\right) \right] - \int_0^1 \Phi(p) d\beta(p, v, F) \right|$$

is zero. But

$$\alpha\left(\frac{k}{v}, v, F\right) - \gamma\left(\frac{k}{v}, \frac{k+1}{v}, v, F\right) = \beta\left(\frac{k+1}{v}, v, F\right) - \beta\left(\frac{k}{v}, v, F\right),$$

so that (2.8) is bounded by

$$\begin{aligned} & \left| \sum_{k=0}^{v-1} \Phi\left(\frac{k}{v}\right) \left[ \beta\left(\frac{k+1}{v}, v, F\right) - \beta\left(\frac{k}{v}, v, F\right) \right] - \int_0^1 \Phi(p) d\beta(p, v, F) \right| \\ & + \left[ \sup_{0 \leq k \leq v-1} \left| \Phi_v\left(\frac{k}{v}\right) - \Phi\left(\frac{k}{v}\right) \right| \right] \cdot \beta(1, v, F), \end{aligned}$$

of which the first term tends to zero by the continuity of  $\Phi$  and monotonicity of  $\beta$ , and the last by the uniformity of the convergence of the  $\Phi_v$ .

To complete the induction using evident abbreviations we write the expression on the left-hand side of (2.6):

$$(2.9) \quad \sum_{k=0}^{v-1} \Phi_v\left(\frac{k}{v}\right) \left[ \prod_{j=1}^n \alpha_j - \prod_{j=1}^n \gamma_j \right] = \sum_{k=0}^{v-1} \left[ \Phi_v\left(\frac{k}{v}\right) \right] \left[ \frac{\alpha_n + \gamma_n}{2} \right] \left[ \prod_{j=1}^{n-1} \alpha_j - \prod_{j=1}^{n-1} \gamma_j \right] + \sum_{k=0}^{v-1} \left[ \Phi_v\left(\frac{k}{v}\right) \right] \left[ \left( \prod_{j=1}^{n-1} \alpha_j + \prod_{j=1}^{n-1} \gamma_j \right) / 2 \right] \left[ \alpha_n - \gamma_n \right].$$

Regarding the first term on the right-hand side of (2.9), set

$$\Phi_v^*(p) \equiv \Phi_v(p) \cdot [(\alpha(p, v, F_n) + \gamma(p, p + 1/v, v, F_n))/2],$$

so that product of the first two braces of the summands of that term may be written  $\Phi_v^*(k/v)$ ; and  $\Phi_v^*(p)$ , in view of the assumed uniform continuity, tends uniformly with  $v$  to

$$(2.10) \quad \Phi^*(p) \equiv \Phi(p)\alpha(p, v, F_n).$$

Hence the induction hypothesis for  $n - 1$  implies that the first term of (2.9) tends with  $v$  to

$$(2.11) \quad \sum_{j=1}^{n-1} \int_0^1 \Phi^*(p) \partial(p, j, n - 1) d\beta(p, v, F_j).$$

Similarly, by using the induction hypothesis for  $n = 1$  and defining

$$\Phi_v^{**}(p) \equiv \Phi_v(p) \cdot \left[ \left( \prod_{j=1}^{n-1} \alpha(p, v, F_j) + \prod_{j=1}^{n-1} \gamma\left(p, p + \frac{1}{v}, v, F_j\right) \right) / 2 \right],$$

$\Phi_v^{**}(p)$  tends uniformly to

$$(2.12) \quad \Phi^{**}(p) \equiv \Phi(p) \prod_{j=1}^{n-1} \alpha(p, v, F_j),$$

and the second term on the right-hand side of (2.9) tends to

$$(2.13) \quad \int_0^1 \Phi^{**}(p) d\beta(p, v, F_n).$$

By using (2.7), (2.10) and (2.12) it is clear that (2.11) and (2.13) add to the desired right-hand side of (2.6).

Substituting  $\Phi_v(p) \equiv 1$  in Lemma 1 now gives the limit required in (2.1) and the conclusion that

$$(2.14) \quad P\{V > v\} = \sum_{j=1}^n \int_0^1 \prod_{i \neq j} \alpha(p, v, F_i) d\beta(p, v, F_j) + \prod_{i=1}^n \alpha(1, v, F_i)$$

when the columns of  $\|X_{ij}\|$  are independently distributed according to distributions  $F_j$  satisfying the continuity conditions of Lemma 1.

Defining

$$W_j(p, v) = \begin{cases} \beta(p, v, F_j) & \text{for } p \in [0, 1), \\ \alpha(1, v, F_j)/n & \text{for } p = 1, \end{cases}$$

we can abbreviate (2.14) to

$$(2.15) \quad P\{V > v\} = \sum_{j=1}^n \int_0^1 \prod_{i \neq j} \alpha(p, v, F_i) dW_j(p, v).$$

Finally, define

$$B^*(p, v): \text{the region } px_1 + (1 - p)x_2 \geq v, x_2 \leq v, \\ \beta^*(p, v, F): P_F\{B^*(p, v)\}.$$

Since

$$\begin{aligned} & \int_0^1 \partial(p, j, n) d\beta^*(p, v, F_j) - \int_0^1 \partial(p, j, n) d\beta(p, v, F_j) \\ &= \int_0^1 \partial(p, j, n) d\alpha(p, v, F_j) \\ &= \prod_{i=1}^n \alpha(1, v, F_i) - \prod_{i=1}^n \alpha(0, v, F_i), \end{aligned}$$

an expression alternate to (2.15) is

$$(2.16) \quad P\{V > v\} = \sum_{j=1}^n \int_0^1 \prod_{i \neq j} \alpha(p, v, F_i) dW_j^*(p, v),$$

where

$$W_j^*(p, v) = \begin{cases} \beta^*(p, v, F_j) & \text{for } p \in (0, 1], \\ \alpha(0, v, F_j)/n & \text{for } p = 0. \end{cases}$$

We denote the right-hand side of (2.14), (2.15) or (2.16) by  $J(v; F_1, \dots, F_n)$ .

**3. Implications.** Consider the linear program (LP):

$$(3.1) \quad \begin{aligned} & \min \{u_1 + u_2\}, \\ & X_{1j}u_1 + X_{2j}u_2 \geq 1, \quad j = 1, \dots, n, \quad u_1, u_2 \geq 0. \end{aligned}$$

It is well known that, for  $yz = 1$ , LP has the optimum  $z > 0$  if and only if the game  $G$  with matrix  $\|X_{ij}\|$  has the value  $y > 0$ . Hence assuming independence among pairs  $(X_{1j}, X_{2j})$  with  $(X_{1j}, X_{2j})$  distributed according to  $F_j$  satisfying the continuity conditions of Lemma 1, we have, for  $t > 0$ ,

$$(3.2) \quad \begin{aligned} P\{\text{optimum of LP} \in (0, t) | F_1, \dots, F_n\} &= P\{V > 1/t | F_1, \dots, F_n\} \\ &\equiv J(1/t; F_1, \dots, F_n). \end{aligned}$$

Consider now the program (1.2). With the change of variables  $u_i = b_i y_i$ , that program may be written:

$$\begin{aligned} & \min \{u_1 + u_2\} \\ & (a_{1j}/d_{1j})u_1 + (a_{2j}/d_{2j})u_2 \geq 1, \quad j = 1, \dots, n, \quad u_1, u_2 \geq 0, \end{aligned}$$



where  $d_{ij} = b_i c_j$  and so is, in fact, essentially of form (3.1). Hence suppose independence among the pairs of coefficients  $(a_{1j}, a_{2j})$  for the program (1.2), and denote the distribution of  $(a_{1j}/d_{1j}, a_{2j}/d_{2j})$  by  $F_j$ . Then, if the  $F_j$  satisfy the continuity conditions of Lemma 1, the probability that the optimum  $M$  of (1.2) is in  $(0, t)$  is given by  $J(1/t; F_1, \dots, F_n)$ . Now the optimum  $M$  must be positive, and the alternative to positive  $M$  is that (1.2) is infeasible. Hence our stochastic analysis of (1.2) is completed by noting that

$$(3.3) \quad P\{(1.2) \text{ infeasible}\} = 1 - \lim_{s \rightarrow 0^+} J(s; F_1, \dots, F_n).$$

Section 2 also provides some information about the geometric behavior of  $P\{V > v\}$  or  $P\{M \in (0, t)\}$ . For example, concerning the former, suppose that all  $F_j$  are the same, say  $F$ , and that  $F$ , besides satisfying the continuity conditions of Lemma 1, also is such that  $\beta(p, v, F)$  is strictly increasing for  $p \in [0, 1]$ . Then by defining

$$\rho = \max_{0 \leq p \leq 1} \alpha(p, v, F),$$

it is clear that  $P\{V > v\}$  is geometric in  $\rho$ , in the sense that, provided  $\rho > 0$ ,

$$\lim_{n \rightarrow \infty} \frac{P\{V > v\}}{(\rho + \varepsilon)^n} = \lim_{n \rightarrow \infty} \frac{(\rho - \varepsilon)^n}{P\{V > v\}} = 0.$$

The values of the two  $1 \times n$  games corresponding respectively to the two rows of  $\|X_{ij}\|$  are of course geometric in  $\alpha(0, v, F)$  and  $\alpha(1, v, F)$  respectively.

Finally, it is of interest to verify the second formula of [2] for  $m = 2$ . This is readily done by choosing the  $F_j$  circular and centered at  $(v, v)$ , in which case (2.14) yields

$$P\{V > v\} = (n)(\frac{1}{2})^{n-1}(\frac{1}{4}) + (\frac{1}{2})^n,$$

as required.

REFERENCES

[1] BERNARD BEREANU, *On stochastic linear programming: The Laplace transform of the distribution of the optimum and applications*, J. Math. Anal. Appl., 15 (1966), pp. 280–294.  
 [2] THOMAS M. COVER, *The probability that a random game is unfair*, Ann. Math. Statist., 37 (1966), pp. 1796–1799.  
 [3] PAUL R. HALMOS, *Measure Theory*, Van Nostrand, New York, 1950.  
 [4] ANDRÁS PRÉKOPA, *On the probability distribution of the optimum of a random linear program*, this Journal, 4 (1966), pp. 211–222.  
 [5] JATI K. SENGUPTA AND GERHARD TINTNER, *The approach of stochastic linear programming: Methods and computation*, submitted to Management Sci.  
 [6] DONALD J. SOULTS, *Asymptotic value distributions for matrix games*, Doctoral thesis, Iowa State University, Ames, 1968.  
 [7] DAVID R. THOMAS AND H. T. DAVID, *Game value distributions. I*, Ann. Math. Statist., 38 (1967), pp. 242–250.  
 [8] DAVID R. THOMAS, *Game value distributions. II*, Ibid., 38 (1967), pp. 251–260.

**ERRATUM:  
 OPTIMAL CONTROLS FOR SYSTEMS WITH TIME LAG\***

A. HALANAY

In § 4, p. 231 the inequality expressing the maximum principle should read

$$\begin{aligned}
 & H(t, \tilde{x}(t), \tilde{x}(\beta_1(t)), \dots, \tilde{x}(\beta_k(t)), u, \tilde{u}(\beta_1(t)), \dots, \tilde{u}(\beta_k(t)), \tilde{b}) \\
 & \quad + H(\gamma_1(t), \tilde{x}(\gamma_1(t)), \tilde{x}(t), \dots, \tilde{x}(\beta_k(\gamma_1(t))), \tilde{u}(\gamma_1(t)), u, \dots, \\
 & \hspace{20em} \tilde{u}(\beta_k(\gamma_1(t))), \tilde{b})\dot{\gamma}_1(t) \\
 & \quad + \dots \\
 & \quad + H(\gamma_k(t), \tilde{x}(\gamma_k(t)), \tilde{x}(\beta_1(\gamma_k(t))), \dots, \tilde{x}(t), \tilde{u}(\gamma_k(t)), \tilde{u}(\beta_1(\gamma_k(t))), \dots, u, \tilde{b})\dot{\gamma}_k(t) \\
 & \leq \text{the same sum where } u \text{ is replaced by } \tilde{u}(t).
 \end{aligned}$$

On p. 233 the coordinate  $k^\rho$  should be

$$\begin{aligned}
 k^\rho = & F_\rho(t, \tilde{x}(t), \dots, \tilde{x}(\beta_k(t)), u, \tilde{u}(\beta_1(t)), \dots, \tilde{u}(\beta_k(t)), \tilde{b}) \\
 & - F_\rho(t, \tilde{x}(t), \dots, \tilde{x}(\beta_k(t)), \tilde{u}(t), \tilde{u}(\beta_1(t)), \dots, \tilde{u}(\beta_k(t)), \tilde{b}) \\
 & + [F_\rho(\gamma_1(t), \tilde{x}(\gamma_1(t)), \tilde{x}(t), \dots, \tilde{x}(\beta_k(\gamma_1(t))), \tilde{u}(\gamma_1(t)), u, \dots, \tilde{u}(\beta_k(\gamma_1(t))), \tilde{b}) \\
 & \quad - F_\rho(\gamma_1(t), \tilde{x}(\gamma_1(t)), \tilde{x}(t), \dots, \tilde{x}(\beta_k(\gamma_1(t))), \tilde{u}(\gamma_1(t)), \tilde{u}(t), \dots, \\
 & \hspace{20em} \tilde{u}(\beta_k(\gamma_1(t))), \tilde{b})]\dot{\gamma}_1(t) \\
 & + \dots \\
 & + [F_\rho(\gamma_k(t), \tilde{x}(\gamma_k(t)), \dots, \tilde{x}(t), \tilde{u}(\gamma_k(t)), \dots, u, \tilde{b}) \\
 & \quad - F_\rho(\gamma_k(t), \tilde{x}(\gamma_k(t)), \dots, \tilde{x}(t), \tilde{u}(\gamma_k(t)), \dots, \tilde{u}(t), \tilde{b})]\dot{\gamma}_k(t).
 \end{aligned}$$

---

\* This Journal, 6 (1968), pp. 215–234.

## OUTPUT CONTROLLABILITY AND SYSTEM SYNTHESIS\*

A. S. MORSE†

**Abstract.** The geometric theory of linear multivariable systems is extended by introducing the concept of a controllable output subspace. Necessary and sufficient conditions for an output subspace to be controllable are given. As an example application, controllable output subspaces are used to solve a generalized state-feedback decoupling problem.

**1. Introduction.** It has recently been shown that the geometric theory of [1] can be successfully applied to a large class of problems in linear systems synthesis. In this article, the theory is extended to encompass systems of the form

$$(1) \quad \dot{x}(t) = Ax(t) + Bu(t),$$

$$(2) \quad y(t) = Cx(t) + Du(t)$$

which differ from those previously considered by the presence of  $D$  in (2). The extended theory is based on the concept of a controllable output subspace which is discussed in the next section. This concept is closely related to the idea of a controllability subspace which plays a key role in a number of problems in linear systems synthesis. As an application of the extended theory, conditions under which the above system can be decoupled with state feedback are found.

With regard to (1) and (2),  $x$  is an  $n$ -vector,  $u$  an  $m$ -vector and  $y$  a  $p$ -vector. All matrices above and below are constant, real-valued and of appropriate size. Script letters are used for real vector spaces;  $\mathcal{X}$ ,  $\mathcal{U}$ ,  $\mathcal{Y}$  are the state, control and output spaces respectively for (1) and (2). Matrices and their maps are denoted by the same symbol. For a linear map  $M: \mathcal{W} \rightarrow \mathcal{Z}$ ,  $\{M\}$  or  $\mathcal{M}$  denotes the range of  $M$ . The notation  $\{A|\mathcal{B}\}$  is defined as

$$\{A|\mathcal{B}\} = \mathcal{B} + A\mathcal{B} + \dots + A^{n-1}\mathcal{B}.$$

It will be convenient to regard  $\mathcal{X}$  and  $\mathcal{Y}$  as subspaces of  $\mathcal{E} \equiv \mathcal{X} \oplus \mathcal{Y}$ . With  $\mathcal{E}$  so defined,  $P$  (respectively  $Q$ ) denotes the projection from  $\mathcal{E}$  onto  $\mathcal{X}$  (respectively  $\mathcal{Y}$ ) along  $\mathcal{Y}$  (respectively  $\mathcal{X}$ ). Use will be made of the extended maps

$$\bar{A}: \mathcal{E} \rightarrow \mathcal{E}, \quad e \mapsto APe;$$

$$\bar{B}: \mathcal{U} \rightarrow \mathcal{E}, \quad u \mapsto Bu;$$

$$\bar{C}: \mathcal{E} \rightarrow \mathcal{E}, \quad e \mapsto CPe;$$

$$\bar{D}: \mathcal{U} \rightarrow \mathcal{E}, \quad u \mapsto Du.$$

Note, in particular, that  $\bar{A} = \bar{A}P = P\bar{A}$ ,  $\bar{C} = Q\bar{C} = \bar{C}P$ ,  $\bar{B} = P\bar{B}$  and  $\bar{D} = Q\bar{D}$ . Below  $A$  (respectively  $B, C, D$ ) and its extension  $\bar{A}$  (respectively  $\bar{B}, \bar{C}, \bar{D}$ ) will be denoted by the same symbol  $A$  (respectively  $B, C, D$ ). The interpretation, in each instance, will be clear from the context.

\* Received by the editors May 26, 1970.

† Office of Control Theory and Application, NASA Electronics Research Center, Cambridge, Massachusetts. Now at Department of Engineering and Applied Science, Yale University, New Haven, Connecticut 06520.

The following development is based on the results of [1] which are assumed to be known.

**2. Controllable output subspaces.** With regard to (1), consider the control

$$(3) \quad u(t) = Fx(t) + Gv(t),$$

where  $v \in \mathcal{V}$  ( $q$ -space),  $G: \mathcal{V} \rightarrow \mathcal{U}$  and

$$(4) \quad F: \mathcal{E} \rightarrow \mathcal{U}.$$

For the extended map  $F$  in (4) to be consistent with its associated matrix in (3), the former must satisfy

$$(5) \quad F = FP.$$

With (3) applied to (1), the largest subspace  $\mathcal{R} \subset \mathcal{X}$  which  $v$  can completely control is given by

$$(6) \quad \mathcal{R} = \{A + BF\{\hat{B}G\}\}.$$

If for fixed  $A, B, \mathcal{R} \subset \mathcal{X}$ , a pair  $(F, G)$  exists for which (6) is true, then  $\mathcal{R}$  is called a controllability subspace of  $(A, B)$ . Existence conditions providing a characterization for controllability subspaces are given in [1].

In [1]–[3], the fundamental role played by controllability subspaces in linear systems synthesis is discussed. However, to study systems with output relations of the form in (2), an additional concept is required. Observe that application of (3) to (1) and (2) affects not only state controllability, but output controllability as well. Thus if  $\mathcal{S}$  is the largest subspace of  $\mathcal{Y}$  which  $v$  can completely control, then clearly

$$(7) \quad \mathcal{S} = (C + DF)\mathcal{R} + \{DG\}$$

with  $\mathcal{R}$  given by (6). This suggests a problem similar to the controllability subspace problem of [1]. Namely, for fixed  $\mathcal{S}, \mathcal{Y}$ , find conditions for the existence of a pair  $(F, G)$  and a subspace  $\mathcal{R} \subset \mathcal{X}$  for which both (6) and (7) are satisfied.

Conditions (6) and (7) may be simplified.

LEMMA 1. Let the maps  $\hat{A}: \mathcal{E} \rightarrow \mathcal{E}$ ,  $\hat{B}: \mathcal{U} \rightarrow \mathcal{E}$  be fixed and suppose  $\hat{A} = \hat{A}P$ . If for fixed  $G$ ,

$$(8) \quad \mathcal{R} \equiv \{P\hat{A}\{P\hat{B}G\}\}, \quad \mathcal{S} \equiv Q\hat{A}\mathcal{R} + Q\{\hat{B}G\},$$

then

$$(9) \quad \mathcal{R} = P\{\hat{A}\hat{\mathcal{B}} \cap (\mathcal{R} + \mathcal{S})\}, \quad \mathcal{S} = Q\{\hat{A}\hat{\mathcal{B}} \cap (\mathcal{R} + \mathcal{S})\}.$$

Conversely, if  $\mathcal{R}$  and  $\mathcal{S}$  are fixed subspaces satisfying (9), there exists a  $G$  for which (8) is true.

*Proof.* From the relation  $\hat{A} = \hat{A}P$  follow the identities

$$(10) \quad \begin{aligned} P\{\hat{A}\{\hat{B}G\}\} &\equiv \{P\hat{A}\{P\hat{B}G\}\}, \\ Q\{\hat{A}\{\hat{B}G\}\} &\equiv Q\hat{A}\{P\hat{A}\{P\hat{B}G\}\} + Q\{\hat{B}G\}. \end{aligned}$$

Write  $\mathcal{W} \equiv \{\hat{A}\{\hat{B}G\}\}$  and  $\mathcal{T} \equiv \{\hat{A}\hat{\mathcal{B}} \cap (\mathcal{R} + \mathcal{S})\}$ .

If (8) is true, (10) provides  $\mathcal{R} = P\mathcal{W}$  and  $\mathcal{S} = Q\mathcal{W}$ . Thus  $\mathcal{W} \subset \mathcal{R} + \mathcal{S}$ . But  $\{\hat{B}G\} \subset \hat{B} \cap \mathcal{W} \subset \hat{B} \cap (\mathcal{R} + \mathcal{S})$ , so  $\mathcal{W} \subset \{\hat{A}|\hat{B} \cap (\mathcal{R} + \mathcal{S})\} = \mathcal{T}$ . Now  $\hat{B} \cap (\mathcal{R} + \mathcal{S}) \subset \mathcal{R} + \mathcal{S}$  and if  $\mathcal{V} \subset \mathcal{R} + \mathcal{S}$ ,  $\hat{A}\mathcal{V} \subset \mathcal{R} + \mathcal{S}$ . Thus  $\mathcal{T} \subset \mathcal{R} + \mathcal{S}$ . There follows  $\mathcal{R} = P\mathcal{W} \subset P\mathcal{T} \subset \mathcal{R}$  and  $\mathcal{S} \subset Q\mathcal{W} \subset Q\mathcal{T} \subset \mathcal{S}$ . Therefore  $\mathcal{R} = P\mathcal{T}$ ,  $\mathcal{S} = Q\mathcal{T}$ , and (9) is true.

Conversely, if (9) holds,  $\mathcal{R} = P\mathcal{T}$  and  $\mathcal{S} = Q\mathcal{T}$ . By Lemma 4.1 of [1], there exists a  $G$  such that  $\mathcal{T} = \{\hat{A}|\{\hat{B}G\}\}$ . It follows from this and (10), that (8) is true. This completes the proof.

Recalling constraint (5) and applying Lemma 1 with

$$\hat{A} \equiv (A + BF) + (C + DF),$$

$$\hat{B} \equiv B + D$$

we can state the existence question as follows.

Let  $A, B, C, D, \mathcal{S} \subset \mathcal{Y}$  be fixed. Find conditions for the existence of a map  $F: \mathcal{U} \rightarrow \mathcal{E}$  and a subspace  $\mathcal{T} \subset \mathcal{E}$  such that

$$(11) \quad F = FP,$$

$$(12) \quad \mathcal{S} = Q\mathcal{T},$$

$$(13) \quad \mathcal{T} = \{A + C + (B + D)F|\{B + D\} \cap (P\mathcal{T} + \mathcal{S})\}.$$

If  $F$  and  $\mathcal{T}$  satisfying (11)–(13) exist,  $\mathcal{S}$  will be called a *controllable output subspace* of  $(A, B, C, D)$ . The subspace  $\mathcal{R} \equiv P\mathcal{T}$  will be called a *generator* of  $\mathcal{S}$ .

It is possible to solve this problem and thus to characterize controllable output subspaces. Write  $\overline{\mathcal{T}}(\mathcal{W})$  for the maximal controllability subspace of  $(A + C, B + D)$  contained in a fixed space  $\mathcal{W} \subset \mathcal{E}$ . A procedure for finding  $\overline{\mathcal{T}}(\mathcal{W})$  is given in [1]. Our main result is now presented.

**THEOREM 1.** Let  $(A, B, C, D), \mathcal{R} \subset \mathcal{X}, \mathcal{S} \subset \mathcal{Y}$  be fixed. Then  $\mathcal{S}$  is a controllable output subspace with generator  $\mathcal{R}$  if and only if

$$(14) \quad \mathcal{S} = Q\overline{\mathcal{T}}(\mathcal{R} + \mathcal{S}), \quad \mathcal{R} = P\overline{\mathcal{T}}(\mathcal{R} + \mathcal{S}).$$

*Proof.* If  $\mathcal{S}$  is a controllable output subspace with generator  $\mathcal{R}$ , then there exist  $F$  and  $\mathcal{T}$  satisfying (11)–(13) with  $\mathcal{R} = P\mathcal{T}$ . Thus,

$$\begin{aligned} (A + C + (B + D)F)(\mathcal{R} + \mathcal{S}) &= (A + C + (B + D)F)P(\mathcal{R} + \mathcal{S}) \\ &= (A + C + (B + D)F)\mathcal{T} \\ &\subset \mathcal{T} \subset \mathcal{R} + \mathcal{S}. \end{aligned}$$

By Theorem 4.3 of [1], the controllability subspace  $\mathcal{T}$  is maximal relative to  $\mathcal{R} + \mathcal{S}$ . That is,  $\mathcal{T} = \overline{\mathcal{T}}(\mathcal{R} + \mathcal{S})$  and (14) is true.

Conversely, if (14) holds,  $\mathcal{S}$  will be a controllable output subspace with  $\mathcal{R}$  as a generator if there exists an  $F$  satisfying (11) for which

$$(15) \quad \overline{\mathcal{T}}(\mathcal{R} + \mathcal{S}) = \{A + C + (B + D)F|\{B + D\} \cap (\mathcal{R} + \mathcal{S})\}.$$

Since  $\overline{\mathcal{T}} \equiv \overline{\mathcal{T}}(\mathcal{R} + \mathcal{S})$  is a controllability subspace, by Theorem 4.1 of [1],  $(A + C)\overline{\mathcal{T}} \subset \{B + D\} + \overline{\mathcal{T}}$ . This and (14) provide

$$(A + C)\mathcal{R} = (A + C)\overline{\mathcal{T}} \subset \{B + D\} + \overline{\mathcal{T}} \subset \{B + D\} + \mathcal{R} + \mathcal{S}.$$

By Lemma 3.2 of [1], there exists an  $\bar{F}$  such that  $(A + C + (B + D)\bar{F})\mathcal{R} \subset \mathcal{R} + \mathcal{S}$ . If  $F \equiv \bar{F}P$ , then

$$(A + C + (B + D)F)(\mathcal{R} + \mathcal{S}) = (A + C + (B + D)\bar{F})\mathcal{R} \subset \mathcal{R} + \mathcal{S}.$$

By Theorem 4.3 of [1], (15) holds for this choice of  $F$ . In addition,  $F$  satisfies (11) as required. This completes the proof.

In general, a fixed controllable output subspace may have many generators. It is useful, therefore, to characterize controllable output subspaces without a hypothesis involving generators. This is easily accomplished.

**COROLLARY 1.** *Let  $(A, B, C, D)$ ,  $\mathcal{S} \subset \mathcal{Y}$  be fixed. Then  $\mathcal{S}$  is a controllable output subspace if and only if there exists a controllability subspace  $\mathcal{T}$  of  $(A + C, B + D)$  such that*

$$(16) \quad \mathcal{S} = Q\mathcal{T}.$$

If (16) holds,  $\mathcal{R} \equiv P\mathcal{T}$  is a generator of  $\mathcal{S}$ .

*Proof.* Necessity is obvious. Take  $\mathcal{T} = \bar{\mathcal{T}}(\mathcal{R} + \mathcal{S})$  and apply Theorem 1.

For sufficiency, note that  $\mathcal{T} \subset \mathcal{R} + \mathcal{S}$ . Thus  $\mathcal{T} \subset \bar{\mathcal{T}}(\mathcal{R} + \mathcal{S})$ , so  $\mathcal{S} = Q\mathcal{T} \subset Q\bar{\mathcal{T}}(\mathcal{R} + \mathcal{S}) \subset \mathcal{S}$  and  $\mathcal{R} = P\mathcal{T} \subset P\bar{\mathcal{T}}(\mathcal{R} + \mathcal{S}) \subset \mathcal{R}$ . It follows that  $\mathcal{R}$  and  $\mathcal{S}$  satisfy Theorem 1 which, in turn, provides the desired result.

*Remark 1.* It is clear from Corollary 1 that  $\mathcal{S}$  is a controllable output subspace if and only if  $\mathcal{S} = Q\bar{\mathcal{T}}(\mathcal{R} + \mathcal{S})$ .

*Remark 2.* For fixed  $\mathcal{R}$  and  $\mathcal{S}$  satisfying (14), let  $\mathbf{F}$  be the class of  $F$  for which (11) and (13) hold with  $\mathcal{T} = \bar{\mathcal{T}}(\mathcal{R} + \mathcal{S})$ . By Theorem 1, this class is non-empty. Starting with the Theorem 4.2 of [1], it is a simple matter to show that the spectrum of  $A + BF$  restricted to  $\mathcal{R}$  can be freely assigned by suitable choice of  $F \in \mathbf{F}$ .

**3. Application to state-feedback decoupling.** Suppose in (2), the system output  $y$  consists of  $k$  subvectors

$$y = \begin{pmatrix} y_1 \\ \vdots \\ y_k \end{pmatrix}$$

which are to be independently controlled. Algebraically, the structure of  $y$  is equivalent to the assumption

$$(17) \quad \mathcal{Y} = \mathcal{Y}_1 \oplus \cdots \oplus \mathcal{Y}_k,$$

where  $\mathcal{Y}_i$  is the output subspace associated with  $y_i$ . Consider a control of the form

$$(18) \quad u(t) = Fx(t) + \sum_{i=1}^k G_i v_i(t).$$

The object of decoupling is to choose  $F$  and the  $G_i$  so that  $v_i$  can control  $y_i$  without influencing  $y_j$  for  $j \neq i$ . If  $\mathcal{R}_i$  is the controllability subspace of  $(A, B)$  associated with  $v_i$ , these requirements dictate that

$$(19) \quad \begin{aligned} \mathcal{R}_i &= \{A + BF | \{BG_i\}\}, & i &= 1, \dots, k, \\ \mathcal{Y}_i &= (C + DF)\mathcal{R}_i + \{DG_i\}, & i &= 1, \dots, k. \end{aligned}$$

Utilizing the results of the last section, we may formulate the decoupling problem as follows.

Find conditions for the existence of a map  $F: \mathcal{E} \rightarrow \mathcal{U}$  and subspaces  $\mathcal{T}_i$ ,  $i = 1, \dots, k$ , such that

$$(20) \quad F = FP,$$

$$(21) \quad \mathcal{T}_i = \{A + C + (B + D)F\{B + D\} \cap \mathcal{T}_i\}, \quad i = 1, \dots, k,$$

$$(22) \quad \mathcal{Y}_i = Q\mathcal{T}_i, \quad i = 1, \dots, k.$$

In short, the  $\mathcal{Y}_i$  must be controllable output subspaces and they must all be constructed with the same  $F$ .

To study this problem, use will be made of the following lemma.

LEMMA 2. Let  $\mathcal{T}_i$ ,  $i = 1, \dots, k$ , be fixed subspaces satisfying (22). Write  $\bar{\mathbf{F}}$  for the class of  $F$  for which (21) holds and assume  $\bar{\mathbf{F}}$  is nonempty. Then there exists an  $F \in \bar{\mathbf{F}}$  for which (20) is true.

*Proof.* Let  $F_0 \in \bar{\mathbf{F}}$ . If  $F_0 = F_0P$ , set  $F = F_0$ . Suppose  $F_0 \neq F_0P$ . Write  $\mathcal{T}_i = \hat{\mathcal{T}}_i \oplus \mathcal{T}_i \cap \mathcal{Y}$ ,  $i = 1, \dots, k$ , where  $\hat{\mathcal{T}}_i$  is an arbitrary completion. Since  $\mathcal{Y}_i = Q\mathcal{T}_i$ ,  $\mathcal{T}_i \subset \mathcal{X} + \mathcal{Y}_i$ . Thus,

$$\begin{aligned} \mathcal{Y} \cap \sum_{i=1}^k \hat{\mathcal{T}}_i &= \mathcal{Y} \cap \sum_{i=1}^k (\mathcal{Y}_i + \mathcal{X}) \cap \hat{\mathcal{T}}_i \\ &= \sum_{i=1}^k \hat{\mathcal{T}}_i \cap \mathcal{Y} = 0. \end{aligned}$$

Therefore the space  $\sum_{i=1}^k \hat{\mathcal{T}}_i$  is isomorphic to its projection on  $\mathcal{X}$  and there exists an  $F_1$  such that

$$F_0 \sum_{i=1}^k \hat{\mathcal{T}}_i = F_1P \sum_{i=1}^k \hat{\mathcal{T}}_i.$$

Write  $F = F_1P$ . Then for  $i = 1, \dots, k$ ,

$$(A + C + (B + D)F)\mathcal{T}_i = (A + C + (B + D)F_0)\hat{\mathcal{T}}_i \subset \mathcal{T}_i.$$

By Lemma 4.3 of [1],  $F \in \bar{\mathbf{F}}$ . In addition,  $F = FP$  as required. This completes the proof.

Using Lemma 2, we can now translate the decoupling problem (20)–(22) into a formulation for which the results of [1] are directly applicable. Define

$$(23) \quad \mathcal{N}_i \equiv \mathcal{X} + \sum_{j \neq i} \mathcal{Y}_j, \quad i = 1, \dots, k.$$

The translated formulation is as follows.

Find conditions for the existence of a map  $F: \mathcal{E} \rightarrow \mathcal{U}$  and subspaces  $\mathcal{T}_i$ ,  $i = 1, \dots, k$ , such that

$$(24) \quad \mathcal{T}_i = \{A + C + (B + D)F\{B + D\} \cap \mathcal{T}_i\}, \quad i = 1, \dots, k,$$

$$(25) \quad \mathcal{T}_i + \mathcal{N}_i = \mathcal{E}, \quad i = 1, \dots, k,$$

$$(26) \quad \mathcal{T}_i \subset \bigcap_{j \neq i} \mathcal{N}_j, \quad i = 1, \dots, k.$$

It may easily be checked that (25) and (26) are equivalent to (22). If this problem has a solution, Lemma 2 insures the existence of an  $F \in \bar{\mathbb{F}}$  for which (20) is true.

Theorem 7.1 of [1] provides the solution to this problem for the case when the number of independent open-loop system inputs equals the number of output subvectors to be controlled. This is equivalent to the assumption

$$(27) \quad \dim(\{B + D\}) = k.$$

We now give the solution to the decoupling problem, stated in notation consistent with (20)–(22).

**THEOREM 2.** *Let (27) hold. The state-feedback decoupling problem has a solution if and only if*

$$Q\bar{\mathcal{T}}_i(\mathcal{Y}_i + \mathcal{X}) = \mathcal{Y}_i, \quad i = 1, \dots, k,$$

and

$$\{B + D\} = \sum_{i=1}^k (\{B + D\} \cap \bar{\mathcal{T}}_i(\mathcal{Y}_i + \mathcal{X})).$$

Furthermore, if  $F, \bar{\mathcal{T}}_1 \dots \bar{\mathcal{T}}_k$  is any solution,

$$\mathcal{T}_i = \bar{\mathcal{T}}_i(\mathcal{Y}_i + \mathcal{X}), \quad i = 1, \dots, k.$$

**Remark 3.** Theorem 2 asserts that if a solution exists the  $\bar{\mathcal{T}}_i$  are unique. It follows that the generators  $\mathcal{R}_i = P\bar{\mathcal{T}}_i$  are also unique.

**4. Conclusion.** The state-feedback decoupling problem is by no means the only problem to which the above ideas may be applied. With some effort, controllable output subspaces may be used to extend the general decoupling problem of [2], the triangular decoupling problem of [3] and perhaps the problem of system inversion [4].

#### REFERENCES

- [1] W. M. WONHAM AND A. S. MORSE, *Decoupling and pole assignment in linear multivariable systems: A geometric approach*, this Journal, 8 (1970), pp. 1–18.
- [2] A. S. MORSE AND W. M. WONHAM, *Decoupling and pole assignment by dynamic compensation*, this Journal, 8 (1970), pp. 317–337.
- [3] ———, *Triangular decoupling of linear multivariable systems*, IEEE Trans. Automatic Control, AC-15 (1970), pp. 447–449.
- [4] L. M. SILVERMAN, *Inversion of multivariable linear systems*, Ibid., AC-14 (1969), pp. 270–277.



## ON INFORMATION STRUCTURES, FEEDBACK AND CAUSALITY\*

H. S. WITSENHAUSEN†

**Abstract.** A finite number of decisions, indexed by  $\alpha \in A$ , are to be taken. Each decision amounts to selecting a point in a measurable space  $(U_\alpha, \mathcal{F}_\alpha)$ . Each decision is based on some information fed back from the system and characterized by a subfield  $\mathcal{I}_\alpha$  of the product space  $(\prod_\alpha U_\alpha, \prod_\alpha \mathcal{F}_\alpha)$ . The decision function for each  $\alpha$  can be any function  $\gamma_\alpha$  measurable from  $\mathcal{I}_\alpha$  to  $\mathcal{F}_\alpha$ .

A property of the  $\{\mathcal{I}_\alpha\}_{\alpha \in A}$  is defined which assures that the setup has a causal interpretation. This property implies that for any combination of choices of the  $\gamma_\alpha$ , the closed loop equations have a unique solution.

The converse implication is false, when  $\text{card } A > 2$ .

**1. Introduction.** In control-oriented works on dynamic games (in particular, stochastic control problems) one usually finds a “dynamic equation” describing the evolution of a “state” in response to decision (control) variables of the players and to random variables. One also finds “output equations” which define output variables for a player as functions of the state, decision and random variables. Then the information structure is defined by allowing each decision variable to be any desired (measurable) function of the output variables generated for that player up to that time.

Such a setup assumes that the time order in which the various decision variables are selected is fixed in advance. It assumes that each player acts as if he had responsibility only for one station. It assumes that this station has perfect memory.

For large complex systems these tacit assumptions are unlikely to hold. In modeling more general situations it is often natural to consider the number of decisions to be programmed as finite, even though the alternatives to be considered can be infinite in number at each decision. But in the larger systems each “player” is in fact some form of organization unified only by its aims and central pre-planning of its policy. Implementation of policy is left for execution by “agents” (which, of course, may be devices). The order in which the various agents of the various organizations will have to act cannot always be predicted, and the information available to different agents, even of the same organization, may be non-comparable in the sense that, of two agents, neither one knows everything his colleague knows.

These difficulties in specifying the information structure of a game were faced and overcome in the early days of game theory [1], [2].

In this paper a different approach is proposed. The decision process is considered as a feedback loop and the game is characterized by its interaction with the policies of the agents, without prejudging questions of chronological order. In this way the relation between causality and the existence of a unique solution for any combination of policies can be brought into focus.

**2. The classical theory.** Von Neumann and Morgenstern [1] introduced a general description of games in extensive form for the case where the number of

---

\* Received by the editors January 15, 1970.

† Bell Telephone Laboratories, Inc., Murray Hill, New Jersey 07974.

decisions is finite and all variables range over finite sets. Their description is in terms of the agents constituting a player's organization. When the umpire calls one of these agents in the course of the play, he provides him with certain data about the situation and the agent must, on this basis, make a decision among a set of alternatives. The agent's decision function is decided in advance by the player. The set of these functions for all his agents constitutes a player's pure strategy.

One weakness of this setup was the assumption that an agent always knows the value of the time variable (move sequence number). Soon afterwards Kuhn [2] introduced a more general description which removes this restriction. A key device of Kuhn is that one separate agent is provided for each possible set of data. Such an agent will select an alternative, according to prior instructions from a player, if and when he is called by the umpire. No data is furnished to the agent. In other words, the decision functions are defined pointwise, with a separate agent for each point. Most information structure theory [3]–[5], [10], [11] is based on Kuhn's model.

The "atomization" of agents in this model is unacceptable when continuous variables are considered for the choice variables. To define the expected payoff for given pure strategies in a stochastic game, measurability on  $\sigma$ -fields smaller than the power sets is essential for the decision functions. For this reason a new formulation of Von Neumann type, with the additional feature of measurability, was constructed by Aumann [6].

For the purposes of the present paper, the Von Neumann and Aumann formulations are inadequate because they assume a fixed known sequencing of the decisions. The Kuhn formulation is inadequate for two reasons. First, the tree describing the game is an expression of the general solution of the closed loop relations. (These relations map information into decisions by the policies, and decisions into information by the rules of the game.) For any combination of policies one can find the corresponding outcome by following the tree along selected branches, and this is an explicit procedure. Thus the difficulties that might arise in solving the loop have been eliminated by defining the game in terms of a general unique solution which must be found before the model can be set up.

Second, as Aumann has pointed out, once decisions are selected from a continuum, measurability conditions must be imposed and this is not conveniently feasible with a tree type model.

The "team theory" of Radner [9] is less general than the Von Neumann description for a single player. Radner assumes that there is an initial chance move and that the players' agents receive information only about this chance move, not about each other's actions. This absence of feedback would be inappropriate in most dynamic situations.

**3. Games as multiple feedback loops.** As is customary, random effects are modeled as due to the mixed (i.e., randomized) strategy of a player called "Chance," who may have several agents. Once this is done the system upon which the agents act is completely deterministic.

Let  $n > 0$  be the finite cardinality of the set  $A$  of all agents in the game. Assume, without loss of generality, that each agent  $\alpha$  in  $A$  will take one decision during the game. If there are circumstances where an agent is never called to action, then one can still let him make a decision which will have no effect. If an agent

were to act more than once, then either the information available on all these occasions is the same, and then these decisions are considered as a single one picked from a product set, or else the information is not the same, and then one provides a separate agent for each occasion.

Agent  $\alpha$  must select his decision from a nonempty set  $U_\alpha$  on which is given a  $\sigma$ -field  $\mathcal{F}_\alpha$ . The policy of the agent will be specified as the dependence of the decision upon information available to him and upon the outputs of random devices that may be used to mix a player's strategy.

The crucial observation is that all variables in the game will be determined by the decisions actually taken, that is, by the element  $u$  of the Cartesian product  $U = \prod_{\alpha \in A} U_\alpha$ . Note that this Cartesian product consists of those maps  $u$  of  $A$  into the union of the  $U_\alpha$  which satisfy  $u(\alpha) \equiv u_\alpha \in U_\alpha$  for all  $\alpha$  in  $A$ . This definition does not require any ordering of the set  $A$ .

For a nonvoid subset  $B$  of  $A$ , denote by  $P_B$  the projection of the product  $\prod_{\alpha \in A} U_\alpha$  upon the product  $\prod_{\alpha \in B} U_\alpha$ . The projection  $P_B(u)$  is just the restriction of  $u$  to the domain  $B$ . Let  $P_\alpha(u) \equiv P_{\{\alpha\}}(u) \equiv u_\alpha$ . Let  $\prod_{\alpha \in B} \mathcal{F}_\alpha$  be the product  $\sigma$ -field on  $\prod_{\alpha \in B} U_\alpha$  and let  $\mathcal{F}(B)$  be the smallest  $\sigma$ -field on  $U$  such that  $P_B$  is measurable. Note that  $B_1 \subset B_2$  implies  $\mathcal{F}(B_1) \subset \mathcal{F}(B_2)$ . The  $\sigma$ -field  $\mathcal{F}(A) = \prod_{\alpha \in A} \mathcal{F}_\alpha$ , also denoted by  $\mathcal{F}$ , is the finest field to be considered on  $U$ .

The outcome of the game, that is, the list of payoffs to the various players, is determined by  $u$ . These payoffs are functions from  $U$  into real numbers, measurable on  $\mathcal{F}$ . Likewise, any data, about the situation, received by agent  $\alpha$ , is information about  $u$  and will determine a  $\sigma$ -field  $\mathcal{I}_\alpha \subset \mathcal{F}$ , the *information field* of agent  $\alpha$ .

For a nonstochastic game and if only pure strategies are considered, the policy that agent  $\alpha$  is assigned by his organization can be any function  $\gamma_\alpha: U \rightarrow U_\alpha$  measurable as a map of the measurable space  $(U, \mathcal{I}_\alpha)$  into  $(U_\alpha, \mathcal{F}_\alpha)$ . Denoting by  $\Gamma_\alpha$  this set of functions, let  $\Gamma = \prod_{\alpha \in A} \Gamma_\alpha$ .

Then for the combination of policies  $\gamma \in \Gamma$  one has the problem of solving the closed loop equations

$$(1) \quad P_\alpha(u) = \gamma_\alpha(u) \quad \text{for all } \alpha \in A.$$

When random effects are present, a sample space  $(\Omega, \mathcal{B}, P)$  for the joint behavior of all random devices is introduced.  $\mathcal{B}$  is a  $\sigma$ -field on the set  $\Omega$  and  $P$  a probability measure on  $(\Omega, \mathcal{B})$ . For each agent a  $\sigma$ -field  $\mathcal{B}_\alpha \subset \mathcal{B}$  is given. The possible policies of agent  $\alpha$  are the functions

$$\gamma_\alpha: U \times \Omega \rightarrow U_\alpha$$

such that  $\gamma_\alpha^{-1}(\mathcal{F}_\alpha) \subset \mathcal{I}_\alpha \times \mathcal{B}_\alpha$ . Again  $\Gamma_\alpha$  will denote this set of functions and  $\Gamma = \prod_{\alpha \in A} \Gamma_\alpha$ .

Then for a given combination of policies  $\gamma \in \Gamma$  the problem is to find, for each  $\omega$  in  $\Omega$ , solutions  $u$  (dependent on  $\omega$ ) to the closed loop equations

$$(2) \quad P_\alpha(u) = \gamma_\alpha(u, \omega) \quad \text{for all } \alpha \in A.$$

One would like that there be one and only one solution for each  $\omega$ , with  $\mathcal{B}$ -measurable  $\omega$ -dependence. One also would like this solution to be obtained by actual play of the extensive form of the game in a way compatible with causality.

None of these properties will hold in general for an arbitrary *information structure* characterized by nonvoid sets  $A, U_\alpha$  and  $\sigma$ -fields  $\mathcal{F}_\alpha, \mathcal{I}_\alpha$ , with the possible additional specification of  $\Omega, \mathcal{B}$  and  $\mathcal{B}_\alpha$ .

For such a structure to become a game one would have to specify a partition of  $A$  into players, a probability measure  $P$  on  $(\Omega, \mathcal{B})$ , the mixed strategy of the chance player's agents and the payoff function for each other player. The selection of  $(\Omega, \mathcal{B}, \mathcal{P})$  and  $\mathcal{B}_\alpha$  must take into account that each player or coalition of players uses an independent random device and that the chance player must be considered as a separate coalition in that respect.

Only the information structure is of concern in this paper.

**4. Some basic questions.** What condition on the fields  $\mathcal{I}_\alpha$  expresses the causality requirement? Since the requirement is a heuristic one, a formal condition (property C below) will be stated and strong arguments given to relate it to the idea of causality.

Define *property S* (solvability) as follows: For each  $\gamma \in \Gamma$  and  $\omega \in \Omega$  there exists one and only one  $u \in U$  satisfying the closed loop equations.

If property S holds, then for each  $\gamma \in \Gamma$  the solution defines a map  $M^\gamma$  from  $\Omega$  into  $U$ . *Property SM* holds when this map is always measurable from  $(\Omega, \mathcal{B})$  into  $(U, \mathcal{F})$ . The importance of this property is that it makes the payoff of each player measurable on the sample space, so that the expected payoff can be defined.

Another important consequence of property SM is that the induced fields  $\mathcal{I}_\alpha^\gamma$ , inverse images under  $M^\gamma$  of  $\mathcal{I}_\alpha$ , are contained in  $\mathcal{B}$ . Thus conditional expectations can be defined with respect to  $\mathcal{I}_\alpha^\gamma$ . Conditional expectations with respect to  $\mathcal{I}_\alpha$  are meaningless because this field is not on the probability space  $\Omega$ . On the other hand, the meaningful conditioning on  $\mathcal{I}_\alpha^\gamma$  yields results that may actually depend on  $\gamma$ . Casual treatment of this point has been a source of confusion in stochastic control theory.

The following relations hold among these properties: C implies SM but, surprisingly, SM does not imply C when more than two agents are involved.

**5. The causality property.** The main difficulty is one of notation. Let  $n$  be the cardinality of the set  $A$  of agents, and for  $1 \leq k \leq n$  denote by  $S_k$  the set of all injections of  $\{1, 2, \dots, k\}$  into  $A$ , i.e., the arrangements without repetition of  $k$  elements from  $A$ . Let  $S^+ = \bigcup_{k=1}^n S_k$ ,  $S_0 = \{\emptyset\}$  and  $S = S_0 \cup S^+$ .

For  $0 \leq i \leq j \leq n$  let  $T_i^j: S_j \rightarrow S_i$  be the truncation map which restricts  $s \in S_j$  to the domain  $\{1, \dots, i\}$  or to  $\emptyset$  for  $i = 0$ .

For  $s \in S$  denote by  $|s|$  the range of  $s$  and by  $\|s\|$  the cardinality of  $|s|$ . Note that  $s \in S_{\|s\|}$ .

It will be assumed that for all  $\alpha$  in  $A$ , and  $x$  in  $U_\alpha$ , the singleton set  $\{x\}$  belongs to  $\mathcal{F}_\alpha$ . As a consequence, the singletons of  $\prod_{\alpha \in B} U_\alpha$  belong to  $\prod_{\alpha \in B} \mathcal{F}_\alpha$ . In particular, the singletons of  $U$  belong to  $\mathcal{F}$ .

Physically, decisions can be considered as events in space-time. Then the causality condition means that, for any play of the game, the actions of the agents can be ordered and that the information available to an agent may depend on decisions of agents acting earlier but cannot depend upon the decisions of agents acting at concurrent or later times. In the language of relativity theory two agents act at concurrent times if the interval separating their actions is of spatial, as

opposed to temporal, character. For different plays of the game the order may turn out to be different because the decisions made were different and these decisions affect the space-time positions of other decisions.

Thus the causality condition requires that the field  $\mathcal{I}_\alpha$ , subject to the condition that the set of agents acting before  $\alpha$  is  $B$ , is contained in the field  $\mathcal{F}(B)$  determined by the decisions of the earlier agents.

A way of stating this requirement precisely is the following.

*Property C.* An information structure  $\{U_\alpha, \mathcal{F}_\alpha, \mathcal{I}_\alpha\}_{\alpha \in A}$  is said to possess property C when there exists a map  $\varphi: U \rightarrow S_n$  such that for all  $k, s, E$  satisfying  $1 \leq k \leq n, s = (\alpha_1, \dots, \alpha_k) \in S_k, E \in \mathcal{I}_{\alpha_k}$ , one has

$$(3) \quad E \cap (T_k^n \circ \varphi)^{-1}(s) \in \mathcal{F}(\|T_{k-1}^k(s)\|).$$

Note that in general a set  $E$  in  $\mathcal{I}_\alpha$  need not belong to  $\mathcal{F}(B)$  for any proper subset  $B$  of  $A - \{\alpha\}$ . Neither is it necessary that the set  $(T_k^n \circ \varphi)^{-1}(s)$  belong to  $\mathcal{I}_\alpha$  for any  $\alpha$ . Also note that for  $k = 1$  one has  $T_0^1(s) = \emptyset$  and  $\mathcal{F}(\emptyset) = \{\emptyset, U\}$ .

In abbreviated notation the property can be stated as the existence of  $\varphi$  such that for  $1 \leq k \leq n$  and  $s = (\alpha_1, \dots, \alpha_k)$  in  $S_k$ , one has

$$\mathcal{I}_{\alpha_k} \cap (T_k^n \circ \varphi)^{-1}(s) \subset \mathcal{F}(\alpha_1, \dots, \alpha_{k-1}).$$

The function  $\varphi$  need not be unique, if only because a transposition of concurrently acting agents cannot have any effect.

LEMMA 1. *If property C holds with  $\varphi$ , then for all  $s \in S^+$ ,*

$$(4) \quad (T_{\|s\|}^n \circ \varphi)^{-1}(s) \in \mathcal{F}(\|T_{\|s\|-1}^{\|s\|}(s)\|).$$

*Proof.* For  $s \in S_k$  one has  $\|s\| = k$  and since  $U$  belongs to  $\mathcal{I}_\alpha$  for all  $\alpha$ , (3) with  $E = U$  yields (4).

This lemma says that after  $k - 1$  agents have taken their decisions, the selection of the next agent to act is thereby determined mathematically. None of the agents, though, need to know this selection. Even the selected agent need not know that he is the  $k$ th one to act nor which  $k - 1$  agents have already acted.

LEMMA 2. *The function  $T_1^n \circ \varphi$  has a constant value  $(\alpha^*)$  and  $\mathcal{I}_{\alpha^*} = \{\emptyset, U\}$ .*

*Proof.* Lemma 1, with  $\|s\| = 1$ , shows that  $T_1^n \circ \varphi$  is measurable on  $\mathcal{F}(\emptyset) = \{\emptyset, U\}$  and is therefore a constant. Denoting this constant by  $(\alpha^*)$ , property C with  $k = 1$  and  $s = (\alpha^*)$  gives for any set  $E$  in  $\mathcal{I}_{\alpha^*}$ ,

$$E \in \mathcal{F}(\emptyset) = \{\emptyset, U\}$$

which shows that  $\mathcal{I}_{\alpha^*} = \{\emptyset, U\}$ , completing the proof.

The interpretation of Lemma 2 is that there has to be at least one ‘‘starting agent’’ and that a starting agent cannot have any information about decisions taken since none have been taken yet. For instance, many games begin with either one of two agents depending on ‘‘the toss.’’ In that case the starting agent is a chance agent who decides according to a fixed mixed strategy the outcome of the toss. The selection of the next agent then depends on this outcome that is on the first agent’s decision.

The same argument as in Lemma 2 shows that the function  $T_2^n \circ \varphi$  is measurable on  $\mathcal{F}(\{\alpha^*\})$ .

**6. The umpire's information fields.** Continuing with the assumption that property C holds with function  $\varphi$ , the possibility for the umpire to organize the play in an orderly manner is to be demonstrated. To this end, introduce the *umpire fields*  $\mathcal{F}^{(k)}$  defined as follows:

$$(5) \quad \begin{aligned} \mathcal{F}^{(0)} &\equiv \{\emptyset, U\}, \\ \mathcal{F}^{(k)} &\equiv \left\{ \bigcup_{s \in S_k} [(T_k^n \circ \varphi)^{-1}(s) \cap F_s] \mid \text{for all } s \in S_k, F_s \in \mathcal{F}(\|s\|) \right\} \end{aligned}$$

for  $1 \leq k \leq n$ .

Heuristically  $\mathcal{F}^{(k)}$  is the field characterizing the information available to the umpire after he has seen  $k$  agents take their decisions.

**LEMMA 3.** *The collections of sets  $\mathcal{F}^{(k)}$  are  $\sigma$ -fields on  $U$ .  $\mathcal{F}^{(1)} = \mathcal{F}(\{\alpha^*\})$ ,  $\mathcal{F}^{(k)} \subset \mathcal{F}$  and  $\mathcal{F}^{(n)} = \mathcal{F}$ .*

*Proof.* That the  $\mathcal{F}^{(k)}$  are indeed  $\sigma$ -fields follows from Auxiliary Lemma 1 (see Appendix) by specialization, noting that  $\{(T_k^n \circ \varphi)^{-1}(s) \mid s \in S_k\}$  is a collection of pairwise disjoint sets covering  $U$ . For  $k = \|s\| = 1$  this collection has only one nonempty member  $(T_1^n \circ \varphi)^{-1}(\alpha^*) = U$ . Thus

$$\mathcal{F}^{(1)} = \{F_{(\alpha^*)} \mid F_{(\alpha^*)} \in \mathcal{F}(\{\alpha^*\})\} = \mathcal{F}(\{\alpha^*\}).$$

Each set  $F_s$  in (5) belongs to a subfield of  $\mathcal{F}$ . The same is true for each set  $(T_k^n \circ \varphi)^{-1}(s)$  by virtue of Lemma 1. Thus each set in  $\mathcal{F}^{(k)}$  belongs to  $\mathcal{F}$ .

For  $k = \|s\| = n$  one has  $|s| = A$  and  $\mathcal{F}(\|s\|) = \mathcal{F}$ . Therefore, given a set  $F$  in  $\mathcal{F}$ , one may take  $F_s = F$  for all  $s \in S_n$  in (5). This represents  $F$  as a set in  $\mathcal{F}^{(n)}$ . Hence  $\mathcal{F} = \mathcal{F}^{(n)}$ .

A most important fact is that the umpire fields are nested.

**LEMMA 4.**  *$\mathcal{F}^{(k)} \subset \mathcal{F}^{(k+1)}$  for  $0 \leq k < n$ .*

*Proof.* The assertion is trivial for  $k = 0$ . For  $0 < k < n$  and  $E \in \mathcal{F}^{(k)}$ , one has

$$E = \bigcup_{s \in S_k} [(T_k^n \circ \varphi)^{-1}(s) \cap F_s]$$

with  $F_s \in \mathcal{F}(\|s\|)$  for all  $s \in S_k$ .

Now since  $T_k^n = T_k^{k+1} \circ T_{k+1}^n$ , one has

$$(T_k^n \circ \varphi)^{-1} = (T_k^{k+1} \circ T_{k+1}^n \circ \varphi)^{-1} = (T_{k+1}^n \circ \varphi)^{-1} \circ (T_k^{k+1})^{-1}$$

for the inverse set mappings. Thus

$$(T_k^n \circ \varphi)^{-1}(s) = \bigcup_{s'} \{(T_{k+1}^n \circ \varphi)^{-1}(s') \mid s' \in S_{k+1}, s = T_k^{k+1}(s')\}.$$

Then

$$\begin{aligned} E &= \bigcup_{s \in S_k} [F_s \cap \bigcup_{s'} \{(T_{k+1}^n \circ \varphi)^{-1}(s') \mid s' \in S_{k+1}, s = T_k^{k+1}(s')\}] \\ &= \bigcup_s \bigcup_{s'} \{F_s \cap (T_{k+1}^n \circ \varphi)^{-1}(s') \mid s \in S_k, s' \in S_{k+1}, s = T_k^{k+1}(s')\} \\ &= \bigcup_{s' \in S_{k+1}} [F_{T_k^{k+1}(s')} \cap (T_{k+1}^n \circ \varphi)^{-1}(s')]. \end{aligned}$$

But if  $s = T_k^{k+1}(s')$  then, as sets,  $|s| \subset |s'|$ , and therefore  $\mathcal{F}(\|s\|) \subset \mathcal{F}(\|s'\|)$  so that

$$F_{T_k^{k+1}(s')} \in \mathcal{F}(\|s'\|).$$

Comparison with (5) shows that  $E \in \mathcal{F}^{(k+1)}$  which proves the lemma.

Heuristically, Lemma 4 says that the umpire enjoys perfect recall. The next lemma says that when the umpire has obtained decisions of  $k - 1$  agents he is able to tell which agent comes into play next.

LEMMA 5. For  $0 < k \leq n$ ,  $T_k^n \circ \varphi$  is measurable on  $\mathcal{F}^{(k-1)}$ .

*Proof.* It suffices to show that for any  $s'$  in  $S_k$  the set  $(T_k^n \circ \varphi)^{-1}(s')$  is of the form

$$(6) \quad \bigcup_{s \in S_{k-1}} [(T_{k-1}^n \circ \varphi)^{-1}(s) \cap F_s]$$

with  $F_s$  in  $\mathcal{F}(|s|)$  for all  $s$  in  $S_{k-1}$ .

Let  $s^* = T_{k-1}^k(s') \in S_{k-1}$  and let  $F_s = (T_k^n \circ \varphi)^{-1}(s')$  for  $s = s^*$ , noting that, by Lemma 1, this set is in  $\mathcal{F}(|s^*|)$ , and  $F_s = \emptyset$  for  $s \in S_{k-1} - \{s^*\}$ .

Then (6) becomes

$$(T_{k-1}^n \circ \varphi)^{-1}(s^*) \cap (T_k^n \circ \varphi)^{-1}(s').$$

Since  $s^* = T_{k-1}^k(s')$  the first of these sets contains the second and the intersection is just the second set, establishing the claim.

**7. The causal solution process.** Assuming that property C holds with function  $\varphi$ , it will now be shown that for  $\gamma \in \Gamma$  the closed loop equations

$$(7) \quad P_\alpha(u) = \gamma_\alpha(u, \omega) \quad \text{for all } \alpha \in A$$

admit, for each  $\omega$ , a unique solution and the dependence of this solution upon  $\omega$  is measurable. Furthermore the solution process can be organized in a recursive way, corresponding to causal play of the game under the direction of the umpire.

LEMMA 6. Equations (7) cannot admit two distinct solutions  $u, \bar{u}$  for the same decision functions  $\gamma$  and the same value of  $\omega$ .

*Proof.* Since  $\omega$  is fixed, one need only consider the section  $\gamma_\alpha(\cdot, \omega)$  of the functions  $\gamma_\alpha$ . These sections are  $\mathcal{I}_\alpha$ -measurable and will be denoted by  $\gamma_\alpha(\cdot)$  for brevity.

Assume then that for some  $u, \bar{u} \in U$  and for all  $\alpha$  in  $A$ , one has both  $P_\alpha(u) = \gamma_\alpha(u)$  and  $P_\alpha(\bar{u}) = \gamma_\alpha(\bar{u})$ .

By Lemma 2, the relations

$$(8) \quad T_k^n(\varphi(u)) = T_k^n(\varphi(\bar{u})) = s = (\alpha_1, \dots, \alpha_k),$$

$$(9) \quad P_{\alpha_i}(u) = P_{\alpha_i}(\bar{u}), \quad i = 1, \dots, k - 1,$$

hold trivially for  $k = 1$ . But when (8), (9) hold for some  $k$  one has

$$(10) \quad P_{\alpha_k}(u) = P_{\alpha_k}(\bar{u}).$$

Indeed, otherwise the sets

$$E = \gamma_{\alpha_k}^{-1}(P_{\alpha_k}(u)) \quad \text{and} \quad \bar{E} = \gamma_{\alpha_k}^{-1}(P_{\alpha_k}(\bar{u}))$$

would be disjoint sets in the field  $\mathcal{I}_{\alpha_k}$ . This is so because of the assumption that the singleton sets  $\{P_{\alpha_k}(u)\}$  and  $\{P_{\alpha_k}(\bar{u})\}$  in  $U_{\alpha_k}$  belong to  $\mathcal{F}_{\alpha_k}$ . Then by property C, the set

$$F = E \cap (T_k^n \circ \varphi)^{-1}(s),$$

which contains  $u$  and the set

$$\bar{F} = \bar{E} \cap (T_k^n \circ \varphi)^{-1}(s),$$

which contains  $\bar{u}$ , would be disjoint sets in the field  $\mathcal{F}(\{\alpha_1, \dots, \alpha_{k-1}\})$ . But since that field is generated by the function  $P_{\{\alpha_1, \dots, \alpha_{k-1}\}}$  which takes the same value at  $u$  and  $\bar{u}$ , a contradiction of Auxiliary Lemma 2 (see Appendix) is obtained.

If  $k = n$ , then (8), (9), (10) state that  $u = \bar{u}$ ; if  $k < n$ , then one will have

$$(11) \quad s' \equiv T_{k+1}^n(\varphi(u)) = T_{k+1}^n(\varphi(\bar{u})) = \bar{s}'.$$

Indeed otherwise  $u$  and  $\bar{u}$  would belong respectively to the disjoint sets  $(T_{k+1}^n \circ \varphi)^{-1}(s')$  and  $(T_{k+1}^n \circ \varphi)^{-1}(\bar{s}')$ . By Lemma 1 these two sets both belong to the field  $\mathcal{F}(\{\alpha_1, \dots, \alpha_k\})$ . This field is determined by the function  $P_{\{\alpha_1, \dots, \alpha_k\}}$ , which by (9), (10) takes the same value at  $u$  and  $\bar{u}$ . Thus auxiliary Lemma 2 would again be contradicted.

Now with the adjunction of (10) and (11), the relations (8), (9) hold with  $k$  increased by unity. Thus Lemma 6 follows by induction.

The successive steps of the recursive solution process, given a set of maps  $\gamma_\alpha: U \times \Omega \rightarrow U_\alpha$ , are described by maps

$$M_k^\gamma: U \times \Omega \rightarrow U,$$

where the superscript  $\gamma$  refers to the system of decision functions and  $k = 1, \dots, n$ . These maps are defined with the help of an arbitrary reference element  $r \in U$ .

The action of  $M_k^\gamma$  is defined as follows:

$$(12) \quad P_\alpha(M_k^\gamma(u, \omega)) = \begin{cases} P_\alpha(u) & \text{for } \alpha \in |T_{k-1}^n(\varphi(u))|, \\ \gamma_\alpha(u, \omega) & \text{for } \alpha = (\varphi(u))_k, \\ P_\alpha(r) & \text{otherwise.} \end{cases}$$

LEMMA 7. *The function  $M_k^\gamma$  is measurable, as a map of the measurable space  $(U \times \Omega, \mathcal{F}^{(k-1)} \times \mathcal{B})$  into the measurable space  $(U, \mathcal{F}^{(k)})$ .*

*Proof.* The set  $U$  can be partitioned into the sets  $(T_k^n \circ \varphi)^{-1}(s)$  with  $s$  ranging over  $S_k$ . By Lemma 5 these sets belong to  $\mathcal{F}^{(k-1)}$  and, a fortiori, to  $\mathcal{F}^{(k)}$  (Lemma 4).

The restriction of  $M_k^\gamma$  to each set in the partition leaves the sequence  $s = (\alpha_1, \dots, \alpha_k)$  unchanged because it preserves the values  $u_{\alpha_1}, \dots, u_{\alpha_{k-1}}$  upon which the restriction of  $T_k^n \circ \varphi$  to the same domain depends.

On each set in the partition, the division of the agents into the three groups of (12) is fixed. By the definition of the umpire fields it thus suffices to check that for each  $s = (\alpha_1, \dots, \alpha_k)$  on  $S_k$  the restriction of  $M_k^\gamma$  to  $(T_k^n \circ \varphi)^{-1}(s) \times \Omega$  is measurable from the field  $\mathcal{F}(\{\alpha_1, \dots, \alpha_{k-1}\}) \times \mathcal{B}$  to the field  $\mathcal{F}(\{\alpha_1, \dots, \alpha_k\})$ . By definition of the latter field, it suffices to check the measurability of  $P_{\alpha_i} \circ M_k^\gamma$ , restricted to the same domain, for  $i = 1, \dots, k$ . Now for  $i = 1, \dots, k - 1$  the restriction is the projection  $P_{\alpha_i}$  defining  $\mathcal{F}(\{\alpha_1, \dots, \alpha_{k-1}\})$ . For the index  $k$  the function  $\gamma_{\alpha_k}$  is obtained. This function is measurable on  $\mathcal{I}_{\alpha_k} \times \mathcal{B}$  and the intersections of sets in  $\mathcal{I}_{\alpha_k}$  with  $(T_k^n \circ \varphi)^{-1}(s)$  are in  $\mathcal{F}(\{\alpha_1, \dots, \alpha_{k-1}\})$  by virtue of property C. This establishes Lemma 7.



Heuristically, the maps  $M_k^\gamma$  correspond to the idea of state transition equations. From that point of view it is the field  $\mathcal{F}^{(k)}$  that represents the state (of knowledge of the umpire).

**THEOREM 1.** *Property C implies property SM.*

*Proof.* Since  $M_1^\gamma$  is measurable on  $\mathcal{F}^{(0)} \times \mathcal{B}$  and  $\mathcal{F}^{(0)} = \{\emptyset, U\}$ , this function may be considered as a measurable map  $\bar{M}_1^\gamma$  of  $(\Omega, \mathcal{B})$  into  $(U, \mathcal{F}^{(1)})$ . For  $1 < k \leq n$  define, recursively,

$$\bar{M}_k^\gamma : (\Omega, \mathcal{B}) \rightarrow (U, \mathcal{F}^{(k)})$$

by

$$(13) \quad \bar{M}_k^\gamma(\omega) = M_k^\gamma(\bar{M}_{k-1}^\gamma(\omega), \omega).$$

Then  $\bar{M}_n^\gamma$  is the desired solution map. Indeed, the procedure (13) establishes in succession and leaves thereafter invariant the following equations:

$$(14) \quad \begin{aligned} \alpha_1 &= (\varphi(u))_1 = \alpha^*, \\ u_{\alpha_1} &= \gamma_{\alpha_1}(u, \omega), \\ &\dots \\ \alpha_n &= (\varphi(u))_n, \\ u_{\alpha_n} &= \gamma_{\alpha_n}(u, \omega). \end{aligned}$$

Therefore, the closed loop equations (7) will be satisfied.

The solution is unique for each  $\omega$  by virtue of Lemma 6.

Finally, by auxiliary Lemma 3, the composition of measurable functions yields measurable functions

$$\bar{M}_k^\gamma : (\Omega, \mathcal{B}) \rightarrow (U, \mathcal{F}^{(k)}).$$

For  $k = n$  one has  $\mathcal{F}^{(k)} = \mathcal{F}$  (Lemma 3), so that the solution map is measurable as claimed.

**8. Noncausal information structures.** In exploring the possible validity of the implication  $S \Rightarrow C$ , only sections  $\gamma_\alpha(\cdot, \omega)$  for fixed  $\omega \in \Omega$  are involved. Therefore the  $\omega$ -dependence of  $\gamma$  will be suppressed.

The precise situation is described by the following theorem.

**THEOREM 2.** *Property S implies property C when  $n = 1$  or  $2$ . The implication is false for  $n > 2$ .*

*Proof.* (i) Suppose  $n = 1$ , and property S holds. Then  $A = \{\alpha\}$ ,  $U_\alpha = U$ ,  $\mathcal{I}_\alpha \subset \mathcal{F}_\alpha = \mathcal{F}$ . If a set  $E$  in  $\mathcal{I}_\alpha$  is nonempty and has nonempty complement  $E^c$ , then choose  $u_0 \in E$ ,  $\bar{u}_0 \in E^c$  and  $\gamma_\alpha(u) = u_0$  for  $u \in E$ ,  $\bar{u}_0$  otherwise. Then  $\gamma_\alpha$  is measurable on  $\mathcal{I}_\alpha$  and the closed loop equation  $u = u_\alpha = \gamma_\alpha(u)$  has two distinct solutions  $u_0$  and  $\bar{u}_0$ . Since this contradicts property S no such set  $E$  can exist, which means that  $\mathcal{I}_\alpha = \{\emptyset, U\}$ . Then property C holds as claimed. By interchanging  $u_0, \bar{u}_0$  in the definition of  $\gamma_\alpha$  a contradiction to existence, instead of to uniqueness, is obtained. The two types of contradictions correspond to the two forms of the classical liars paradox.

(ii) Suppose  $n = 2$  and S holds. Then  $A = \{\alpha, \beta\}$ ,  $U = U_\alpha \times U_\beta$ ,  $\mathcal{F} = \mathcal{F}_\alpha \times \mathcal{F}_\beta$ ,  $\mathcal{I}_\alpha \subset \mathcal{F}$ ,  $\mathcal{I}_\beta \subset \mathcal{F}$ .

If for some fixed  $u_\beta^* \in U_\beta$  one had a set  $E$  in  $\mathcal{I}_\alpha$  such that both  $E$  and its complement  $E^c$  have a nonempty intersection with the set  $\{u|P_\beta(u) = u_\beta^*\}$ , then a contradiction to S is obtained by taking  $\gamma_\beta(u) = u_\beta^*$  and constructing  $\gamma_\alpha$  as in (i). Hence all sets  $E$  in  $\mathcal{I}_\alpha$  are cylindrical, in the sense that  $(u_\alpha, u_\beta) \in E$  implies  $(\bar{u}_\alpha, u_\beta) \in E$  for all  $\bar{u}_\alpha$  in  $U_\alpha$ . This means that  $\mathcal{I}_\alpha \subset \mathcal{F}(\{\beta\})$ . Likewise  $\mathcal{I}_\beta \subset \mathcal{F}(\{\alpha\})$ . Therefore property C will hold if and only if at least one of  $\mathcal{I}_\alpha, \mathcal{I}_\beta$  is the trivial field  $\{\emptyset, U\}$ . Suppose this is not the case. Then there are sets  $E \in \mathcal{F}_\alpha, F \in \mathcal{F}_\beta$  such that  $U_\alpha = E + E^c, U_\beta = F + F^c$  are proper partitions and  $E \times U_\beta \in \mathcal{I}_\beta, F \times U_\alpha \in \mathcal{I}_\alpha$ . Select  $u_\alpha^0 \in E, \bar{u}_\alpha^0 \in E^c, u_\beta^0 \in F, \bar{u}_\beta^0 \in F^c$ . Define  $\gamma_\alpha(u) = u_\alpha^0$  for  $u \in F \times U_\alpha, \bar{u}_\alpha^0$  otherwise,  $\gamma_\beta(u) = u_\beta^0$  for  $u \in E \times U_\beta, \bar{u}_\beta^0$  otherwise. These functions are measurable on the respective information fields and the closed loop equations admit the two distinct solutions  $u^0 = (u_\alpha^0, u_\beta^0)$  and  $\bar{u}^0 = (\bar{u}_\alpha^0, \bar{u}_\beta^0)$ . This contradiction proves the claim.

(iii) That the implication does not hold for  $n = 3$  (and, a fortiori, for larger  $n$ ) is shown by the following counterexample which is adapted from a counterexample of R. L. Graham [7] to the corresponding conjecture in combinatorial set theory.

With  $A = \{\alpha, \beta, \gamma\}$  let  $U_\alpha = U_\beta = U_\gamma = \{\text{true, false}\}$  so that  $\text{card } U = 8$  and  $\mathcal{F} = 2^U$ . The notation  $x\bar{y}$  will denote the Boolean expression “ $x$  and (not  $y$ ).” Let  $\mathcal{I}_\alpha$  be the field generated by  $u_\beta\bar{u}_\gamma$ . It consists of  $\emptyset, U, \{u|u_\beta\bar{u}_\gamma = \text{true}\}$ , and  $\{u|u_\beta\bar{u}_\gamma = \text{false}\}$ . Likewise let  $\mathcal{I}_\beta$  be generated by  $u, \bar{u}_\alpha$  and  $\mathcal{I}_\gamma$  by  $u_\alpha\bar{u}_\beta$ . Then each agent has a choice of four possible decision functions giving a total of 64 possible combinations. It is easily verified that for each of these combinations the closed loop equations are satisfied by exactly one of the eight elements of  $U$ . Since none of the information fields is trivial, property C does not hold while property S does.

Thus the normal form of a game with such a structure is perfectly well-defined. The strategies determine the outcome, but the extensive form of the game cannot be played without some form of precognition.

Since the counterexample used in Theorem 2, part (iii), is nonrandom it trivially satisfies property SM. This shows that property SM does not imply property C.

The last question is a purely technical one: when C is not assumed, does property S imply SM? In most special cases of interest the answer is yes, but in the generality of the present paper this assertion appears to be false, though a counterexample is not on paper.

**9. Conclusions.** The facts brought to light above are contributions towards a systematic information structure theory. They are entirely compatible with the more usual state equation–output equation formulation, as long as it is realized that (i) the classical information structure is a very special one, and (ii) output equations matter only by virtue of the fields they determine: different outputs determining the same fields are equivalent.

The information fields  $\mathcal{I}_\alpha$  belong to the lattice of all subfields of  $\mathcal{F}$ . When games differ only as to information, relations between values and information can be derived accordingly.

Various important notions can easily be made precise within the given framework. For example, define a *station* as a set of agents, belonging to the same organization, such that the set of information fields of these agents is totally

ordered under inclusion (i.e., nested). A player (or coalition) enjoys *perfect recall* when his whole organization forms a single station. Then it is clear that the derivation of Aumann's theorem on behavior strategies in [6] holds for games satisfying property C when the player under consideration enjoys perfect recall as just defined.

A physical interpretation of the noncausal information structures with property SM, which exist by Theorem 2, appears rather difficult [12].

**Appendix.** A measurable space  $(X, \mathcal{B})$  consists of a set  $X$  and a  $\sigma$ -field  $\mathcal{B}$  of subsets of  $X$ . A function  $f$  from measurable space  $(X, \mathcal{B})$  into measurable space  $(Y, \mathcal{C})$  is called measurable when the inverse image under  $f$  of any set in  $\mathcal{C}$  belongs to  $\mathcal{B}$ .

The following facts are simple exercises in Boolean algebra.

**AUXILIARY LEMMA 1.** *Let  $S$  be a nonempty set equal to the disjoint union of the sets  $\{E_i | i \in I\}$ , where  $I$  is a nonempty index set.<sup>1</sup> For each  $i$  in  $I$  let  $\mathcal{F}_i$  be a  $\sigma$ -field on  $S$  (more generally an  $\mathcal{S}$ -class in the sense of Loève [8, p. 59]). Then the collection of sets*

$$\mathcal{C} = \left\{ \bigcup_{i \in I} (E_i \cap F_i) \mid \text{for all } i \in I, F_i \in \mathcal{F}_i \right\}$$

*is a  $\sigma$ -field (respectively,  $\mathcal{S}$ -class) on  $S$ .*

**AUXILIARY LEMMA 2.** *Let  $f$  be a function from set  $X$  into measurable space  $(Y, \mathcal{C})$  and let  $\mathcal{B} = f^{-1}(\mathcal{C})$  be the  $\sigma$ -field generated on  $X$  by  $f$ . Assume that for some  $x_1, x_2$  in  $X$ , one has  $f(x_1) = f(x_2)$ . Then there is no set in  $\mathcal{B}$  that contains  $x_1$  but not  $x_2$ .*

**AUXILIARY LEMMA 3.** *Let  $(X, \mathcal{B})$ ,  $(Y_1, \mathcal{C}_1)$ ,  $(Y_2, \mathcal{C}_2)$ ,  $(Z, \mathcal{D})$  be measurable spaces and let  $(Y_1 \times Y_2, \mathcal{C}_1 \times \mathcal{C}_2)$  denote the product space with the product field, the smallest field such that projections are measurable. Suppose the functions  $f_1: (X, \mathcal{B}) \rightarrow (Y_1, \mathcal{C}_1)$ ,  $f_2: (X, \mathcal{B}) \rightarrow (Y_2, \mathcal{C}_2)$  and  $g: (Y_1 \times Y_2, \mathcal{C}_1 \times \mathcal{C}_2) \rightarrow (Z, \mathcal{D})$  are measurable. Then the function  $h: (X, \mathcal{B}) \rightarrow (Z, \mathcal{D})$  defined by  $h(x) = g(f_1(x), f_2(x))$  is measurable.*

**Acknowledgment.** This paper was based on an invited contribution to the first International Conference on the Theory and Application of Differential Games, Amherst, Massachusetts, September 29, 1969.

#### REFERENCES

- [1] J. VON NEUMANN AND O. MORGENSTERN, *The Theory of Games and Economic Behavior*, Princeton Univ. Press, Princeton, N.J., 1944, Chap. II.
- [2] H. W. KUHN, *Extensive games and the problem of information*, Contributions to the Theory of Games, vol. 2, Annals of Mathematics Studies, no. 28, Princeton Univ. Press, Princeton, N.J., 1953, pp. 193–216.
- [3] N. DALKEY, *Equivalence of information patterns and essentially determinate games*, *Ibid.*, pp. 217–243.
- [4] G. L. THOMPSON, *Signaling strategies in  $n$ -person games*, *Ibid.*, pp. 267–277.
- [5] B. J. BIRCH, *On games with almost complete information*, Proc. Cambridge Philos. Soc., 51 (1955), pp. 275–287.

<sup>1</sup> Some of the  $E_i$  may be empty.

- [6] R. J. AUMANN, *Mixed and behavior strategies in infinite extensive games*, Advances in Game Theory, Princeton Univ. Press, Princeton, N.J., 1964, pp. 627–650.
- [7] R. L. GRAHAM, Private communication, 1969.
- [8] M. LOEVE, *Probability Theory*, 3rd ed., Van Nostrand, Princeton, N.J., 1963.
- [9] R. RADNER, *Team decision problems*, Ann. Math. Statist., 33 (1962), pp. 857–881.
- [10] W. H. CLINGMAN, *Permanently optimal strategies in extensive games*, Proc. Amer. Math. Soc., 17 (1966), pp. 156–161.
- [11] E. MARCHI, *A note on temporal games*, Z. Wahrscheinlichkeitstheorie Verw. Gebiete, 12 (1969), pp. 98–106.
- [12] R. G. NEWTON, *Particles that travel faster than light?*, Science, 167 (1970), pp. 1569–1574.

## NECESSARY CONDITIONS FOR JOINING OPTIMAL SINGULAR AND NONSINGULAR SUBARCS\*

J. P. MCDANELL† AND W. F. POWERS‡

**Abstract.** Necessary conditions for the optimality of junctions between singular and nonsingular subarcs are developed for singular optimal control problems. Previously known necessary conditions concerning the continuity and smoothness of a piecewise analytic optimal control at a junction are clarified and extended. The main result is that the sum of the order of the singular arc and the lowest order time derivative of the control which is discontinuous at the junction must be an odd integer when the strengthened generalized Legendre-Clebsch condition is satisfied. Also, new necessary conditions which do not require an analyticity assumption are developed. These aid in characterizing problems which may possess nonanalytic junctions.

**1. Introduction.** Optimal control problems in which the control variables appear only linearly admit the possibility of the occurrence of singular extremals. The analysis of such problems is complicated by the fact that the solution, in general, consists of some combination of singular and nonsingular subarcs, the number and sequence of which are not known a priori. If the solution is totally singular, recent results [1], [2] are available to prove optimality in a large number of cases. If the solution is totally nonsingular, it is the familiar bang-bang control generated by a switching function with isolated zeros, as determined by the minimum principle. However, the mathematical characterization of optimal controls which contain both singular and nonsingular subarcs is far from complete.

This paper is concerned with the problem of characterizing the continuity and smoothness properties of the optimal control at a junction between singular and nonsingular subarcs. The analysis was motivated by the preliminary results obtained in this direction by Kelley, Kopp and Moyer [3] and Johnson [4]. We shall comment on their results in a later section.

**2. Problem statement.** The class of problems to which this analysis applies is the following. Determine the scalar control  $u^*(t)$ ,  $t \in [t_0, t_f]$ , which minimizes the functional

$$(2.1) \quad J(u) = G(t_f, x(t_f)) + \int_{t_0}^{t_f} [L_0(t, x) + L_u(t, x)u] dt,$$

where the system equation is

$$(2.2) \quad \dot{x} = f_0(t, x) + f_u(t, x)u$$

subject to the constraints

$$(2.3) \quad |u(t)| \leq K(t), \quad t \in [t_0, t_f],$$

$$(2.4) \quad \{t_0, x(t_0), t_f, x(t_f)\} \in S.$$

---

\* Received by the editors May 26, 1970, and in revised form August 26, 1970. This research was supported by the National Science Foundation under Grant GK-4990 and by the National Aeronautics and Space Administration under Grant NGR-23-005-329.

† Computer, Information, and Control Engineering Program, University of Michigan, Ann Arbor, Michigan 48104.

‡ Department of Aerospace Engineering, University of Michigan, Ann Arbor, Michigan 48104.

Here  $x$  is an  $n$ -vector and  $S$  is a closed subset of  $R^{2n+2}$ . The functions  $f_0, f_u, L_0, L_u$  are assumed to be analytic in both arguments in a suitable domain;  $K(t)$  is assumed to be analytic in a neighborhood of each junction and  $|u(t)| < K(t)$  almost everywhere on the singular subarcs. Of course, the usual case  $|u| \leq K$  with  $K = \text{const.}$  is included as a special case. We restrict attention to a scalar control in order to simplify notation. A similar analysis holds for each component of a vector control.

Clearly, the Hamiltonian for this problem is linear in the control, i.e.,

$$(2.5) \quad H(t, x, \lambda, u) = \lambda^T f_0(t, x) + L_0(t, x) + [\lambda^T f_u(t, x) + L_u(t, x)]u.$$

The multiplier equations are then given by

$$(2.6) \quad \dot{\lambda} = -H_x(t, x, \lambda, u),$$

where  $H_x$  is also linear in  $u$ . The coefficient of  $u$  in (2.5) is called the *switching function*, which we shall designate as  $\phi(t)$ , i.e.,

$$(2.7) \quad \phi(t) \equiv H_u(t, x(t), \lambda(t)).$$

The minimum principle (i.e., Pontryagin's maximum principle in a minimum form) states that for almost every  $t \in [t_0, t_f]$  and each  $u$  satisfying  $|u| \leq K(t)$ , the optimal control  $u^*(t)$  must satisfy

$$(2.8) \quad H(t, x(t), \lambda(t), u^*(t)) \leq H(t, x(t), \lambda(t), u).$$

Therefore, as is well known, on each open subinterval of  $[t_0, t_f]$  there are two distinct possibilities for  $u^*$ . Either

$$(2.9) \quad u^*(t) = -K(t) \text{sgn } \phi(t)$$

or

$$(2.10) \quad \phi(t) \equiv 0.$$

Equations (2.9) and (2.10) define, respectively, the nonsingular and singular subarcs of the optimal control.

The class of problems defined above will be called *singular control problems*, even though only a portion of the total trajectory may be singular.

**3. Notation and definitions.** The following definitions will clarify the terminology used in this paper.

**DEFINITION 1.** A real-valued function  $g$  is said to be *piecewise analytic* on an interval  $(a, b)$  if for each  $t_c \in (a, b)$  there exist  $t_1 \in (a, t_c)$  and  $t_2 \in (t_c, b)$  such that  $g$  is analytic on the open subintervals  $(t_1, t_c)$  and  $(t_c, t_2)$ .

**DEFINITION 2.** A junction between singular and nonsingular subarcs of the control is said to be a *nonanalytic junction* if the control is not piecewise analytic in any neighborhood of the junction.

**DEFINITION 3.** Let  $u$  be an optimal singular control on the interval  $[t_1, t_2]$ , and let  $(d^{2q}/dt^{2q})[H_u(t, x, \lambda)]$  be the lowest order total derivative of  $H_u$  in which  $u$  appears explicitly with a coefficient which is not identically zero on  $[t_1, t_2]$ . Then the integer  $q$  is called the *order of the singular arc*.

Implicit in Definition 3 is the property that  $u$  first appears explicitly in an even order derivative of  $H_u$ ; i.e., it is correct to refer to  $q$  as an integer. For a proof of this property see Robbins [5].

We also need the well-known generalized Legendre-Clebsch necessary condition for optimality of singular subarcs [3].

**THEOREM** (Generalized Legendre-Clebsch condition). *On an optimal singular subarc of order  $q$ , it is necessary that*

$$(3.1) \quad (-1)^q \frac{\partial}{\partial u} \left[ \frac{d^{2q}}{dt^{2q}} H_u \right] \geq 0.$$

Condition (3.1) hereafter will be called the *GLC condition*. By the strengthened GLC condition we mean that strict inequality holds in (3.1).

In this paper it will be convenient to consider the lowest order derivative of a function to be its zeroth derivative, by which we mean the function itself. We shall use the notation

$$g^{(0)} \equiv g, \quad g^{(i)} \equiv d^i g / dt^i, \quad i = 1, 2, \dots$$

Also, where the context makes the meaning clear, we shall use  $u$  to designate the optimal control instead of  $u^*$ .

**4. The junction theorems.** As indicated in § 1, the theory for totally singular and totally nonsingular optimal controls is rather well developed. The main difficulty with singular control problems occurs when both singular and nonsingular subarcs are present. Since a useful sufficient condition for such problems is not available, one is led naturally to the study of necessary conditions which are valid in the neighborhood of a junction between singular and nonsingular subarcs. It is expected that such conditions can be used to eliminate candidate extremals and/or predict beforehand the way in which singular and nonsingular subarcs must be joined, e.g., whether the optimal control is continuous or discontinuous at a junction.

If the optimal control is well-behaved in a neighborhood of a junction, then the following property must hold.

**THEOREM 1.** *Let  $t_c$  be a point at which singular and nonsingular subarcs of an optimal control  $u$  are joined, and let  $q$  be the order of the singular arc. Suppose the strengthened GLC condition is satisfied at  $t_c$ , i.e.,  $(-1)^q (\partial/\partial u) H_u^{(2q)} > 0$ , and assume that the control is piecewise analytic in a neighborhood of  $t_c$ . Let  $u^{(r)}$  ( $r \geq 0$ ) be the lowest order derivative of  $u$  which is discontinuous at  $t_c$ . Then  $q + r$  is an odd integer.*

*Proof.* Since  $H_u^{(2q)}$  is the lowest order time derivative of  $H_u$  which contains  $u$  explicitly, from (2.2), (2.5) and (2.6) we see that it must have the form

$$(4.1) \quad H_u^{(2q)}(t, x, \lambda, u) \equiv A(t, x, \lambda) + B(t, x, \lambda)u.$$

Define the functions  $\alpha$  and  $\beta$  as follows :

$$(4.2) \quad \alpha(t) \equiv A(t, x(t), \lambda(t)),$$

$$(4.3) \quad \beta(t) \equiv B(t, x(t), \lambda(t)).$$

From the hypotheses it is clear that  $\alpha$  and  $\beta$  are continuous and have at least  $r$  continuous derivatives at  $t_c$ . The switching function  $\phi$  as defined by (2.7) has exactly  $2q + r - 1$  continuous derivatives at  $t_c$ .

Let  $\varepsilon$  be a nonzero real number of arbitrarily small magnitude such that  $t_c + \varepsilon$  is a point on the nonsingular side of  $t_c$  and  $t_c - \varepsilon$  is a point on the singular

side. Let  $u_n$  and  $u_s$  designate the control  $u$  on the nonsingular and singular sides of  $t_c$ , respectively. By  $u_n^{(i)}(t_c)$  and  $u_s^{(i)}(t_c)$  we mean the limit as  $\varepsilon \rightarrow 0$  of  $u^{(i)}(t_c + \varepsilon)$  and  $u^{(i)}(t_c - \varepsilon)$ , respectively.

We wish to expand  $\phi(t_c + \varepsilon)$  in a Taylor series about  $t_c$ . Let  $k = 2q + r$ . Then  $\phi^{(k)}$  will be the lowest order derivative of the switching function  $\phi$  which is discontinuous at  $t_c$ , and since  $\phi \equiv 0$  on the singular side of  $t_c$ , the first nonzero term of the Taylor series will be the term containing  $\phi^{(k)}$ . Noting that

$$(4.4) \quad \phi^{(k)} \equiv \frac{d^r}{dt^r}[\alpha + \beta u]$$

we can write

$$(4.5) \quad \phi(t_c + \varepsilon) = \frac{\varepsilon^k}{k!} \left[ \alpha^{(r)}(t_c) + \sum_{i=0}^r \binom{r}{i} \beta^{(r-i)}(t_c) u_n^{(i)}(t_c) \right] + o(\varepsilon^k),$$

where Leibniz's formula [6, p. 184] for differentiation of a product has been used to differentiate  $\beta u_n$ .

On the singular arc,

$$(4.6) \quad \phi^{(2q)} = \alpha + \beta u_s \equiv 0.$$

Therefore,  $\alpha \equiv -\beta u_s$ , and

$$(4.7) \quad \alpha^{(r)} = \frac{d^r}{dt^r}[-\beta u_s] = -\sum_{i=0}^r \binom{r}{i} \beta^{(r-i)} u_s^{(i)}.$$

Substituting from (4.7) into (4.5), we have

$$(4.8) \quad \phi(t_c + \varepsilon) = \frac{\varepsilon^k}{k!} \sum_{i=0}^r \binom{r}{i} \beta^{(r-i)}(t_c) [u_n^{(i)}(t_c) - u_s^{(i)}(t_c)] + o(\varepsilon^k).$$

If  $r > 0$ ,

$$(4.9) \quad u_n^{(i)}(t_c) = u_s^{(i)}(t_c), \quad i = 0, \dots, r - 1.$$

Therefore, (4.8) becomes

$$(4.10) \quad \phi(t_c + \varepsilon) = \frac{\varepsilon^k}{k!} \beta(t_c) [u_n^{(r)}(t_c) - u_s^{(r)}(t_c)] + o(\varepsilon^k).$$

Let  $\sigma = -\text{sgn } \phi(t_c + \varepsilon)$  so that  $u_n(t) = \sigma K(t)$ . Then recalling that  $u_n^{(i)}(t_c) \equiv \lim_{\varepsilon \rightarrow 0} u_n^{(i)}(t_c + \varepsilon)$  we have

$$(4.11) \quad u_n^{(i)}(t_c) \equiv \sigma K^{(i)}(t_c), \quad i = 0, \dots, r.$$

Now consider the following series expansion on the singular arc:

$$(4.12) \quad \sigma K(t_c - \varepsilon) - u(t_c - \varepsilon) = \sum_{i=0}^r \frac{(-\varepsilon)^i}{i!} [\sigma K^{(i)}(t_c) - u_s^{(i)}(t_c)] + o(\varepsilon^r).$$

The right-hand side of (4.12) can be simplified using (4.9) and (4.11) to obtain

$$(4.13) \quad \sigma K(t_c - \varepsilon) - u(t_c - \varepsilon) = \frac{(-1)^r \varepsilon^r}{r!} [u_n^{(r)}(t_c) - u_s^{(r)}(t_c)] + o(\varepsilon^r).$$



Substituting from (4.13) into (4.10) and recalling that  $k = 2q + r$ , we obtain

$$(4.14) \quad \phi(t_c + \varepsilon) = (-1)^r \frac{\varepsilon^{2q} r!}{k!} \beta(t_c) [\sigma K(t_c - \varepsilon) - u(t_c - \varepsilon)] + o(\varepsilon^k).$$

From the application of the minimum principle on the nonsingular subarc (see (2.9)) we have  $\sigma = 1$  if  $\phi(t_c + \varepsilon) < 0$  and  $\sigma = -1$  if  $\phi(t_c + \varepsilon) > 0$ . Therefore, the following inequality must hold:

$$(4.15) \quad (-1)^r \varepsilon^{2q} \beta(t_c) [K(t_c - \varepsilon) \pm u(t_c - \varepsilon)] < 0.$$

From the GLC condition we have

$$(4.16) \quad (-1)^q \beta(t_c) > 0.$$

Multiplying the left-hand side of (4.15) by the positive quantity in (4.16) we obtain

$$(4.17) \quad (-1)^{q+r} \varepsilon^{2q} \beta^2(t_c) [K(t_c - \varepsilon) \pm u(t_c - \varepsilon)] < 0.$$

Since  $|u(t)| \leq K(t)$  for all  $t \in [t_0, t_f]$ , and the singular arc is assumed to be interior almost everywhere, the bracketed quantity in (4.17) is strictly positive, regardless of the choice of sign on  $\pm u(t_c - \varepsilon)$ . Also  $\varepsilon^{2q} > 0$  regardless of the sign of  $\varepsilon$ . Therefore, condition (4.17) reduces to

$$(4.18) \quad (-1)^{q+r} < 0$$

from which it is clear that  $q + r$  is an odd integer. This completes the proof.

Theorem 1 implies the following important corollaries.

**COROLLARY 1.** *In  $q$ -even problems, assuming  $u$  is piecewise analytic and the strengthened GLC condition is satisfied, the optimal control is continuous at each junction.*

**COROLLARY 2.** *In  $q$ -odd problems, assuming  $u$  is piecewise analytic and the strengthened GLC condition is satisfied, the optimal control either has a jump discontinuity at each junction, or else the singular control joins the boundary smoothly, i.e., with continuous first derivative.*

In the corollaries above, especially Corollary 1, which applies to the  $q$ -even case, the assumption that  $u$  is piecewise analytic is not to be taken lightly. In fact, the authors have not seen or been able to produce a  $q$ -even example with a continuous junction, i.e., the junction is usually nonanalytic if  $q$  is even.

The conclusions reached by Kelley, Kopp and Moyer [3] are consistent with those stated in Corollaries 1 and 2, with one important exception—they ruled out the possibility of a continuous junction for  $q$ -odd problems. This erroneous conclusion resulted from the claim that continuity of  $u$  implies  $(\partial/\partial u)H_u^{(2q)} > 0$  which is not true in general, as can be seen from (4.15) (in which  $\beta \equiv (\partial/\partial u)H_u^{(2q)}$ ). That such a junction is realizable will be demonstrated by means of a simple example in a later section.

Theorem 1 requires that the *strengthened* GLC condition be satisfied at the junction point. While this is the usual and most important case, the possibility exists that the GLC condition is satisfied with equality. To treat this case, note from Definition 3 that for a  $q$ th order singular arc the GLC expression  $(\partial/\partial u)H_u^{(2q)}$  (i.e.,  $\beta$ ) cannot be identically zero on the singular subarc. Therefore, in view of our

analyticity assumptions, a derivative of some order must be nonzero at the junction point  $t_c$  even if  $\beta(t_c) = 0$ . This leads to the following theorem, which is a generalization of Theorem 1, but is stated separately to avoid obscuring the result for the important case covered by Theorem 1.

**THEOREM 2.** *Let  $t_c$  be a point at which singular and nonsingular subarcs of an optimal control  $u$  are joined, and let  $q$  be the order of the singular arc. Assume that the control is piecewise analytic in a neighborhood of  $t_c$ . Let  $u^{(r)}$  ( $r \geq 0$ ) be the lowest order derivative of  $u$  which is discontinuous at  $t_c$ , and let  $\beta^{(m)}$  ( $m \geq 0$ ) be the lowest order derivative of the GLC expression  $(\partial/\partial u)H_u^{(2q)} \equiv \beta$  which is nonzero at  $t_c$ . Then, (i) if  $m \leq r$ ,  $q + r + m$  is an odd integer; (ii) if  $m > r$ ,  $-\text{sgn} [\beta^{(m)}(t_c^+) \beta^{(m)}(t_c^-)] = (-1)^{q+r+m}$ .*

*Proof outline.* The proof is similar to that for Theorem 1; however, in order to obtain a nontrivial term in the Taylor series expansion for  $\phi(t_c + \varepsilon)$ , one must consider higher order terms with the result that (4.15) is replaced by

$$(4.19) \quad (-1)^r \varepsilon^{2q+m} \beta_n^{(m)}(t_c) [K(t_c - \varepsilon) \pm u_s(t_c - \varepsilon)] < 0.$$

A Taylor series expansion for  $\beta(t)$  on the singular arc yields

$$(4.20) \quad \beta(t_c - \varepsilon) = \frac{(-\varepsilon)^m}{m!} \beta_s^{(m)}(t_c) + o(\varepsilon^m),$$

where the subscripts  $s$  and  $n$  on  $\beta^{(m)}(t_c)$  indicate the limit at  $t_c$  on the singular and nonsingular sides, respectively. Since  $\beta_s^{(m)}(t_c) \neq 0$ , from the GLC condition and (4.20) we have

$$(4.21) \quad (-1)^{q+m} \varepsilon^m \beta_s^{(m)}(t_c) > 0.$$

From (4.19) and (4.21) it follows that

$$(4.22) \quad (-1)^{q+r+m} \beta^{(m)}(t_c^+) \beta^{(m)}(t_c^-) < 0.$$

If  $m \leq r$ ,  $\beta^{(m)}$  is continuous at  $t_c$ , in which case (4.22) implies that  $q + r + m$  is an odd integer. If  $m > r$ ,  $\beta^{(m)}$  may not be continuous at  $t_c$ , and the conclusion of Theorem 2 for this case follows.

The main restriction in Theorems 1 and 2 is the assumption that the control is piecewise analytic in a neighborhood of the junction. This hypothesis is usually satisfied on the singular subarc, but not always on the nonsingular subarc. Thus, we are led to consider properties which do not require the assumption of analyticity, as stated in the following theorem. The functions  $A$  and  $B$  in this theorem are those defined by the identity (4.1).

**THEOREM 3.** *Let  $u$  be an optimal control which contains both nonsingular subarcs and piecewise continuous,  $q$ -th order singular subarcs.*

- (i) *If  $H_u^{(2q)} \neq 0$  on the nonsingular side of a junction, then the control is discontinuous.*
- (ii) *If  $A = 0$ ,  $B \neq 0$ , and  $K \neq 0$  at a junction, then the control is discontinuous.*
- (iii) *If  $u$  is piecewise continuous on the nonsingular subarc,  $H_u^{(2q)} = 0$  on the nonsingular side of a junction, and  $B \neq 0$  at the junction, then the control is continuous.*

*Proof.* Using the same notation as in the proof of Theorem 4.1, and recalling that  $H_u^{(2q)} \equiv 0$  on the singular subarc, we have for case (i),

$$(4.23) \quad \alpha(t_c) \pm \beta(t_c)K(t_c) \neq 0 = \alpha(t_c) + \beta(t_c)u_s(t_c)$$

from which we obtain  $|u_s(t_c)| \neq K(t_c)$ . Therefore,  $u$  is discontinuous.

For case (ii),  $\alpha(t_c) = 0$  and  $\beta(t_c) \neq 0$  imply  $u_s(t_c) = 0$ , and since  $K(t_c) \neq 0$ , the control is discontinuous.

For case (iii) we have

$$(4.24) \quad \alpha(t_c) + \beta(t_c)u_n(t_c) = 0 = \alpha(t_c) + \beta(t_c)u_s(t_c).$$

Since  $\beta(t_c) \neq 0$ , we must have  $u_n(t_c) = u_s(t_c)$ .

Case (ii) of Theorem 4.3 may appear to be a rather special case, but it occurs frequently enough to be of interest. Note that this result is independent of an even or odd assumption. Because of this, we can couple (ii) with the previous result for  $q$ -even problems to obtain the following interesting property.

**COROLLARY 3.** *If  $q$  is even,  $A(t, x, \lambda) \equiv 0$ ,  $K(t_c) \neq 0$ , and  $B(t_c, x(t_c), \lambda(t_c)) \neq 0$ , where  $t_c$  is a junction point between optimal singular and nonsingular subarcs, then the junction is nonanalytic.*

*Proof.* Assume the contrary, i.e., that the optimal control is piecewise analytic in a neighborhood of  $t_c$ . Then by Corollary 1 the control is continuous at  $t_c$ , but by Theorem 3 (ii) the control is discontinuous, which supplies the necessary contradiction.

In the next section this corollary will be used to predict the nonanalytic junction in the well-known Fuller problem.

**5. Example of a nonanalytic junction.** Consider the Fuller problem [7], which is to minimize

$$(5.1) \quad J = \frac{1}{2} \int_0^T x_1^2 dt$$

subject to

$$(5.2) \quad \begin{aligned} \dot{x}_1 &= x_2, & x_1(0) &= \xi_1 \neq 0, \\ \dot{x}_2 &= u, & x_2(0) &= \xi_2, \end{aligned}$$

$$(5.3) \quad |u| \leq 1,$$

where  $T$  is fixed. The Hamiltonian, the multiplier equations and the switching function are given by

$$(5.4) \quad H = \lambda_1 x_2 + \lambda_2 u + \frac{1}{2} x_1^2,$$

$$(5.5) \quad \dot{\lambda}_1 = -x_1, \quad \lambda_1(T) = 0,$$

$$\dot{\lambda}_2 = -\lambda_1, \quad \lambda_2(T) = 0,$$

$$(5.6) \quad \phi \equiv H_u = \lambda_2.$$

The lowest order derivative of  $H_u$  which contains  $u$  explicitly is

$$(5.7) \quad H_u^{(4)} = u$$

from which we see that the order of the singular arc is even, namely  $q = 2$ . Also, the strengthened GLC condition holds, and  $A(t, x, \lambda) \equiv 0$ . Thus, we have precisely the conditions of Corollary 3, indicating that any junctions which occur must be nonanalytic junctions.

This problem has been studied thoroughly by Fuller [7] and Johansen [8], and the result is well known. The singular arc is given by

$$(5.8) \quad u_s = x_1 = x_2 = 0.$$

Since  $\xi_1 \neq 0$ , the initial control must be nonsingular. The nonsingular arc is characterized by the nonlinear differential equation

$$(5.9) \quad \phi^{(4)} = -\operatorname{sgn} \phi.$$

The solution of (5.9) yields a switching function with an infinite number of zeros such that the ratio of the lengths of successive intervals between zeros is a constant. If  $T$  is sufficiently large, the resulting nonsingular (bang-bang) control drives the state to the origin in a finite time  $t_c$ , with an infinite number of switches occurring in a neighborhood of  $t_c$ , at which point the optimal control becomes singular. The control is clearly discontinuous at the junction point  $t_c$ , as it must be according to (iii) of Theorem 3. Even though  $q$  is even, Corollary 1 is not violated because the control is not piecewise analytic in a neighborhood of the junction.

The predicted behavior at the junction is useful knowledge for numerical computational schemes; e.g., Jacobson [9] was able to successfully compute bang-bang solutions for this problem with  $T$  sufficiently small so that the singular arc did not occur. However, for large  $T$  the nonanalytic junction came into play. After computing about ten switches, the method became unstable [10].

Note that the optimal control for this "innocent looking," second order example is measurable but not piecewise continuous. Aside from its physical applicability, the existence of such examples is useful for motivating the assumption of measurable controls in the proof of the minimum principle.

**6. Example of a smooth junction with  $q$  odd.** This example demonstrates not only the realizability of a smooth junction with  $q$  odd, but also the dependence of junction phenomena upon boundary conditions. For this case we consider the performance index

$$(6.1) \quad J = \frac{1}{2} \int_0^T (x_2^2 - x_1^2) dt,$$

where  $T = 2.985$ . The equations of motion and constraints are given by

$$(6.2) \quad \begin{aligned} \dot{x}_1 &= x_2, & x_1(0) &= 0, & x_1(T) &= \sigma_1, \\ \dot{x}_2 &= u, & x_2(0) &= 1, & x_2(T) &= \sigma_2, \end{aligned}$$

$$(6.3) \quad |u| \leq 1.$$

The Hamiltonian and the multiplier equations are given by

$$(6.4) \quad H = \lambda_1 x_2 + \lambda_2 u + \frac{1}{2} x_2^2 - \frac{1}{2} x_1^2,$$

$$(6.5) \quad \dot{\lambda}_1 = x_1, \quad \dot{\lambda}_2 = -\lambda_1 - x_2.$$

The switching function and its second derivative are

$$(6.6) \quad \phi \equiv H_u = \lambda_2,$$

$$(6.7) \quad \ddot{H}_u = -x_1 - u$$

so the singular arc is first order; i.e.,  $q$  is odd. From (6.7) we see that the strengthened GLC condition is satisfied. Setting the right-hand side of (6.7) equal to zero and substituting in the equations of motion (6.2), we can readily verify that a singular arc emanating from the initial state  $(0, 1)$  is given by

$$(6.8) \quad u_s = -\sin t, \quad x_1 = \sin t, \quad x_2 = \cos t.$$

If the terminal state  $(\sigma_1, \sigma_2)$  is chosen to lie on the trajectory (6.8), the solution is totally singular, as can be shown by the sufficient conditions in [1] and [2]. However, in this paper we are concerned with junctions. Consider the case where  $\sigma_1 = 0, \sigma_2 = -\sqrt{2}$ . For this case we propose as a candidate for the optimal control

$$(6.9) \quad u = \begin{cases} -\sin t, & t \in [0, \pi/2), \\ -1, & t \in [\pi/2, T]. \end{cases}$$

This control is admissible and satisfies all the necessary conditions for optimality, including that of Theorem 1. There is a junction at  $t_c = \pi/2$ . The control and its first derivative are continuous at  $t_c$ , but the second derivative is discontinuous, so we have  $r = 2, q = 1$ , and  $q + r$  is an odd integer.

The authors are unaware of any workable sufficient conditions in the literature which are applicable to this particular type of problem, i.e., nonconvex and containing both singular and nonsingular subarcs. Consequently, we employed a gradient-type numerical method to justify that the candidate control is indeed optimal, within the bounds of a numerical justification. The modified conjugate gradient method of Pagurek and Woodside [11] was used with penalty functions to enforce the terminal constraints. The result is shown in Fig. 1. The control is continuous and smooth at the junction as expected.

It is apparent that the fortuitous occurrence of this smooth junction is a direct result of our judicious choice of the terminal boundary conditions. In fact, to generate this phenomenon, the form of the candidate control was first selected on the basis of intuition; then a convenient point on the resulting trajectory was selected as the fixed terminal state, and finally the corresponding time was taken to be the explicit final time.

By changing the terminal state, we were able to generate discontinuous controls, which undoubtedly are the usual case. These are shown in Figs. 2 and 3. For these cases  $r = 0$ , and the condition of Theorem 1 is satisfied again. To further emphasize the special character of the smooth junction, the phase plane trajectories for the controls in Figs. 1–2 are given in Fig. 4.

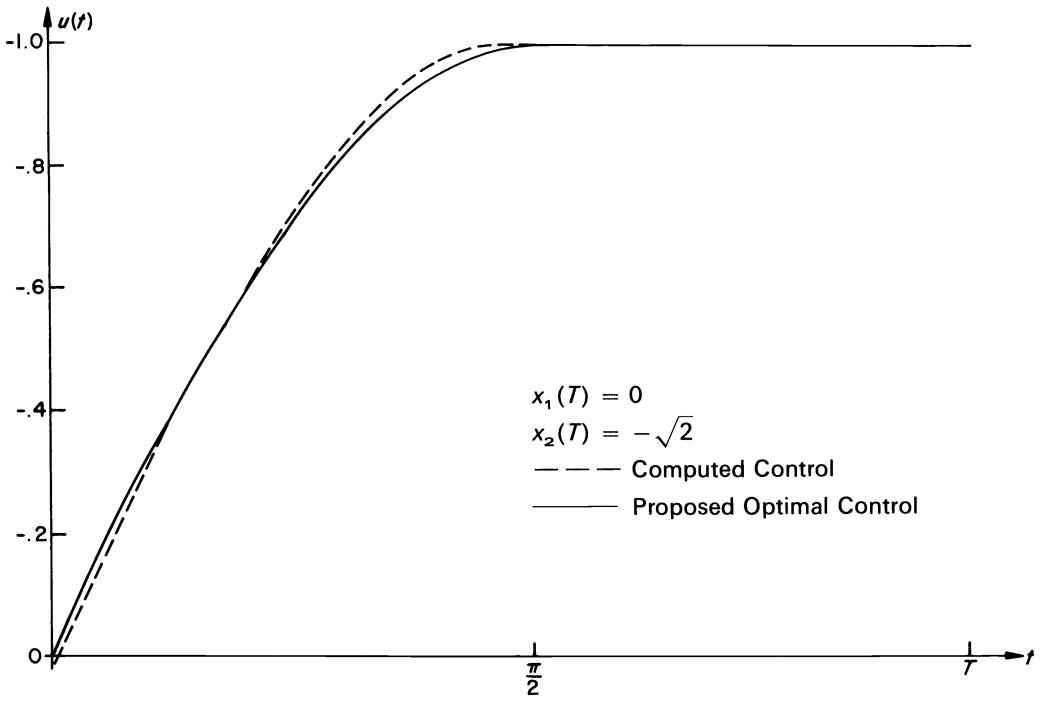


FIG. 1. Example of an optimal control which is continuous at a singular-to-nonsingular junction with  $q$  odd

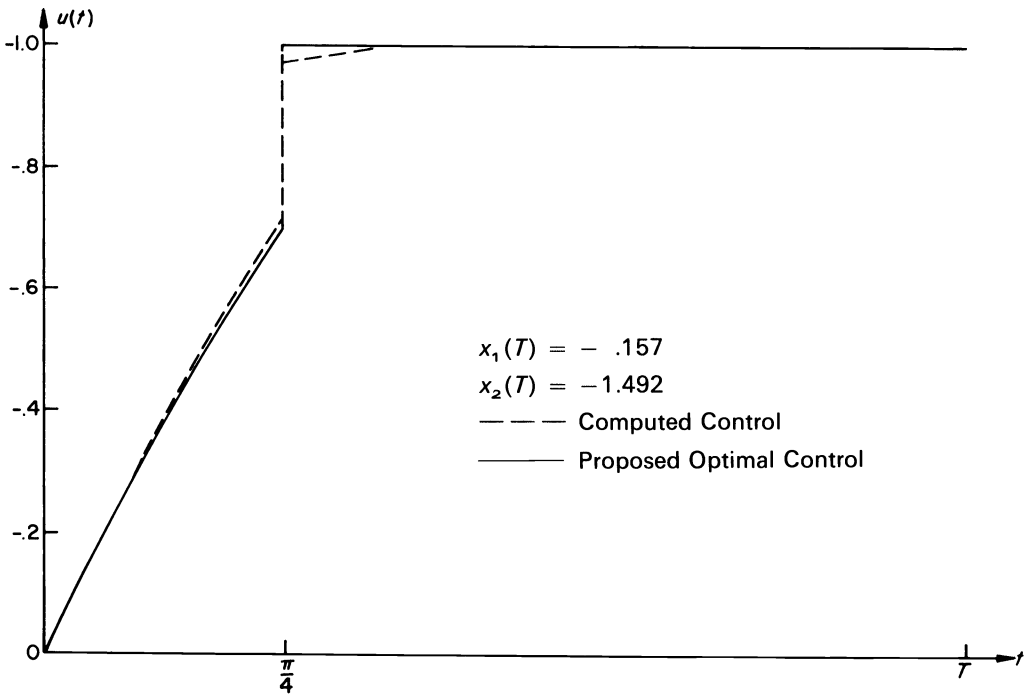


FIG. 2. A typical discontinuous junction

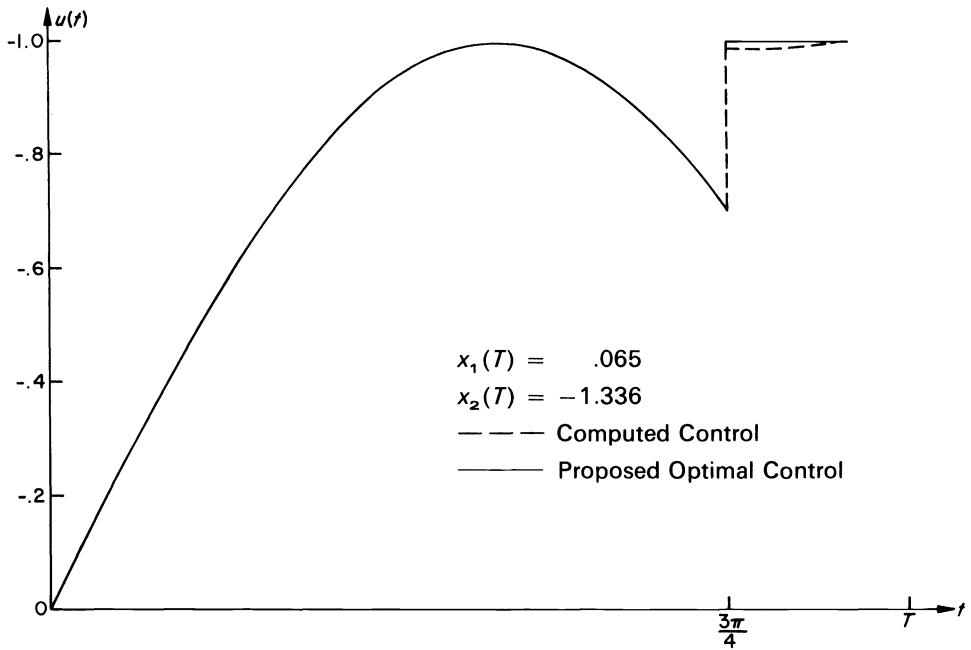


FIG. 3. A predominantly singular control with discontinuous junction to a short nonsingular subarc

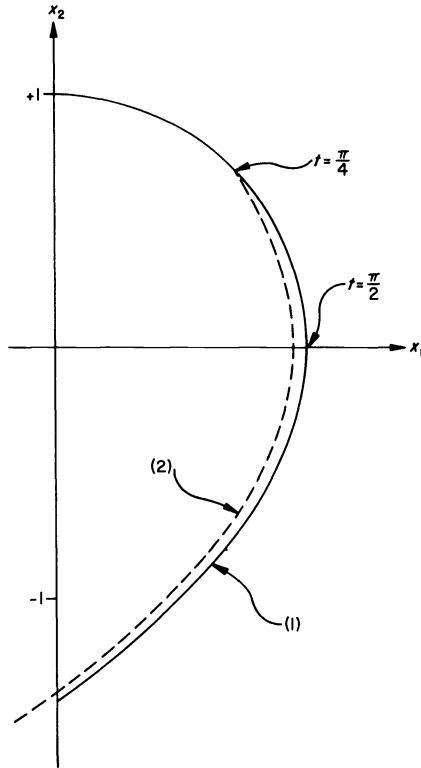


FIG. 4. Phase plane trajectories for the controls in Figs. 1 and 2

**7. Conclusion.** Necessary conditions for the optimality of junctions between singular and nonsingular subarcs in singular optimal control problems have been developed. Necessary conditions developed previously by Kelley, Kopp and Moyer [3], which involve an analyticity assumption, have been clarified and extended. The main result in this direction is that the sum of the order of the singular arc and the order of the lowest time derivative of the control which is discontinuous at the junction must be an odd integer when the strengthened generalized Legendre-Clebsch condition is satisfied. Also, new necessary conditions which do not involve an analyticity assumption have been developed. These conditions aid mainly in characterizing problems which may possess nonanalytic junctions.

It should be emphasized that these are *local* necessary conditions for optimality. Yet, as indicated by the example in § 6, the point at which a junction occurs is determined mainly by initial and terminal boundary conditions, i.e., by essentially nonlocal information. This means that any junction theory which, for example, might be used to establish criteria for switching between singular and nonsingular arcs in an indirect computational scheme will have to take such nonlocal information into account.

It is becoming increasingly apparent that a close relationship exists between singular problems and bounded state problems [12], [13]. In this regard it is interesting to note that the result of Theorem 1 bears some similarity to a result of Jacobson, Lele and Speyer [13] which identifies certain properties of optimal trajectories associated with odd order state space constraints. Such similarities suggest the possibility of a duality between these two classes of problems which might be profitably exploited.

#### REFERENCES

- [1] J. P. MCDANELL AND W. F. POWERS, *New Jacobi-type necessary and sufficient conditions for singular optimization problems*, AIAA J., 8 (1970), pp. 1416–1420.
- [2] J. L. SPEYER AND D. H. JACOBSON, *Necessary and sufficient conditions for optimality for singular control problems: A transformation approach*, Rep. 69-24, Analytical Mechanics Associates, Inc., Cambridge, Mass., 1969; J. Math. Anal. Appl., to appear.
- [3] H. J. KELLEY, R. E. KOPP AND H. G. MOYER, *Singular extremals*, Topics in Optimization, G. Leitmann, ed., Academic Press, New York, 1967.
- [4] C. D. JOHANSEN, *Singular Solutions in Problems of Optimal Control*, Advances in Control Systems, vol. 2, C. T. Leondes, ed., Academic Press, New York, 1965.
- [5] H. M. ROBBINS, *A generalized Legendre-Clebsch condition for the singular cases of optimal control*, IBM J. Res. Develop., 11 (1967), pp. 361–372.
- [6] J. DIEUDONNÉ, *Foundations of Modern Analysis*, Academic Press, New York, 1960.
- [7] A. T. FULLER, *Study of an optimum nonlinear control system*, J. Electronics Control, 15 (1963), pp. 63–71.
- [8] D. E. JOHANSEN, *Solution of a linear mean square estimation problem when process statistics are undefined*, Joint Automatic Control Conference, Troy, New York, 1965.
- [9] D. H. JACOBSON, *Differential dynamic programming methods for solving bang-bang control problems*, IEEE Trans. Automatic Control, AC-13 (1968), pp. 661–675.
- [10] ———, personal communication.
- [11] B. PAGUREK AND C. M. WOODSIDE, *The conjugate gradient method for optimal control problems with bounded control variables*, Automatica, 4 (1968), pp. 337–349.



- [12] D. H. JACOBSON AND M. M. LELE, *A transformation technique for optimal control problems with a state variable inequality constraint*, Tech. Rep. 574, Harvard Univ., Cambridge, Mass., 1968.
- [13] D. H. JACOBSON, M. M. LELE AND J. L. SPEYER, *New necessary conditions of optimality for control problems with state variable inequality constraints*, Tech. Rep. 597, Harvard Univ., Cambridge, Mass., 1969.

## DISCRETE SPLINES VIA MATHEMATICAL PROGRAMMING\*

O. L. MANGASARIAN† AND L. L. SCHUMAKER‡

**Abstract.** Existence, uniqueness and characterizing properties are given for a class of constrained minimization problems in real Euclidean space. These problems are the discrete analogues of minimization problems in Banach space whose solutions are generalized splines. Solutions of these discrete problems, which are called discrete splines, can be obtained by algorithms of mathematical programming.

**1. Introduction.** The purpose of this paper is to investigate some constrained minimization problems in real Euclidean space which are suggested by certain minimization problems in Banach space whose solutions are generalized splines (see [5], [9]). Because of the analogy with classical splines, and because the solutions of the discrete problems exhibit a spline-like structure, we call them discrete splines. Our aim here is to obtain existence, uniqueness and characterizing properties of the discrete splines. We include two detailed examples and a section with remarks indicating application to computing continuously constrained spline interpolants as well as some problems for further study.

**2. Definition of discrete splines.** Let  $L$  be a general forward difference operator of order  $m$  of the form

$$(2.1) \quad (Ly)_j = \sum_{v=0}^m a_v y_{v+j},$$

$$a_m \neq 0, \quad j = 0, 1, \dots, N - m.$$

$L$  maps a point  $y = (y_0, y_1, \dots, y_N)$  in  $R^{N+1}$  into a point in  $R^{N-m+1}$ . Denote by  $N_0$  the null space of  $L$ .

For  $i = 1, 2, \dots, k$ , let  $\mathcal{M}_i$  be forward difference operators of order  $m_i$  mapping points of  $R^{N+1}$  into  $R^{N-m_i+1}$ , where  $0 \leq m_i \leq m$  and

$$(2.2) \quad (\mathcal{M}_i y)_j = \sum_{v=0}^{m_i} a_{iv} y_{v+j},$$

$$a_{im_i} \neq 0, \quad j = 0, 1, \dots, N - m_i.$$

Let  $\alpha_i, \beta_i, i = 1, 2, \dots, k$ , be points in  $R^{N-m_i+1}$  such that  $\alpha_i \leq \beta_i$ ; that is, each component of  $\alpha_i$  is less than or equal to the corresponding component of  $\beta_i$ :

$$\alpha_{ij} \leq \beta_{ij}, \quad j = 0, 1, \dots, N - m_i.$$

Define

$$(2.3) \quad U = \{y \in R^{N+1} : \alpha_i \leq \mathcal{M}_i y \leq \beta_i, i = 1, 2, \dots, k\}.$$

---

\* Received by the editors March 19, 1970, and in revised form August 7, 1970.

† Computer Sciences Department, University of Wisconsin, Madison, Wisconsin 53706. The work of this author was supported by the National Science Foundation under Grant GJ-362 and the United States Army under Contract DA-31-124-ARO-D-462.

‡ Department of Mathematics, University of Texas, Austin, Texas 78712. The work of this author was supported by USAFOSR-69-1812.

Finally, let  $1 \leq p \leq \infty$  and let  $g \in R^{N-m+1}$  be prescribed. Throughout this paper we shall be concerned with the discrete minimization problems

$$(2.4) \quad \begin{aligned} & \underset{y \in U}{\text{minimize}} (\|Ly - g\|_p)^p, & 1 \leq p < \infty, \\ & \underset{y \in U}{\text{minimize}} \|Ly - g\|_\infty, & p = \infty, \end{aligned}$$

where

$$\|x\|_p = \left( \sum_{j=0}^{N-m} |x_j|^p \right)^{1/p}, \quad 1 \leq p < \infty,$$

and

$$\|x\|_\infty = \max_{0 \leq j \leq N-m} |x_j|.$$

We shall call the solutions of (2.4) *discrete splines*. The reason for this terminology will become more apparent in later sections.

**3. Existence and uniqueness of discrete splines.**

**THEOREM 3.1.**

- (i) A solution of (2.4) exists if  $U$  is nonempty.
- (ii) For  $1 < p < \infty$ , any two solutions of (2.4) differ by an element of  $N_0$ .
- (iii) If  $1 < p < \infty$  and  $s \in U$  is a solution of (2.4), then it is unique if and only if  $N_0 \cap U(s) = \{0\}$ , where  $U(s) = \{f - s : f \in U\}$ .

*Proof of Theorem 3.1(i).* This theorem is a consequence of the following more general theorem. (For  $p = 1$  or  $\infty$ , (2.4) is a bounded linear program and hence has a solution [3, Theorem 2, p. 134].)

**THEOREM 3.2.** Let  $A$  be an  $l \times n$  matrix,  $L$  a  $k \times n$  matrix,  $b$  an  $l$ -vector,  $g$  a  $k$ -vector,

$$\begin{aligned} X &= \{x : x \in R^n, Ax \leq b\}, \\ S &= \{Lx - g\} = \{z : z \in R^k, z = Lx - g, x \in X\}, \end{aligned}$$

and let  $\varphi$  be a lower semicontinuous function on  $S$  satisfying some growth condition such as this: There exist a  $\hat{z} \in S$  and  $\rho \geq 0$  such that for any  $z \in S$ ,

$$(3.1) \quad \|z\|_\infty > \rho \Rightarrow \varphi(z) > \varphi(\hat{z}).$$

Then the minimization problem

$$(3.2) \quad \min \{\varphi(Lx - g) : x \in R^n, Ax \leq b\}$$

has a solution provided  $X$  is nonempty.

*Proof of Theorem 3.2.* Because of the growth condition (3.1) the minimization problem (3.2) is equivalent to

$$\min \{\varphi(z) : z \in S, \|z\|_\infty \leq \rho\}.$$

Because  $\varphi$  is lower semicontinuous on  $S$ , and because  $\{z : \|z\|_\infty \leq \rho\}$  is compact, this problem has a solution provided that  $S$  is closed. We show now that  $S$  is closed.<sup>1</sup>

<sup>1</sup> A linear map of a closed convex set is *not* in general closed. For example the linear map  $f(X)$  of the closed convex set  $X = \{x : x \in R^2, x_2 \geq e^{-x_1}\}$ , where  $f(x) = x_2$ , is the open interval  $(0, \infty)$ . What we show above is that a linear map of a closed *polyhedral* set is closed.

Let  $\bar{z}$  be a point of closure of  $S$ . Then

$$(3.3) \quad \inf \{ \varepsilon : \varepsilon \in \mathbb{R}, -e\varepsilon \leq z - \bar{z} \leq e\varepsilon, z \in S \} = 0,$$

where  $e$  is a  $k$ -vector of ones. But since (3.3) is a linear program, the infimum 0 is attained at some  $z^* \in S$  (see [3], [8]). Hence  $z^* - \bar{z} = 0, \bar{z} = z^* \in S$ , and  $S$  is closed. (This proof is similar to that of Laurent [7] except that we use linear programming to establish the fact that  $S$  is closed, whereas Laurent uses a theorem of Dieudonné concerning the closure of the algebraic difference of two sets in a Hausdorff topological vector space.)

Theorem 3.1(i) follows from Theorem 3.2 by making the identifications:  $x = y, z = Lx - g, k = N - m + 1, n = N + 1, l = 2 \sum_{i=1}^k (N - m_i + 1), \varphi(z) = (\|z\|_p)^p, 1 \leq p < \infty$ , and

$$A = \begin{bmatrix} \tilde{M}_1 \\ \vdots \\ \tilde{M}_k \\ -\tilde{M}_1 \\ \vdots \\ -\tilde{M}_k \end{bmatrix}_{l \times N+1}, \quad b = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_k \\ -\alpha_1 \\ \vdots \\ -\alpha_k \end{bmatrix}_{l \times 1},$$

$$\tilde{M}_i = \begin{bmatrix} a_{i0} & \cdots & a_{im_i} & 0 & \cdots & 0 \\ 0 & a_{i0} & \cdots & a_{im_i} & 0 & \cdots & 0 \\ \vdots & & & & & & \vdots \\ 0 & & & & 0 & a_{i0} & \cdots & a_{im_i} \end{bmatrix}_{N-m_i+1 \times N+1},$$

observing that all vector norms are continuous functions [4], and that the growth condition (3.1) is satisfied as follows by  $\varphi(z) = (\|z\|_p)^p, 1 \leq p < \infty$ . Let  $\hat{z}$  be any point in  $S$ . Then for any  $z \in S$ ,

$$\begin{aligned} \varphi(z) \leq \varphi(\hat{z}) &\Rightarrow \|z\|_p \leq \|\hat{z}\|_p \\ &\Rightarrow \alpha_p \|z\|_\infty \leq \|\hat{z}\|_p \quad \text{for some } \alpha_p > 0 \\ &\quad \text{(since } \alpha_p \|z\|_\infty \leq \|z\|_p, \alpha_p > 0 \text{ (see [4]))} \\ &\Rightarrow \|z\|_\infty \leq \rho = \|\hat{z}\|_p / \alpha_p. \end{aligned}$$

*Proof of Theorem 3.1(ii).* Since for  $1 < p < \infty, (\|z - g\|_p)^p$  is a strictly convex function of  $z$ , it follows for any minimum solution  $y, z = Ly$  is unique. Hence for any two solutions  $y^1$  and  $y^2$  we have that  $Ly^1 = Ly^2$  or  $L(y^1 - y^2) = 0$  and hence differ by an element of the null space  $N_0$  of  $L$ .

*Proof of Theorem 3.1(iii).* Let  $s$  be a solution of (2.4). If  $N_0 \cap U(s) \neq \{0\}$ , then we may add any nontrivial element  $\hat{s} \in N_0 \cap U(s)$  to  $s$  and obtain  $s + \hat{s} \in U$  and  $L(s + \hat{s}) = Ls + L\hat{s} = Ls$ , and hence  $s + \hat{s}$  is another distinct solution.

Conversely, if  $s_1$  and  $s_2$  are two solutions of (2.4), then by part (ii) of this theorem  $L(s_1 - s_2) = 0$ , and hence  $s_1 - s_2 \in N_0$ . But  $s_1 - s_2 \in U(s_2)$ . Hence

$s_1 - s_2 \in N_0 \cap U(s_2)$ , which implies that  $s_1 = s_2$  whenever this set consists of the zero element only.

The existence assertion (i) of Theorem 3.1 could also be established without the use of programming theory by following the lines of the proof of a similar theorem in an appropriate Sobolev space for the continuous analogue of (2.4) (see [9]).

**4. Necessary and sufficient characterization of discrete splines.**

**THEOREM 4.1.** *For  $1 < p < \infty$  a point  $y \in U$  is a solution of (2.4) if and only if for  $i = 1, 2, \dots, k$  there exist vectors  $\lambda_i \geq 0, \mu_i \geq 0$  in  $R^{N-m_i+1}$  such that*

$$(4.1) \quad H_1 = \sum_{i=1}^k [\lambda_i(-\mathcal{M}_i y + \alpha_i) + \mu_i(\mathcal{M}_i y - \beta_i)] = 0$$

and

$$(4.2) \quad H = (\|Ly - g\|_p)^p + H_1$$

satisfy

$$(4.3) \quad \nabla_j H = \partial H / \partial y_j = 0, \quad j = 0, 1, \dots, N,$$

where

$$\lambda_i(-\mathcal{M}_i y + \alpha_i) = \sum_{j=0}^{N-m_i} \lambda_{ij}[-(\mathcal{M}_i y)_j + \alpha_{ij}].$$

*Proof.* The conditions (4.1)–(4.3) follow by a direct application of the Kuhn–Tucker optimality theorem of nonlinear programming [6], [8]. That they are necessary follows from the fact that the set  $U$  is polyhedral, and hence no constraint qualification is needed. That the conditions are sufficient follows from the convexity of  $U$  and the objective function  $(\|Ly - g\|_p)^p$ .

These conditions completely characterize the discrete splines, although in this generality they do not provide a very good idea of the structure of such splines. In the next section we shall specialize problem (2.4) and thereby define what we will call discrete polynomial splines, whose structure can then be more explicitly delineated.

For convenience we note that for  $p = 2$ ,

$$\begin{aligned} \nabla_j (\|Ly\|_2)^2 &= \nabla_j \sum_{i=0}^{N-m} \left( \sum_{v=0}^m a_v y_{v+i} \right)^2 \\ (4.4) \quad &\begin{cases} 2 \sum_{l=0}^j a_{j-l} \left( \sum_{v=0}^m a_v y_{v+l} \right), & j = 0, 1, \dots, m-1, \\ (4.5) \quad = \begin{cases} 2 \sum_{l=0}^m a_{m-l} \left( \sum_{v=0}^m a_v y_{v+l+j-m} \right), & j = m, \dots, N-m, \\ (4.6) \quad \begin{cases} 2 \sum_{l=0}^{N-j} a_{m-l} \left( \sum_{v=0}^m a_v y_{v+l+j-m} \right), & j = N-m+1, \dots, N. \end{cases} \end{cases} \end{cases} \end{aligned}$$

**5. Discrete natural polynomial splines.** Given integers  $N, m$ , a set of integers  $\Lambda = \{0 \leq i_1 < i_2 < \dots < i_k \leq N\}$ , and real numbers  $\gamma = \{\gamma_j\}_1^k$ , we define a discrete natural polynomial spline of degree  $2m - 1$  with knots  $i_1 < \dots < i_k$  as a vector  $y = (y_0, \dots, y_N)$  which solves the minimization problem

$$(5.1) \quad \underset{y \in U_N(\gamma)}{\text{minimize}} \frac{1}{h^{2m-1}} \sum_{i=0}^{N-m} (\Delta^m y_i)^2,$$

where

$$(5.2) \quad U_N(\gamma) = \{y = (y_0, \dots, y_N) : y_{i_j} = \gamma_j, j = 1, 2, \dots, k\}$$

and

$$\Delta^m y_i = \sum_{v=0}^m (-1)^{m-v} \binom{m}{m-v} y_{i+v}$$

is the  $m$ th forward difference with  $h = 1/N$ .

The characterization of discrete natural polynomial splines is given by the following theorem.

**THEOREM 5.1.** *A vector  $y \in U_N(\gamma)$  solves (5.1) if and only if, for  $j \notin \Lambda$ ,*

$$(5.3a) \quad \Delta^{2m} y_{j-m} = 0, \quad j = m, \dots, N - m,$$

$$(5.3b) \quad \Delta^m y_j = 0, \quad j = 0, 1, \dots, i_1 - 1,$$

$$(5.3c) \quad \Delta^m y_j = 0, \quad j = i_k + 1, \dots, N - m.$$

*Proof.* Let  $\Lambda_1 = \{m, \dots, N - m\}$ ,  $\Lambda_2 = \{0, \dots, m - 1\}$  and  $\Lambda_3 = \{N - m + 1, \dots, N\}$ . By setting  $a_v = h^{-(2m-1)} (-1)^{m-v} \binom{m}{m-v}$ , we can write the following with the help of (4.4)–(4.6):

$$\begin{aligned} & \frac{\partial}{\partial y_j} \left( \frac{1}{h^{2m-1}} \sum_{i=0}^{N-m} (\Delta^m y_i)^2 \right) \\ &= \begin{cases} \frac{2}{h^{4m-2}} \sum_{l=0}^m (-1)^l \binom{m}{l} \sum_{v=0}^m (-1)^{m-v} \binom{m}{m-v} y_{v+l+j-m}, & j \in \Lambda_1, \\ \frac{2}{h^{4m-2}} \sum_{l=0}^j (-1)^{m-j+l} \binom{m}{m-j+l} \sum_{v=0}^m (-1)^{m-v} \binom{m}{m-v} y_{v+l}, & j \in \Lambda_2, \\ \frac{2}{h^{4m-2}} \sum_{l=0}^{N-j} (-1)^l \binom{m}{l} \sum_{v=0}^m (-1)^{m-v} \binom{m}{m-v} y_{v+l+j-m}, & j \in \Lambda_3, \end{cases} \end{aligned}$$

$$= \begin{cases} \frac{2}{h^{4m-2}} \sum_{l=0}^m (-1)^l \binom{m}{l} \Delta^m y_{l+j-m} = \frac{2}{h^{4m-2}} \Delta^{2m} y_{j-m}, & j \in \Lambda_1, \\ \frac{2}{h^{4m-2}} \sum_{l=0}^j (-1)^{m-j+l} \binom{m}{j-l} \Delta^m y_l, & j \in \Lambda_2, \\ \frac{2}{h^{4m-2}} \sum_{l=0}^{N-j} (-1)^l \binom{m}{l} \Delta^m y_{l+j-m}, & j \in \Lambda_3. \end{cases}$$

Now by Theorem 4.1 we have for  $j \notin \Lambda$ ,

$$(5.4a) \quad \Delta^{2m} y_{j-m} = 0, \quad j \in \Lambda_1,$$

$$(5.4b) \quad (\varphi y)_j = \sum_{l=0}^j (-1)^{m+j+l} \binom{m}{j-l} \Delta^m y_l = 0, \quad j \in \Lambda_2,$$

$$(5.4c) \quad (\psi y)_j = \sum_{l=0}^{N-j} (-1)^l \binom{m}{l} \Delta^m y_{l+j-m} = 0, \quad j \in \Lambda_3.$$

By coupling (5.4a) for  $j < i_1$  with (5.4b) we obtain (5.3b). Similarly, (5.3c) follows from (5.4c) and (5.4a) for  $j > i_k$ .

Property (5.3a) asserts that the values  $y_i$  lie on polynomials of degree  $m - 1$  (since  $\Delta^m y_i = 0$ ) in the intervals  $[0, i_1 - 1]$  and  $[i_k + 1, N]$ , while between the knots the  $y_i$  lie on polynomials of degree  $2m - 1$  ( $\Delta^{2m} y_i = 0$ ). It is not difficult to determine how these polynomials tie together at the knots. For example, let  $j$  be a knot of  $y$  separated from the other knots by at least 1; i.e., if  $j_1, j_2$  are the knots preceding and following  $j$ , respectively, then  $j_1 < j - 1$  and  $j + 1 < j_2$ . Let  $P_L(t)$  and  $P_R(t)$  be the polynomials of degree  $2m - 1$  on which the  $y_i$  lie to the left and right of  $j$ . Since  $\Delta^{2m} y_{j-m-1} = 0$ , the points  $\{y_i\}_{j-m-1}^{j+m-1}$  lie on  $P_L$ . Similarly, since  $\Delta^{2m} y_{j-m+1} = 0$ , the points  $\{y_i\}_{j-m+1}^{j+m+1}$  lie on  $P_R$ . Thus we conclude that the values  $\{y_i\}_{j-m+1}^{j+m-1}$  are common to both  $P_L$  and  $P_R$ . Consequently,  $\Delta^v P_L(j - m + 1) = \Delta^v P_R(j - m + 1) = \Delta^v y_{j-m+1}$  for  $v = 0, 1, \dots, 2m - 2$ . If these differences divided by  $h^v$  converge as  $N \rightarrow \infty$ , then they approach the corresponding derivatives of  $P_L$  and  $P_R$ , and it follows that  $P_L$  and  $P_R$  are tied together of class  $C^{2m-2}$  in the limit.

It is interesting to compare the discrete natural polynomial spline with the classical natural polynomial spline. We recall that a natural polynomial spline of degree  $2m - 1$  with knots  $x_1 < x_2 < \dots < x_k$  is a function  $s(x) \in C^{2m-2}$  such that

$$(5.5a) \quad s^{(2m)}(x) = 0 \quad \text{for } x \in (x_i, x_{i+1}), \quad i = 1, 2, \dots, k - 1,$$

$$(5.5b) \quad s^{(m)}(x) = 0 \quad \text{for } x < x_1, \quad x > x_k.$$

These splines are the solutions of

$$(5.6) \quad \underset{f \in U(\gamma)}{\text{minimize}} \|f^{(m)}\|_2,$$

where  $U(\gamma) = \{f \in H_2^m : f(x_i) = \gamma_i, \quad i = 1, 2, \dots, k\}$  and  $H_2^m = \{f \in C^{(m-1)} : f^{(m)}$  absolutely continuous,  $f^{(m)} \in L_2\}$ .

It is now clear how the discrete polynomial splines resemble their classical analogues in the continuous version (5.6), while (2.4) was obtained as the dis-

cretization of the minimization problem defining very general splines in [9]. While in general the discrete natural splines are not simply the discretized values of their continuous counterparts, it should not be surprising that they are closely connected, and in §7 we shall make some remarks about convergence in an appropriate sense.

**6. Examples.** To further motivate the study of the minimization problem (2.4) and to further illustrate the nature of discrete splines as well as their relationship with the well-known continuous splines, we devote this section to some rather simple examples.

*Example 6.1.* The sequence  $\{s_i = s(ih)\}_{i=-2N}^{2N}$ , where  $h = 1/N$  and

$$(6.1) \quad s(t) = \begin{cases} 1, & -2 \leq t < -1, \\ -t, & -1 \leq t < 0, \\ 2t, & 0 \leq t < 1, \\ 2, & 1 \leq t \leq 2, \end{cases}$$

is the unique solution of the problem :

$$\text{minimize } \frac{1}{h} \sum_{-2N}^{2N-1} (\Delta y_i)^2$$

subject to  $1 \leq y_{-N} \leq \frac{3}{2}, \quad -\frac{1}{4} \leq y_0 \leq 0, \quad 2 \leq y_N \leq \frac{5}{2},$

where  $\Delta y_i = -y_i + y_{i+1}$ . We note that the function  $s(t)$  above uniquely minimizes  $\int_{-2}^2 (f')^2 dt$  subject to  $1 \leq f(-1) \leq \frac{3}{2}, \quad -\frac{1}{4} \leq f(0) \leq 0,$  and  $2 \leq f(1) \leq \frac{5}{2}$ . Hence the sequence solving the discretized problem is obtained by discretizing the continuous solution of the continuous problem. (As noted above, this does not hold in general.)

*Discussion of Example 6.1.* By Theorem 4.1,  $y = (y_{-2N}, y_{-2N+1}, \dots, y_{2N})$  satisfying the constraints is a solution of Example 6.1 if and only if

$$\Delta y_{-2N} = 0,$$

$$\Delta^2 y_i = 0, \quad i = -2N, -2N + 1, \dots, 2N - 2; \quad i \neq -N - 1, -1, N - 1,$$

$$(2/h)\Delta^2 y_{i-1} = \mu_i - \lambda_i, \quad i = -N, 0, N,$$

$$\Delta y_{2N-1} = 0, \quad \lambda_{-N}, \mu_{-N}, \lambda_0, \mu_0, \lambda_N, \mu_N \geq 0,$$

$$(1 - y_{-N})\lambda_{-N} + (y_{-N} - \frac{3}{2})\mu_{-N} = 0, \quad (-\frac{1}{4} - y_0)\lambda_0 + y_0\mu_0 = 0,$$

$$(2 - y_N)\lambda_N + (y_N - \frac{5}{2})\mu_N = 0.$$

It is easily checked that the above conditions are satisfied by the sequence  $\{s_i\}_{i=-2N}^{2N}$  defined in Example 6.1. Indeed, since  $s_i$  hits the lower, upper, and lower constraints respectively at  $-N, 0, N$ , we have  $\mu_{-N} = \lambda_0 = \mu_N = 0$  and

$$(2/h)\Delta^2 s_{-N-1} = -2 = -\lambda_{-N} \leq 0,$$

$$(2/h)\Delta^2 s_{-1} = 4 = \mu_0 \geq 0,$$

$$(2/h)\Delta^2 s_{N-1} = -4 = -\lambda_N \leq 0.$$



The uniqueness of the solution  $\{s_i\}_{-2N}^{2N}$  is established by Theorem 3.1(iii). We have here

$$N_0 = \{y: y_i = \text{const.}, i = -2N, \dots, 2N\}.$$

We also have that

$$y = (y_{-2N}, \dots, y_{2N}) \in U(s) = \{f - s: f \in U\}$$

implies that  $0 \leq y_{-N} \leq \frac{1}{2}, -\frac{1}{4} \leq y_0 \leq 0, 0 \leq y_N \leq \frac{1}{2}$ . Thus  $N_0 \cap U(s) = \{y: y_i = 0, i = -2N, \dots, 2N\}$  and hence  $\{s_i\}_{-2N}^{2N}$  is unique.

*Example 6.2.* Let  $0 < \gamma \leq 1$ . The sequence

$$s_i = \begin{cases} (i/N) - (i/N)^2, & i = 0, 1, \dots, \bar{i}, \\ \alpha(i - N) + \gamma, & i = \bar{i} + 1, \dots, N, \end{cases}$$

where  $\bar{i}$  is the integer satisfying

$$1 + \frac{1}{N} - \frac{2(\bar{i} + 1)}{N} + \frac{(\bar{i} + 1)^2 - (\bar{i} + 1)}{N^2} < \gamma \leq 1 + \frac{1}{N} - \frac{2\bar{i}}{N} + \frac{\bar{i}^2 - \bar{i}}{N^2},$$

and

$$\alpha = \frac{1}{N - \bar{i}} \left( \gamma - \frac{\bar{i}}{N} + \frac{\bar{i}^2}{N^2} \right),$$

is the unique solution of the problem :

$$\text{minimize } \frac{1}{h} \sum_{i=0}^{N-1} (\Delta y_i)^2$$

subject to

$$i/N - (i/N)^2 \leq y_i \leq i/N, \quad i = 0, 1, \dots, N - 1, \\ y_N = \gamma.$$

As  $N \rightarrow \infty$ , the point  $\bar{i}/N$ , which is the last point at which  $s_i$  is on the curve  $t - t^2$ , converges to the value  $1 - \sqrt{\gamma}$ ; and the  $s_i$  become the discretized values of the function

$$s(t) = \begin{cases} t - t^2, & 0 \leq t < 1 - \sqrt{\gamma}, \\ (2\sqrt{\gamma} - 1)t + (1 + \gamma - 2\sqrt{\gamma}), & 1 - \sqrt{\gamma} \leq t \leq 1. \end{cases}$$

The function  $s(t)$  above is known (see [9, Example 2.2]) to be the unique solution of the constrained variational problem :

$$\text{minimize} \\ (6.2) \quad \int_0^1 (f')^2 dt \quad \text{subject to } t - t^2 \leq f(t) \leq t \quad \text{and } f(1) = \gamma.$$

*Discussion of Example 6.2.* By Theorem 4.1, a sequence  $\{s_i\}_0^N$  satisfying the constraints is a solution of Example 6.2 if and only if

$$\begin{aligned} (2/h)\Delta s_0 &= \mu_0 - \lambda_0, \\ (2/h)\Delta^2 s_{i-1} &= \mu_i - \lambda_i, \quad i = 1, \dots, N - 1, \\ (2/h)\Delta s_{N-1} &= -\mu_N + \lambda_N, \\ \left(\frac{i}{N} - \left(\frac{i}{N}\right)^2 - s_i\right)\lambda_i + \left(s_i - \frac{i}{N}\right)\mu_i &= 0, \quad i = 0, \dots, N - 1, \\ (-s_N + \gamma)\lambda_N + (s_N - \gamma)\mu_N &= 0. \end{aligned}$$

It is easily checked that these conditions are satisfied by the sequence  $\{s_i\}_0^N$  given in Example 6.2. Since for  $i = 1, \dots, \bar{i}$ ,  $s_i = i/N - (i/N)^2$  (that is,  $s_i$  is on the lower constraint), we have  $\mu_i = 0$ ,  $i = 1, \dots, \bar{i}$ . For  $i = \bar{i} + 1, \dots, N - 1$ ,  $s_i$  is not on any constraint and so  $\lambda_i = \mu_i = 0$  for these  $i$ . We have then

$$\begin{aligned} \frac{2}{h}\Delta^2 s_{i-1} &= \frac{2}{h} \left[ \frac{N^2\gamma - i^2 - N^2 - N + i + 2Ni}{N^2(N - i)} \right] \leq 0, \\ \frac{2}{h}\Delta s_0 &= 2(1 - h) = \mu_0 - \lambda_0, \\ \frac{2}{h}\Delta^2 s_{i-1} &= -\frac{4}{N} = -\lambda_i \leq 0, \quad i = 1, \dots, \bar{i} - 1, \\ \frac{2}{h}\Delta^2 s_{i-1} &= 0, \quad i = \bar{i} + 1, \dots, N - 1. \end{aligned}$$

This example was also solved numerically on a *Borroughs 5500* computer for  $\gamma = \frac{1}{8}$  using *Rosen's* gradient projection algorithm [10] of nonlinear programming. Computer time for the case of  $N = 10$  was 10 seconds and for  $N = 20$  was 40 seconds. The numerical results agreed with the known solution.

It is interesting to note that if we use the 1-norm instead of the 2-norm in this example, that is,

$$\text{minimize } \frac{1}{h} \sum_{i=0}^{N-1} |\Delta y_i|,$$

then for  $\frac{1}{4} \leq \gamma \leq 1$ , any nondecreasing sequence  $\{y_i\}_0^N$  satisfying the constraints solves the problem, and the minimum value is  $\gamma$ . For  $0 \leq \gamma \leq \frac{1}{4}$ , any sequence  $\{y_i\}_0^N$  satisfying the constraints such that  $\{y_i\}_{i \leq N/2}$  is nondecreasing with  $y_{(N-1)/2}$  or  $y_{N/2}$  on the lower constraint and  $\{y_i\}_{i \geq N/2}$  is nonincreasing solves the problem. If we use the  $\infty$ -norm, that is,

$$\text{minimize } \max_{0 \leq i \leq N-1} (1/h)|\Delta y_i|,$$

then for  $\gamma \leq 1 - h$ , any sequence  $\{y_i\}_0^N$  satisfying the constraints with  $|\Delta y_i| \leq |\Delta y_0| = h - h^2$  is a solution. If  $1 - h < \gamma \leq 1$ , then  $y_i = i\gamma/N$  uniquely solves the problem. As  $N \rightarrow \infty$ ,  $\Delta y_0/h \rightarrow 1$  and  $|\Delta y_i| \leq h$  for  $i = 1, \dots, N - 1$ . Both of these problems using the 1-norm and  $\infty$ -norm can be formulated as linear programs. Numerical experiments verifying the results in these cases were also performed.

**7. Applications and remarks.** (a) At present we see the principal application of discrete splines as a means to obtaining reasonable approximations to the interpolation problems with side constraints discussed in [9]. Example 6.2 discussed above illustrates how a rather difficult constrained minimization problem (6.2) for determining a certain spline interpolant can be approximately solved by computing a discrete spline. The computation of discrete splines is not difficult in view of the many available nonlinear programming algorithms (cf. [10], [12]). By using the artifice of an enlarged constraint set, Daniel [2] has shown that the discrete splines converge in an appropriate sense to the spline solutions of the continuous problems. Although rates of convergence are lacking, our experience with some simple examples indicates the practicality of this computational approach.

(b) As a possible further application, we mention the theory of best quadrature formulas in the sense of Sard and its close connection with natural polynomial splines (see, e.g., [5], [11]). Certain summation formulas may have similar connections with discrete splines.

(c) The theory of optimal control was convenient in the study of the continuous analogue of (2.4) in [9]. Another promising direction of investigation would be to examine the use of discrete optimal control for (2.4).

(d) Atteia [1] and Laurent [7] have studied rather general splines in Hilbert space settings which can be considered to subsume problem (2.4) in the special case of  $p = 2$ . The present investigation differs from theirs in that we consider some very specific problems in a Euclidean space, and thus can get more explicit characterization results.

**Acknowledgments.** We wish to acknowledge useful discussions with J. W. Daniel and K. Gehner, suggestions by M. Golomb and P. Laurent, and the computational programming of D. Kuba.

#### REFERENCES

- [1] M. ATTEIA, *Fonctions spline définies sur un ensemble convexe*, Numer. Math., 12 (1968), pp. 192–210.
- [2] J. W. DANIEL, *Convergence of a discretization for constrained spline function problems*, this Journal, 9 (1971).
- [3] G. B. DANTZIG, *Linear Programming and Extensions*, Princeton Univ. Press, Princeton, N.J., 1963.
- [4] E. ISAACSON AND H. B. KELLER, *Analysis of Numerical Methods*, John Wiley, New York, 1966.
- [5] J. JEROME AND L. L. SCHUMAKER, *On Lg-splines*, J. Approx. Theory, 2 (1969), pp. 29–49.
- [6] H. W. KUHN AND A. W. TUCKER, *Nonlinear programming*, Proc. Second Berkeley Symposium on Mathematical Statistics and Probability, J. Neyman, ed., Univ. of California Press, Berkeley, 1951, pp. 481–492.
- [7] P. J. LAURENT, *Construction of spline functions in a convex set*, Approximations with Special Emphasis on Spline Functions, I. J. Schoenberg, ed., Academic Press, New York, 1969, pp. 415–446.
- [8] O. L. MANGASARIAN, *Nonlinear Programming*, McGraw-Hill, New York, 1969.
- [9] O. L. MANGASARIAN AND L. L. SCHUMAKER, *Splines via optimal control*, Approximations with Special Emphasis on Spline Functions, I. J. Schoenberg, ed., Academic Press, New York, 1969, pp. 119–156.
- [10] J. B. ROSEN, *The gradient projection method for nonlinear programming*, SIAM J. Appl. Math., 8 (1960), pp. 181–217.
- [11] I. J. SCHOENBERG, *On best approximations of linear operators*, Nederl. Akad. Wetensch. Proc. Ser. A, 67 (1964), pp. 155–163.
- [12] W. I. ZANGWILL, *Nonlinear Programming, A Unified Approach*, Prentice-Hall, Englewood Cliffs, N.J., 1969.

## OPTIMAL STATIONARY CONTROL WITH STATE AND CONTROL DEPENDENT NOISE\*

U. G. HAUSSMANN†

**Abstract.** The steady state optimal linear regulator with state and control dependent noise is analyzed in a manner similar to that developed by Wonham [1]. By state dependent noise we mean Gaussian white noise with coefficient linear in the state variable, and similarly for control dependent noise. Using a Lyapunov criterion for the existence of stationary probability distributions due to Zakai, it is possible to treat equations leading to diffusion processes with degenerate differential generators. It is found that if the noise is sufficiently small, then an optimal control exists. Further analysis, again using Lyapunov methods, yields conditions under which an optimal control exists no matter how large the noise is.

**1. Introduction.** Consider the linear control system described formally by

$$(1.1) \quad \dot{x} = Ax - Bu - C(u)\dot{\omega}_1 + D(x)\dot{\omega}_2 + E\dot{\omega}_3.$$

Here  $u$  is the control,  $\dot{\omega}_1$ ,  $\dot{\omega}_2$  and  $\dot{\omega}_3$  are independent Gaussian white noise disturbances, and  $C$  and  $D$  are linear in their arguments.  $D(x)\dot{\omega}_2$  and  $C(u)\dot{\omega}_1$  can represent wideband random perturbations of the system matrix  $A$  and of the input matrix  $B$  respectively.

The problem of interest is to choose a control  $u = \varphi(x)$  which minimizes, in the steady state, the expected quadratic cost

$$(1.2) \quad \mathcal{E}\{x'Mx + u'Nu\}.$$

In [1] and [2] it is shown that if  $EE'$  is positive definite then an optimal control exists, provided the control and state dependent noise is sufficiently small. Using a recent result due to Zakai [3], we are able to remove the restriction on  $E$ . Furthermore, we show that if the control dependent noise as well as the state dependent noise due to the stable modes affects only the stable modes, then an optimal control exists no matter how large the control dependent noise  $C(u)\dot{\omega}_1$  is. Further conditions which imply the existence of an optimal control independently of the size of the state dependent noise are also given.

We state the problem precisely in § 2. In § 3 we remove the restriction on  $E$ , and § 4 contains the proof that existence is independent of the noise level under certain conditions. In § 5 we apply the method to a stability problem and in § 6 we conclude with an example.

**2. The problem.** We consider the stochastic differential equation

$$(2.1) \quad dx = Ax dt - Bu dt - C(u) d\omega_1 + D(x) d\omega_2 + E d\omega_3, \quad t \geq 0,$$

where<sup>2</sup>  $x \in R^n$ ,  $u \in R^m$ , and  $\omega_1$ ,  $\omega_2$  and  $\omega_3$  are independent Wiener processes of

---

\* Received by the editors March 31, 1970.

† Department of Mathematics, University of British Columbia, Vancouver, British Columbia, Canada. This work was supported by the National Research Council of Canada under Contracts 67-4058 and 67-3105.

<sup>1</sup> ' denotes the transpose.

<sup>2</sup>  $R^n$  denotes real Euclidean  $n$ -space with norm given by  $\|x\| = \{\sum_{i=1}^n |x_i|^2\}^{1/2}$ .

dimension  $d_1, d_2$  and  $d_3$ , respectively.  $C$  and  $D$  are given by

$$(2.2) \quad C(u) = \sum_{i=1}^m C^i u_i,$$

$$(2.3) \quad D(x) = \sum_{i=1}^n D^i x_i,$$

respectively, and  $A, B, C^i, D^i, E$  are real constant matrices of corresponding dimensions. The pair  $(A, B)$  is assumed to be stabilizable.

If  $u$  in (2.1) has the form  $u = \varphi(x)$  where

$$(2.4) \quad \|\varphi(x) - \varphi(y)\| \leq k\|x - y\|, \quad x, y \in R^n,$$

then (2.1) is an equation of Itô's type. Hence if the random variable  $x(0)$  is independent of  $\omega_i, i = 1, 2, 3$ , then (2.1) determines a diffusion process

$$X_\varphi = \{x(t) : t \geq 0\}$$

(see [2], [4]). Of interest to us is the case where  $X_\varphi$  has an invariant probability measure  $\mu_\varphi$  defined on the Borel sets of  $R^n$ ; i.e., if  $x(0)$  has the probability distribution  $\mu_\varphi$ , then so does  $x(t), t > 0$ . As has been shown in [5] and [6], if  $EE' > 0$  (i.e., positive definite) then suitable sufficient conditions can be given so that an invariant  $\mu_\varphi$  exists. In this case  $\mu_\varphi$  is unique. Unfortunately for most problems in control theory one does not have  $EE' > 0$ . For example if the system is given (formally) by

$$(2.5) \quad \begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ a_1 & a_2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} (\dot{\omega}_3 + u)$$

so that it corresponds to

$$\frac{d^2 y}{dt^2} = a_2 \frac{dy}{dt} + a_1 y + \frac{d\omega_3}{dt} + u,$$

then  $EE'$  is not positive definite (so that the differential generator of the associated diffusion process is not fully elliptic) even though the system is controllable. Recently Zakai [3] has given alternate conditions which guarantee that an invariant  $\mu_\varphi$  exists. However it is in general not unique. (For systems of the form (2.5) it is unique [3].)

We define  $\Phi$ , the class of admissible control laws, to be the set of functions  $\varphi(\cdot)$  such that:

- (i)  $\varphi$  satisfies (2.4) for some constant  $k$ ;
- (ii) an invariant probability measure  $\mu_\varphi$  exists;
- (iii) for any invariant probability measure  $\mu$ ,

$$(2.6) \quad \mathcal{E}_\mu\{\|x\|^2\} = \int_{R^n} \|x\|^2 \mu(dx) < \infty.$$

Let  $L(x, u)$  be given by

$$(2.7) \quad L(x, u) = x'Mx + u'Nu,$$

where  $M$  and  $N$  are constant symmetric positive definite matrices of dimension

$n \times n, m \times m$ , respectively. For  $\varphi$  in  $\Phi$  let  $\mathcal{E}_\varphi\{L(x, \varphi)\}$  denote the infimum of  $\mathcal{E}_\mu\{L(x, \varphi)\}$  taken over all probability measures invariant with respect to  $X_\varphi$ . In [1], [2] Wonham considered the problem (with the extra restriction on  $E$ ) of finding conditions on  $C$  and  $D$  such that a control  $\varphi^\circ$  exists, optimal in the sense that

$$\mathcal{E}_{\varphi^\circ}\{L(x, \varphi^\circ)\} = \inf[\mathcal{E}_\varphi\{L(x, \varphi)\} : \varphi \in \Phi].$$

We shall refine these conditions.

Let us point out here that taking the infimum of  $\mathcal{E}_\mu$  in the definition of  $\mathcal{E}_\varphi$  is irrelevant, for under certain assumptions which we shall always make, it follows that for the optimal control  $\varphi^\circ$ , all invariant probability measures yield the same expected cost, i.e.,

$$(2.8) \quad \mathcal{E}_{\varphi^\circ}\{L(x, \varphi^\circ)\} = \mathcal{E}_\mu\{L(x, \varphi^\circ)\}$$

for all invariant probability measures  $\mu$  (see Theorem 3.2).

**3. Admissible and optimal controls.** First we must show that  $\Phi$  is nonvoid. Let  $\mathcal{L}_u$  be the differential operator given by

$$(3.1) \quad \mathcal{L}_u V(x) = \frac{1}{2} \text{tr} \{C(u)' V_{xx} C(u) + D(x)' V_{xx} D(x) + E' V_{xx} E\} + (Ax - Bu)' V_x,$$

where  $V_x$  denotes the vector  $\partial V / \partial x_i$  and  $V_{xx}$  denotes the matrix  $\partial^2 V / \partial x_i \partial x_j$ . We assume from now on that  $u = \varphi(x)$  in (2.1) with  $\varphi$  satisfying (2.4). Then  $\mathcal{L}_\varphi$  is the differential generator of  $X_\varphi$ , also written simply as  $\mathcal{L}$ .

Let us consider the following control law  $\varphi$  and function  $V$ :

$$(3.2) \quad \varphi(x) = Kx, \quad V(x) = x'Px,$$

where  $K$  and  $P$  are two constant matrices,  $P \geq 0$  (i.e.,  $x'Px \geq 0, x \in R^n$ ), to be determined such that

$$(3.3) \quad \mathcal{L}_\varphi V(x) = \lambda - L(x, \varphi(x))$$

for some  $\lambda > 0$ . Unfortunately such a  $V$  does not satisfy certain conditions imposed in [3]<sup>3</sup> and so results of [3] do not apply directly.

**THEOREM 3.1.** *Assume  $P$  and  $K$  exist ( $P \geq 0$ ) such that (3.2) and (3.3) are satisfied. Then  $\varphi$  is an admissible control.*

*Proof.* For  $n = 1, 2, \dots$  define

$$(3.4) \quad D_n(x) = \begin{cases} D(x) & \text{if } \|x\| \leq n, \\ D[x(n/\|x\|)] & \text{if } \|x\| > n, \end{cases}$$

$$\tilde{C}_n(x) = \begin{cases} C(Kx) & \text{if } \|x\| \leq n, \\ C[(n/\|x\|)Kx] & \text{if } \|x\| > n. \end{cases}$$

<sup>3</sup> Zakai calls this condition (B). It requires that  $\mathcal{E}\{f(t)|x(0) = a\}$  be bounded in  $t$  on any finite  $t$  interval, for any  $a \in R^n$ , where

$$f(t) = \|V_x(x(t))[C(K(x(t))), D(x(t)), E]\|^2.$$

We can only conclude that  $f(t) \leq k[1 + \|x(t)\|^4]$  for some  $k > 0$ .

Then<sup>4</sup>

$$(3.5) \quad \begin{aligned} \|D(x) - D_n(x)\| &\leq \begin{cases} 0 & \text{if } \|x\| \leq n, \\ k\|x\| & \text{if } \|x\| > n, \end{cases} \\ \|C(Kx) - \tilde{C}_n(x)\| &\leq \begin{cases} 0 & \text{if } \|x\| \leq n, \\ k\|x\| & \text{if } \|x\| > n, \end{cases} \end{aligned}$$

where

$$k = \max \left\{ \|K\| \sum_{i=1}^n \|C^i\|, \sum_{i=1}^n \|D^i\|, \|A - BK\| \right\}.$$

Also  $D_n(\cdot)$  and  $\tilde{C}_n(\cdot)$  are Lipschitz continuous with constant  $k$  (as are  $D(\cdot)$  and  $C(\varphi(\cdot))$ ). Now let  $x^n$  be the unique solution of

$$(3.6) \quad dx = (A - BK)x dt - \tilde{C}_n(x) d\omega_1 + D_n(x) d\omega_2 + E d\omega_3, \quad x(0) = \xi.$$

We let  $\mathcal{L}_n$  be the corresponding differential generator so that

$$(3.7) \quad \mathcal{L}_n V(x) = \begin{cases} x' \{ \Delta(P) + K' \Gamma(P) K + \hat{A}' P + P \hat{A} \} x + \text{tr} \{ E' P E \}, & \|x\| \leq n, \\ x' \left\{ \frac{n^2}{\|x\|^2} [ \Delta(P) + K' \Gamma(P) K ] + \hat{A}' P + \hat{P} A \right\} x + \text{tr} \{ E' P E \}, & \|x\| > n, \end{cases}$$

where

$$(\Delta(P))_{ij} = \text{tr} \{ (D^i)' P D^j \}, \quad (\Gamma(P))_{ij} = \text{tr} \{ (C^i)' P C^j \}, \quad \hat{A} = A - BK.$$

Then  $\mathcal{L}_n V(x) = \mathcal{L} V(x)$  if  $\|x\| \leq n$ , and  $\mathcal{L}_n V(x) \leq \mathcal{L} V(x)$  if  $\|x\| > n$ , because  $\Delta(P) + K' \Gamma(P) K \geq 0$  as  $P \geq 0$ .

We shall require the following three properties of  $V$ :

$$(3.8a) \quad \mathcal{E}_x \{ V(x(t)) \} = \mathcal{E} \{ V(x(t)) | x(0) = x \} \leq c \mathcal{E}_x \{ \|x(t)\|^2 \}, \quad c < \infty,$$

$$(3.8b) \quad \mathcal{E}_x \{ | \mathcal{L}_n V(x(t)) | \} \leq c \mathcal{E}_x \{ (1 + \|x(t)\|^2) \}, \quad n = 1, 2, \dots,$$

$$(3.8c) \quad \mathcal{E}_x \{ \| V_x(x(t))' [ \tilde{C}_n(x(t)), D_n(x(t)), E ] \|^2 \} \leq c \mathcal{E}_x \{ \|x(t)\|^2 \}.$$

Thus each of the three left-hand expressions is bounded in  $t$  over finite  $t$  intervals, for each  $x$  in  $R^n$ .

As  $L(x, u) = x' M x + u' N u$  then there is an  $R_0 < \infty$  such that  $\lambda - L(x, Kx) < -\eta$  for some  $\eta > 0$  if  $\|x\| > R_0$ . Now the proof of [3, Theorem 1] shows that

$$(3.9) \quad -\frac{1}{t} \mathcal{E} \{ V(\xi) \} + \eta \leq (\rho + \eta) \frac{1}{t} \int_0^t \text{Pr} \{ \|x^n(s)\| \leq R_0 \} ds$$

for any  $n$  if  $x^n(0) = \xi$  has the distribution  $\mu^+$ , where  $\mu^+$  is any Borel probability measure on  $R^n$  with compact support. Here  $\rho$  is the maximum of  $\mathcal{L} V(x)$  over  $\|x\| \leq R_0$ , so that  $\mathcal{L}_n V(x) \leq \rho$  for all  $n$ .

<sup>4</sup> If  $F$  is a matrix mapping  $R^n$  into  $R^m$ , then  $\|F\|$  denotes the usual operator norm.

In the Appendix we show that  $\mathcal{E}\{\|x_t^n - x_t\|^2\} \rightarrow 0$ . Hence a subsequence, which we call  $x_t^n$  again, converges to  $x_t$  with probability one, i.e.,

$$\Pr \left\{ \lim_{n \rightarrow \infty} \|x_t^n - x_t\| = 0 \right\} = 1.$$

By Egoroff's theorem  $x_t^n(\omega)$  converges to  $x_t(\omega)$  almost uniformly ( $\omega$ ). Hence for any  $\delta > 0$  there is an  $n(s)$  such that for  $n > n(s)$ ,

$$\Pr \{ \|x_s\| \leq R_0 + 1 \} \geq \Pr \{ \|x_s^n\| \leq R_0 \} - \delta$$

and so<sup>5</sup>

$$\Pr \{ \|x_s\| \leq R_0 + 1 \} \geq \lim_{n \rightarrow \infty} \Pr \{ \|x_s^n\| \leq R_0 \}.$$

By (3.9) and the dominated convergence theorem ( $t < \infty$ ) we have

$$-\frac{1}{t} \mathcal{E}\{V(\xi)\} + \eta \leq (\rho + \eta) \frac{1}{t} \int_0^t \Pr \{ \|x(s)\| \leq R_0 + 1 \} ds,$$

and so (see [3]) the process  $X_\varphi$  has an invariant measure.

Now  $|L(x_s^n) - L(x_s)| \leq k_1 \|x_s^n - x_s\|^2$ , so that

$$\left| \mathcal{E}_x \left\{ \int_0^t L(x_s^n) ds - \int_0^t L(x_s) ds \right\} \right| \leq k_1 \mathcal{E}_x \left\{ \int_0^t \|x_s^n - x_s\|^2 ds \right\}.$$

As  $t < \infty$ ,  $\mathcal{E}_x\{\|x_s^n\|^2\} < \infty$ ,  $\mathcal{E}_x\{\|x_s\|^2\} < \infty$  and as  $x(s, \omega)$  is measurable in  $(s, \omega)$ , then

$$\begin{aligned} \lim_{n \rightarrow \infty} \mathcal{E}_x \left\{ \int_0^t \|x_s^n - x_s\|^2 ds \right\} &= \lim_{n \rightarrow \infty} \int_0^t \mathcal{E}_x \{ \|x_s^n - x_s\|^2 \} ds \\ &= \int_0^t \lim_{n \rightarrow \infty} \mathcal{E}_x \{ \|x_s^n - x_s\|^2 \} ds. \end{aligned}$$

By Itô's formula and (3.8),

$$\begin{aligned} 0 &\leq \mathcal{E}_x \{ V(x_t^n) \} = V(x) + \mathcal{E}_x \left\{ \int_0^t \mathcal{L}_n V(x_s^n) ds \right\} \\ &\leq V(x) + \lambda t - \mathcal{E}_x \left\{ \int_0^t L(x_s^n) ds \right\}. \end{aligned}$$

It follows that  $\mathcal{E}_x \{ \int_0^t L(x_s) ds \} \leq V(x) + \lambda t$ . Now the argument given in [3, Theorem 2] shows that

$$\int_{\mathbb{R}^n} L(x) \mu(dx) \leq \lambda$$

for any invariant probability measure  $\mu$ . Hence  $\mathcal{E}_\mu\{\|x\|^2\} < \infty$  and  $\varphi(x) = Kx$  is admissible. The theorem is established.

<sup>5</sup> This limit can be shown to exist.



COROLLARY. Under the condition of the above theorem it follows that  $X_\varphi$  is stable insofar as

$$(3.10) \quad \mathcal{E}_x\{\tau\} \leq V(x)/\eta, \quad \|x\| > R_0,$$

where  $\tau$  is the first passage time of  $x_t$  to the set  $\{x: \|x\| \leq R_0\}$ .

This follows because if  $\tau_n$  denotes the corresponding random variable for  $x_t^n$ , then

$$0 \leq \mathcal{E}_x\{V(x^n(t \wedge \tau_n))\} = V(x) + \mathcal{E}_x\left\{\int_0^{t \wedge \tau_n} \mathcal{L}_n V(x^n(s)) ds\right\},$$

i.e.,

$$0 \leq V(x) - \eta \mathcal{E}_x\{t \wedge \tau_n\}.$$

But  $t \wedge \tau_n(\omega) \leq t$ , and so

$$0 \leq V(x) - \eta \mathcal{E}_x\{t \wedge \lim_{n \rightarrow \infty} \tau_n\} = V(x) - \eta \mathcal{E}_x\{t \wedge \tau\}.$$

Now (3.10) follows by monotone convergence. Relation (3.10) implies that the process is recurrent (and positive if  $\mathcal{L}$  is fully elliptic) [5].

Evidently (3.2) yields a solution of (3.3) if and only if

$$(3.11) \quad \lambda = \text{tr}\{E'PE\}$$

and

$$(3.12) \quad \Delta(P) + K'\Gamma(P)K + (A - BK)'P + P(A - BK) + M + K'NK = 0.$$

If  $(A, B)$  is stabilizable, then  $K$  can be chosen so that  $A - BK$  is stable, i.e., all the eigenvalues have negative real parts. The following lemma shows that if  $C^i, D^i$  are sufficiently small then (3.12) has a unique positive definite solution. This and Theorem 3.1 imply that  $\varphi \in \Phi$ .

LEMMA 3.1. If  $Q > 0$  and  $A$  is stable, then

$$(3.13) \quad \Delta(P) + K'\Gamma(P)K + A'P + PA + Q = 0$$

has a unique solution  $P > 0$ , provided

$$(3.14) \quad \left\| \int_0^\infty e^{tA'} [K'\Gamma(I)K + \Delta(I)] e^{tA} dt \right\| < 1.$$

For the proof see [1], [2].

Thus an admissible control exists if

$$(3.15) \quad \inf_K \left\| \int_0^\infty e^{t(A-BK)'} [K'\Gamma(I)K + \Delta(I)] e^{t(A-BK)} dt \right\| < 1,$$

where the infimum is taken only over  $K$  such that  $A - BK$  is stable. In [1], [2] the following theorem is proved.

THEOREM 3.2. If  $(A, B)$  is stabilizable, if  $M > 0, N > 0$ , and if  $\Gamma, \Delta$  satisfy (3.15), then an optimal control exists, of the form

$$(3.16) \quad \varphi^\circ(x) = K^\circ x.$$

Furthermore,

$$(3.17) \quad K^\circ = [\Gamma(P^\circ) + N]^{-1} B' P^\circ,$$

where  $P^\circ$  is the unique solution, in the class of positive semidefinite matrices, of the equation

$$(3.18) \quad M + A'P + PA + \Delta(P) - PB[\Gamma(P) + N]^{-1}B'P = 0.$$

The minimum cost is

$$(3.19) \quad \mathcal{E}_{\phi^\circ} \{L(x, \phi^\circ(x))\} = \text{tr} \{E'P^\circ E\}.$$

From the same proof it follows that (2.8) holds.

**4. Dependence on noise.** Theorem 3.2 shows that the existence of an optimal control depends on  $C$  and  $D$  only through (3.15), although of course the actual value of the control always depends on  $C$  and  $D$ . We now investigate (3.15) in an effort to determine if this condition can be relaxed. Let  $\alpha(\lambda)$  be the minimal polynomial of  $A$ . It factors as follows:

$$(4.1) \quad \alpha(\lambda) = \alpha^+(\lambda)\alpha^-(\lambda),$$

where all the zeros of  $\alpha^+$  lie in the closed right half-plane and those of  $\alpha^-$  in the open left half-plane. Also define

$$R_+^n = \{x \in R^n : \alpha^+(A)x = 0\}, \quad R_-^n = \{x \in R^n : \alpha^-(A)x = 0\}.$$

As  $\alpha^+$  and  $\alpha^-$  are coprime it follows that  $R^n = R_+^n \oplus R_-^n$ , although the direct sum is not necessarily orthogonal. Let  $P^+$  and  $P^-$  be the projections of  $R^n$  onto  $R_+^n$  and  $R_-^n$  along  $R_-^n$  and  $R_+^n$ , respectively.

The next two lemmas are useful for computation. The proofs are straightforward and so are omitted. However first we need some notation. Let  $g_1, \dots, g_r$ ,  $r \leq d_2$ , be a basis of the range space of the  $n^2 \times d_2$  matrix

$$(4.2) \quad \begin{bmatrix} D^1 \\ D^2 \\ \vdots \\ D^n \end{bmatrix},$$

and let  $G^i$  be the matrix composed of  $n$  columns whose direct sum is  $g_i$ ; i.e., if  $G_j^i$  is the  $j$ th column of  $G^i$ , then

$$(4.3) \quad g_i = \begin{bmatrix} G_1^i \\ \vdots \\ G_n^i \end{bmatrix}, \quad i = 1, 2, \dots, r.$$

**LEMMA 4.1.** Let  $H(\omega) = [h_1, h_2, \dots, h_n]$ , where  $h_i(\omega) = D^i\omega$ . Then  $R_-^n$  is invariant with respect to  $H(\omega)$  for all  $\omega \in R^{d_2}$  if and only if  $R_-^n$  is invariant with respect to  $G^i$ ,  $i = 1, \dots, r$ .

**LEMMA 4.2.** Assume  $\{e_i\}_1^n$  is the usual basis in  $R^n$ . If  $\{e_i\}_{i \in S}$  spans  $R_-^n$ , then  $R_-^n$  is invariant with respect to  $H(\omega)$  for all  $\omega$  if and only if  $\{D^i\} \subseteq R_-^n$ ,  $i \in S$ , where  $\{D^i\}$  denotes the range space of  $D^i$ .

The following lemma expresses an algebraic condition as an analytic one and is useful for computation. Krasovskii stated a related sufficient condition [7]. Let  $s$  be an eigenvector of  $A$ , i.e., for some  $\lambda$ , possibly complex,

$$As = \lambda s.$$

Let  $s^+$  denote any eigenvector  $s$  for which the corresponding eigenvalue has non-negative real part. If  $C$  is an  $n \times m$  matrix, then define  $R_C[A]$  by

$$(4.4) \quad R_C[A] = \{[C, AC, A^2C, \dots, A^{n-1}C]\}.$$

Recall that  $\{\cdot\}$  denotes the range space. Suppose the dimension of this space is  $q$ . Then  $C_A$  denotes an  $n \times q$  matrix whose columns are a basis of  $R_C[A]$ ,  $n \geq q$ .

LEMMA 4.3. *A necessary and sufficient condition that  $\{C\} \subseteq R_-^n$  is that*

$$(4.5) \quad \text{rank}(A - \lambda I)C_A = q$$

for each  $\lambda$  with  $\text{Re } \lambda \geq 0$ .

*Proof.* We first note that  $\{C\} \subseteq R_-^n$  if and only if no  $s^+$  lies in  $R_C[A]$ . This follows readily by considering  $A$  in Jordan form and noting that in that case  $R_+^n$  and  $R_-^n$  are orthogonal. Also note that  $\{C\} \subseteq R_-^n$  if and only if  $c_i \in R_-^n$ ,  $i = 1, \dots, m$ , where  $c_i$  is the  $i$ th column of  $C$ , and that  $R_C[A] = \text{span}\{R_{c_i}[A]: i = 1, \dots, m\}$ .

If  $c_A^i$  is the  $i$ th column of  $C_A$ , then any  $s \in R_C[A]$  is of the form

$$s = \sum_{i=1}^q \sigma_i c_A^i.$$

It follows that some  $s^+$  is in  $R_C[A]$  if and only if

$$(A - \lambda I)C_A \sigma = 0$$

has a nontrivial solution  $\sigma$  for some  $\lambda$  with  $\text{Re } \lambda \geq 0$ . A simple argument now completes the proof.

We now wish to split our system (2.1) into two parts, one on  $R_+^n$  and one on  $R_-^n$ . If  $R_+^n$  has dimension  $l_1$ , we can choose a map  $T_+$  of  $R^n$  onto  $R^{l_1}$  so that  $T_+$  restricted to  $R_+^n$  is one-to-one and onto, and  $T_+ R_-^n = \{0\}$ . Then  $T_+$  is an  $l_1 \times n$  matrix. We denote its generalized inverse by  $T_+^\dagger$  so that

$$T_+ T_+^\dagger = I, \quad T_+^\dagger T_+ = P^+, \quad P^+ T_+^\dagger = T_+^\dagger, \quad T_+ P^+ = T_+.$$

Similarly define  $T_-$ , an  $l_2 \times n$  matrix corresponding to  $R_-^n$ . Also define

$$(4.6) \quad \tilde{D}_+^j = T_+ \sum_{i=1}^n D^i (T_+^\dagger)_{ij}, \quad j = 1, 2, \dots, l_1,$$

and

$$(4.7) \quad \Delta_+(R)_{ij} = \text{tr}\{(\tilde{D}_+^i)' R (\tilde{D}_+^j)\}, \quad i, j = 1, \dots, l_1.$$

If we let  $A_+ = T_+ A T_+^\dagger$  and define  $A_-$ ,  $\Delta_-$  and  $\tilde{D}_-^j$  likewise, then  $A_-$  is a stable  $l_2 \times l_2$  matrix.

THEOREM 4.1. *If  $(A, B)$  is stabilizable (i.e.,  $R_+^n \subseteq R_B[A]$ , [8]), and if  $R_-^n$  is invariant with respect to  $H(\omega)$  for all  $\omega$ , then the system (2.1) has an admissible control*

$\varphi$ , provided that

$$(4.8) \quad \left\| \int_0^\infty e^{tA'} \Delta_-(I) e^{tA} dt \right\| = \sup \left\{ \int_0^\infty x' e^{tA'} \Delta(T_- T_-) e^{tA} x dt : x \in R_-^n, \|T_- x\| \leq 1 \right\} < 1,$$

and

$$(4.9) \quad \inf_{K_+} \left\| \int_0^\infty e^{t(A_+ - T_+ BK_+)' } [\Delta_+(I) + K_+ \Gamma(T_+ T_+) K_+] e^{t(A_+ - T_+ BK_+)} dt \right\| = \inf_K \sup \left\{ \int_0^\infty x' e^{t(A - P^+ BK)' } [\Delta(T_+ T_+) + K' \Gamma(T_+ T_+) K] \cdot e^{t(A - P^+ BK)} x dt : x \in R_+^n, \|T_+ x\| \leq 1 \right\} < 1.$$

Moreover  $\varphi$  can be chosen linear in  $x$ , independent of  $P^- x$ . Also (4.8) is satisfied if

$$(4.10) \quad \sum_{ijk} |(T_- D^k)_{ij}|^2 < \left[ \sup \left\{ \int_0^\infty \|e^{tA} x\|^2 dt : x \in R_-^n, \|T_- x\| \leq 1 \right\} \right]^{-1}.$$

*Proof.* As  $R_-^n$  is invariant under  $H(\omega)$  then

$$(4.11) \quad T_+ D(x)\omega = T_+ P^+ D(P^+ x)\omega = T_+ D(T_+^\dagger T_+ x)\omega = \tilde{D}_+(T_+ x)\omega,$$

where  $\tilde{D}_+(y) = \sum_{i=1}^{l_1} \tilde{D}_+^i y_i$ . We let  $T_+ x = y$ ,  $T_- x = z$ . Then (2.1) becomes

$$(4.12) \quad dy = (A_+ y - T_+ B u) dt - T_+ C(u) d\omega_1 + \tilde{D}_+(y) d\omega_2 + T_+ E d\omega_3,$$

$$(4.13) \quad dz = (A_- z - T_- B u) dt - T_- C(u) d\omega_1 + \tilde{D}_-(z) d\omega_2 + T_- E d\omega_3 + T_- D(T_+^\dagger y) d\omega_2,$$

where  $\tilde{D}_-(z) = \sum_{i=1}^{l_2} \tilde{D}_-^i z_i = T_- D(T_+^\dagger z)$ . Consider (4.13) with  $y = 0$ ,  $u = 0$ , i.e.,

$$(4.14) \quad dz = A_- z dt + \tilde{D}_-(z) d\omega_2 + T_- E d\omega_3.$$

As (4.14) has a stable matrix, then by (4.8) and Lemma 3.1 there is a symmetric positive definite  $l_2 \times l_2$  matrix  $P_2$  such that

$$(4.15) \quad \Delta_-(P_2) + A'_- P_2 + P_2 A_- + I = 0.$$

As  $R_+^n \subseteq R_B[A]$  then  $R^{l_1} = T_+ R_+^n \subseteq R_{T_+ B}[A_+]$  so that  $(A_+, T_+ B)$  is controllable [8]. Now (4.9) and Lemma 3.1 show that there are an  $l_1 \times m$  matrix  $K_+$  and an  $l_1 \times l_1$  matrix  $P_1 > 0$ , such that

$$(4.16) \quad K_+ \Gamma(T_+ P_1 T_+) K_+ + \Delta_+(P_1) + (A_+ - T_+ BK_+)' P_1 + P_1 (A_+ - T_+ BK_+) + I = 0.$$

Let  $\lambda_1 = \text{tr} \{(T_+ E)' P_1 (T_+ E)\}$ ,  $\lambda_2 = \text{tr} \{(T_- E)' P_2 (T_- E)\}$ , and let  $K = K_+ T_+$ .

To show that  $\varphi(x) = Kx$  is an admissible control we shall use Theorem 3.1. Consider  $V(x) = x' P x$ , where  $P = T_+ P_1 T_+ + \beta T_- P_2 T_- > 0$  and  $\beta > 0$  is to be

chosen so that for some  $\lambda > 0, \alpha > 0$ ,

$$(4.17) \quad \mathcal{L}_\phi V(x) \leq \lambda - \alpha \|x\|^2.$$

This condition, rather than (3.3), suffices in Theorem 3.1. Using (4.11), (4.15) and (4.16) we have

$$(4.18) \quad \begin{aligned} \mathcal{L}_\phi V(x) &= y'K'_+\Gamma(T'_+P_1T_+)K_+y + \beta y'K'_+\Gamma(T'_-P_2T_-)K_+y \\ &\quad + y'\Delta_+(P_1)y + \beta z'\Delta_-(P_2)z + \beta\{y'(T'_+\dagger)\Delta(T'_-P_2T_-)T'_+\dagger z \\ &\quad + z'(T'_+\dagger)\Delta(T'_-P_2T_-)T'_+\dagger y + y'(T'_+\dagger)\Delta(T'_-P_2T_-)T'_+\dagger y\} \\ &\quad + \lambda_1 + \beta\lambda_2 + y'(A_+ - T_+BK_+)P_1y + y'P_1(A_+ - T_+BK_+)y \\ &\quad + \beta z'A'_-P_2z + \beta z'P_2A_-z - 2\beta y'(T_-BK_+)P_2z \\ &= \lambda_1 + \beta\lambda_2 - \|y\|^2 - \beta\|z\|^2 + \beta\{y'K'_+\Gamma(T'_-P_2T_-)K_+y \\ &\quad + 2z'(T'_+\dagger)\Delta(T'_-P_2T_-)T'_+\dagger y + y'(T'_+\dagger)\Delta(T'_-P_2T_-)T'_+\dagger y \\ &\quad - 2z'P_2T_-BK_+y\} \\ &\leq \lambda_1 + \beta\lambda_2 - (\beta/(4\gamma))\|x\|^2 \end{aligned}$$

if  $0 < \beta$  is sufficiently small and  $\gamma = \max\{\|T'_+\dagger\|^2, \|T'_-\dagger\|^2\}$ . This last assertion is proved as follows. Let  $f(y, z)$  be given by

$$\begin{aligned} f(y, z) &= y'\{K'_+\Gamma(T'_-P_2T_-)K_+ + (T'_+\dagger)\Delta(T'_-P_2T_-)T'_+\dagger\}y \\ &\quad + 2z'\{(T'_-\dagger)\Delta(T'_-P_2T_-)T'_-\dagger - P_2T_-BK_+\}y. \end{aligned}$$

Consider

$$(4.19) \quad \|y\|^2(1 - \frac{1}{2}\beta) + \frac{1}{2}\beta\|z\|^2 \geq \beta c_1\|y\|^2 + 2\beta c_2\|y\|\|z\|,$$

where  $c_1 \geq 0, c_2 \geq 0$ . If  $\|z\| = 0$ , then (4.19) holds for all  $y$  if  $0 < \beta \leq (c_1 + \frac{1}{2})^{-1}$ . Otherwise let  $\|y\| = k\|z\|, k$  variable in  $[0, \infty)$ . Then (4.19) holds for all  $k$  if

$$(4.20) \quad g(k) = k^2(1 - \frac{1}{2}\beta) + \frac{1}{2}\beta - \beta c_1 k^2 - 2\beta c_2 k \geq 0$$

for all  $k \in [0, \infty)$ . But if  $0 < \beta < (c_1 + \frac{1}{2})^{-1}$  and if  $\beta < (2c_2^2 + c_1 + \frac{1}{2})^{-1}$ , then (4.20) holds for all  $k \in [0, \infty)$ . If

$$\begin{aligned} c_1 &= \|K'_+\Gamma(T'_-P_2T_-)K_+ + (T'_+\dagger)\Delta(T'_-P_2T_-)T'_+\dagger\|, \\ c_2 &= \|(T'_-\dagger)\Delta(T'_-P_2T_-)T'_-\dagger - P_2T_-BK_+\|, \end{aligned}$$

then

$$(4.21) \quad f(y, z) \leq c_1\|y\|^2 + 2c_2\|y\|\|z\|.$$

As  $\|x\|^2 \leq 2 \max\{\|T'_+\dagger\|^2, \|T'_-\dagger\|^2\}[\|y\|^2 + \|z\|^2]$ , then

$$\begin{aligned} \lambda_1 + \beta\lambda_2 - \beta/(4\gamma)\|x\|^2 &\geq \lambda_1 + \beta\lambda_2 - \frac{1}{2}\beta[\|y\|^2 + \|z\|^2] \\ &\geq \lambda_1 + \beta\lambda_2 - \|y\|^2 - \beta\|z\|^2 + \beta f(y, z) \end{aligned}$$

by (4.21) and (4.19), and (4.18) is established. Now Theorem 3.1 shows that  $\phi(x) = Kx$  is in  $\Phi$ .

**COROLLARY 1.** *If  $\{C^i\} \subseteq R^n$  for  $i \in S \subseteq \{1, 2, \dots, m\}$  in addition to the assumptions of the theorem, then an admissible control exists no matter how large  $C^i$ ,  $i \in S$ , is.*

*Proof.* Existence depends on  $C^i$  only through  $\Gamma$  in (4.9). As  $\Gamma(T'_+ T_+)_ij = \text{tr} \{(T_+ C^i)'(T_+ C^j)\}$  and as  $P^+ C^i = 0$  for  $i \in S$ , then  $\Gamma(T'_+ T_+)$  is independent of  $C^i$ ,  $i \in S$ .

We recall that if  $\{e_i\}_{i \in S}$  spans  $R^n_-$ , then by hypothesis and Lemma 4.2,  $\{D^i\} \subseteq R^n_-$ ,  $i \in S$ , so that  $\Delta_+(I)$  is independent of  $D^i$ ,  $i \in S$ .

**COROLLARY 2.** *If  $R^n_+ \subseteq R_B[A]$ ,  $\{C^i\} \subseteq R^n_-$ ,  $i = 1, \dots, m$ , and if  $R^n_- \subseteq \bigcap_{i=1}^r N[G^i]$  and  $\{D^i\} \subseteq R^n_-$ ,  $i = 1, \dots, n$ , then an admissible control exists.<sup>6</sup>*

*Proof.* As  $R^n_- \subseteq N[H(\omega)]$  by hypothesis then  $D(P^- x) = 0$  so that  $\Delta_-(I) = 0$ . Moreover  $\Gamma(T'_+ T_+) = 0$  and  $\Delta(T'_+ T_+) = 0$  so the conclusion follows.

Let us sum up what this last corollary says. We assume that at least the unstable modes of  $A$  are controllable, that the state and control dependent noise affects only the stable modes and that the state dependent noise is due only to the unstable modes; then an admissible control exists.

Let us now assume that the conditions of Theorem 4.1 are satisfied so that we have an admissible control  $\varphi(x) = Kx$  and a function  $V(x) = x'Px$  such that

$$\mathcal{L}_\varphi V(x) \leq \lambda - \alpha \|x\|^2, \quad \lambda > 0, \quad \alpha > 0.$$

In fact  $\mathcal{L}_\varphi V(x)$  is of the form

$$\mathcal{L}_\varphi V(x) = l(P) + f_\varphi(x, P),$$

where  $l(\cdot)$  and  $f_\varphi(x, \cdot)$  are linear and  $f_\varphi(x, P) \leq -\alpha \|x\|^2$ ,  $\alpha > 0$ . We can choose  $\eta < \infty$  such that

$$\alpha \eta \|x\|^2 \geq L(x, Kx),$$

and so

$$\mathcal{L}_\varphi(x'Qx) \leq l(Q) - L(x, \varphi(x))$$

if  $Q = \eta P$ . Hence

$$0 \geq \Delta(Q) + K'\Gamma(Q)K + (A - BK)'Q + Q(A - BK) + M + K'NK.$$

This condition suffices to make Wonham's argument [1] proving the existence of an optimal control valid.

**THEOREM 4.2.** *Under the assumptions of Theorem 4.1 or of one of its corollaries, an optimal control exists and (3.16), (3.17), (3.18) and (3.19) are valid.*

We point out that the following may be useful with regard to (4.9). Consider the system

$$(4.22) \quad dx/dt = Ax - Bu, \quad x(0) = x$$

and

$$W(x) = \int_0^\infty x(t)' \Delta x(t) + \alpha u(t)' u(t) dt,$$

<sup>6</sup>  $N[A]$  denotes the null space of  $A$ .

where  $(A, B)$  are controllable,  $A, B, \Delta$  are real  $n \times n$ ,  $n \times m$  and  $n \times n$  matrices, respectively,  $\Delta > 0$ , and  $\alpha$  is a positive scalar.

From the theory of the linear regulator it is known that  $\inf_{u=Kx} W(x) = x'P(\alpha)x$ , where  $P(\alpha)$  is the unique positive solution of

$$(4.23) \quad A'P + PA - PBB'P/\alpha + \Delta = 0.$$

It can be shown that  $P(\alpha) - P(\beta) \geq 0$  if  $\alpha > \beta > 0$ , and hence that  $\lim_{\alpha \rightarrow 0} P(\alpha) = P_0$  exists. Then

$$\|P_0\| = \inf_K \left\| \int_0^\infty e^{t(A-BK)'} \Delta e^{t(A-BK)} dt \right\|.$$

LEMMA 4.4.  $\|P_0\| = 0$  if and only if  $\text{rank } B = n$ .

*Proof.* By (4.23),

$$A'P_0 + P_0A - \lim_{\alpha \rightarrow 0} P(\alpha) \frac{BB'}{\alpha} P(\alpha) + \Delta = 0,$$

and so  $P_0BB'P_0 = 0$ . If  $\text{rank } B = n$ , then  $P_0 = 0$ . Conversely if  $P_0 = 0$ , then

$$\lim_{\alpha \rightarrow 0} \frac{P(\alpha)}{\sqrt{\alpha}} BB' \frac{P(\alpha)}{\sqrt{\alpha}} = \Delta.$$

Let  $p_\alpha^i$  be the  $i$ th column vector of  $P(\alpha)$  and let  $q_\alpha^i = (1/\sqrt{\alpha})B'p_\alpha^i$ . Thus  $q_\alpha^i \rightarrow q^i$  and  $(q^i)'q^j = \Delta_{ij}$ . As  $\Delta > 0$  then the  $q^i$  are linearly independent. Hence  $\text{rank } B = n$ .

Now we can replace the condition  $\{D^i\} \subseteq R_-^n$ ,  $i = 1, \dots, n$ , in Corollary 2 of Theorem 4.1 by  $R_+^n \subseteq \{B\}$ .

The method of this section yields the following alternate form of Corollary 2. Suppose  $C = 0$  and  $W$  is a subspace such that  $\{D^i\} \subseteq W \subseteq \bigcup_{j=1}^r N[G^j]$ ,  $i = 1, \dots, n$ . Suppose there is a  $K$  such that  $\hat{A} = A - BK$  has stable modes  $\hat{R}_-^n$  and unstable modes  $\hat{R}_+^n$ , where  $W \subseteq \hat{R}_-^n$ . We also assume that there is a space  $W^c$  complementary to  $W$  such that  $W^c \subseteq R_B[A]$ , and that  $W$  and  $W^c$  completely reduce  $\hat{A}$ . Then an admissible control exists independently of  $D$ . If  $W^c \subseteq \{B\}$ , then  $\{D^i\} \subseteq W$  is not required.

**5. An application to stability.** If we set  $E = 0$  in (2.1) and if  $u = \varphi(x) = Kx$ , then  $x(t) \equiv 0$  is the unique solution of (2.1) with  $x(0) = 0$ ; i.e., 0 is an equilibrium point. In this case the degenerate probability measure  $\mu$  with support at 0 is invariant and  $\mathcal{E}_\mu\{x'x\} = 0$ .

Let us consider the problem of stability of the second mean. We have

$$\begin{aligned} \frac{d}{dt} \mathcal{E}_x\{x_t'Px_t\} &= \mathcal{E}_x\{\mathcal{L}_\varphi(x_t'Px_t)\} \\ (5.1) \quad &= \mathcal{E}_x\{x_t'[\Delta(P) + K'\Gamma(P) + A'P + PA]x_t\} \\ &\leq \alpha\lambda(P)\mathcal{E}_x\{x_t'Px_t\} \end{aligned}$$

if  $P > 0$ , where  $\alpha > 0$  and  $\lambda(P)$  is the largest eigenvalue of  $\Delta(P) + K'\Gamma(P)K + A'P + PA$ . If  $\lambda(P) < 0$ , then there exist  $\beta, \gamma > 0$  such that

$$(5.2) \quad \mathcal{E}_x\{x_t'Px_t\} \leq \beta(x'Px) \exp\{-\gamma t\},$$

and so  $\mathcal{E}\{x_t'Px_t\}$  is globally exponentially stable. It follows that if there is a  $P > 0$  such that  $\lambda(P) < 0$ , then the second mean is exponentially stable; i.e., (5.2) holds with  $P = I$ . On the other hand, Theorem 4.1 and its corollaries give conditions such that  $P > 0$  exists and

$$\mathcal{L}_\phi(x'Px) \leq -\eta x'Px$$

with  $\eta > 0$ . Hence if the conditions of the theorem or one of its corollaries are satisfied, then we have exponential stability of the second mean. This implies [9] that the zero solution is asymptotically stable with probability one. More complete results have been obtained in [10] for systems generated by linear  $n$ th order differential equations whose constant coefficients are distorted by white noise.

As we mentioned, the invariant measure is degenerate, and so the stationary optimal control is of no interest. However, now the problem of minimizing

$$(5.3) \quad \mathcal{E}_x \left\{ \int_0^\infty [x_t'Mx_t + u_t'Nu_t] dt \right\}$$

has a solution under exactly the conditions of Theorem 4.1 (see [1]).

**6. Example.** We consider only one elementary system which does, however, show that the results of § 4 are an improvement over Theorem 3.2. Let

$$A = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}, \quad B = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad C(u) = \begin{bmatrix} 0 \\ \alpha \end{bmatrix} u, \quad D(x) = \begin{bmatrix} \beta \\ \gamma \end{bmatrix} x_1.$$

Then

$$R_-^2 = \left\{ \begin{bmatrix} 0 \\ \lambda \end{bmatrix} : \lambda \in R^1 \right\} \quad \text{and} \quad \{C^1\} \subseteq R_-^2, \quad D(R_-^2) = 0.$$

Also  $(A, B)$  is stabilizable and  $R_+^2 \subseteq \{B\}$ . Hence an optimal control exists for all values of  $\alpha, \beta, \gamma$ .

On the other hand,

$$\Delta(I) = (\beta^2 + \gamma^2) \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \quad \text{and} \quad \Gamma(I) = \alpha^2.$$

Then

$$(6.1) \quad \begin{aligned} \alpha^2 \inf_K Q(K) &\equiv \alpha^2 \inf_K \left\| \int_0^\infty e^{t(A-BK)'} K' K e^{t(A-BK)} dt \right\| \\ &\leq \inf_K \left\| \int_0^\infty e^{t(A-BK)'} [\Delta(I) + \alpha^2 K' K] e^{t(A-BK)} dt \right\|. \end{aligned}$$

Computation shows that  $K = [2, 0]$  minimizes  $Q(K)$  in the class of  $K$ 's such that  $A - BK$  is stable. Also  $Q([2, 0]) = 2$  so that Theorem 3.2 requires at least  $\alpha^2 < \frac{1}{2}$ .

**Appendix.** We shall prove the following lemma.

LEMMA. If  $x(0) = x^n(0) = \xi$  with  $\mathcal{E}\{\|\xi\|^2\} < \infty$ , and if  $T < \infty$ , then

$$(A.1) \quad \lim_{n \rightarrow \infty} \mathcal{E}\{\|x^n(t) - x(t)\|^2\} = 0$$



uniformly in  $t$  for  $0 \leq t < T$ . Here  $x(t)$  satisfies (2.1) and  $x^n(t)$  satisfies (3.6).

*Proof.* We shall write  $x(t)$  as  $x_t$ . Let

$$y_t = \int_0^t D(x_s) d\omega_s, \quad y_t^n = \int_0^t D_n(x_s) d\omega_s.$$

As  $\mathcal{E}\{\|\xi\|^2\} < \infty$ , then  $\mathcal{E}\{\max_{0 \leq t \leq T} \|x_t\|^2\} < \infty$  so

$$\mathcal{E}\{\|y_t - y_t^n\|^2\} = \mathcal{E}\left\{\int_0^t \|D(x_s) - D_n(x_s^n)\|^2 ds\right\}.$$

Now

$$\|D(x_s) - D_n(x_s^n)\|^2 \leq 2\|D(x_s) - D_n(x_s)\|^2 + 2\|D_n(x_s) - D_n(x_s^n)\|^2.$$

If we let  $\mathcal{M}_t^n = [0, t] \cap \{s: \|x_s\| > n\}$ , then

$$\begin{aligned} \mathcal{E}\{\|y_t - y_t^n\|^2\} &\leq 2k^2 \mathcal{E}\left\{\int_{\mathcal{M}_t^n} \|x_s\|^2 ds\right\} + 2k^2 \mathcal{E}\left\{\int_0^t \|x_s - x_s^n\|^2 ds\right\} \\ &= 2k^2(\alpha_1 + \alpha_2) \end{aligned}$$

if we set  $\alpha_1 = \mathcal{E}\{\int_{\mathcal{M}_t^n} \|x_s\|^2 ds\}$  and  $\alpha_2 = \mathcal{E}\{\int_0^t \|x_s^n - x_s\|^2 ds\}$ . The same result holds for  $y_t = \int_0^t C(Kx_s) d\omega_s$  and  $y_t^n = \int_0^t \tilde{C}_n(x_s) d\omega_s$ . Moreover,

$$\mathcal{E}\left\{\left\|\int_0^t (A - BK)(x_s - x_s^n) ds\right\|^2\right\} \leq 2tk^2\alpha_2.$$

It follows that

$$M_t^n \equiv \mathcal{E}\{\|x_t^n - x_t\|^2\} \leq 6tk^2\alpha_2 + 12k^2(\alpha_1 + \alpha_2),$$

i.e.,

$$M_t^n \leq 6k^2(2 + t) \int_0^t M_s^n ds + 12k^2\alpha_1,$$

i.e.,

$$M_t^n \leq 12k^2\alpha_1 \exp\{6k^2(2 + t)t\} \leq c_1(T)\alpha_1.$$

However,

$$\begin{aligned} \alpha_1 &= \mathcal{E}\left\{\int_{\mathcal{M}_t^n} \|x_s\|^2 ds\right\} = \int_0^t \mathcal{E}\{\|x_s\|^2 1_{\|x_s\| > n}(s)\} ds \\ (A.2) \quad &= \int_0^t \int_{\|y\| > n} \|y\|^2 P_s(dy) ds, \end{aligned}$$

where  $P_s(G) = \Pr\{x_s \in G\}$ . As  $\mathcal{E}\{\max_{0 \leq s \leq T} \|x_s\|^2\} < \infty$  then the inner integral of the last term in (A.2) converges to zero uniformly in  $s$ . This establishes (A.1).

**Acknowledgment.** I should like to thank Dr. W. M. Wonham for his suggestions during the course of this work.

## REFERENCES

- [1] W. M. WONHAM, *Optimal stationary control of a linear system with state-dependent noise*, this Journal, 5 (1967), pp. 486–500.
- [2] ———, *Random differential equations in control theory*, Probabilistic Methods in Applied Mathematics, vol. II, A. T. Bharucha-Reid, ed., Academic Press, New York, 1969.
- [3] M. ZAKAI, *A Liapunov criterion for existence of stationary probability distributions for systems with noise*, this Journal, 7 (1969), pp. 390–397.
- [4] E. B. DYNKIN, *Markov Processes*, vol. I, Academic Press, New York, 1965.
- [5] W. M. WONHAM, *Liapunov criteria for weak stochastic stability*, J. Differential Equations, 2 (1966), pp. 195–207.
- [6] ———, *A Liapunov method for the estimation of statistical averages*, *Ibid.*, 2 (1966), pp. 365–377.
- [7] N. N. KRASOVSKII, *Stabilization of systems in which noise is dependent on the value of the control signal*, Engng. Cybernetics, 2 (1965), pp. 94–102.
- [8] W. M. WONHAM, *On pole assignment in multi-input controllable linear systems*, IEEE Trans. Automatic Control, AC-12 (1967), pp. 660–665.
- [9] F. KOZIN, *On almost sure sample asymptotic properties of diffusion processes defined by stochastic differential equations*, J. Math. Kyoto Univ., 4 (1965), pp. 515–528.
- [10] M. B. NEVEL'SON AND R. Z. KHAS'MINSKII, *Stability of a linear system with random disturbances of its parameters*, J. Appl. Math. Mech., 30 (1966), pp. 487–493.

## INPUT-OUTPUT STRUCTURE OF LINEAR SYSTEMS WITH APPLICATION TO THE DECOUPLING PROBLEM\*

L. M. SILVERMAN AND H. J. PAYNE†

**1. Introduction.** Let  $\mathcal{S} = (A, B, C, D)$  be a linear time-invariant system representation described by the equations

$$(1.1) \quad \dot{x}(t) = Ax(t) + Bu(t),$$

$$(1.2) \quad y(t) = Cx(t) + Du(t),$$

where  $x(t) \in \Sigma = R^n$ ,  $u(t) \in R^r$ ,  $y(t) \in R^m$  and the constant matrices  $A$ ,  $B$ ,  $C$  and  $D$  are  $n \times n$ ,  $n \times r$ ,  $m \times n$  and  $m \times r$ , respectively. For each initial state  $x_0 \in \Sigma$ ,  $\mathcal{S}$  defines a mapping  $\mathcal{S}_{x_0}: \mathcal{U} \rightarrow \mathcal{Y}$ , where  $\mathcal{U}$  denotes the input function space over  $[0, \infty)$  and  $\mathcal{Y}$  the corresponding output space. For simplicity,  $\mathcal{U}$  will be taken to be  $\mathcal{L}_r^0$ , the space of  $r$ -vector functions continuous on  $[0, \infty)$ .

In the first half of this paper (§§ 2-5) two basic problems related to the input-output structure of the system  $\mathcal{S}$  are examined. The first is that of explicitly characterizing the functional range of  $\mathcal{S}$ , while the second is that of describing the set of inputs which generate given elements of the range. Of particular interest for the second problem is the set of inputs which will generate the zero output. These problems have been considered by several authors in the single-input, single-output case [1]-[4] with a complete solution being given in [4]. In the multivariable case attention, for the most part, has centered on situations in which  $r = m$  and  $\mathcal{S}$  is one-to-one [5]-[8], with a complete solution being given in [8]. However, for many applications, particularly the decoupling problem, this restriction is quite unrealistic. Hence we consider the general problem here. The basis of our approach to the problem is the structure algorithm developed in [8]. In § 2 this algorithm is described and generalized to the case  $m \neq r$ , and in § 3 several basic properties of the algorithm pertinent to the problems considered are derived. A complete solution to the range and input characterization problems is given in § 4. The main results of this section are Theorems 4.1 and 4.2. Theorem 4.1 gives the necessary and sufficient conditions that a function must satisfy in order that it be in the range of  $\mathcal{S}_{x_0}$  while Theorem 4.2 specifies an "inverse system" representation which can be utilized to generate the set of all inputs  $u$  such that  $w = \mathcal{S}_{x_0}[u]$  for any  $w$  in the range of  $\mathcal{S}_{x_0}$ . Corollaries of these theorems include necessary and sufficient conditions for  $\mathcal{S}_{x_0}$  to be left or right invertible and a measure of the rank of the transfer function matrix of  $\mathcal{S}$  in terms of quantities defined by the structure algorithm. In § 5 the results are specialized to the problem of generating the zero output. Two equivalent characterizations of the set of "zeroing" inputs are given, the first being an open loop specification as the set of outputs of a fixed dynamic system, and the second, a feedback law of the type

---

\* Received by the editors April 2, 1970.

† Department of Electrical Engineering, University of Southern California, Los Angeles, California 90007. This work was supported by the Joint Services Electronics Program (U.S. Army, U.S. Navy, and U.S. Air Force) under Grant AFOSR-69-1622A.

$u = Fx + Gv$ . A complete characterization of the class of all pairs  $(F, G)$  which can be used for this purpose is also given. Application is also made in this section to the problem of input isolation discussed by Wonham and Morse [9]. It is shown that the structure algorithm very simply classifies the set of inputs which can be isolated from the output by state feedback.

In the second part of the paper (§§ 6–9) application of the structure theory developed in the first part is made to the problem of decoupling  $\mathcal{S}$  with respect to a specified partition of the outputs into  $p$  subsets either by static feedback (of the type  $u = Fx + Gv$  for  $F$  and  $G$  constant matrices) or dynamic compensation. The static decoupling problem for the case  $p = m = r$  and  $D = 0$  was first defined by Morgan [10], [11], with a complete solution given by Falb and Wolovich [12], [13], [14] and Gilbert [15]. Partial results were also obtained by Rekasius [16]. The more general static decoupling problem (with  $D = 0$ ) was first defined by Wonham and Morse [9] who also provided solutions for several special cases. The approach of Wonham and Morse differs considerably from that of previous work on the problem with their emphasis being on geometric characterizations. As indicated in their introduction, although they solve two large classes of static decoupling problems their results do not lend themselves directly to computer implementation. The approach taken here is a completely general but algebraic one and is well-suited to computer implementations. In this respect it is similar to Gilbert's [15] procedure when specialized to the case he considers.

The case  $D \neq 0$  is handled with no additional complication and leads to results more general than those of [12] and [15] even for the case  $p = m = r$ .<sup>1</sup> Moreover, for the case  $D = 0$ , a broader class of decoupling problems is solved than previously reported [9]. An outline of §§ 6–8 follows.

The general static decoupling problem is defined in § 6. The definition given is essentially the same as that given in [9] for the case  $D = 0$ . In § 7 an explicit algebraic characterization of the class of all feedback pairs  $(F, G)$  which decouple  $\mathcal{S}$  preserving output controllability is given in terms of matrices defined by the structure algorithm. This characterization consists of a set of (nonlinear) algebraic equations, the set of whose solutions coincides with the set of all decoupling pairs, together with a matrix rank condition which insures output controllability. As yet, a finite algorithm for determining existence of a solution to these equations has not been found. For many important special cases, however, solutions have been found and a number of these are given in § 8. Included in the results of this section are the cases solved (by quite different methods) by Wonham and Morse [9]. A complete solution to the static decoupling problem when  $G$  is constrained to be nonsingular is also given. This case includes the now standard result of Falb and Wolovich [12] and Gilbert [15]. Emphasis throughout this section is on constructive proofs which can be directly implemented on a digital computer.

The dynamic decoupling problem is treated in § 9. This problem was first considered in full generality (for  $D = 0$ ) by Morse and Wonham [18], and a complete solution was given for the case in which state feedback is allowed. A simple constructive solution to this problem is also given here. Moreover, the method used here leads quite directly to a complete solution of the more practical

---

<sup>1</sup> Recently, Morse [17] extended some, but not all, of the results of [9] for the case  $D \neq 0$ .

problem—that of decoupling with dynamic compensation and *output* feedback alone, while preserving output controllability and stability of the closed loop system (see Theorem 9.3). This result completely solves the outstanding “classical” problem of decoupling a system specified by a transfer function matrix which originally motivated Morgan’s state variable approach to the problem (see [11] for references to the early literature on the problem).

**2. The structure algorithm.** Let  $q_0 = \text{rank } D$  and let  $\bar{D}_0$  be the submatrix formed from the first  $q_0$  independent rows of  $D$ . Then there exists an  $m \times m$  nonsingular matrix  $S_0$  such that

$$S_0 D = \begin{bmatrix} \bar{D}_0 \\ 0 \end{bmatrix}.$$

Using  $S_0$  as an output transformation yields a new system representation  $\mathcal{T}_0$  defined by the equations

$$(2.1) \quad \dot{x}(t) = Ax(t) + Bu(t),$$

$$(2.2) \quad y_0(t) = C_0 x(t) + D_0 u(t),$$

where  $y_0(t) = S_0 y(t)$ ,  $C_0 = S_0 C$  and  $D_0 = S_0 D$ . It will prove convenient to partition  $y_0$  and  $C_0$  conformably with  $D_0$  as

$$y_0 = \begin{bmatrix} \bar{y}_0 \\ \tilde{y}_0 \end{bmatrix}, \quad C_0 = \begin{bmatrix} \bar{C}_0 \\ \tilde{C}_0 \end{bmatrix}$$

so that  $\bar{y}_0 = \bar{C}_0 x + \bar{D}_0 u$  and  $\tilde{y}_0 = \tilde{C}_0 x$ . As for the case  $r = m$  (see [8]), the remainder of the sequence  $\mathcal{T}_i$  can be defined inductively. We assume that  $\mathcal{T}_k$  has the form

$$(2.3) \quad \dot{x}(t) = Ax(t) + Bu(t),$$

$$(2.4) \quad y_k(t) = C_k x(t) + D_k u(t),$$

with the partitioning

$$y_k(t) = \begin{bmatrix} \bar{y}_k(t) \\ \tilde{y}_k(t) \end{bmatrix}, \quad C_k = \begin{bmatrix} \bar{C}_k \\ \tilde{C}_k \end{bmatrix}, \quad D_k = \begin{bmatrix} \bar{D}_k \\ 0 \end{bmatrix},$$

where  $\bar{D}_k$  has  $q_k$  rows and rank  $q_k$ ,  $\bar{y}_k$  and  $\bar{C}_k$  have  $q_k$  rows, and  $\tilde{y}_k$  and  $\tilde{C}_k$  have  $m - q_k$  rows. Observe that if

$$M_k = \begin{bmatrix} I_{q_k} & 0 \\ 0 & I_{m-q_k}(d/dt) \end{bmatrix},$$

then

$$M_k y_k(t) = \begin{bmatrix} \bar{C}_k \\ \tilde{C}_k A \end{bmatrix} x(t) + \begin{bmatrix} \bar{D}_k \\ \tilde{C}_k B \end{bmatrix} u(t).$$

Note that  $\tilde{y}_k = \tilde{C}_k x(t)$ , hence  $\tilde{y}_k$  is continuously differentiable so that  $M_k y_k(t)$  always exists. In the expression for  $M_k y_k(t)$ , let  $H_{k+1}$  and  $J_{k+1}$  denote the matrices multiplying  $x(t)$  and  $u(t)$ , respectively, and let  $q_{k+1} = \text{rank } J_{k+1}$ . If  $\bar{D}_{k+1}$  is the

matrix formed from the first  $q_{k+1}$  independent rows of  $J_{k+1}$ , then there exists an  $m \times m$  nonsingular matrix  $S_{k+1}$  such that

$$S_{k+1}J_{k+1} = \begin{bmatrix} \bar{D}_{k+1} \\ 0 \end{bmatrix}.$$

$\mathcal{T}_{k+1}$  is then defined by the equations

$$(2.5) \quad \dot{x}(t) = Ax(t) + Bu(t),$$

$$(2.6) \quad y_{k+1}(t) = C_{k+1}x(t) + D_{k+1}u(t),$$

where  $y_{k+1} = S_{k+1}M_k y_k(t)$ ,  $C_{k+1} = S_{k+1}H_{k+1}$  and  $D_{k+1} = S_{k+1}J_{k+1}$ . It follows from the above that  $y_k(t) = N_k y(t)$ , where<sup>2</sup>

$$(2.7) \quad N_k = \prod_{i=0}^k S_{k-i}M_{k-i-1}, \quad k = 0, 1, \dots \quad (M_{-1} \equiv I),$$

is a sequence of nonsingular matrix differential operators. Furthermore,  $\bar{N}_k$  and  $\tilde{N}_k$  are defined by  $\bar{y}_k = \bar{N}_k y(t)$  and  $\tilde{y}_k = \tilde{N}_k y(t)$ .

It is clear that the matrices  $S_i$  defined above are not unique in general. The following method of constructing the matrices proves to be most convenient and will be utilized in the remainder of the paper.

Let  $S_{k+1}^*$  be the (unique) permutation matrix such that  $S_{k+1}^*J_{k+1}$  has the following structure. Its first  $q_{k+1}$  rows are the first  $q_{k+1}$  independent rows of  $J_{k+1}$  with the relative order maintained, and its last  $m - q_{k+1}$  rows are the remaining rows of  $J_{k+1}$ , also with the relative order maintained. It is clear from the form of  $J_{k+1}$  that  $S_{k+1}^*$  has the form

$$(2.8) \quad S_{k+1}^* = \begin{bmatrix} I_{q_k} & 0 \\ 0 & \bar{R}_{k+1} \\ 0 & \tilde{R}_{k+1} \end{bmatrix},$$

where  $\bar{R}_{k+1}$  has  $q_{k+1} - q_k$  rows and  $\tilde{R}_{k+1}$  has  $m - q_{k+1}$  rows. Furthermore,

$$R_{k+1} = \begin{bmatrix} \bar{R}_{k+1} \\ \tilde{R}_{k+1} \end{bmatrix}$$

is a permutation matrix. Setting

$$(2.9) \quad \bar{D}_{k+1} = \begin{bmatrix} \bar{D}_k \\ \bar{R}_{k+1} \tilde{C}_k B \end{bmatrix}$$

and

$$(2.10) \quad K_{k+1}^* = \tilde{R}_{k+1} \tilde{C}_k B \bar{D}_{k+1}^\dagger,$$

where

$$\bar{D}_{k+1}^\dagger = \bar{D}'_{k+1} (\bar{D}_{k+1} \bar{D}'_{k+1})^{-1},$$

<sup>2</sup>  $\prod_{i=j}^k A_i \equiv A_k A_{k-1} \dots A_j$ . The empty product is taken to be the identity.

we have that  $S_{k+1}$  can be chosen as

$$(2.11) \quad S_{k+1} = \begin{bmatrix} I_{q_{k+1}} & 0 \\ -K_{k+1}^* & I_{m-q_{k+1}} \end{bmatrix} S_{k+1}^*$$

so that

$$(2.12) \quad \bar{D}_{k+1} = \begin{bmatrix} \bar{D}_k \\ \bar{R}_{k+1} \tilde{C}_k B \end{bmatrix}, \quad \bar{C}_{k+1} = \begin{bmatrix} \bar{C}_k \\ \bar{R}_{k+1} \tilde{C}_k A \end{bmatrix}$$

and

$$(2.13) \quad \tilde{C}_{k+1} = \tilde{R}_{k+1} \tilde{C}_k A - K_{k+1}^* \bar{C}_{k+1}.$$

This unique explicit representation of the algorithm will be assumed without further comment in the remainder of the paper.

Note that if  $a$  is the first integer such that  $q_a = q_n$  and

$$(2.14) \quad K_{k+1} = \tilde{R}_{k+1} \tilde{C}_k B \bar{D}_a^\dagger, \quad 0 \leq k \leq n-1,$$

then (2.13) can be expressed as

$$(2.15) \quad \tilde{C}_{k+1} = \tilde{R}_{k+1} \tilde{C}_k A - K_{k+1} \bar{C}_a, \quad 0 \leq k \leq n-1,$$

or

$$(2.16) \quad \tilde{C}_{k+1} = \tilde{R}_{k+1} \tilde{C}_k (A - B \bar{D}_a^\dagger \bar{C}_a), \quad 0 \leq k \leq n-1.$$

This follows from the obvious nesting properties of the sequences  $\bar{C}_k$  and  $\bar{D}_k$ ,  $k = 0, 1, \dots, n$  (see [8]).

**3. Properties of the structure algorithm.** We shall first consider the effect of state feedback on the structure algorithm. Let  $u$  be a feedback law of the form

$$(3.1) \quad u = Fx + v.$$

Then the closed loop system is represented by the quadruple  $(A + BF, B, C + DF, D)$ . If  $P$  is any operator associated with  $\mathcal{S}$  and  $Q$  is the corresponding operator of  $(A + BF, B, C + DF, D)$ , we shall say that  $P \xrightarrow{F} Q$ . For example, it is clear that  $A \xrightarrow{F} A + BF, C \xrightarrow{F} C + DF$  while  $B$  and  $D$  are invariant under this type of feedback.

LEMMA 3.1. *The operators  $N_i$  and the matrices  $D_i$  and  $\tilde{C}_i$  are invariant under feedback of the type (3.1), and*

$$\bar{C}_i \xrightarrow{F} \bar{C}_i + \bar{D}_i F$$

for  $i = 0, 1, \dots$ .

*Proof.* The proof is by induction on the steps of the algorithm. For  $i = 0$ , the result is obvious. Assuming that the result holds for  $i = k$  then we clearly have that

$$H_{k+1} \xrightarrow{F} \begin{bmatrix} \bar{C}_k + \bar{D}_k F \\ \tilde{C}_k A + \tilde{C}_k BF \end{bmatrix} = H_{k+1} + J_{k+1} F$$

and

$$J_{k+1} \xrightarrow{F} \begin{bmatrix} \bar{D}_k \\ \tilde{C}_k B \end{bmatrix} = J_{k+1}.$$

Hence,  $S_{k+1} \xrightarrow{F} S_{k+1}$  which in turn implies that  $N_{k+1} \xrightarrow{F} N_{k+1}$ ,  $D_{k+1} \xrightarrow{F} D_{k+1}$  and

$$C_{k+1} \xrightarrow{F} S_{k+1}(H_{k+1} + J_{k+1}F) = C_{k+1} + D_{k+1}F,$$

which completes the proof.

*Remark.* Lemma 3.1 yields as a special case the  $F$ -invariants observed by Gilbert [15]. His results obtain by setting  $D = 0$  and  $m = 1$ .

Let the matrix  $L_k$  be defined as

$$(3.2) \quad L'_k = [\tilde{C}'_0 : \tilde{C}'_1 : \cdots : \tilde{C}'_{k-1}].$$

The matrices  $L_k$  have many properties in common with observability matrices. In fact,  $L_k$  identifies a subspace onto which the projection of the state of  $\mathcal{S}$  is determined by the output alone. This follows from the relationship

$$(3.3) \quad \tilde{Y}'_k(t) = L_k x(t),$$

where

$$(3.4) \quad \tilde{Y}'_k(t) = [\tilde{y}'_0 : \tilde{y}'_1 : \cdots : \tilde{y}'_{k-1}].$$

LEMMA 3.2. *If  $\text{rank } L_k = \text{rank } L_{k+1}$ , then  $\text{rank } L_i = \text{rank } L_k$  for all  $i > k$ .*

*Proof.* Let  $a$  be the first integer such that  $q_a = q_n$ . Then by construction,  $\bar{D}_i = \bar{D}_a$  and  $\bar{C}_i = \bar{C}_a$  for all  $i \geq a$  and  $\bar{D}_i$  and  $\bar{C}_i$  are leading submatrices of  $\bar{D}_a$  and  $\bar{C}_a$ , respectively, for  $i < a$ . Hence, by setting  $F = -\bar{D}_a^\dagger \bar{C}_a$ , where  $\bar{D}_a^\dagger = \bar{D}'_a(\bar{D}_a \bar{D}'_a)^{-1}$ , it follows from Lemma 3.1 that  $\bar{C}_i \xrightarrow{F} 0$  for all  $i$ . Also by Lemma 3.1,  $L_i \xrightarrow{F} L_i$  so that there is no loss of generality for purposes of the proof in assuming that  $\bar{C}_i = 0$  for all  $i$ . It follows from (2.12) and (2.13) that with  $\bar{C}_k = 0$ ,

$$(3.5) \quad \tilde{C}_k A = R_{k+1}^{-1} \begin{bmatrix} 0 \\ \tilde{C}_{k+1} \end{bmatrix}.$$

Suppose now that  $\text{rank } L_k = \text{rank } L_{k+1}$ . This implies that there exist matrices  $P_i$  such that

$$(3.6) \quad \tilde{C}_k = \sum_{i=0}^{k-1} P_i \tilde{C}_i.$$

By (2.13) with  $\bar{C}_k = 0$  we have

$$(3.7) \quad \tilde{C}_{k+1} = \tilde{R}_{k+1} \tilde{C}_k A = \sum_{i=0}^{k-1} \tilde{R}_{k+1} P_i \tilde{C}_i A.$$

Therefore, by (3.5),

$$(3.8) \quad \tilde{C}_{k+1} = \sum_{i=0}^{k-1} P_i^1 \tilde{C}_{i+1} = \sum_{i=1}^k P_{i-1}^1 \tilde{C}_i,$$

where  $P_i^1$  is an appropriately defined matrix. Substituting (3.6) into (3.8), therefore,



determines matrices  $P_i^2$  such that

$$\tilde{C}_{k+1} = \sum_{i=0}^{k-1} P_i^2 \tilde{C}_i$$

so that  $\text{rank } L_{k+2} = \text{rank } L_k$ . The result for  $i > k + 2$  follows by an obvious induction argument.

**THEOREM 3.1.** *Let  $a$  be the first integer such that  $q_a = q_n$ , let  $\beta$  be the first integer such that  $\text{rank } L_\beta = \text{rank } L_{\beta+1}$  and let  $\gamma = \text{rank } \tilde{C}_0$ . Then*

$$(3.9) \quad a \leq \beta \leq n - \gamma + 1.$$

*Proof.* The upper bound on  $\beta$  follows from the observation that  $L_1$  has rank  $\gamma$  and the fact that

$$\text{rank } L_i - \text{rank } L_{i-1} \geq 1 \quad \text{for } i \leq \beta$$

by Lemma 3.2.

The lower bound on  $\beta$  is clear since  $\text{rank } D_i = \text{rank } D_\beta$  for  $i > \beta$ , again by Lemma 3.2.

*Remark.* Theorem 3.1 provides a stopping rule for the algorithm.

**LEMMA 3.3.** *There exist an integer  $\delta$ ,  $\beta \leq \delta \leq n$ , and matrices  $P_i$ ,  $i = 0, \dots, \delta - 1$ , such that*

$$(3.10) \quad \tilde{C}_\delta = \sum_{i=0}^{\delta-1} P_i \left( \prod_{j=i+1}^{\delta} \tilde{R}_j \right) \tilde{C}_i.$$

*Proof.* As in the proof of Lemma 3.2, it may be assumed without loss of generality that  $\tilde{C}_i = 0$  for all  $i$ . With this assumption  $\tilde{C}_{k+1} = \tilde{R}_{k+1} \tilde{C}_k A$  so that

$$(3.11) \quad \tilde{C}_k = \left( \prod_{i=1}^k \tilde{R}_i \right) \tilde{C}_0 A^k$$

and

$$(3.12) \quad \left( \prod_{i=k+1}^{\delta} \tilde{R}_i \right) \tilde{C}_k = \left( \prod_{i=1}^{\delta} \tilde{R}_i \right) \tilde{C}_0 A^k,$$

where  $\delta$  is as yet undetermined. Let  $W_\delta = \left( \prod_{i=1}^{\delta} \tilde{R}_i \right) \tilde{C}_0$ . Then for a fixed  $\delta$ , it follows by standard observability matrix arguments (of which the proof of Lemma 3.2 is a generalization) that there exists an integer  $\gamma \leq n$  such that  $(W_\delta A^\gamma)'$  is in the range of

$$[W_\delta' \vdots \dots \vdots (W_\delta A^{\gamma-1})'].$$

Hence, there exists a first integer  $\delta \leq n$  for which  $\delta = \gamma$  so that

$$W_\delta A^\delta = \sum_{i=0}^{\delta-1} P_i W_\delta A^i$$

for appropriate matrices  $P_i$ . Consequently, since by (3.11),  $\tilde{C}_\delta = W_\delta A^\delta$ , equation (3.10) follows.

The lower bound on  $\delta$  derives from the observation that (3.10) implies that  $\text{rank } L_{\delta+1} = \text{rank } L_\delta$ .

*Remark.* The integer  $\delta$  is not as simply defined as  $a$  and  $\beta$ . However, only the form of (3.10) is of importance below and  $\delta$  need not be calculated for the main results.

For the case  $q_a = m$  the results of this section can be considerably strengthened.

**THEOREM 3.2.** *If  $q_a = m$ , then*

(i) *the rows of  $L_a$  are linearly independent, and*

(ii)  *$a = \beta = \delta$ .*

*Proof.* The proof of (i) is essentially the same as that of Theorem 3 in [8] and need not be repeated here. That  $\beta = a$  follows from (i) and the fact that if  $q_a = m$ , then  $L_j = L_a$  for  $j \geq a$ . That  $\delta = a$  follows from the fact that  $\tilde{C}_a = 0$  so that (3.10) is satisfied trivially with  $a = \delta$  and  $P_i = 0, i = 0, \dots, a - 1$ .

A final property of the structure algorithm that will be used in the sequel is the following lemma.

**LEMMA 3.4.** *Let  $\lambda \in R^r$ . Then  $\bar{D}_a \lambda = 0$  if and only if  $\begin{bmatrix} D \\ L_\beta B \end{bmatrix} \lambda = 0$ .*

*Proof.* From the structure algorithm,

$$\bar{R}_{k+1} \tilde{C}_k B = [0 : I : 0] \bar{D}_a$$

and

$$\tilde{R}_{k+1} \tilde{C}_k B = K_{k+1} \bar{D}_a$$

so that

$$\tilde{C}_k B = R_{k+1}^{-1} \begin{bmatrix} 0 : I : 0 \\ K_{k+1} \end{bmatrix} \bar{D}_a.$$

Hence  $D_a \lambda = 0$  implies  $L_\beta B \lambda = 0$ . By the nesting property of  $\bar{D}_k$  and the definition of  $D_0$ , one also has  $D \lambda = 0$ .

For sufficiency, if  $\bar{D}_0 \lambda = 0$  and  $\tilde{C}_0 B \lambda = 0$ , then it is clear that  $\bar{D}_1 \lambda = 0$ . Proceeding by induction, suppose  $\bar{D}_k \lambda = 0$ . Then if  $\tilde{C}_k B \lambda = 0$ , it follows that

$$\bar{D}_{k+1} \lambda = \begin{bmatrix} \bar{D}_k \\ \bar{R}_{k+1} \tilde{C}_k B \end{bmatrix} \lambda = 0.$$

This completes the proof.

**4. System range.** In this section it will be shown that the structure algorithm completely determines the functional range of  $\mathcal{S}$  and provides a method for generating elements of the range. Two operators derived from the structure algorithm will be utilized for this purpose:

$$(4.1) \quad M_1 = \left( \frac{d^\delta}{dt^\delta} - \sum_{i=0}^{\delta-1} P_i \frac{d^i}{dt^i} \right) \prod_{j=0}^a \tilde{R}_j,$$

$$(4.2) \quad M_2 = \sum_{j=0}^{\delta} \left( \prod_{l=j+1}^a \tilde{R}_l \right) K_j \frac{d^{\delta-j}}{dt^{\delta-j}} - \sum_{j=0}^{\delta-1} P_j \sum_{k=0}^j \left( \prod_{l=k+1}^a \tilde{R}_l \right) K_k \frac{d^{j-k}}{dt^{j-k}}.$$

Here and below,  $\delta$  and  $P_i, i = 0, \dots, \delta - 1$ , are the quantities defined in Lemma 3.3;  $a$  and  $\beta$  are as defined in Theorem 3.1. The matrices defined in § 2 will be utilized without further comment.

Let  $\mathcal{Y}_k$  be the space of  $m$ -vector functions for which  $N_k y$  is defined and continuous. Then we have the following preliminary result.

LEMMA 4.1. *Let  $y \in \mathcal{Y}_\delta$  and let  $y_j = N_j y, j = 0, \dots, \delta - 1$ . If  $M_1 y(t) \equiv 0$  on  $[0, \infty)$ ,  $\bar{y}_a(t) \equiv 0$  on  $[0, \infty)$  and  $\tilde{y}_j(0) = 0, j = 0, \dots, \beta - 1$ , then  $y(t) \equiv 0$  on  $[0, \infty)$ .*

*Proof.* First we show that  $\bar{y}_a(t) \equiv 0$  on the interval implies  $\tilde{y}_j(t) = (\prod_{i=0}^j \tilde{R}_i) y^{(j)}(t)$ . To see this, we proceed by induction. By the algorithm,

$$\bar{y}_0(t) = \bar{R}_0 y(t)$$

and

$$\tilde{y}_0(t) = \tilde{R}_0 y(t) - K_0^* \bar{y}_0(t).$$

But  $\bar{y}_a(t) \equiv 0$  implies that  $\bar{y}_j(t) \equiv 0, j = 0, 1, \dots, a$ . Therefore,  $\tilde{y}_0(t) = \tilde{R}_0 y(t)$ . Let the desired relation be true for  $j = k - 1$ . Then by the algorithm,

$$\bar{y}_k(t) = \begin{bmatrix} \bar{y}_{k-1}(t) \\ \bar{R}_k \tilde{y}_{k-1}^{(1)}(t) \end{bmatrix}$$

and

$$\tilde{y}_k(t) = \tilde{R}_k \tilde{y}_{k-1}^{(1)}(t) - K_k^* \bar{y}_k(t).$$

Since  $\bar{y}_k(t) \equiv 0, \tilde{y}_k(t) = \tilde{R}_k \tilde{y}_{k-1}^{(1)}(t)$ , which establishes the suggested relation.

From the differential equation  $M_1 y(t) = 0$ , the initial conditions  $\tilde{y}_j(0) = 0, j = 0, 1, \dots, \beta - 1$ , and the result established above, it follows that  $(\prod_{j=0}^a \tilde{R}_j) y(t) \equiv 0$ , hence  $\tilde{y}_a(t) \equiv 0$ . We proceed by induction to show that  $\tilde{y}_k(t) \equiv 0$  for  $k < a$ . Suppose  $\tilde{y}_k(t) \equiv 0$ . Then  $\bar{y}_k(t) = 0$  implies  $\tilde{R}_k \tilde{y}_{k-1}^{(1)}(t) \equiv 0$ , and  $\bar{R}_k \tilde{y}_{k-1}^{(1)}(t) \equiv 0$ . Integrating and employing the initial condition, we have  $\tilde{R}_k \tilde{y}_{k-1}(t) \equiv 0$  and  $\bar{R}_k \tilde{y}_{k-1}(t) \equiv 0$ . Then, since  $[\bar{R}'_k : \tilde{R}'_k]$  is nonsingular,  $\tilde{y}_{k-1}(t) \equiv 0$ . By induction,  $\tilde{y}_0(t) \equiv 0$ . Finally, since  $[\bar{R}'_0 : \tilde{R}'_0]' y(t) = [\bar{y}'_0 : \tilde{y}'_0]', y(t) \equiv 0$ .

The following theorem gives an explicit characterization of the range of a linear system under continuous inputs. A partial result of this nature was first given in [8, Theorem 5], where the history of the problem is also detailed.

THEOREM 4.1 (Range theorem). *An  $m$ -vector function  $w$  defined on  $[0, \infty)$  is in the range of  $\mathcal{L}_{x_0}$  if and only if*

- (i)  $w \in \mathcal{Y}_\delta$ ,
- (ii)  $\tilde{N}_i w(t)|_{t=0} = \tilde{C}_i x_0, i = 0, \dots, \beta - 1$ , and
- (iii)  $(M_1 - M_2 \bar{N}_a) w(t) = 0$  for all  $t \in [0, \infty)$ .

*Proof.* For necessity, observe that if  $w$  is the output of  $\mathcal{L}$  for an initial state  $x_0$  and continuous input  $v$  then

$$(4.3) \quad \bar{N}_k w(t) = \bar{C}_k x(t) + \bar{D}_k v(t),$$

$$(4.4) \quad \tilde{N}_k w(t) = \tilde{C}_k x(t).$$

Conditions (i) and (ii) follow immediately from (4.3) and (4.4). In fact, it is clear that the two conditions are satisfied for all  $i = 0, 1, \dots$ . Setting  $\bar{w}_k(t) = \bar{N}_k w(t)$  and  $\tilde{w}_k(t) = \tilde{N}_k w(t)$ , it follows from the structure algorithm that

$$(4.5) \quad \tilde{w}_k(t) = \tilde{R}_k \tilde{w}_{k-1}^{(1)}(t) - K_k \bar{w}_k(t).$$

Then, iterating, one finds

$$(4.6) \quad \tilde{w}_k(t) = \left( \sum_{j=0}^k \tilde{R}_j \right) w^{(k)}(t) - \sum_{j=0}^k \left( \prod_{l=j+1}^k \tilde{R}_l \right) K_j \bar{w}_a^{(k-j)}(t).$$

Employing Lemma 3.3, and noting that  $\tilde{R}_j = I$  for  $j > a$ , we clearly see that

$$(4.7) \quad \tilde{w}_\delta(t) = \sum_{i=0}^{\delta-1} P_i \left( \prod_{j=i+1}^a \tilde{R}_j \right) \tilde{w}_i(t).$$

Substituting for  $\tilde{w}_\delta(t)$  and  $\tilde{w}_i(t)$  from (4.6), one finds that  $w$  must satisfy condition (iii).

For sufficiency, assume that  $w$  is such that conditions (i)–(iii) hold. Let  $K$  be a matrix whose columns form a basis for the null space of  $\bar{D}_a$ . Then, since  $\bar{D}_a$  has full row rank,

$$(4.8) \quad T = [\bar{D}_a^\dagger : K]$$

is nonsingular. Now let

$$u(t) = \bar{D}_a^\dagger v_1(t) + K v_2(t)$$

and set

$$(4.9) \quad v_1(t) = \bar{w}_a(t) - \bar{C}_a z(t),$$

where  $z$  is the solution of

$$(4.10) \quad \dot{z}(t) = (A - B\bar{D}_a^\dagger \bar{C}_a)z(t) + B\bar{D}_a^\dagger \bar{w}_a(t) + BK v_2(t)$$

with initial condition  $z(0) = x_0$ . It is clear that, for any  $v_2$ , the output of  $\mathcal{L}_{x_0}$  with this input has the property that  $\bar{y}_a(t) \equiv \bar{w}_a(t)$  and  $y$  satisfies properties (i)–(iii). Defining  $r(t) = y(t) - w(t)$  and  $r_j(t) = N_j r(t)$  we then have that  $r(t)$  satisfies (iii), and  $\bar{r}_a(t) \equiv 0$ . Hence,  $M_1 r(t) \equiv 0$  and since  $\tilde{r}_j(0) = 0, j = 0, \dots, \beta - 1$ , Lemma 4.1 implies  $r(t) \equiv 0$  which completes the proof.

*Remark.* In the case  $q_a = m$ , it is easily seen that  $\tilde{R}_a = K_a = 0$  so that  $M_1 = M_2 = 0$ . Thus only conditions (i) and (ii) need be satisfied to insure that  $w$  is in the range of  $\mathcal{L}_{x_0}$ . For the subcase  $q_a = m = r$ , therefore, the range theorem specializes to Theorem 5 of [8]. In general, condition (iii) can be expressed in the form

$$M_1 w(t) = M_2 \bar{w}_a(t).$$

This equation shows that  $w(t)$  in the range of  $\mathcal{L}_{x_0}$  is uniquely determined by  $\bar{w}_a(t)$  and the initial conditions.

**COROLLARY 4.1.** *The range of  $\mathcal{S}$  is invariant under feedback of the type (3.1).*

*Proof.* By Lemma 3.1,  $\mathcal{Y}_i, \tilde{N}_i, \tilde{C}_i$  and  $(M_1 - M_2 \bar{N}_a)$  are all invariant under this type of feedback. Hence, the result is immediate from Theorem 4.1.

The following theorem gives a complete characterization of the class of inputs which will generate a particular function  $w$  in the range of  $\mathcal{L}_{x_0}$ .

**THEOREM 4.2.** *Let  $w$  be in the range of  $\mathcal{L}_{x_0}$ . Then  $\mathcal{L}_{x_0}[u] = w$  if and only if  $u$  is the output of the system*

$$(4.11) \quad \dot{z}(t) = (A - B\bar{D}_a^\dagger \bar{C}_a)z(t) + B\bar{D}_a^\dagger \bar{N}_a w(t) + BK v(t),$$

$$(4.12) \quad u(t) = -\bar{D}_a^\dagger \bar{C}_a z(t) + \bar{D}_a^\dagger \bar{N}_a w(t) + K v(t)$$

with  $z(0) = x_0$  for some  $v$ .

*Proof.* Sufficiency was established in the proof of Theorem 4.1. To establish necessity, let  $u(t) = \bar{D}_a^\dagger v_1(t) + K v_2(t)$  be any input such that  $\mathcal{L}_{x_0}[u] = w$ . Then from the structure algorithm we have

$$\dot{x}(t) = Ax(t) + B\bar{D}_a^\dagger v_1(t) + BK v_2(t)$$

and

$$\bar{N}_a w(t) = \bar{C}_a x(t) + v_1(t)$$

with  $x(0) = x_0$ . Solving the second of these equations for  $v_1$  and substituting into the first yields the pair of equations (4.11), (4.12) with  $z = x$  and  $v = v_2$ .

*Remark.* The system representation (4.11)–(4.12) is a generalization of the inverse system representation discussed in [8] for the case  $r = m = q_a$ . For this case observe that  $K = 0$  so that a unique  $u$  exists satisfying  $\mathcal{L}_{x_0}[u] = w$  for a given  $x_0$  and  $w$ . This is true more generally when a *left inverse* for  $\mathcal{L}_{x_0}$  exists. The conditions for left invertibility are given in the following corollary.

**COROLLARY 4.2.** *There exists a left inverse  $\mathcal{L}_{x_0}^L : \mathcal{Y} \rightarrow \mathcal{L}_r^0$  such that  $\mathcal{L}_{x_0}^L \mathcal{L}_{x_0} = I$ , the identity operator, if and only if  $q_a = r$ .*

*Proof.* When  $q_a = r$ ,  $K = 0$  so that from (4.11)–(4.12) it is clear that the system

$$(4.13) \quad \dot{z}(t) = (A - B\bar{D}_a^\dagger \bar{C}_a)z(t) + B\bar{D}_a^\dagger \bar{N}_a w(t),$$

$$(4.14) \quad u(t) = -\bar{D}_a^\dagger \bar{C}_a z(t) + \bar{D}_a^\dagger \bar{N}_a w(t)$$

with  $z(0) = x_0$  is a left inverse for  $\mathcal{L}_{x_0}$ .

Let  $\mathcal{Y}_{a,x_0} = \{w \in \mathcal{Y}_a : \tilde{N}_i w(t)|_{t=0} = \tilde{C}_i x_0, i = 0, \dots, \beta - 1\}$ . Then from the remarks following Theorem 4.1,  $\mathcal{Y}_{a,x_0}$  is precisely the range of  $\mathcal{L}_{x_0}$  when  $q_a = m$ . Moreover, it is clear by Theorem 4.2 that the system representation (4.13)–(4.14) is a right inverse for  $\mathcal{L}_{x_0}$  (there may be others). Hence, we have the following corollary.

**COROLLARY 4.3.** *There exists a right inverse  $\mathcal{L}_{x_0}^R : \mathcal{Y}_{a,x_0} \rightarrow \mathcal{L}_r^0$  such that  $\mathcal{L}_{x_0} \mathcal{L}_{x_0}^R = I$  if and only if  $q_a = m$ .*

*Remark.* Both the left and right inverse of  $\mathcal{L}_{x_0}$ , when they exist, may be represented in precisely the same way—as a dynamical system,  $(A - B\bar{D}_a^\dagger \bar{C}_a, B\bar{D}_a^\dagger, -\bar{D}_a^\dagger \bar{C}_a, \bar{D}_a^\dagger)$  following a bank of differentiators realizing  $N_a$ . When  $q_a = m = r$  this is seen to be the inverse system representation defined in [8].

It should also be noted that in general a representation of the “inverse system” of lower dynamic order than that of (4.11)–(4.12) can be found. Using the procedure outlined in the proof of Theorem 4 in [8] we can show that there exists a representation equivalent to (4.11)–(4.12) whose dynamic portion has dimension  $n - \text{rank } L_\beta$ .

Let  $H(s)$  be the transfer function matrix of  $(A, B, C, D)$ . Then another generalization of the results of [8] is the following theorem.

**THEOREM 4.3.**  $\text{rank } H(s) = q_a$  for almost all  $s$ .

*Proof.* Let  $\mathcal{L}$  denote the Laplace transform operator. Then it is clear that for zero initial conditions,

$$\mathcal{L}[N_a y(t)] = [C_a(sI - A)^{-1}B + D_a]\hat{u}(s) = \hat{N}_a(s)\hat{y}(s),$$

where  $\hat{N}_a(s)$  results from replacing the operator  $d/dt$  by  $s$  in  $N_a$ ,  $\hat{u}(s) = \mathcal{L}[u]$ , and  $\hat{y}(s) = \mathcal{L}[y]$ . Since  $N_a$  is a nonsingular operator,

$$(4.15) \quad H(s) = \hat{N}_a^{-1}(s)[C_a(sI - A)^{-1}B + D_a].$$

It follows immediately from (4.15) that  $\text{rank } H(s) \geq q_a = \text{rank } \bar{D}_a$  for almost all  $s$ .

Also note that  $\text{rank } H(s) = \text{rank } H(s)T$ , where  $T$  is defined by (4.8) in the proof of Theorem 4.1. It is clear from that proof that there exists  $v_1$  such that for arbitrary  $v_2$ ,

$$H(s)T \begin{bmatrix} \hat{v}_1(s) \\ \hat{v}_2(s) \end{bmatrix} = 0.$$

Hence, for all  $s$  the null space of  $H(s)T$  has dimension at least  $r - q_a$ . Hence

$$\text{rank } H(s) = \text{rank } H(s)T \leq r - (r - q_a) = q_a$$

for all  $s$ , which completes the proof.

**5. Zeroing the output and input isolation.** A problem which is important in many areas of control [2]–[4] and one which is central to the decoupling problem is that of “zeroing the output” of a system either by feedback or open loop control. A complete solution to the open loop problem follows directly from the range theorem and its corollaries.

**THEOREM 5.1.**

(i) *There exists an input  $u \in \mathcal{U}$  such that  $\mathcal{L}_{x_0}[u]$  is identically zero on  $[0, \infty)$  if and only if  $x_0$  is in the null space of  $L_\beta$ .*

(ii) *Let  $x_0$  be in the null space of  $L_\beta$ . Then the output  $y$  of  $\mathcal{S}_{x_0}$  is identically zero on  $[0, \infty)$  if and only if  $\bar{y}_a(t)$  is identically zero on  $[0, \infty)$ .*

(iii) *Let  $x_0$  be in the null space of  $L_\beta$ . Then  $\mathcal{L}_{x_0}[u]$  is identically zero on  $[0, \infty)$  if and only if  $u$  can be expressed as the output of the system representation*

$$(5.1) \quad \dot{z}(t) = (A - B\bar{D}_a^\dagger \bar{C}_a)z(t) + BKv(t),$$

$$(5.2) \quad u(t) = -\bar{D}_a^\dagger \bar{C}_a z(t) + Kv(t)$$

for some  $v$  and  $z(0) = x_0$ .

*Proof.* Part (i) follows immediately from Theorem 4.1. Part (ii) follows from the sufficiency proof of Theorem 4.1. Part (iii) follows from Theorem 4.2.

*Remark.* Let  $\mathcal{Z} = ((A - B\bar{D}_a^\dagger \bar{C}_a), BK, -\bar{D}_a^\dagger \bar{C}_a, K)$ . It is seen from part (iii) of Theorem 5.1 that the system representation  $\mathcal{Z}$  completely characterizes the set of inputs which zero the output of  $\mathcal{S}_{x_0}$  in the sense that any such input must be in the range of  $\mathcal{L}_{x_0}$ . Clearly then, any system representation having the same range as  $\mathcal{L}_{x_0}$  will generate the same set of inputs. In particular, for the case  $x_0 = 0$ ,  $\mathcal{Z}$  can be replaced by any system which is zero state equivalent to it. Hence, by constructing a minimal such equivalent  $(\hat{A}, \hat{B}, \hat{C}, \hat{D})$  we can always find a completely controllable and observable system representation whose range for the zero initial state is the set of all inputs which will zero the output of  $\mathcal{S}_0$ . Moreover, since system range is invariant under state feedback, by Corollary 4.1, we can replace  $(\hat{A}, \hat{B}, \hat{C}, \hat{D})$  by  $(\hat{A} + \hat{B}F, \hat{B}, \hat{C} + \hat{D}F, \hat{D})$  for any  $F$ . Since  $(\hat{A}, \hat{B}, \hat{C}, \hat{D})$  is controllable by construction,  $\hat{A} + \hat{B}F$  can have any distribution of eigenvalues by appropriate choice of  $F$  [19]. We summarize these remarks in the following corollary.

**COROLLARY 5.1.** *There exists a system representation  $\hat{\mathcal{L}} = (\hat{A}, \hat{B}, \hat{C}, \hat{D})$  having any desired set of eigenvalues with the following property:*

$$\mathcal{S}_0[u] = 0 \quad \text{if and only if} \quad u = \hat{\mathcal{L}}_0[v] \quad \text{for some } v.$$

The set of inputs which zero the output of  $\mathcal{S}_0$  also has a feedback representation. It is clear from the form of (5.1)–(5.2) that the same set of inputs can be generated by setting

$$(5.3) \quad u(t) = -\bar{D}_a^\dagger \bar{C}_a x(t) + Kv(t).$$

Hence we have the next corollary.

**COROLLARY 5.2.** *Let  $x_0$  be in the null space of  $L_\beta$ . Then  $\mathcal{S}_{x_0}[u]$  is identically zero on  $[0, \infty)$  if and only if  $u$  can be expressed in the feedback form (5.3) for some  $v$ .*

In the remainder of this section we shall be concerned with zero state properties of  $\mathcal{S}$ . The notation  $(A, B, C, D) \sim (\hat{A}, \hat{B}, \hat{C}, \hat{D})$  will be used if  $\mathcal{S}_0[u] = \hat{\mathcal{S}}_0[u]$  for all  $u$ . If  $\mathcal{S}_0[u] = 0$  for all  $u$ , we say that  $(A, B, C, D) \sim 0$ .

Closely related to the problem of zeroing the output of a system is that of isolating the output from a subset of the inputs via a control law of the type (3.1). More precisely, suppose

$$(5.4) \quad \dot{x} = Ax + Bu + Ew,$$

$$(5.5) \quad y = Cx + Du,$$

where  $w$  is a “disturbance input” to the original system representation (2.1)–(2.2). When does there exist a control law of the type (3.1) such that in the closed loop system the output is unaffected by variations of  $w$ ? Equivalently, when does there exist a matrix  $F$  such that  $(A + BF, E, C + DF) \sim 0$ ?

This problem was first stated in essentially this form (with  $D = 0$ ) by Wonham and Morse [9] who also provided a solution in terms of a recursively defined subspace. It is shown below that the solution of this problem follows quite easily from the range theorem and that the relevant subspace ( $\mathcal{V}$  of [9]) determined by Wonham and Morse is actually the null space of  $L_\beta$ . Moreover, use of the structure algorithm gives an explicit characterization of the entire family of feedback matrices  $F$  which will isolate a given input set.

**THEOREM 5.2.** *There exists a matrix  $F$  such that  $(A + BF, E, C + DF) \sim 0$  if and only if  $L_\beta E = 0$ .*

*Proof.* To establish sufficiency, it is first noted that the inversion algorithm is invariant with respect to all  $E$  in the null space of  $L_\beta$ . The proof of this fact follows by a simple induction argument based on the observation that

$$J_{i+1} = \begin{bmatrix} \bar{D}_i \\ \tilde{C}_i(B + E) \end{bmatrix} = \begin{bmatrix} \bar{D}_i \\ \tilde{C}_i B \end{bmatrix}.$$

Since  $\bar{C}_a \xrightarrow{F} \bar{C}_a + D_a F$  it follows by Corollary 5.2 that if  $F = -\bar{D}_a^\dagger \bar{C}_a$  and  $u = Fx$ , the closed loop system  $(A + BF, E, C + DF) \sim 0$  which completes the sufficiency proof.

Suppose now that the zero state response of

$$\dot{x}(t) = (A + BF)x(t) + Ew(t),$$

$$y(t) = (C + DF)x(t)$$

is identically zero for all  $w$ . Then certainly,  $N_i y(t) \equiv 0$  for all  $w$  and  $x(0) = 0$ . Observe that  $\tilde{N}_0 y(t) = \tilde{C}_0 x(t)$  so that

$$\frac{d}{dt} \tilde{N}_0 y(t) = \tilde{C}_0 (A + BF)x(t) + \tilde{C}_0 Ew(t) \equiv 0$$

which for  $t = 0$  and  $x(0) = 0$  implies  $\tilde{C}_0 E = 0$ , since  $w(0)$  is unconstrained. By Lemma 3.1 it follows that  $\tilde{N}_1 y(t) = C_1 x(t)$ . Suppose now that  $\tilde{C}_j E = 0$  and  $\tilde{N}_{j+1} y(t) = \tilde{C}_{j+1} x(t)$ . Then repeating the preceding arguments leads to  $\tilde{C}_{j+1} E = 0$  and  $\tilde{N}_{j+2} y(t) = \tilde{C}_{j+2} x(t)$ . Consequently it follows by induction that  $\tilde{C}_j E = 0$  for all  $j = 0, 1, \dots$ , which completes the proof.

Since the matrix  $[\bar{D}_a^\dagger : K]$  is nonsingular, any feedback matrix  $F$  can be expressed in the form

$$(5.6) \quad F = \bar{D}_a^\dagger F_1 + KF_2,$$

where  $F_1$  is a  $q_a \times n$  matrix and  $F_2$  is an  $(r - q_a) \times n$  matrix. With this notation, a useful characterization of the class of matrices which will isolate a given disturbance matrix  $E$  can be given.

**THEOREM 5.3.** *Let  $E$  be such that  $L_\beta E = 0$  and let  $F$  be expanded in the form (5.6). Then  $(A + BF, E, C + DF) \sim 0$  if and only if*

$$(5.7) \quad (\bar{C}_a + F_1)(A - B\bar{D}_a^\dagger \bar{C}_a + BKF_2)^j E = 0$$

for  $j = 0, 1, \dots, n - 1$ .

*Proof.* By Theorem 5.1, part (ii), if  $x_0 = 0$ , the output  $y$  of the system  $(A + BF, E, C + DF)$  is identically zero on  $[0, \infty)$  if and only if  $\bar{y}_a(t)$  is identically zero on the interval. Now

$$\bar{y}_a(t) = (\bar{C}_a + \bar{D}_a F)x(t),$$

where

$$\dot{x}(t) = (A + BF)x(t) + Ew(t).$$

Hence  $\bar{y}_a(t) \equiv 0$  for all  $w$  if and only if

$$(5.8) \quad (\bar{C}_a + \bar{D}_a F)(A + BF)^j E = 0, \quad j = 0, \dots, n - 1.$$

With  $F$  represented by (5.6) this equation becomes

$$(5.9) \quad (\bar{C}_a + F_1)(A + BF)^j E = 0, \quad j = 0, \dots, n - 1.$$

It remains to show that (5.7) and (5.9) are equivalent. It is clear they are equivalent for  $j = 0$ .

Suppose (5.9) holds for  $j = 0, 1, \dots, n - 1$  and (5.7) holds for  $j = 0, 1, \dots, k$ , and consider

$$(5.10) \quad (\bar{C}_a + F_1)(A + BF)^{k+1} E = (\bar{C}_a + F_1)(A + BF)(A + BF)^k E.$$

By (5.6),

$$(5.11) \quad BF = B\bar{D}_a^\dagger F_1 + BKF_2,$$

and using (5.9) with  $j = k$  we have

$$(\bar{C}_a + F_1)(A + BF)^{k+1} E = (\bar{C}_a + F_1)(A - B\bar{D}_a^\dagger \bar{C}_a + BKF_2)(A + BF)^k E.$$



Continuing in this way by using (5.9) and (5.11) with  $j = k - 1, k - 2, \dots, 1$ , one finds that

$$(\bar{C}_a + F_1)(A + BF)^{k+1}E = (\bar{C}_a + F_1)(A - B\bar{D}_a^\dagger\bar{C}_a + BKF_2)^{k+1}E.$$

Hence by induction, (5.9) implies (5.7). That (5.9) follows from (5.7) is easily established by reversing the steps of the proof.

When  $E$  is contained in the column range of  $B$ , i.e.,  $E = BG$  for some matrix  $G$ , it follows from Lemma 3.4 that if  $\bar{D}_a G = 0$ , then  $L_\beta E = 0$  (the reverse implication is also true if  $D = 0$ ). This observation together with Theorems 5.2 and 5.3 yields a characterization of the class of all feedback laws of the type

$$(5.12) \quad u = Fx + Gv$$

which will zero the output of  $\mathcal{S}_0$ .

**THEOREM 5.4.** *Let  $F$  be expanded in the form (5.6). Then*

$$(5.13) \quad (A + BF, BG, C + DF, DG) \sim 0$$

*if and only if*

$$(5.14) \quad \bar{D}_a G = 0$$

*and*

$$(5.15) \quad (\bar{C}_a + F_1)(A - B\bar{D}_a^\dagger\bar{C}_a + BKF_2)^j BG = 0$$

*for  $j = 0, 1, \dots, n - 1$ .*

*Remark.* The class of all pairs  $(F, G)$  which zero the output of  $\mathcal{S}_0$  can be determined explicitly from (5.14) and (5.15); each  $G$  satisfying (5.14) and each  $F_2$  yield a set of matrices  $F_1$  which satisfy the linear equations

$$(5.16) \quad (\bar{C}_a + F_1)Q(F_2, G) = 0,$$

where

$$Q(F_2, G) = [BG \ : \ \dots \ : \ (A - B\bar{D}_a^\dagger\bar{C}_a + BKF_2)^{n-1}BG].$$

Note that if  $G = K$ , the range of  $Q(F_2, K)$  is the same as that of  $Q(0, K)$ . Hence, we have the following special case of Theorem 5.4.

**COROLLARY 5.3.**

$$(5.17) \quad (A + BF, GK, C + DF, DK) \sim 0$$

*if and only if*

$$(5.18) \quad (\bar{C}_a + \bar{D}_a F)Q(0, K) = 0.$$

Note that in contrast to (5.16), equation (5.18) is linear in  $F$ .

Since left invertibility of a system is of particular importance in many applications, we shall conclude our general study of input-output properties with a summary of several equivalent characterizations of left invertibility in Theorem 5.5 below. Verification is straightforward and will be left to the reader.<sup>3</sup>

<sup>3</sup>  $\mathcal{R}(B)$  denotes the column range of  $B$  and  $\eta(L_\beta)$  denotes the null space of  $L_\beta$ .

**THEOREM 5.5.** *The following statements are equivalent :*

- (i)  $\mathcal{L}_{x_0}$  is left invertible for all  $x_0$ ;
- (ii)  $\text{rank } H(s) = r$  for almost all  $s$ ;
- (iii)  $\text{rank } \bar{D}_a = r$ ;
- (iv)  $\text{rank} \begin{bmatrix} D \\ L_\beta B \end{bmatrix} = r$ .

*If  $D = 0$ , then the following statements are also equivalent to the above :*

- (v)  $\eta(L_\beta) \cap \mathcal{R}(B) = 0$ ;
- (vi) if  $E \neq 0$  is contained in  $\mathcal{R}(B)$ , then there exists no matrix  $F$  such that  $(A + BF, E, C) \sim 0$ ;
- (vii) if  $x_0 \neq 0$  is contained in  $\mathcal{R}(B)$ , then there exists no input  $u \in \mathcal{U}$  such that  $\mathcal{L}_{x_0}[u] = 0$  on  $[0, \infty)$ .

**6. State feedback decoupling.** Consider a partition of the output of  $(A, B, C, D)$  into  $p$  nonempty subsets of components

$$(6.1) \quad y(t) = \begin{bmatrix} y^1(t) \\ \vdots \\ y^p(t) \end{bmatrix},$$

and let  $m_i$  denote the size of the subvector  $y^i$  ( $0 < m_i$  and  $\sum_{i=1}^p m_i = m$ ). This partition induces a corresponding partition of  $C$  and  $D$ :

$$(6.2) \quad C = \begin{bmatrix} C^1 \\ \vdots \\ C^p \end{bmatrix}, \quad D = \begin{bmatrix} D^1 \\ \vdots \\ D^p \end{bmatrix},$$

where  $C^i$  and  $D^i$  each have  $m_i$  rows.

The basic decoupling problem we shall consider is that of finding a control law of the type

$$u(t) = Fx(t) + Gv(t)$$

(denoted as  $(F, G)$ ), where

$$v(t) = \begin{bmatrix} v^1(t) \\ \vdots \\ v^p(t) \end{bmatrix}$$

is such that the  $i$ th input set  $v^i$  affects only the  $i$ th output set  $y^i$ . Let  $G$  be partitioned conformably with  $v$  as

$$(6.3) \quad G = [G^1 \vdots \dots \vdots G^p].$$

Then a more formal definition is the following.

**DEFINITION 6.1.** The feedback pair  $(F, G)$  decouples  $(A, B, C, D)$  (relative to the output partition (6.2)) if and only if

$$(6.4) \quad (A + BF, BG^i, C^j + D^jF, D^jG^i) \sim 0, \quad j \neq i.$$

The state feedback decoupling problem then consists of finding a pair  $(F, G)$  satisfying (6.4) for a specified output partition. It is clear that the pair  $(F, 0)$  will always decouple  $(A, B, C, D)$  so that without any further constraints the state feedback decoupling problem always has a solution. However, if the ultimate purpose of decoupling is to control the output of the original system in some non-trivial way, additional constraints are necessary. One such constraint is output controllability [20] of the closed loop system.

Defining

$$(6.5) \quad Q(F, G) = [BG \vdots \cdots \vdots (A + BF)^{n-1}BG]$$

and

$$(6.6) \quad P(C, D, F, G) = [(C + DF)Q(F, G) \vdots DG],$$

we have that the well-known necessary and sufficient condition for output controllability of  $(A + BF, BG, C + DF, DG)$  is (see [20])

$$(6.7) \quad \text{rank } P(C, D, F, G) = m.$$

If the closed loop system is decoupled, then this condition is clearly equivalent to saying that

$$(6.8) \quad \text{rank } P(C^i, D^i, F, G^i) = m_i, \quad i = 1, \dots, p.$$

**DEFINITION 6.2.** The feedback pair  $(F, G)$  output controllably decouples  $(A, B, C, D)$  if and only if (6.4) and (6.8) hold.

For the case  $D = 0$  the definition of the decoupling problem is equivalent to that of Wonham and Morse [9]. Other constraints on the decoupled system may also be desired such as the state controllability, observability, stability, etc., but the output controllability constraint appears to be the weakest meaningful one.

In the following sections we shall develop criteria and algorithms for decoupling based on the results of §§ 4 and 5. To do this, however, some new notation is required. Let  $C$  and  $D$  be partitioned as in (6.2). Then  $\Gamma^i$  is defined to be the matrix formed by deleting  $C^i$  from  $C$  and  $\Delta^i$  the matrix formed by deleting  $D^i$  from  $D$ . Also with  $G$  partitioned as in (6.3), define  $\Omega^i$  to be the matrix formed by deleting  $G^i$  from  $G$ .

We shall have occasion to perform the structure algorithm described in § 2 on the subsystems  $(A, B, C^i, D^i)$  and  $(A, B, \Gamma^i, \Delta^i)$ . The matrices appearing at the  $j$ th stage of the algorithm for  $(A, B, C^i, D^i)$  will be denoted as

$$(6.9) \quad C_j^i = \begin{bmatrix} \bar{C}_j^i \\ \dots \\ \check{C}_j^i \end{bmatrix}, \quad D_j^i = \begin{bmatrix} \bar{D}_j^i \\ \dots \\ 0 \end{bmatrix},$$

where  $q_j^i = \text{rank } \bar{D}_j^i$  is equal to the number of rows in  $\bar{D}_j^i$  and  $\bar{C}_j^i$ . Also, define  $a_i$  to be the first integer such that  $\text{rank } \bar{D}_{a_i}^i = \text{rank } \bar{D}_n^i$ . For notational simplicity, we write  $\bar{D}_{a_i}^i = \bar{D}_{a_i}$  and  $C_{a_i}^i = C_{a_i}$ .

Similarly, for  $(A, B, \Gamma^i, \Delta^i)$  the appropriate matrices appearing at the  $j$ th step of the algorithm will be denoted as

$$(6.10) \quad \Gamma_j^i = \begin{bmatrix} \bar{\Gamma}_j^i \\ \dots \\ \check{\Gamma}_j^i \end{bmatrix}, \quad \Delta_j^i = \begin{bmatrix} \bar{\Delta}_j^i \\ \dots \\ 0 \end{bmatrix},$$

where  $\xi_j^i = \text{rank } \bar{\Delta}_j^i$  is equal to the number of rows in  $\bar{\Delta}_j^i$  and  $\bar{\Gamma}_j^i$ . Furthermore,  $\alpha_i$  is defined to be the first integer such that  $\text{rank } \Delta_{\alpha_i}^i = \text{rank } \Delta_n^i$ . Again for simplicity, write  $\bar{\Delta}_{\alpha_i}^i = \bar{\Delta}_{\alpha_i}$  and  $\bar{\Gamma}_{\alpha_i}^i = \bar{\Gamma}_{\alpha_i}$ .

With the above notation, observe that condition (6.4) can be restated in two equivalent ways:

$$(6.11) \quad (A + BF, BG^i, \Gamma^i + \Delta^i F, \Delta^i G^i) \sim 0, \quad i = 1, \dots, p,$$

$$(6.12) \quad (A + BF, B\Omega^i, C^i + D^i F, D^i \Omega^i) \sim 0, \quad i = 1, \dots, p,$$

Furthermore, the output controllability constraint (6.8) has the equivalent formulation

$$(6.13) \quad \text{rank } P(\Gamma^i, \Delta^i, F, \Omega^i) = m - m_i, \quad i = 1, \dots, p.$$

**7. Decoupling characterizations.** In this section we shall provide two alternate characterizations of the class of state feedback decoupling pairs. It will be shown in the following section that both characterizations are useful in that one leads to necessary conditions and the other to sufficient conditions for decoupling.

First, define  $\mathcal{F}(G^i)$  to be the family of feedback matrices  $F$  such that the pair  $(F, G^i)$  zeros all but the  $i$ th set of outputs, i.e.,

$$(7.1) \quad \mathcal{F}(G^i) = \{F : (A + BF, BG^i, \Gamma^i + \Delta^i F, \Delta^i G^i) \sim 0\}.$$

A complete characterization of  $\mathcal{F}(G^i)$  is provided by Theorem 5.4.

Following the development of § 5 we first note that any feedback matrix  $F$  can be uniquely represented in the form

$$(7.2) \quad F = \bar{\Delta}_{\alpha_i}^\dagger F_{i1} + K_i F_{i2},$$

where  $K_i$  is a fixed matrix whose columns form a basis for the right null space of  $\Delta_{\alpha_i}$ . Further, define

$$(7.3) \quad Q_i(F_{i2}, G^i) = [BG^i \vdots \dots \vdots (A - B\bar{\Delta}_{\alpha_i}^\dagger \bar{\Gamma}_{\alpha_i} + BK_i F_{i2})^{n-1} BG^i].$$

Then the following is an immediate corollary of Theorem 5.4.

**COROLLARY 7.1.**  $F \in \mathcal{F}(G^i)$  if and only if

$$(7.4) \quad (\bar{\Gamma}_{\alpha_i} + F_{i1})Q_i(F_{i2}, G^i) = 0.$$

Similarly, denoting by  $\mathcal{F}^*(\Omega^i)$  the family of all matrices  $F$  for which

$$(A + BF, \Omega^i, C^i + D^i F, D^i \Omega^i) \sim 0$$

and representing  $F$  as

$$(7.5) \quad F = \bar{D}_{\alpha_i}^\dagger F_{i1}^* + K_i^* F_{i2}^*$$

where  $K_i^*$  is a fixed matrix whose columns form a basis for the right null space of  $\bar{D}_{\alpha_i}$ , we may define

$$(7.6) \quad Q_i^*(F_{i2}^*, \Omega^i) = [B\Omega^i \vdots \dots \vdots (A - B\bar{D}_{\alpha_i}^\dagger \bar{C}_{\alpha_i} + BK_i^* F_{i2}^*)^{n-1} B\Omega^i].$$

Again by Theorem 5.4 we have the following corollary.

**COROLLARY 7.2.**  $F \in \mathcal{F}^*(\Omega^i)$  if and only if

$$(7.7) \quad (\bar{C}_{\alpha_i} + F_{i1}^*)Q_i^*(F_{i2}^*, \Omega^i) = 0.$$

The output controllability criterion (6.8) admits of considerable simplification when  $F$  is restricted to  $\mathcal{F}(G^i)$ . Let

$$(7.8) \quad P_i(F_{i2}, G^i) = [(C^i - D^i \bar{\Delta}_{\alpha_i}^\dagger \bar{\Gamma}_{\alpha_i} + D^i K_i F_{i2}) Q_i(F_{i2}, G^i) \vdots D^i G^i].$$

LEMMA 7.1. *If  $F \in \mathcal{F}(G^i)$ , then*

$$(7.9) \quad Q(F, G^i) = Q_i(F_{i2}, G^i),$$

$$(7.10) \quad P(C^i, D^i, F, G^i) = P_i(F_{i2}, G^i).$$

*Proof.* Let  $F_{i0} = \bar{\Gamma}_{\alpha_i} + F_{i1}$ . Then by direct substitution,

$$(A + BF)^j B G^i = (A - B \bar{\Delta}_{\alpha_i}^\dagger \bar{\Gamma}_{\alpha_i} + B K_i F_{i2} + B \bar{\Delta}_{\alpha_i}^\dagger \bar{\Gamma}_{\alpha_i} F_{i0})^j B G^i.$$

Hence, by (7.4) and a straightforward induction argument it follows that

$$(A + BF)^j B G^i = (A - B \bar{\Delta}_{\alpha_i}^\dagger \bar{\Gamma}_{\alpha_i} + B K_i F_{i2})^j B G^i, \quad j = 0, 1, \dots,$$

which establishes (7.9). Equation (7.10) follows by a similar argument.

Similarly, we define

$$(7.11) \quad P_i^*(F_{i2}^*, \Omega^i) = [(\Gamma^i - \Delta^i \bar{D}_{\alpha_i}^\dagger \bar{C}_{\alpha_i} + \Delta^i K_i^* F_{i2}^*) \Omega^i \vdots \Delta^i \Omega^i].$$

A parallel proof to the above establishes the next lemma.

LEMMA 7.2. *If  $F \in \mathcal{F}^*(\Omega^i)$ , then*

$$(7.12) \quad Q^*(F, \Omega^i) = Q_i^*(F_{i2}^*, \Omega^i),$$

$$(7.13) \quad P(\Gamma^i, \Delta^i, F, \Omega^i) = P_i^*(F_{i2}^*, \Omega^i).$$

The first major characterization of decoupling pairs is given by the following theorem.

THEOREM 7.1. *The pair  $(F, G)$  output controllably decouples  $(A, B, C, D)$  if and only if*

$$(7.14) \quad \bar{\Delta}_{\alpha_i} G^i = 0, \quad i = 1, \dots, p,$$

$$(7.15) \quad (\bar{\Gamma}_{\alpha_i} + F_{i1}) Q_i(F_{i2}, G^i) = 0, \quad i = 1, \dots, p,$$

$$(7.16) \quad \text{rank } P_i(F_{i2}, G^i) = m_i, \quad i = 1, \dots, p.$$

*Proof.* It follows by Theorem 5.4 and Corollary 7.1 that (7.14) and (7.15) are necessary and sufficient for  $(F, G)$  to decouple  $(A, B, C, D)$ . In light of Lemma 7.1, equation (7.16) is just a restatement of the output controllability criterion.

By a parallel proof, we also have the following characterization of decoupling pairs.

THEOREM 7.2. *The pair  $(F, G)$  output controllably decouples  $(A, B, C, D)$  if and only if*

$$(7.17) \quad \bar{D}_{\alpha_i} \Omega^i = 0, \quad i = 1, \dots, p,$$

$$(7.18) \quad (\bar{C}_{\alpha_i} + F_{i1}^*) Q_i^*(F_{i2}^*, \Omega^i) = 0, \quad i = 1, \dots, p,$$

$$(7.19) \quad \text{rank } P_i^*(F_{i2}^*, \Omega^i) = m - m_i, \quad i = 1, \dots, p.$$

A considerable simplification of the decoupling characterization results if attention is restricted to the class of decoupling pairs for which  $G$  is square and nonsingular. Such decoupling pairs which preserve the original number of independent inputs to  $(A, B, C, D)$  can be justified in many problems either from physical or mathematical constraints. An example of the latter is the important special case in which the number of inputs to  $(A, B, C, D)$  equals the number of output partitions (in particular, when  $D = 0$ ,  $m_i = 1$  and  $m = r$  we have the problem of Falb and Wolovich). The simplification obtained when  $G$  is nonsingular is that output controllability is always preserved independently of  $F$ .

**LEMMA 7.3.** *If  $G$  is a square nonsingular matrix, then  $(A, B, C, D)$  is output controllable if and only if  $(A + BF, BF, C + DF, DG)$  is output controllable for any  $F$ .*

The proof of this lemma is essentially the same as for the well-known results for state controllability [2], and will be omitted.

Following immediately from Lemma 7.3 and Theorems 7.1 and 7.2, respectively, are the two simplified characterizations.

**THEOREM 7.3.** *The pair  $(F, G)$ , with  $G$  square and nonsingular, output controllably decouples  $(A, B, C, D)$  if and only if (7.14) and (7.15) hold and  $(A, B, C, D)$  is output controllable.*

**THEOREM 7.4.** *The pair  $(F, G)$ , with  $G$  square and nonsingular, output controllably decouples  $(A, B, C, D)$  if and only if (7.17) and (7.18) hold and  $(A, B, C, D)$  is output controllable.*

**8. Decoupling criteria.** In this section, the decoupling characterizations as expressed in Theorems 7.1–7.4 of the previous section will be employed to develop explicit criteria for the existence of decoupling pairs  $(F, G)$ . The results presented, while not representing a complete resolution of this question, do include within them all known previous results [9]–[15], as well as several important extensions.

In general, the determination of decoupling pairs  $(F, G)$  is difficult since  $F$  and  $G$  are interrelated in a highly complicated way through the equations (e.g., (7.14)–(7.16)) characterizing output controllable decoupling. With additional constraints on  $G$ , however, the problem becomes quite tractable.

The first group of results in this section derive from the choice  $G = K$ , where

$$(8.1) \quad K = [K_1 \ \cdots \ K_p]$$

with the  $K_i$  as defined in the previous section. Theorems 8.3–8.5 present sufficient conditions for the existence of a matrix  $F$  to complete a decoupling pair of the form  $(F, K)$ .

The remaining results of this section pertain to those situations in which  $G$  is chosen to be (or must be) nonsingular. Necessary and sufficient conditions are derived for a decoupling pair of this type to exist.

The column space of a matrix  $A$  will be denoted as  $\mathcal{R}[A]$ . Let  $\bar{Q}_i$  be a matrix whose columns form a basis for  $\mathcal{R}[Q_i(0, K_i)]$  and let  $\bar{Q}_i^*$  be a matrix whose columns form a basis for  $\mathcal{R}[Q_i^*(0, K_i^*)]$ .  $\bar{P}_i$  and  $\bar{P}_i^*$  are defined similarly.

The following two lemmas are basic to the development of the necessary condition expressed in Theorems 8.1 and 8.2 and the sufficient conditions expressed in Theorems 8.3–8.5.

LEMMA 8.1. *If  $\bar{\Delta}_{\alpha_i} G^i = 0$ , then*

$$(8.2) \quad \mathcal{R}[Q_i(F_{i2}, G^i)] \subseteq \mathcal{R}[\bar{Q}_i],$$

$$(8.3) \quad \mathcal{R}[P_i(F_{i2}, G^i)] \subseteq \mathcal{R}[\bar{P}_i].$$

*These relationships hold with equality if  $G^i = K_i$ .*

*Proof.* To establish (8.2), it need only be shown that there exists a matrix  $L$  such that  $Q_i(F_{i2}, G^i) = Q_i(0, K_i)L$ . Since  $\bar{\Delta}_{\alpha_i} G^i = 0$ ,  $G^i = K_i M$  for some matrix  $M$ . Then  $Q_i(F_{i2}, G^i)$  consists of blocks of the form  $(A - B\bar{\Delta}_{\alpha_i}^\dagger \bar{\Gamma}_{\alpha_i} + BK_i F_{i2})^k BK_i M$ ,  $k = 0, 1, \dots, n - 1$ , so that it is sufficient to establish the existence of a matrix  $L$  (not the same for each block) satisfying a similar relationship for each block. Now, if each such block is expanded, each resultant term is of the form  $BK_i N$  or  $(A - B\bar{\Delta}_{\alpha_i}^\dagger \bar{\Gamma}_{\alpha_i})^j N$ ,  $j \leq k$ , for some matrix  $N$  (not the same for each term). Hence each term is expressible in the form  $Q_i(0, K_i)P$  for some matrix  $P$  (not the same for each term), which establishes (8.2). That equality holds when  $G^i = K_i$  follows from the invariance of the controllability space of the pair  $(A - B\bar{\Delta}_{\alpha_i}^\dagger \bar{\Gamma}_{\alpha_i}, BK_i)$  under feedback.

To establish (8.3), define matrices  $M$  and  $N$  by  $G^i = K_i M$ ,  $Q_i(F_{i2}, G^i) = Q_i(0, K_i)N$ . Then, employing an obvious partitioning,

$$P_i(F_{i2}, G^i) = P_i(0, K_i) \begin{bmatrix} N & 0 \\ F_{i2} Q_i(F_{i2}, G^i) & M \end{bmatrix}.$$

The condition of equality when  $G^i = K_i$  follows from invariance of controllability matrices under feedback.

A parallel form of this lemma is established in an identical fashion.

LEMMA 8.2. *If  $\bar{D}_{\alpha_i} \Omega^i = 0$ , then*

$$(8.4) \quad \mathcal{R}[Q_i^*(F_{i2}^*, \Omega^i)] \subseteq \mathcal{R}[\bar{Q}_i^*],$$

$$(8.5) \quad \mathcal{R}[P_i^*(F_{i2}^*, \Omega^i)] \subseteq \mathcal{R}[\bar{P}_i^*].$$

*These relationships hold with equality if  $\Omega^i = K_i^*$ .*

A parallel can be drawn here with the work of Wonham and Morse [9]. It is clear from the various definitions that for the case  $D = 0$ ,  $\mathcal{R}[\bar{Q}_i]$  is the maximum dimension controllability subspace which is simultaneously in the null spaces of  $C^j$ ,  $j \neq i$ . For the more general situation it is clear that  $\mathcal{R}[\bar{P}_i]$  is the largest output controllability subspace in  $R^{m_i}$  which can be generated while simultaneously zeroing  $y^j(t)$ ,  $j \neq i$ . These remarks lead to the following necessary condition.

THEOREM 8.1. *There exists a pair  $(F, G)$  which output controllably decouples  $(A, B, C, D)$  only if*

$$(8.6) \quad \text{rank } \bar{P}_i = m_i, \quad i = 1, \dots, p.$$

A similar argument establishes the following theorem.

THEOREM 8.2. *There exists a pair  $(F, G)$  which output controllably decouples  $(A, B, C, D)$  only if*

$$(8.7) \quad \text{rank } \bar{P}_i^* = m - m_i, \quad i = 1, \dots, p.$$

Theorem 8.1 provides a useful point of departure for the development of sufficient conditions for decoupling. For if (8.6) is satisfied, output controllability

is guaranteed by the choice  $G = K$  and, as will be seen, the further conditions to be imposed on  $F$  are considerably simplified. On the other hand, a parallel approach starting with the choice  $\Omega^i = K_i^*$  is not in general possible since there may not exist a  $G$  to simultaneously realize all the  $K_i^*$ .

As suggested, considerable simplification is achieved if attention is restricted to the choice  $G = K$ .

It follows from Theorem 7.1 and Lemma 8.1 that if condition (8.6) is satisfied, the pair  $(F, K)$  output controllably decouples  $\mathcal{S}$  if and only if the equations

$$(8.8) \quad (\bar{\Gamma}_{\alpha_i} + \bar{\Delta}_{\alpha_i}F)\bar{Q}_i = 0, \quad i = 1, \dots, p,$$

have a solution for  $F$ . (We have used the fact that  $F_{i1} = \bar{\Delta}_{\alpha_i}F$ .) Since  $\bar{Q}_i$  does not depend on  $F$ , this set of equations is linear in the elements of  $F$ . Hence standard techniques can be utilized to determine if a solution exists. If no solution exists, however, it cannot be concluded that a solution to the original decoupling problem does not exist, since a solution may still be possible with matrices  $G_i$  whose columns do not span the null space of  $\bar{\Delta}_{\alpha_i}$ .

Let

$$(8.9) \quad \hat{\Gamma} = \begin{bmatrix} \bar{\Gamma}_{\alpha_1} \\ \vdots \\ \bar{\Gamma}_{\alpha_p} \end{bmatrix}, \quad \hat{\Delta} = \begin{bmatrix} \bar{\Delta}_{\alpha_1} \\ \vdots \\ \bar{\Delta}_{\alpha_p} \end{bmatrix}.$$

Then a special case for which a solution to (8.8) exists is given in the next theorem.

**THEOREM 8.3.** *If rank  $P_i = m_i$ ,  $i = 1, \dots, p$ , and  $\mathcal{R}[\hat{\Gamma}] \subset \mathcal{R}[\hat{\Delta}]$ , then there exists a pair  $(F, K)$  which output controllably decouples  $(A, B, C, D)$ .*

*Proof.* Since  $\mathcal{R}[\hat{\Gamma}] \subset \mathcal{R}[\hat{\Delta}]$ , there exists a matrix  $F$  such that  $\hat{\Delta}F = -\hat{\Gamma}$ . This choice satisfies (8.8) which completes the proof.

Although it is not possible to develop useful sufficient conditions along similar lines by considering the choice  $\Omega^i = K_i^*$ , the following lemma provides a link so that the characterization of Theorem 7.2 may be employed for a useful parallel result.

For a fixed partition of  $G$  in the form (6.3), we have the next lemma.

**LEMMA 8.3.** *If  $\bar{\Delta}_{\alpha_i}G^i = 0$ ,  $i = 1, \dots, p$ , then  $\bar{D}_{\alpha_i}\Omega^i = 0$ ,  $i = 1, \dots, p$ .*

*Proof.* From the structure algorithm developed in § 2, it is clear that for any  $i$ ,  $i = 1, \dots, p$ , the rows of  $\bar{D}_{\alpha_i}$  can be expressed as a linear combination of the rows of  $\bar{\Delta}_{\alpha_j}$  for any  $j \neq i$ . Hence, if  $\bar{\Delta}_{\alpha_j}G^j = 0$ ,  $j \neq i$ ,  $\bar{D}_{\alpha_i}G^j = 0$ ,  $j \neq i$ , i.e.,  $\bar{D}_{\alpha_i}\Omega^i = 0$ .

Lemma 8.3 allows one to consider the choice  $G = K$  in connection with the first two conditions, (7.17) and (7.18), of Theorem 7.2. The satisfaction of these two guarantee decoupling, while (8.6) guarantees output controllability. Using Lemmas 8.2 and 8.3, one now sees that it is sufficient for decoupling to require that a solution for  $F$  exist for the set of equations

$$(8.10) \quad (\bar{C}_{\alpha_i} + \bar{D}_{\alpha_i}F)\bar{Q}_i^* = 0, \quad i = 1, \dots, p,$$

with  $F$  given by (7.5).



Let

$$(8.11) \quad \hat{C} = \begin{bmatrix} \bar{C}_{a_1} \\ \vdots \\ \bar{C}_{a_p} \end{bmatrix}, \quad \hat{D} = \begin{bmatrix} \bar{D}_{a_1} \\ \vdots \\ \bar{D}_{a_p} \end{bmatrix}.$$

Then, by a proof analogous to that for Theorem 8.3, we have the following theorem.

**THEOREM 8.4.** *If rank  $\bar{P}_i = m_i$ ,  $i = 1, \dots, p$ , and  $\mathcal{R}[\hat{C}] \subset \mathcal{R}[\hat{D}]$ , then there exists a pair  $(F, G)$  which output controllably decouples  $(A, B, C, D)$ .*

A special case of this result is the following.

**COROLLARY 8.1.** *If rank  $\bar{P}_i = m_i$ ,  $i = 1, \dots, p$ , and the rows of  $\hat{D}$  are linearly independent, then there exists a pair  $(F, G)$  which output controllably decouples  $(A, B, C, D)$ .*

The sufficiency conditions expressed in this corollary, while only a specialization of the conditions of Theorem 8.3, will be seen to be necessary as well if it is required that  $G$  be nonsingular.

Under the conditions of Corollary 8.1, a class of decoupling pairs can be explicitly displayed. The independence of the rows of  $\hat{D}$  implies the existence of a right inverse,  $\hat{D}^\dagger$ . Hence, the pair  $(-\hat{D}^\dagger \hat{C}, K)$  output controllably decouples  $(A, B, C, D)$ .

Let

$$(8.12) \quad \bar{Q} = [\bar{Q}_1 \ \vdots \ \dots \ \vdots \ \bar{Q}_p].$$

Then another special case for which a solution to (8.8) is guaranteed is given in the next theorem.

**THEOREM 8.5.** *If rank  $\bar{P}_i = m_i$ ,  $i = 1, \dots, p$ , and the columns of  $\bar{Q}$  are linearly independent, then there exists a pair  $(F, G)$  which output controllably decouples  $(A, B, C, D)$ . Moreover, under these conditions, the pair  $(F_1, K)$ , with*

$$(8.13) \quad F_1 = -[\bar{\Delta}_{a_1}^\dagger \bar{\Gamma}_{a_1} \bar{Q}_1 \ \vdots \ \dots \ \vdots \ \bar{\Delta}_{a_p}^\dagger \bar{\Gamma}_{a_p} \bar{Q}_p] \bar{Q}^\dagger,$$

output controllably decouples  $(A, B, C, D)$ .

*Proof.* Equation (8.8) can be rewritten as

$$\bar{\Delta}_{\alpha_i} (F \bar{Q}_i) = -\bar{\Gamma}_{\alpha_i} \bar{Q}_i, \quad i = 1, \dots, p.$$

Since the rows of  $\bar{\Delta}_{\alpha_i}$  are linearly independent, a solution for  $F \bar{Q}_i$  for each of these equations exists. Specifically,

$$F \bar{Q}_i = -\bar{\Delta}_{\alpha_i}^\dagger \bar{\Gamma}_{\alpha_i} \bar{Q}_i, \quad i = 1, \dots, p.$$

Collecting these equations, we have

$$F \bar{Q} = -[\bar{\Delta}_{\alpha_1}^\dagger \bar{\Gamma}_{\alpha_1} \bar{Q}_1 \ \vdots \ \dots \ \vdots \ \bar{\Delta}_{\alpha_p}^\dagger \bar{\Gamma}_{\alpha_p} \bar{Q}_p].$$

Since the columns of  $\bar{Q}$  are linearly independent,  $F = F_1$  is a solution of (8.8) which, together with  $G = K$ , output controllably decouples  $(A, B, C, D)$ .

*Remark.* The general solution of (8.8) is obtained by noting that the complete solution for  $F \bar{Q}$  is

$$F \bar{Q} = F_1 \bar{Q} + K F_2,$$

where  $F_2 = \text{diag}(F_{i2})$  and  $F_{i2}$ ,  $i = 1, \dots, p$ , are arbitrary. ( $F_{i2}$  has the same number

of rows as  $K_i$  has columns, and the same number of columns as  $\bar{Q}_i$  has columns.) This in turn yields

$$(8.14) \quad F = F_1 + KF_2\bar{Q}^\dagger.$$

This general form for  $F$  will now be utilized to show that when the conditions of Theorem 8.5 are satisfied and  $\text{rank } \bar{Q} = n$ , decoupling can be achieved with  $A + BF$  strictly stable (all eigenvalues having negative real part).

To see this, first note that by Lemma 8.1 the columns of  $\bar{Q}_i$  form a basis for  $\mathcal{R}[Q(F, K_i)]$ . Performing the coordinate transformation  $z = \bar{Q}^{-1}x$ , where  $x$  is the state of the decoupled system, one can then readily verify that

$$\hat{A} = \bar{Q}^{-1}(A + BF_1)\bar{Q} + \bar{Q}^{-1}(BF_2\bar{Q}^{-1})\bar{Q}$$

has the block diagonal form

$$\hat{A} = \begin{bmatrix} \hat{A}_1 & & 0 \\ & \ddots & \\ 0 & & \hat{A}_p \end{bmatrix} + \begin{bmatrix} \hat{B}_1 F_{12} & & 0 \\ & \ddots & \\ 0 & & \hat{B}_p F_{p2} \end{bmatrix}$$

and

$$\hat{B} = \bar{Q}^{-1}BK = \begin{bmatrix} \hat{B}_1 & & 0 \\ & \ddots & \\ 0 & & \hat{B}_p \end{bmatrix}.$$

The pairs  $(\hat{A}_i, \hat{B}_i)$  are clearly controllable, so that  $(\hat{A}_i + \hat{B}_i F_{i2})$  can be assigned any desired spectrum by choice of  $F_{i2}$  (see [18]). Hence,  $\hat{A}$  and therefore  $(A + BF)$  can be assigned any spectrum with an  $F$  of the form (8.14).

A corollary of Theorem 8.5 is Theorem 4.1 of Wonham and Morse [9] (stated somewhat differently).

**COROLLARY 8.2.** *Suppose  $D = 0$  and  $m = n$ , with  $\text{rank } C = n$ . Then there exists a pair  $(F, G)$  which output controllably decouples  $(A, B, C, D)$  if and only if*

$$(8.15) \quad \text{rank } \bar{P}_i = m_i, \quad i = 1, \dots, p.$$

*Proof.* Since the number of rows of  $\bar{P}_i$  is  $m_i$ , necessity follows from Theorem 7.1. To prove sufficiency, consider the matrix  $C\bar{Q}$ . By construction,  $C^i\bar{Q}_j = 0, i \neq j$ , and  $C^j\bar{Q}_j = \bar{P}_j$ , i.e.,  $C\bar{Q} = \text{diag}(\bar{P}_i)$ , so that  $\text{rank } C\bar{Q} = \sum_{i=1}^p \text{rank } \bar{P}_i = n$ . Since  $\text{rank } C = n$  it follows that  $\text{rank } \bar{Q} = n$ . Also, since  $\Gamma^j\bar{Q}_j = 0, j = 1, \dots, p$ , and  $\text{rank } \Gamma^j = n - m_j, \text{rank } \bar{Q}_j \leq m_j$ . Since each  $\bar{Q}_j$  has full column rank, it follows that  $\bar{Q}$  has  $n$  linearly independent columns. Then the desired result follows immediately from Theorem 8.5.

If the necessary condition of Theorem 8.1 is satisfied but no solution to (8.8) exists, the next logical step is to choose a  $G$  of smaller dimension than  $K$  and re-examine the decoupling conditions. A guideline for a second choice of  $G$  is available from the characterization theorems, Theorems 7.1 and 7.2. From the definition of  $K_i$ , it is evident that  $G$  must be chosen in the form  $G^i = K_i N_i$  in order to satisfy (7.14). Hence, a reasonable approach is to select a set  $N_i, i = 1, \dots, p$ , then examine (7.15) and (7.17) for the possibility of a solution for  $F$ . Equation (7.17)

should be examined first to determine if a set  $F_{i_2}$ ,  $i = 1, \dots, p$ , exists to achieve output controllability. Then if this proves successful, (7.15) can be examined by standard techniques, since for  $F_{i_2}$  fixed, this equation is linear in the elements of  $F_{i_1}$ . As yet, no way of generally terminating this algorithm after a finite number of steps has been found.

At this point, we shall return to the condition utilized in Corollary 8.1—that  $\hat{D}$  have linearly independent rows. It will be shown that for a large class of systems this condition is both necessary and sufficient for decoupling.

LEMMA 8.4. *If there exists a nonsingular matrix  $G$  such that for some  $F$  the pair  $(F, G)$  decouples  $(A, B, C, D)$ , then the rows of  $\hat{D}$  are linearly independent.*

*Proof.* If  $(F, G)$  decouples  $(A, B, C, D)$ , then from (7.17),  $\bar{D}_{a_i}G^j = 0$ ,  $i \neq j$ ; hence  $\hat{D}G$  must be block diagonal. Also, since the rows of  $\bar{D}_{a_i}$  are linearly independent and  $G$  is nonsingular, the rows of  $\bar{D}_{a_i}G$  are linearly independent. Due to its block diagonal structure, it follows that the rows of  $\hat{D}G$  are linearly independent. Then since  $G$  is nonsingular, the result follows.

THEOREM 8.6. *There exists a pair  $(F, G)$  with  $G$  nonsingular which output controllably decouples  $(A, B, C, D)$  if and only if*

- (i)  *$(A, B, C, D)$  is output controllable, and*
- (ii) *the rows of  $\hat{D}$  are linearly independent.*

*Moreover, if conditions (i) and (ii) are satisfied, the pair*

$$(8.16) \quad \hat{F} = -\hat{D}^\dagger \hat{C}, \quad \hat{G} = [\hat{D}^\dagger : K^*],$$

where the columns of  $K^*$  form a basis for the null space of  $\hat{D}$ , output controllably decouples  $(A, B, C, D)$ .

*Proof.* Necessity follows by the preceding lemma and Theorem 7.4. Sufficiency follows by observing that  $\hat{G}$  and  $\hat{F}$  satisfy conditions (7.17) and (7.18), respectively.

The following corollary covers a case considered by Wonham and Morse [9, Theorem 5.1].

COROLLARY 8.3. *If  $r = p$ , then there exists a pair  $(F, G)$  which output controllably decouples  $(A, B, C, D)$  if and only if*

- (i)  *$(A, B, C, D)$  is output controllable, and*
- (ii) *the rows of  $\hat{D}$  are linearly independent.*

*Proof.* It is clear that if  $r = p$ , the  $G$  matrix of any pair which output controllably decouples must have rank  $r$  so that the result follows by the preceding theorem.

*Remark.* If (ii) holds and  $r = p$ , then it is clear that  $\bar{D}_{a_i}$  is a row vector and rank  $\hat{D} = r$ . Hence, if (i) and (ii) hold,  $(-\hat{D}^{-1}\hat{C}, \hat{D}^{-1})$  output controllably decouples  $(A, B, C, D)$ . For the case  $D = 0$ , Wonham and Morse [9] also give a necessary and sufficient condition for the existence of a decoupling pair when  $r = p$ . Their condition is quite different in form from that given above. One advantage of the present formulation is that it reduces to the “standard” result given by Falb and Wolovich [12] for the subcase  $r = p = m$ . It is clear from the structure algorithm and the above remark that with  $D = 0$  and  $r = p$  that if (ii) holds, there is a row  $r^i$  of  $C^i$  such that

$$\bar{D}_{a_i} = r^i A^{a_i - 1} B$$

and

$$\bar{C}_{a_i} = r^i A^{a_i},$$

where  $a_i$  is the first positive integer  $j$  such that  $r^j A^{j-1} B \neq 0$ . Hence, if  $C^i$  contains only one row for all  $i$ ,  $\hat{C} = A^*$  and  $\hat{D} = B^*$  in the notation of Falb and Wolovich [12], which makes the equivalence with their results apparent.

In our notation, therefore, the Falb-Wolovich result takes the following form.

**COROLLARY 8.4.** *If  $m = p = r$ , there exists a pair  $(F, G)$  which output controllably decouples  $(A, B, C, D)$  if and only if*

$$(8.17) \quad \text{rank } \hat{D} = m.$$

Moreover, if (8.15) is satisfied,  $(-\hat{D}^{-1}\hat{C}, \hat{D}^{-1})$  is a decoupling pair.

*Proof.* Condition (i) of Corollary 8.3 is redundant in this case since  $\text{rank } \hat{D} = m$  implies the decoupled system is invertible, a stronger condition than output controllability.

When the conditions of Corollary 8.3 are satisfied, it is also possible to give a characterization of the complete class of decoupling pairs.

**LEMMA 8.5.** *Suppose  $p = r$ ,  $(A, B, C, D)$  is output controllable, and  $\hat{D}$  has linearly independent rows. Then there exists a matrix  $F$  such that the pair  $(F, G)$  output controllably decouples  $(A, B, C, D)$  if and only if  $G$  has the form*

$$(8.18) \quad G = \hat{D}^{-1}\Lambda,$$

where  $\Lambda = \text{diag}(\lambda_i)$  is any nonsingular diagonal matrix.

*Proof.* Necessity follows immediately from (7.17). For sufficiency, note that  $G$  of the form (8.18) satisfies (7.17). By noting that  $F = -\hat{D}^{-1}\hat{C}$  satisfies (7.18) and applying Theorem 7.4, the proof is completed.

In order to delineate the class of matrices  $F$  which together with a  $G$  of the form (8.18) output controllably decouple  $(A, B, C, D)$ , some additional notation will be needed. Let  $\Sigma_i$  be the identity matrix with the  $i$ th column deleted, let  $\rho_i = n - \text{rank } Q_i^*(0, \hat{D}^{-1}\Sigma_i)$  and let  $H_i$  be any  $\rho_i \times n$  matrix having rank  $\rho_i$  such that

$$(8.19) \quad H_i Q_i^*(0, \hat{D}^{-1}\Sigma_i) = 0, \quad i = 1, \dots, p.$$

Also, let

$$(8.20) \quad H = \begin{bmatrix} H_1 \\ \vdots \\ H_p \end{bmatrix}.$$

Then we have the following theorem.

**THEOREM 8.7.** *Suppose  $p = r$ ,  $(A, B, C, D)$  is output controllable, and  $\hat{D}$  has linearly independent rows. Then  $(F, G)$  output controllably decouples  $(A, B, C, D)$  if and only if  $G$  has the form (8.18) and  $F$  has the form*

$$(8.21) \quad F = \hat{D}^{-1}(-\hat{C} + FH),$$

where  $\hat{F} = \text{diag}(f_i)$  and  $f_i, i = 1, \dots, p$ , is an arbitrary  $1 \times \rho_i$  matrix.

*Proof.* Lemma 8.5 established the form of  $G$ . The necessary and sufficient conditions for  $(F, G)$  to output controllably decouple  $(A, B, C, D)$  then reduces to the requirement that  $F$  satisfy (7.15) here rewritten in the form

$$(8.22) \quad (\bar{C}_{a_i} + \bar{D}_{a_i}F)Q_i^*(F_{i2}^*, \Omega^i) = 0, \quad i = 1, \dots, p.$$

Note that  $\Omega^i = \hat{D}^{-1} \Lambda \Sigma_i$ , hence  $\bar{D}_a \Omega^i = 0$ . Therefore, since  $\text{rank } \Omega^i = r - 1$ , one can take  $K_i^* = \Omega^i$ , and, by applying Lemma 8.2,  $\mathcal{R}[Q_i^*(F_{i2}^*, \Omega^i)] = \mathcal{R}[Q_i^*(0, \hat{D}^{-1} \Lambda \Sigma_i)]$ . As  $\Lambda$  is a nonsingular, diagonal matrix, it can be seen that  $H_i$  satisfies (8.19) if and only if  $H_i Q_i^*(F_{i2}^*, \Omega^i) = 0$ . Then  $F$  in the form (8.21) is seen to satisfy (8.22). For necessity, note that any  $F$  can be expressed as

$$(8.23) \quad F = -\hat{D}^{-1} \hat{C} + \hat{D}^{-1} M$$

for some  $M$ . Substitution of (8.23) into (8.22) yields the requirement that

$$M_i Q_i^*(0, \hat{D}^{-1} \Sigma_i) = 0, \quad i = 1, \dots, p,$$

where

$$M = \begin{bmatrix} M_1 \\ \vdots \\ M_p \end{bmatrix}.$$

The form (8.21) is then seen to include all solutions  $F$ .

This characterization of decoupling pairs appears to be new even for the subcase  $r = p = m$ . It is, however, quite similar to Gilbert's characterization [15, Theorem 5]. Moreover, the above characterization can be utilized to very simply resolve the question of eigenvalue placement in the decoupled system. This method is also similar to, but somewhat more direct than, Gilbert's procedure.

The following lemma is equivalent to the related result of Gilbert [15, Lemma 1] so that the proof can be omitted.

LEMMA 8.6. *Let  $H$  be the matrix defined by (8.19) and (8.20). Then*

- (i) *the rows of  $H$  are linearly independent,*
- (ii) *there exists a matrix  $A_i$  such that*

$$H_i(A - B\hat{D}^{-1}\hat{C}) = A_i H_i$$

and

- (iii)  $\tilde{C}_j^i$  for  $j = 0, 1, \dots$  is in the row range of  $H_i$ .

Remark. As a consequence of (iii) above, it can be assumed without loss of generality that if

$$(8.24) \quad L_j^i = \begin{bmatrix} \tilde{C}_0^i \\ \vdots \\ \tilde{C}_{j-1}^i \end{bmatrix}$$

and  $\beta_i$  is the first integer such that  $\text{rank } L_{\beta_i}^i = \text{rank } L_{\beta_i+1}^i = \gamma_i$ , then the first  $\gamma_i$  rows of  $H_i$  can be chosen to be the first  $\gamma_i$  independent rows of  $L_{\beta_i}^i$ .

Following from Lemma 8.6 and Corollary 8.3 is a simple representation of the class of all decoupled systems for which the conditions of Theorem 8.7 hold. First note that (8.19) implies

$$(8.25) \quad HB\hat{D}^{-1} = \begin{bmatrix} b_1 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & b_p \end{bmatrix},$$

where  $b_i$  is a  $\rho_i \times 1$  matrix. It is clear from Lemma 8.6 and (8.25) that

$$H(A + BF)H^\dagger = \begin{bmatrix} (A_1 - b_1 f_1) & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & (A_p - b_p f_p) \end{bmatrix},$$

where  $F$  is any matrix of the form (8.22) and  $H^\dagger = H'(HH')^{-1}$ .

If  $J$  is any matrix whose columns form a basis for the right null space of  $H$  and

$$T = \begin{bmatrix} H \\ J^\dagger \end{bmatrix}, \quad J^\dagger = (J'J)^{-1}J',$$

is used as a coordinate transformation, it can then be verified that the class of decoupled systems  $(T(A + BF)T^{-1}, TBG, (C + DF)T^{-1}, DG)$  has the following explicit form (assume without loss of generality that  $D^i = D_0^i$  and  $C^i = C_0^i$ ):

$$T(A + BF)T^{-1} = \left[ \begin{array}{ccc|c} (A_1 - b_1 f_1) & \cdots & 0 & \\ \vdots & & \vdots & 0 \\ 0 & \cdots & (A_p - b_p f_p) & \\ \hline J^\dagger(A - B\hat{D}^{-1}\hat{C} + B\hat{D}^{-1}\hat{F}H)H^\dagger & & & J^\dagger(A - B\hat{D}^{-1}\hat{C})J \end{array} \right],$$

$$TBG = \begin{bmatrix} b_1 \lambda_1 & \cdots & 0 \\ \vdots & & \vdots \\ 0 & \cdots & b_p \lambda_p \\ J^\dagger B \hat{D}^{-1} \Lambda \end{bmatrix}, \quad DG = \begin{bmatrix} D^1 G \\ \vdots \\ D^p G \end{bmatrix},$$

where

$$D^i G = \begin{cases} \begin{bmatrix} e_i \lambda_i \\ 0 \end{bmatrix} & \text{if } D^i \neq 0, \\ 0 & \text{if } D^i = 0, \end{cases}$$

with  $e_i$  the  $i$ th row of  $I_r$ , and

$$(C + DF)T^{-1} = \begin{bmatrix} (C^1 + D^1 F)T^{-1} \\ \vdots \\ (C^p + D^p F)T^{-1} \end{bmatrix},$$

where

$$(C^i + D^i F)T^{-1} = \begin{cases} \begin{bmatrix} \tilde{J}_i \\ E_{m_i-1}^i \end{bmatrix} & \text{if } D^i \neq 0, \\ E_{m_i}^i & \text{if } D^i = 0, \end{cases}$$

and  $E_i^i$  is an  $l \times n$  matrix whose  $k$ th row is the  $(\sum_{j=1}^{i-1} \rho_j + k)$ th row of  $I_n$  and  $f_i$  is

the row vector  $\vec{f}_i$  with all entries replaced by zero except those indexed  $\sum_{j=1}^{i-1} \rho_j + k, k = 1, \dots, \rho_i$ .

Several results follow directly from this characterization. The most important is the following which generalizes and reformulates the related results of Gilbert [15] and Wonham and Morse [9].

**THEOREM 8.8.** *Suppose  $p = r, (A, B, C, D)$  is state controllable and the rows of  $\hat{D}$  are linearly independent. Then there exists a pair  $(F, G)$  such that  $(A + BF, BG, C + DF, DG)$  is output controllably decoupled and strictly stable if and only if  $J^\dagger(A - B\hat{D}^{-1}\hat{C})J$  is strictly stable.*

*Proof.* Controllability of  $(A, B, C, D)$  implies that the pairs  $(A_i, b_i)$  are controllable so that  $(A_i + b_i f_i)$  can be assigned any desired spectrum by choice of  $f_i$ . It is also clear that the eigenvalues of  $J^\dagger(A - B\hat{D}^{-1}\hat{C})J$  are invariant under decoupling feedback.

*Remark.* Under the hypothesis of Theorem 8.8,  $(A, B, C, D)$  is left invertible and the dynamic part of the left inverse has the representation (see (4.13)–(4.14))

$$(A - B\hat{D}^{-1}\hat{C}, B\hat{D}^{-1}, -\hat{D}^{-1}\hat{C}, \hat{D}^{-1}).$$

It is clear, therefore, that a sufficient condition for stable decoupling when  $\hat{D}$  is nonsingular is that  $(A, B, C, D)$  be “minimum phase” in the sense that the eigenvalues of the inverse system have negative real part.

It is also clear from the characterization of decoupled systems that the transfer function matrix  $Z(s, F, G)$  has the form

$$Z(s, F, G) = \text{diag} (Z_i(s, F, G)),$$

where

$$Z_i(s, F, G) = \begin{cases} \begin{bmatrix} \vec{f}_i \\ E_{m_i-1}^i \end{bmatrix} (sI - A_i + b_i f_i)^{-1} b_i \lambda_i + e_i \lambda_i & \text{if } D^i \neq 0, \\ E_{m_i}^i (sI - A_i + b_i f_i)^{-1} b_i \lambda_i & \text{if } D^i = 0. \end{cases}$$

Hence, the transfer function matrix of the decoupled system can be assigned any desired pole pattern.

**9. Open loop and dynamic decoupling.** As seen in the previous section, a complete solution to the state feedback decoupling problem is not yet possible. Moreover, even for the class of systems in which the decoupling problem is resolved, stable decoupling by state feedback is not always possible. This has motivated the examination of more general decoupling laws [15], [18], [20]. Morse and Wonham [18] have shown that if the class of decoupling laws is sufficiently enlarged, necessary and sufficient conditions can be given both for decoupling and stable decoupling (the conditions turn out to be identical). This problem will also be considered here. It is shown below that the open loop characterization of zeroing inputs given by Corollary 5.1 leads to a very simple proof of the general decoupling criterion as well as an explicit realization of a dynamic decoupling law. Moreover, this approach leads to a considerably stronger result—a necessary and sufficient condition for stable, dynamic decoupling with *output* feedback only.

We shall first consider the open loop decoupling problem.

Let  $\mathcal{U}_i$  denote a subset of the input space  $\mathcal{U}$  of  $\mathcal{S}$  and let  $\mathcal{S}_{x_0}^i$  denote the system defined by the quadruple  $(A, B, C^i, D^i)$  with initial state  $x_0$ . Then the following definition of an open loop decoupling set of inputs will be adopted. It is essentially equivalent to the open loop decoupling definition given in [18].

DEFINITION 9.1. A family of input sets  $\{\mathcal{U}_i\}_{i=1}^p$  output controllably decouples  $\mathcal{S}$  if and only if:

- (i)  $\mathcal{S}_0^j[u] = 0$  for  $j \neq i$  and  $u \in \mathcal{U}_i$ ; and
- (ii) for every  $\lambda \in R^{m_i}$  there is a  $\hat{u} \in \mathcal{U}_i$  such that  $\hat{y}^i(1) = \lambda$ , where  $\hat{y}^i(1)$  is the response of  $\mathcal{S}_0^i$  to  $\hat{u}$  evaluated at  $t = 1$ .

The following theorem gives a necessary and sufficient condition for the existence of a decoupling family. This result is equivalent (for the case  $D = 0$ ) to the corresponding result of Morse and Wonham [18, Theorem 6.2]. The proof is believed to be simpler and more constructive.

THEOREM 9.1. *There exists a family of input sets  $\{\mathcal{U}_i\}_{i=1}^p$  which output controllably decouples  $\mathcal{S}$  if and only if*

$$(9.1) \quad \text{rank } \bar{P}_i = m_i, \quad i = 1, \dots, p.$$

*Proof. Necessity.* Suppose  $\{\mathcal{U}_i\}_{i=1}^p$  output controllably decouples  $\mathcal{S}$ . By Corollary 5.2 the set  $\bar{\mathcal{U}}_i$  of all inputs  $u$  for which  $\mathcal{S}_0^j[u] = 0$  for  $j \neq i$  can be generated by the feedback law

$$(9.2) \quad u = -\bar{\Delta}_{\alpha_i}^\dagger \bar{\Gamma}_{\alpha_i} x + K_i v.$$

Moreover,  $\bar{P}_i$  is the output controllability matrix of  $\mathcal{S}^i$  subject to feedback of the form (9.2). Hence if condition (ii) of Definition 9.1 holds,  $\bar{P}_i$  must have rank  $m_i$  since by condition (i),  $\mathcal{U}_i \subset \bar{\mathcal{U}}_i$ .

*Sufficiency.* If  $\text{rank } \bar{P}_i = m_i$ , then by definition, the set of input functions  $\bar{\mathcal{U}}_i$  generated by the feedback law (9.2) with  $x_0 = 0$  satisfies conditions (i) and (ii) of Definition 9.1.

*Remark.* As seen in the previous section, it is not generally true that all the  $\bar{\mathcal{U}}_i$  can be generated simultaneously with a single feedback law. However, as was shown in Theorem 5.1, the set  $\mathcal{U}_i$  can be generated in an open loop manner by the system representation  $\mathcal{Z}^i$ ,

$$(9.3) \quad \dot{z}^i = (A - B\bar{\Delta}_{\alpha_i}^\dagger \bar{\Gamma}_{\alpha_i})z^i + BK_i v_i,$$

$$(9.4) \quad u^i = -\bar{\Delta}_{\alpha_i}^\dagger \bar{\Gamma}_{\alpha_i} z^i + K_i v_i$$

with  $z^i(0) = 0$ . It is clear, therefore, if we set

$$u = \sum_{i=1}^p u^i,$$

where  $u^i$  is the output of (9.3)–(9.4), decoupling is achieved with the dynamic control law defined by the direct sum of the system representations  $\mathcal{Z}^i$ .

More generally one can define a general dynamic control law by a pair of functions characterizing a dynamical system  $\mathcal{T}$  coupled from  $\mathcal{S}$ :

$$(9.5) \quad z(t) = \mathcal{O}(t, t_0, z(t_0), x_{[t_0,t]}, v_{[t_0,t]}),$$

$$(9.6) \quad u(t) = \mathcal{B}(t, z(t), x(t), v(t)),$$



where  $\mathcal{O}$  is the state transition function [22] of  $\mathcal{T}$ ,  $z(t) \in \Sigma'$ , the state space of  $\mathcal{T}$ ,  $v \in \mathcal{V}$ , the input space of the closed loop system and  $\mathcal{B}$  is the readout function [22] of  $\mathcal{T}$ . The overall response of  $\mathcal{T}$  can be represented by a single operator

$$(9.7) \quad u(t) = \mathcal{T}_{x_0, z_0}[t, v],$$

where it is understood that  $t_0 = 0$ ,  $x(0) = x_0$ ,  $z(0) = z_0$  and  $v = v_{[0, t]}$ . Let  $v$  be partitioned as in (6.3) and let  $\mathcal{T}_{x_0, z_0}^i[t, v^i] = \mathcal{T}_{x_0, z_0}[t, v]$  when  $v^j = 0, j \neq i$ . The closed loop system will be denoted as  $\mathcal{S}[\mathcal{T}]$ .

A general definition of dynamic decoupling can now be given.

DEFINITION 9.2. The control law (9.7) *output controllably decouples*  $\mathcal{S}$  if and only if

(i)  $\mathcal{S}_0^j[\mathcal{T}_{0,0}^i[t, v^i]] = 0$  for  $j \neq i$ ,

and

(ii) the closed loop system  $\mathcal{S}[\mathcal{T}]$  is output controllable.

The next theorem follows immediately from Theorem 5.1 and the subsequent remarks.

THEOREM 9.2. *There exists a control law of the type (9.7) which output controllably decouples  $\mathcal{S}$  if and only if*

$$(9.8) \quad \text{rank } \bar{P}_i = m_i, \quad i = 1, \dots, p.$$

Dynamic decoupling can always be achieved with a strictly stable closed loop system if (9.8) is satisfied and  $\mathcal{S}$  is (state) controllable. To see this, first note that there is no loss of generality in assuming that  $A$  is strictly stable since if  $\mathcal{S}$  is controllable a preliminary feedback law of the type  $u = Fx + w$  can be used to give  $A + BF$  any desired set of eigenvalues [19]. Also, it follows from Corollary 5.1 that the set of inputs  $\bar{u}_i$  can be generated by a system

$$\mathcal{X}^i = (\hat{A}^i, \hat{B}^i, \hat{C}^i, \hat{D}^i)$$

with  $\hat{A}^i$  having any desired set of eigenvalues. Consequently, the combined system  $(\hat{A}, \hat{B}, \hat{C}, \hat{D})$  can have any desired set of eigenvalues since

$$\hat{A} = \begin{bmatrix} A & B\hat{C}^1 & B\hat{C}^2 & \dots & B\hat{C}^p \\ 0 & \hat{A}^1 & 0 & \dots & 0 \\ \vdots & \vdots & & \ddots & \\ \vdots & \vdots & & & \\ 0 & 0 & 0 & \dots & \hat{A}^p \end{bmatrix}.$$

For completeness note that

$$\hat{B} = \begin{bmatrix} B\hat{D}^1 & B\hat{D}^2 & \dots & B\hat{D}^p \\ \hat{B}^1 & 0 & \dots & 0 \\ \vdots & & & \\ \vdots & & & \\ 0 & 0 & \dots & \hat{B}^p \end{bmatrix},$$

$$\hat{C} = [C : D\hat{C}^1 : \dots : D\hat{C}^p], \quad \hat{D} = [D\hat{D}^1 : D\hat{D}^2 : \dots : D\hat{D}^p].$$

It is of interest to observe that if the original  $A$  matrix is stable, no state feedback is required to achieve stable decoupling when condition (9.8) holds. In

many applications the state of  $\mathcal{S}$  cannot be measured directly so that this type of decoupling law is quite important. These ideas can be extended by incorporating an observer [23], [24] into the dynamic decoupling law. It is well known that if  $\mathcal{S}$  is controllable and observable, an observer of the type

$$(9.9) \quad \dot{z} = \check{A}z + \check{B}y$$

with

$$(9.10) \quad u = \check{F}z + w$$

can be found so that the closed loop system has any desired set of eigenvalues. It is clear, therefore, that if we can show that  $\text{rank } \bar{P}_i$  is invariant under this type of feedback a necessary and sufficient condition for stable dynamic decoupling which utilizes only output measurements will result. To show that this is indeed the case, a generalization of Lemma 3.1 will be given.

Observe first that under feedback of the type (9.9)–(9.10), the closed loop system  $\check{\mathcal{S}}$  is represented by the equations

$$(9.11) \quad \begin{bmatrix} \dot{x} \\ \dot{z} \end{bmatrix} = \begin{bmatrix} A & B\check{F} \\ \check{B}C & (\check{A} + \check{B}D\check{F}) \end{bmatrix} \begin{bmatrix} x \\ z \end{bmatrix} + \begin{bmatrix} B \\ \check{B}D \end{bmatrix} w,$$

$$(9.12) \quad y = [C \ ; \ D\check{F}] \begin{bmatrix} x \\ z \end{bmatrix} + Dw.$$

If  $P$  is any operator associated with  $\mathcal{S}$  and  $Q$  is the corresponding operator of  $\check{\mathcal{S}}$ , we shall say that

$$P \xrightarrow{(\check{A}, \check{B}, \check{F})} Q.$$

LEMMA 9.1. *The operators  $N_i$  and the matrices  $D_i$  of the structure algorithm are invariant under dynamic output feedback of the type (9.9)–(9.10) and*

$$(9.13) \quad \bar{C}_i \xrightarrow{(\check{A}, \check{B}, \check{F})} [\bar{C}_i \ ; \ \bar{D}_i\check{F}],$$

$$(9.14) \quad \check{C}_i \xrightarrow{(\check{A}, \check{B}, \check{F})} [\check{C}_i \ ; \ 0].$$

The proof of this result follows in a manner essentially the same as that used to establish Lemma 3.1 and is therefore omitted.

*Remark.* Lemma 9.1 establishes the invariants of the structure algorithm applied to the whole output of  $\mathcal{S}$ . It is clear that similar results hold for the algorithm applied to a subset of the output. In particular,  $\bar{\Delta}_{\alpha_i}$  is invariant so that  $K_i$  is also invariant and

$$(9.15) \quad \bar{\Gamma}_{\alpha_i} \xrightarrow{(\check{A}, \check{B}, \check{F})} [\bar{\Gamma}_{\alpha_i} \ ; \ \bar{\Delta}_{\alpha_i}\check{F}]$$

and

$$(9.16) \quad \tilde{\Gamma}_{\alpha_i} \xrightarrow{(\check{A}, \check{B}, \check{F})} [\tilde{\Gamma}_{\alpha_i} \ ; \ 0].$$

LEMMA 9.2.  *$\text{rank } \bar{P}_i$  is invariant under dynamic output feedback of the type (9.9)–(9.10).*

*Proof.* Using (9.11), (9.12) and (9.15), (9.16) we can see that

$$(9.17) \quad A - B\bar{\Delta}_{\alpha_i}\bar{\Gamma}_{\alpha_i} \xrightarrow{(\check{A}, \check{B}, \check{F})} \left[ \begin{array}{c|c} A - B\bar{\Delta}_{\alpha_i}^\dagger\bar{\Gamma}_{\alpha_i} & B(I - \bar{\Delta}_{\alpha_i}^\dagger\bar{\Delta}_{\alpha_i})\check{F} \\ \hline \check{B}(C - D\bar{\Delta}_{\alpha_i}^\dagger\bar{\Gamma}_{\alpha_i}) & A + \check{B}D(I - \bar{\Delta}_{\alpha_i}^\dagger\bar{\Delta}_{\alpha_i})\check{F} \end{array} \right]$$

and

$$(9.18) \quad C^i - D^i\bar{\Delta}_{\alpha_i}\bar{\Gamma}_{\alpha_i} \xrightarrow{(\check{A}, \check{B}, \check{F})} [C^i - D^i\bar{\Delta}_{\alpha_i}^\dagger\bar{\Gamma}_{\alpha_i} : D^i(I - \bar{\Delta}_{\alpha_i}^\dagger\bar{\Delta}_{\alpha_i})\check{F}].$$

Since the columns of  $I - \bar{\Delta}_{\alpha_i}^\dagger\bar{\Delta}_{\alpha_i}$  span the null space of  $\bar{\Delta}_{\alpha_i}$  there exists a matrix  $M_i$  such that

$$(9.19) \quad I - \bar{\Delta}_{\alpha_i}^\dagger\bar{\Delta}_{\alpha_i} = K_i M_i.$$

Let  $\check{P}_i(0, K_i)$  be defined by

$$(9.20) \quad P_i(0, K_i) \xrightarrow{(\check{A}, \check{B}, \check{F})} \check{P}_i(0, K_i),$$

and for convenience reorder the columns of  $P_i(0, K_i)$  so that

$$(9.21) \quad P_i(0, K_i) = [D^i K_i : (C^i - D^i\bar{\Delta}_{\alpha_i}^\dagger\bar{\Gamma}_{\alpha_i})Q_i(0, K_i)]$$

with  $Q_i(0, K_i)$  given by (7.3). Using (9.17), (9.18) and (9.19) one can then see that  $\check{P}_i(0, K_i)$  has the form

$$\check{P}_i(0, K_i) = P_i(0, K_i) \begin{bmatrix} I_{r_1} & J_{11} & \cdots & J_{1n} \\ 0 & I_{r_2} & \cdots & J_{2,n-1} \\ \vdots & \vdots & & \vdots \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & I_{r_p} \end{bmatrix},$$

where  $r_i = \text{rank } K_i$  ( $J_{11} = M_{11}\check{F}\check{B}D$ , etc.). Since the matrix postmultiplying  $P_i(0, K_i)$  is nonsingular, the desired result follows.

It is now clear by Lemma 9.2 and Theorem 9.2 that if  $\text{rank } \bar{P}_i = m_i$  for  $i = 1, \dots, p$ , and  $\mathcal{S}$  is controllable and observable, then  $\mathcal{S}$  can be output controllably decoupled by dynamic compensation and output feedback alone with the resulting closed loop system having any desired eigenvalues. The procedure to accomplish this is first to construct an observer of the type (9.9)–(9.10) to shift the eigenvalues of  $A$  (the eigenvalues of the observer can be placed arbitrarily). The closed loop system  $\check{\mathcal{S}}$  resulting is then decoupled by the procedure given in the remarks following Theorem 9.1. The combined compensation for decoupling can be illustrated in block diagonal form as shown in Fig. 1, where  $\mathcal{O}$  represents the observer and  $\mathcal{P}$  the precompensator for decoupling.

The above results can be summarized by the following theorem.

**THEOREM 9.3.** *Let  $\mathcal{S}$  be controllable and observable. Then there exists a dynamic control law of the type*

$$u = \hat{C}_1 z_1 + \hat{C}_2 z_2 + \hat{D}v,$$

where

$$\dot{z}_1 = \hat{A}_1 z_1 + \hat{B}_1 y \quad (\text{observer})$$

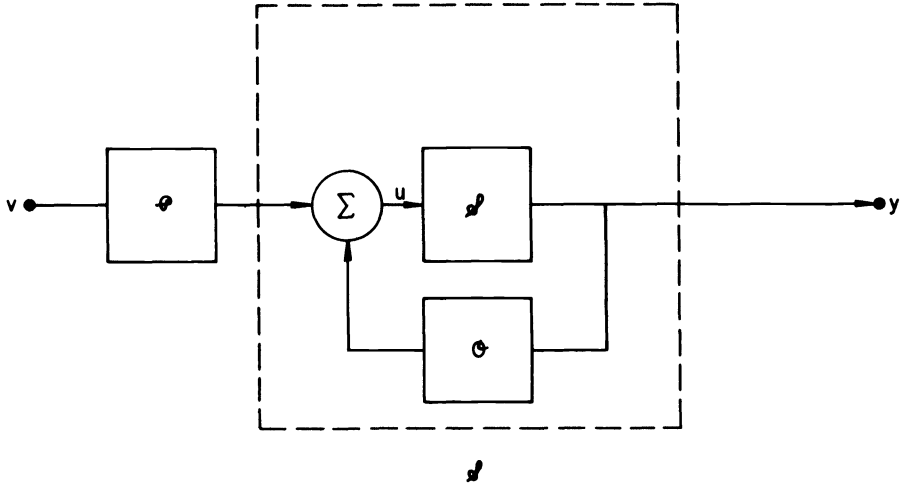


FIG. 1

and

$$\dot{z}_2 = \hat{A}_2 z_2 + \hat{B}_2 v \quad (\text{decoupler})$$

such that the closed loop system

$$\begin{bmatrix} \dot{x} \\ \dot{z}_1 \\ \dot{z}_2 \end{bmatrix} = \begin{bmatrix} A & BC_1 & BC_2 \\ \hat{B}_1 C & (\hat{A}_1 + \hat{B}_1 DC_1) & \hat{B}_1 DC_2 \\ 0 & 0 & \hat{A}_2 \end{bmatrix} \begin{bmatrix} x \\ z_1 \\ z_2 \end{bmatrix} + \begin{bmatrix} B\hat{D} \\ \hat{B}_1 D\hat{D} \\ B_2 \end{bmatrix} v,$$

$$y = [C : DC_1 : DC_2] \begin{bmatrix} x \\ z_1 \\ z_2 \end{bmatrix} + D\hat{D}v$$

is output controllably decoupled and strictly stable if and only if  $\text{rank } \bar{P}_i = m_i$  for  $i = 1, \dots, p$ .

Theorem 9.3 completely solves the “classical” decoupling problem of decoupling a system which is specified by its transfer function matrix (see Morgan [11] for references to the early literature) since a transfer function matrix is only a faithful representation of its minimal (controllable and observable) realizations. Note that for the case  $r = m = p$  the condition on  $\text{rank } \bar{P}_i$  is equivalent to invertibility-nonsingularity of the transfer function matrix.

The question of minimizing the compensating dynamical system for decoupling is still open. For the case where state feedback is allowed, some progress has been made by Morse and Wonham [18].

**10. Concluding remarks.** Several problems of interest related to input-output structure and decoupling are still outstanding. The most obvious is that of generalizing the results given here to the time-variable case. An indication of how such generalizations can be obtained is given in [8]. The problem of determining, by a

finite algorithm, whether or not a given system can be decoupled is still not resolved in general. The characterization theorems of § 7 should be useful in this respect. The related problem of decoupling with a minimal order compensator and output feedback is also open.

**Acknowledgment.** The authors would like to thank E. G. Gilbert, A. S. Morse and W. M. Wonham for several stimulating discussions on the decoupling problem.

## REFERENCES

- [1] L. WEISS, *On a question related to the control of linear systems*, IEEE Trans. Automatic Control, AC-9 (1964), pp. 176-177.
- [2] R. W. BROCKETT, *Poles, zeros and feedback: State space interpretation*, Ibid., AC-10 (1965), pp. 129-135.
- [3] R. SIVAN, *On zeroing the output and maintaining it zero*, Ibid., AC-10 (1965), pp. 193-194.
- [4] L. M. SILVERMAN, *Properties and application of inverse systems*, Ibid., AC-13 (1968), pp. 436-437.
- [5] R. W. BROCKETT AND M. D. MESAROVIC, *The reproducibility of multivariable control systems*, J. Math. Anal. Appl., 11 (1965), pp. 548-563.
- [6] D. C. YOULA AND P. DORATO, *On the inverse of linear dynamical systems*, Electrophysics Memo. PIBMR I-1319-66, Polytechnic Institute of Brooklyn, New York, 1966.
- [7] M. K. SAIN AND J. L. MASSEY, *Invertibility of linear time-invariant systems*, IEEE Trans. Automatic Control, AC-14 (1969), pp. 141-149.
- [8] L. M. SILVERMAN, *Inversion of multivariable linear systems*, Ibid., AC-14 (1969), pp. 270-276.
- [9] W. M. WONHAM AND A. S. MORSE, *Decoupling and pole assignment in linear multivariable systems: A geometric approach*, this Journal, 8 (1970), pp. 1-18.
- [10] B. S. MORGAN, JR., *The synthesis of linear multivariable systems by state variable feedback*, Proc. 1964 JACC, Stanford, Calif., pp. 468-472.
- [11] ———, *Multivariable systems*, 1965 IEEE International Convention Record, Part 6, pp. 87-95.
- [12] P. L. FALB AND W. A. WOLOVICH, *On the decoupling of multivariable systems*, Proc. 1967 JACC, Philadelphia, pp. 791-796.
- [13] ———, *Decoupling in the design of multivariable control systems*, IEEE Trans. Automatic Control, AC-12 (1967), pp. 651-659.
- [14] W. A. WOLOVICH AND P. L. FALB, *On the structure of multivariable systems*, this Journal, 7 (1969), pp. 437-451.
- [15] E. G. GILBERT, *The decoupling of multivariable systems by state feedback*, this Journal, 7 (1969), pp. 50-63.
- [16] Z. V. REKASIUS, *Decoupling of multivariable systems by means of state variable feedback*, Proc. Third Allerton Conference on Circuit and System Theory, Monticello, Ill., 1965, pp. 439-448.
- [17] A. S. MORSE, *Output controllability and system synthesis*, this Journal, 9 (1971), pp. 143-148.
- [18] A. S. MORSE AND W. M. WONHAM, *Decoupling and pole assignment by dynamic compensation*, this Journal, 8 (1970), pp. 317-337.
- [19] W. M. WONHAM, *On pole assignment in multi-input controllable linear systems*, IEEE Trans. Automatic Control, AC-12 (1967), pp. 660-665.
- [20] E. KREINDLER AND P. E. SARACHIK, *On the concepts of controllability and observability of linear systems*, Ibid., AC-9 (1964), pp. 129-136.
- [21] L. M. SILVERMAN, *Decoupling with state feedback and precompensation*, Ibid., to appear.
- [22] R. E. KALMAN, P. L. FALB AND M. ARBIB, *Topics in Mathematical System Theory*, McGraw-Hill, New York, 1969.
- [23] D. G. LUENBERGER, *Observers for multivariable systems*, IEEE Trans. Automatic Control, AC-11 (1966), pp. 190-197.
- [24] J. J. BONGIORNO, JR. AND D. C. YOULA, *On observers in multivariable control systems*, Internat. J. Control, 8 (1968), pp. 221-243.

## EQUILIBRIUM FEEDBACK CONTROL IN LINEAR GAMES WITH QUADRATIC COSTS\*

D. L. LUKES†

**Abstract.** Problems concerning equilibrium strategies for linear games are treated in a Hilbert space setting. First a nonlinear equation is derived for the equilibrium solutions in the player's allowed operator feedback spaces. The results obtained from the study of this equation are then used to compute the solution to a class of differential games and to establish the uniqueness of the solution.

### 1. An $n$ -player game in Hilbert space.

**1.1. Introduction.** This paper studies the class of games modeled by the linear state equation

$$(1.1) \quad x = x_0 + \sum_1^n \mathcal{B}_j u_j$$

with *state variables*  $x, x_0$  in a Hilbert space  $H$ ; *control variables*  $u_j$  in a Hilbert space  $H_j, j = 1, 2, \dots, n$ , and  $\mathcal{B}_j: H_j \rightarrow H$  specified bounded linear transformations. The  $n$  players compute nonzero sum *costs* by means of (1.1) and prescribed quadratic functionals (1.2').

The existence, uniqueness and stability of equilibrium solutions in  $\sum_1^n \oplus H_j$  were investigated by Russell and Lukes [1] for games defined in terms of (1.1), (1.2'). That treatment includes the games in which (1.1) is an ordinary differential equation (a differential game). Varaiya [2] has developed an existence theory for the differential game with costs satisfying certain convexity conditions and with constraints upon the  $u_i$ .

The feedback synthesis of the equilibrium solution obtained for the differential game in [1] and further studied in [3] motivated the author's search for solutions to the game (1.1), (1.2') in subspaces of  $[H, \sum_1^n \oplus H_j]$ , the Banach space of all bounded linear transformations from  $H$  into  $\sum_1^n \oplus H_j$ , rather than in  $\sum_1^n \oplus H_j$  itself. The subsequent solutions presented below are probably more realistic in terms of their game-theoretic interpretation.

A nonlinear operator equation whose solutions are the equilibriums is derived in Theorem 1.1. Sufficient conditions for this equation to have a solution are then presented by Theorems 1.3, 1.5. To demonstrate the applicability of the results, the differential game is treated as an example, with operator subspaces taken to be matrix operators. In the proof of the local playability of the resulting game, Theorem 2.3, a method emerges for estimating the time interval over which the game has a solution. With the local playability established, an equilibrium matrix solution is thereby obtained, computed in terms of the solution to a system of quadratic matrix differential equations. This solution agrees with the results derived by Case [4] and others who have studied differential games

---

\* Received by the editors July 23, 1970, and in revised form October 5, 1970.

† Department of Applied Mathematics and Computer Science, University of Virginia, Charlottesville, Virginia 22901. This research was supported by the National Aeronautics and Space Administration under Grant NGR 47-005-029.

by means of a Hamilton–Jacobi approach. The theory of the present article moreover establishes the uniqueness of those solutions.

The following abstract approach is intended to illuminate the underlying structure common to a general class of games including dynamic games of the above type.

**1.2. The cost functionals.** Player  $P_i$ , who commands control variable  $u_i$

functional

$$(1.2') \quad \mathcal{C}_i = |x - \bar{x}_i|_{\mathcal{W}_i}^2 + |u_i|_{\mathcal{U}_i}^2, \quad i = 1, 2, \dots, n,$$

in which target states  $\bar{x}_i \in H$  and bounded self-adjoint linear operators  $\mathcal{W}_i, \mathcal{U}_i$  are given. Of the associated quadratic forms we require  $|x|_{\mathcal{W}_i}^2 = (x, \mathcal{W}_i x) \geq 0$  for all  $x \in H$  and  $|u_i|_{\mathcal{U}_i}^2 = (u_i, \mathcal{U}_i u_i) > 0$  for all nonzero  $u_i \in H_i$ . By making the obvious preliminary change of control variables we can take the  $\mathcal{U}_i$  to be identity operators and hereafter taking  $\bar{x}_i = 0$  deal with the less cumbersome forms

$$(1.2) \quad \mathcal{C}_i = |x|_{\mathcal{W}_i}^2 + |u_i|^2, \quad i = 1, 2, \dots, n.$$

**1.3. Equilibrium feedback control.** *Feedback controls* are defined to be bounded linear transformations  $\mathcal{L}_i: H \rightarrow H_i, i = 1, 2, \dots, n$ , which give rise to the *abstract variation of parameters formulas*

$$(1.3) \quad x_i = \left( \mathcal{I} - \sum_{j \neq i}^n \mathcal{B}_j \mathcal{L}_j \right)^{-1} (x_0 + \mathcal{B}_i u_i)$$

with  $x_0 \in H, u_i \in H_i, i = 1, 2, \dots, n$ , and

$$(1.4) \quad x = \left( \mathcal{I} - \sum_1^n \mathcal{B}_j \mathcal{L}_j \right)^{-1} x_0, \quad x_0 \in H.$$

Hence if player  $P_i$  plays  $u_i \in H_i$  in (1.1) while the remaining players play feedback controls  $u_j = \mathcal{L}_j x$ , then  $P_i$  incurs the cost computed from (1.2),

$$(1.5) \quad \begin{aligned} \mathcal{C}_i(\mathcal{L}_1, \mathcal{L}_2, \dots, u_i, \dots, \mathcal{L}_n) &= |x_i|_{\mathcal{W}_i}^2 + |u_i|^2 \\ &= |x_0 + \mathcal{B}_i u_i|_{\Omega_i}^2 + |u_i|^2, \end{aligned}$$

where

$$(1.6) \quad \Omega_i = \left( \mathcal{I} - \sum_{j \neq i}^n \mathcal{B}_j \mathcal{L}_j \right)^{-1*} \mathcal{W}_i \left( \mathcal{I} - \sum_{j \neq i}^n \mathcal{B}_j \mathcal{L}_j \right)^{-1}.$$

We shall restrict player  $P_i$ 's feedback controls  $\mathcal{L}_i \in L_i = L_i[H, H_i]$  where  $L_i$  is a closed subspace of  $[H, H_i]$ , the Banach space of all bounded linear transformations from  $H$  into  $H_i, i = 1, 2, \dots, n$ . A *feedback system*  $\mathcal{L} = (\mathcal{L}_1, \mathcal{L}_2, \dots, \mathcal{L}_n)$  will refer to an element,  $\mathcal{L} = \sum_1^n \oplus \mathcal{L}_i$  with  $\mathcal{L}_i \in L_i$ , in the space of transformations  $\sum_1^n \oplus L_i$  from  $\sum_1^n \oplus H$  into  $\sum_1^n \oplus H_i$ . We call a feedback system  $(\hat{\mathcal{L}}_1, \hat{\mathcal{L}}_2, \dots, \hat{\mathcal{L}}_n)$  an *equilibrium feedback system* (or simply an *equilibrium*) if it produces state and control responses defined by the equations

$$(1.7) \quad \hat{x} = x_0 + \sum_1^n \mathcal{B}_j(\hat{\mathcal{L}}_j \hat{x}),$$

$$(1.8) \quad \hat{u}_i = \hat{\mathcal{L}}_i \hat{x}, \quad i = 1, 2, \dots, n,$$

such that

$$(1.9) \quad \mathcal{C}_i(\hat{\mathcal{L}}_1, \hat{\mathcal{L}}_2, \dots, \hat{u}_i, \dots, \hat{\mathcal{L}}_n) \leq \mathcal{C}_i(\hat{\mathcal{L}}_1, \hat{\mathcal{L}}_2, \dots, u_i, \dots, \hat{\mathcal{L}}_n)$$

for all  $u_i \in H_i, i = 1, 2, \dots, n$ .

**THEOREM 1.1.** *A system of feedback controls  $(\mathcal{L}_1, \mathcal{L}_2, \dots, \mathcal{L}_n)$  is an equilibrium system if and only if it satisfies the equations*

$$(1.10) \quad \left[ \mathcal{L}_i + \mathcal{B}_i^* \left( \mathcal{I} - \sum_{j \neq i}^n \mathcal{B}_j \mathcal{L}_j \right)^{-1*} \mathcal{W}_i \right] \left( \mathcal{I} - \sum_1^n \mathcal{B}_j \mathcal{L}_j \right)^{-1} x_0 = 0,$$

$i = 1, 2, \dots, n$ . The corresponding control response is given by the formula

$$(1.11) \quad u_i = -(\mathcal{I} + \mathcal{B}_i^* \Omega_i \mathcal{B}_i)^{-1} \mathcal{B}_i^* \Omega_i x_0,$$

where  $\Omega_i$  is computed from the equilibrium system according to (1.6).

*Proof.* Assume  $(\mathcal{L}_1, \mathcal{L}_2, \dots, \mathcal{L}_n)$  is a feedback system satisfying (1.10). Define  $\Omega_i$  by (1.6) and consider

$$(1.12) \quad \hat{u}_i = -(\mathcal{I} + \mathcal{B}_i^* \Omega_i \mathcal{B}_i)^{-1} \mathcal{B}_i^* \Omega_i x_0,$$

$i = 1, 2, \dots, n$ . Note that the inverse operator in (1.12) exists. Now let  $u_i \in H_i, i = 1, 2, \dots, n$ , be arbitrary and compute

$$(1.13) \quad \begin{aligned} & \mathcal{C}_i(\mathcal{L}_1, \mathcal{L}_2, \dots, u_i, \dots, \mathcal{L}_n) - \mathcal{C}_i(\mathcal{L}_1, \mathcal{L}_2, \dots, \hat{u}_i, \dots, \mathcal{L}_n) - |u_i - \hat{u}_i|_{\mathcal{I} + \mathcal{B}_i^* \Omega_i \mathcal{B}_i}^2 \\ &= |x_0 + \mathcal{B}_i u_i|_{\Omega_i}^2 + |u_i|^2 - |x_0 + \mathcal{B}_i \hat{u}_i|_{\Omega_i}^2 - |\hat{u}_i|^2 \\ &\quad - |u_i - \hat{u}_i|_{\mathcal{I} + \mathcal{B}_i^* \Omega_i \mathcal{B}_i}^2 \\ &= 2([\mathcal{B}_i^* \Omega_i x_0 + (\mathcal{I} + \mathcal{B}_i^* \Omega_i \mathcal{B}_i) \hat{u}_i], u_i - \hat{u}_i) = 0, \end{aligned}$$

the last equality following from (1.12). Hence with (1.13) holding for all  $u_i \in H_i, i = 1, 2, \dots, n$ , the inequalities (1.9) are satisfied. The completion of the argument involves showing  $\hat{u}_i, i = 1, 2, \dots, n$ , is the feedback response to  $(\mathcal{L}_1, \mathcal{L}_2, \dots, \mathcal{L}_n)$ . We compute

$$(1.14) \quad \begin{aligned} & (\mathcal{I} + \mathcal{B}_i^* \Omega_i \mathcal{B}_i) \left[ \hat{u}_i - \mathcal{L}_i \left( \mathcal{I} - \sum_1^n \mathcal{B}_j \mathcal{L}_j \right)^{-1} x_0 \right] \\ &= - \left[ \mathcal{B}_i^* \Omega_i + (\mathcal{I} + \mathcal{B}_i^* \Omega_i \mathcal{B}_i) \mathcal{L}_i \left( \mathcal{I} - \sum_1^n \mathcal{B}_j \mathcal{L}_j \right)^{-1} \right] x_0 \\ &= - \left[ \mathcal{B}_i^* \Omega_i \left( \mathcal{I} - \sum_1^n \mathcal{B}_j \mathcal{L}_j \right) + (\mathcal{I} + \mathcal{B}_i^* \Omega_i \mathcal{B}_i) \mathcal{L}_i \right] \\ &\quad \cdot \left( \mathcal{I} - \sum_1^n \mathcal{B}_j \mathcal{L}_j \right)^{-1} x_0 \\ &= - \left[ \mathcal{L}_i + \mathcal{B}_i^* \Omega_i \left( \mathcal{I} - \sum_{j \neq i}^n \mathcal{B}_j \mathcal{L}_j \right) \right] \left( \mathcal{I} - \sum_1^n \mathcal{B}_j \mathcal{L}_j \right)^{-1} x_0. \end{aligned}$$



Substitution of (1.6) into (1.14) and application of (1.10) gives

$$(1.15) \quad (\mathcal{I} + \mathcal{B}_i^* \Omega_i \mathcal{B}_i) \left[ \hat{u}_i - \mathcal{L}_i \left( \mathcal{I} - \sum_1^n \mathcal{B}_j \mathcal{L}_j \right)^{-1} x_0 \right] = 0.$$

The equation obtained upon cancellation of the nonsingular operator factor in (1.15) shows that  $(\mathcal{L}_1, \mathcal{L}_2, \dots, \mathcal{L}_n)$  meets the definition of an equilibrium given by (1.7)–(1.9).

To prove the converse let  $\mathcal{L} = (\mathcal{L}_1, \mathcal{L}_2, \dots, \mathcal{L}_n)$  satisfy (1.7)–(1.9). Since (1.9) holds for all  $u_i \in H_i$  it holds in particular for the unique  $u_i$  which minimizes the right-hand cost term of (1.9). Using formula (1.5) we can easily compute the minimizing control to be the control given by formula (1.11). With (1.5) we may rewrite (1.9) as

$$|x_0 + \mathcal{B}_i \hat{u}_i|_{\Omega_i}^2 - |x_0 + \mathcal{B}_i u_i|_{\Omega_i}^2 + |\hat{u}_i|^2 - |u_i|^2 \leq 0.$$

Then by substituting out  $u_i$  with (1.11) and  $\hat{u}_i = \mathcal{L}_i (\mathcal{I} - \sum_1^n \mathcal{B}_j \mathcal{L}_j)^{-1} x_0$  which is the solution of (1.7)–(1.8), a rather long calculation based upon the expansion of the two differences of squares transforms the previous inequality into

$$\left[ \mathcal{L}_i \left( \mathcal{I} - \sum_1^n \mathcal{B}_j \mathcal{L}_j \right)^{-1} + (\mathcal{I} + \mathcal{B}_i^* \Omega_i \mathcal{B}_i)^{-1} \mathcal{B}_i^* \Omega_i \right] x_0 \Big|_{(\mathcal{I} + \mathcal{B}_i^* \Omega_i \mathcal{B}_i)}^2 \leq 0.$$

Of course this implies that

$$\left[ \mathcal{L}_i \left( \mathcal{I} - \sum_1^n \mathcal{B}_j \mathcal{L}_j \right)^{-1} + (\mathcal{I} + \mathcal{B}_i^* \Omega_i \mathcal{B}_i)^{-1} \mathcal{B}_i^* \Omega_i \right] x_0 = 0.$$

Elimination of  $\Omega_i$  by (1.6) from this equation leads directly to (1.10). But  $\mathcal{L}$  was an arbitrary equilibrium. This completes the proof that the equilibrium points of the game coincide with the solutions of (1.10).

*Remark.* We note (1.10) says that in order for (1.1) to be in an equilibrium feedback configuration each  $\mathcal{L}_i$  must act upon the closed loop output to  $x_0$  exactly as the adjoint of  $P_i$ 's transfer function composed with  $-\mathcal{W}_i$ .

Thus the quest for an equilibrium leads to the problem of solving (1.10) for  $\mathcal{L} = (\mathcal{L}_1, \mathcal{L}_2, \dots, \mathcal{L}_n)$ . The following fixed-point theorem [5, p. 216] will be useful in deriving sufficient conditions for the existence of a solution.

**THEOREM 1.2.** *Let  $F$  be a Banach space,  $V$  an open ball in  $F$  of center  $y_0$  and radius  $\gamma$ . Let  $v$  be a mapping of  $V$  into  $F$  such that  $|v(y_1) - v(y_2)| \leq k|y_1 - y_2|$  for any pair of points  $y_1, y_2$  of  $V$ , where  $k$  is a constant such that  $0 \leq k < 1$ . Then, if  $|v(y_0) - y_0| < \gamma(1 - k)$ , there is one and only one point  $z \in V$  such that  $z = v(z)$ .*

*Notation.* To prevent any confusion arising from the variety of norms to be dealt with, a word about our notation seems in order. Norms of vectors induced by inner products on  $R^m$ ,  $H$  and  $H_i$  will be denoted as  $|y|$ ,  $|x|$  and  $|u_i|$ , etc. The norms of matrices and in general operators between Hilbert spaces will be denoted by double bars; e.g.  $\|\mathcal{B}_i\| = \sup_{|u_i|=1} |\mathcal{B}_i u_i|$ . However, for sums of spaces and operators we use norms such as  $|z| = \sum_1^m |z_k|$  and  $\|\mathcal{B}\| = \sum_1^n \|\mathcal{B}_i\|$  for  $z = \sum_1^m \oplus z_k \in \sum_1^m \oplus H$  or  $\mathcal{B} = \sum_1^n \oplus \mathcal{B}_i$ , an operator from  $\sum_1^n \oplus H$  into  $\sum_1^n \oplus H_i$ .

In some applications, (1.10) can be solved by setting the first operator factor equal to zero and solving

$$(1.16) \quad \mathcal{L}_i + \mathcal{B}_i^* \left( \mathcal{I} - \sum_{j \neq i}^n \mathcal{B}_j \mathcal{L}_j \right)^{-1*} \mathcal{W}_i = 0,$$

$i = 1, 2, \dots, n$ . Static games in which all operators are matrices often fall into this category. Other applications, e.g., differential games, will require that all factors in (1.10) be retained.

**THEOREM 1.3.** *A ball  $\|\mathcal{L}\| < \lambda$  in  $\sum_1^n \oplus L_i$  for which  $\mathcal{B}_i^*(\mathcal{I} - \sum_{j \neq i}^n \mathcal{B}_j \mathcal{L}_j)^{-1*} \cdot \mathcal{W}_i \in L_i, i = 1, 2, \dots, n$ , contains a solution to (1.16) and hence to (1.10) if*

$$(1.17) \quad \|\mathcal{B}\| \|\mathcal{W}\|^{1/2} < (1 - \|\mathcal{B}\| \lambda) \left( 1 - \frac{\|\mathcal{B}\| \|\mathcal{W}\|}{\lambda} \right)^{1/2}.$$

*Proof.* Consider  $v(\mathcal{L}) = (v_1, v_2, \dots, v_n)(\mathcal{L})$ , where

$$(1.18) \quad v_i(\mathcal{L}) = -\mathcal{B}_i^* \left( \mathcal{I} - \sum_{j \neq i}^n \mathcal{B}_j \mathcal{L}_j \right)^{-1*} \mathcal{W}_i.$$

For  $\mathcal{L}^1, \mathcal{L}^2$  in the ball  $\|\mathcal{L}\| < \lambda$ ,

$$(1.19) \quad \begin{aligned} &v_i(\mathcal{L}^1) - v_i(\mathcal{L}^2) \\ &= -\mathcal{B}_i^* \left( \mathcal{I} - \sum_{j \neq i}^n \mathcal{B}_j \mathcal{L}_j^1 \right)^{-1*} \sum_{j \neq i}^n (\mathcal{L}_j^1 - \mathcal{L}_j^2) \mathcal{B}_j^* \left( \mathcal{I} - \sum_{j \neq i}^n \mathcal{B}_j \mathcal{L}_j^2 \right)^{-1*} \mathcal{W}_i. \end{aligned}$$

Inequality (1.17) implies that  $\|\mathcal{B}\| \lambda < 1$  which in turn we use to estimate (1.19) as

$$(1.20) \quad \begin{aligned} \|v_i(\mathcal{L}^1) - v_i(\mathcal{L}^2)\| &\leq \frac{\|\mathcal{B}_i\| \|\mathcal{W}_i\| \sum_1^n \|\mathcal{L}_j^1 - \mathcal{L}_j^2\| \|\mathcal{B}_j\|}{(1 - \sum_1^n \|\mathcal{B}_j\| \|\mathcal{L}_j^1\|)(1 - \sum_1^n \|\mathcal{B}_j\| \|\mathcal{L}_j^2\|)} \\ &\leq \frac{\|\mathcal{B}_i\| \|\mathcal{W}_i\| \|\mathcal{B}\| \|\mathcal{L}^1 - \mathcal{L}^2\|}{(1 - \|\mathcal{B}\| \lambda)^2}. \end{aligned}$$

Summing (1.20) we have

$$(1.21) \quad \begin{aligned} \|v(\mathcal{L}^1) - v(\mathcal{L}^2)\| &\leq \sum_1^n \|\mathcal{B}_i\| \|\mathcal{W}_i\| \|\mathcal{B}\| \|\mathcal{L}^1 - \mathcal{L}^2\| / (1 - \|\mathcal{B}\| \lambda)^2 \\ &\leq \frac{\|\mathcal{B}\|^2 \|\mathcal{W}\|}{(1 - \|\mathcal{B}\| \lambda)^2} \|\mathcal{L}^1 - \mathcal{L}^2\| = k \|\mathcal{L}^1 - \mathcal{L}^2\|, \end{aligned}$$

where we note from (1.17) that

$$(1.22) \quad k = \frac{\|\mathcal{B}\|^2 \|\mathcal{W}\|}{(1 - \|\mathcal{B}\| \lambda)^2} < 1.$$

Thus we have shown

$$(1.23) \quad \|v(\mathcal{L}^1) - v(\mathcal{L}^2)\| \leq k \|\mathcal{L}^1 - \mathcal{L}^2\|$$

for  $\|\mathcal{L}^k\| < \lambda, k = 1, 2$ . Computing  $v_i(0) = -\mathcal{B}_i^* \mathcal{W}_i$  we find  $\|v_i(0)\| \leq \|\mathcal{B}_i\| \|\mathcal{W}_i\|$ , and hence

$$(1.24) \quad \|v(0)\| \leq \|\mathcal{B}\| \|\mathcal{W}\|.$$

Observe that (1.17) can be rewritten as

$$(1.25) \quad \|\mathcal{B}\| \|\mathcal{W}\| < \lambda \left[ 1 - \frac{\|\mathcal{B}\|^2 \|\mathcal{W}\|}{(1 - \|\mathcal{B}\| \lambda)^2} \right].$$

Combining (1.22), (1.24) and (1.25) shows

$$(1.26) \quad \|v(0) - 0\| \leq \lambda(1 - k).$$

With (1.23) and (1.26) established, the proof can now be completed by direct application of Theorem 1.2.

*Remark.* We add that Theorem 1.2 further implies that the solution satisfying the conditions of Theorem 1.3 is unique. However, in general, the solutions to (1.10) and (1.16) are not unique. For example, the scalar system  $\mathcal{B}_1 = \mathcal{W}_1 = 1$ ,  $\mathcal{B}_2 = \mathcal{W}_2 = 2$  gives rise to equations (1.16),

$$l_1 + \frac{1}{1 - 2l_2} = 0, \quad l_2 + \frac{4}{1 - l_1} = 0$$

having two solutions

$$l_1 = 4 \pm \sqrt{17}, \quad l_2 = -3/2 \pm \sqrt{17}/2.$$

Looking at (1.16) we would expect to find a solution near  $\mathcal{L} = 0$  if  $\|\mathcal{B}\|$  or  $\|\mathcal{W}\|$  is small. This is indeed verified by Theorem 1.3 since (1.17) will then have a solution  $\lambda > 0$ . Theorem 1.5 likewise represents an application of Theorem 1.2 with  $y_0 = 0$ . If we have prior knowledge about the location of a solution near a point other than the origin, then more appropriate results would be obtained using  $y_0 \neq 0$ .

LEMMA 1.4. *The operator function  $v(\mathcal{B}, \mathcal{L}) = (v_1, v_2, \dots, v_n)(\mathcal{B}, \mathcal{L})$  defined on the region  $\|\mathcal{B}\| \|\mathcal{L}\| < 1$  by*

$$v_i(\mathcal{B}, \mathcal{L}) = - \left[ \mathcal{L}_i \sum_1^n \mathcal{B}_j \mathcal{L}_j + \mathcal{B}_i^* \left( \mathcal{I} - \sum_{j \neq i}^n \mathcal{B}_j \mathcal{L}_j \right)^{-1*} \mathcal{W}_i \right] \left( \mathcal{I} - \sum_1^n \mathcal{B}_j \mathcal{L}_j \right)^{-1}$$

*satisfies the estimate*

$$(1.27) \quad \|v(\mathcal{B}, \mathcal{L}^1) - v(\mathcal{B}, \mathcal{L}^2)\| < \alpha \|\mathcal{B}\| \|\mathcal{L}^1 - \mathcal{L}^2\|$$

*for  $\|\mathcal{B}\| \leq \beta$ ,  $\|\mathcal{L}^1\| \leq \lambda$ ,  $\|\mathcal{L}^2\| < \lambda$ , where  $\beta, \lambda$  are any positive numbers such that  $\beta\lambda < 1$  and  $\alpha$  is the constant*

$$\alpha = \frac{3\lambda + 2\|\mathcal{W}\|\beta}{(1 - \beta\lambda)^3}.$$

*Proof.* As a notational convenience let  $\mathcal{R}_k = (\mathcal{I} - \sum_1^n \mathcal{B}_j \mathcal{L}_j^k)^{-1}$  and  $\mathcal{R}_{i,k} = (\mathcal{I} - \sum_{j \neq i}^n \mathcal{B}_j \mathcal{L}_j^k)^{-1}$  for an operator  $\mathcal{L}^k \in L$ ,  $i = 1, 2, \dots, n$ . The estimate

$$\begin{aligned} \left\| \left( \mathcal{I} - \sum_{j \neq i}^n \mathcal{B}_j \mathcal{L}_j \right)^{-1} \right\| &\leq \left( 1 - \sum_{j \neq i}^n \|\mathcal{B}_j\| \|\mathcal{L}_j\| \right)^{-1} \\ &\leq \left( 1 - \sum_1^n \|\mathcal{B}_j\| \|\mathcal{L}_j\| \right)^{-1} \leq (1 - \|\mathcal{B}\| \|\mathcal{L}\|)^{-1} \end{aligned}$$

implies that  $v$  is defined as an operator function on the region  $\|\mathcal{B}\|\|\mathcal{L}\| < 1$ , and in view of our assumptions that  $\|\mathcal{B}\| \leq \beta$ ,  $\|\mathcal{L}^1\| \leq \lambda$ ,  $\|\mathcal{L}^2\| \leq \lambda$ ,  $\beta\lambda < 1$  leads to the bounds

$$(1.28) \quad \|\mathcal{R}_{i,k}\| \leq (1 - \beta\lambda)^{-1}, \quad \|\mathcal{R}_k\| \leq (1 - \beta\lambda)^{-1}$$

for  $k = 1, 2; i = 1, 2, \dots, n$ .

Recalling the definitions of  $v_i$ ,  $\mathcal{R}_{i,k}$ ,  $\mathcal{R}_k$  we add and subtract to get

$$(1.29) \quad \begin{aligned} v_i(\mathcal{B}, \mathcal{L}^1) - v_i(\mathcal{B}, \mathcal{L}^2) &= \left[ \mathcal{L}_i^2 \sum_1^n \mathcal{B}_j \mathcal{L}_j^2 + \mathcal{B}_i^* \mathcal{R}_{i,2}^* \mathcal{W}_i \right] \mathcal{R}_2 \\ &\quad - \left[ \mathcal{L}_i^2 \sum_1^n \mathcal{B}_j \mathcal{L}_j^2 + \mathcal{B}_i^* \mathcal{R}_{i,2}^* \mathcal{W}_i \right] \mathcal{R}_1 \\ &\quad + \left[ \mathcal{L}_i^2 \sum_1^n \mathcal{B}_j \mathcal{L}_j^2 + \mathcal{B}_i^* \mathcal{R}_{i,2}^* \mathcal{W}_i \right] \mathcal{R}_1 \\ &\quad - \left[ \mathcal{L}_i^1 \sum_1^n \mathcal{B}_j \mathcal{L}_j^1 + \mathcal{B}_i^* \mathcal{R}_{i,1}^* \mathcal{W}_i \right] \mathcal{R}_1 \\ &= \left[ \mathcal{L}_i^2 \sum_1^n \mathcal{B}_j \mathcal{L}_j^2 + \mathcal{B}_i^* \mathcal{R}_{i,2}^* \mathcal{W}_i \right] [\mathcal{R}_2 - \mathcal{R}_1] \\ &\quad + \left[ (\mathcal{L}_i^2 - \mathcal{L}_i^1) \sum_1^n \mathcal{B}_j \mathcal{L}_j^2 + \mathcal{L}_i^1 \sum_1^n \mathcal{B}_j (\mathcal{L}_j^2 - \mathcal{L}_j^1) \right] \mathcal{R}_1 \\ &\quad + \mathcal{B}_i^* [\mathcal{R}_{i,2}^* - \mathcal{R}_{i,1}^*] \mathcal{W}_i \mathcal{R}_1. \end{aligned}$$

Rewriting the differences, we have

$$(1.30) \quad \mathcal{R}_2 - \mathcal{R}_1 = \mathcal{R}_2 \sum_1^n \mathcal{B}_j (\mathcal{L}_j^2 - \mathcal{L}_j^1) \mathcal{R}_1,$$

$$(1.31) \quad \mathcal{R}_{i,2}^* - \mathcal{R}_{i,1}^* = \mathcal{R}_{i,2}^* \sum_{j \neq i}^n (\mathcal{L}_j^2 - \mathcal{L}_j^1) \mathcal{B}_j^* \mathcal{R}_{i,1}^*$$

in (1.29), and using (1.28) gives the upper estimate of  $\|v_i(\mathcal{B}, \mathcal{L}^1) - v_i(\mathcal{B}, \mathcal{L}^2)\|$ ,

$$(1.32) \quad \begin{aligned} &\left[ \|\mathcal{L}_i^2\| \beta\lambda + \|\mathcal{B}_i\| \left( \frac{1}{1 - \beta\lambda} \right) \|\mathcal{W}\| \right] \frac{\|\mathcal{B}\| \|\mathcal{L}^1 - \mathcal{L}^2\|}{(1 - \beta\lambda)^3} \\ &+ [\|\mathcal{L}_i^2 - \mathcal{L}_i^1\| \|\mathcal{B}\| \lambda + \|\mathcal{L}_i^1\| \|\mathcal{B}\| \|\mathcal{L}^2 - \mathcal{L}^1\|] \left( \frac{1}{1 - \beta\lambda} \right) \\ &+ \|\mathcal{B}_i\| \left( \frac{1}{1 - \beta\lambda} \right)^3 \|\mathcal{B}\| \|\mathcal{L}^2 - \mathcal{L}^1\| \|\mathcal{W}_i\|. \end{aligned}$$

But  $\beta\lambda < 1$  and hence  $1/(1 - \beta\lambda) > 1$ . Consequently (1.32) produces the summed

upper estimate of  $\|v(\mathcal{B}, \mathcal{L}^1) - v(\mathcal{B}, \mathcal{L}^2)\|$ ,

$$(1.33) \quad \left[ \frac{\lambda(1 - \beta\lambda) + \beta\|\mathcal{W}\| + 2\lambda + \beta\|\mathcal{W}\|}{(1 - \beta\lambda)^3} \right] \|\mathcal{B}\| \|\mathcal{L}^1 - \mathcal{L}^2\| \\ \leq \left[ \frac{3\lambda + 2\|\mathcal{W}\|\beta}{(1 - \beta\lambda)^3} \right] \|\mathcal{B}\| \|\mathcal{L}^1 - \mathcal{L}^2\|$$

which completes the proof.

With each vector  $Z = \sum_1^m \oplus z_k$  in  $\sum_1^m \oplus H$  ( $m$  any positive integer) we associate a linear transformation  $\mathcal{T}_Z(\cdot): L \rightarrow \sum_{i=1}^n \oplus \sum_{k=1}^m \oplus H_i$  given by

$$(1.34) \quad \mathcal{T}_Z(\mathcal{L}) = \sum_1^n \oplus \mathcal{T}_Z^i(\mathcal{L}) \quad \text{with} \quad \mathcal{T}_Z^i(\mathcal{L}) = \sum_{k=1}^m \oplus \mathcal{L}_i z_k.$$

In the next theorem we again consider the operator function

$$(1.35) \quad v(\mathcal{B}, \mathcal{L}) = \sum_1^n \oplus v_i(\mathcal{B}, \mathcal{L}),$$

where

$$v_i(\mathcal{B}, \mathcal{L}) = - \left[ \mathcal{L}_i \sum_1^n \mathcal{B}_j \mathcal{L}_j + \mathcal{B}_i^* \left( \mathcal{I} - \sum_{j \neq i} \mathcal{B}_j \mathcal{L}_j \right)^{-1*} \mathcal{W}_i \right] \left( \mathcal{I} - \sum_1^n \mathcal{B}_j \mathcal{L}_j \right)^{-1}.$$

**THEOREM 1.5.** *Suppose  $\mathcal{T}_Z(\cdot)$  is a homeomorphism whose range  $\mathcal{T}_Z(L)$  contains the vectors  $v(\mathcal{B}, \mathcal{L})Z$  for all  $\mathcal{L}$  in the ball  $\|\mathcal{L}\| < \lambda$ . If  $\|\mathcal{B}\| < \beta$  where both*

$$(1.36) \quad \beta\lambda < 1$$

and

$$(1.37) \quad \beta|Z| \|\mathcal{T}_Z^{-1}\| \{ (1 - \beta\lambda)^3 \|\mathcal{W}\| + \lambda[3\lambda + 2\|\mathcal{W}\|\beta] \} < \lambda(1 - \beta\lambda)^3,$$

then there exists a solution  $\mathcal{L}$  in  $L$  to (1.10) with  $\|\mathcal{L}\| < \lambda$ , holding for all  $x_0 \in \text{span} \{z_1, z_2, \dots, z_m\}$ .

*Proof.* By adding and subtracting the term  $\mathcal{L}_i$  and simplifying, we can rewrite (1.10) in the equivalent form

$$(1.38) \quad \mathcal{L}_i x_0 = v_i(\mathcal{B}, \mathcal{L}) x_0.$$

Since we are assuming that  $x_0 \in \text{span} \{z_1, z_2, \dots, z_m\}$ ,  $x_0 = \sum_1^m c_k z_k$  for scalars  $c_k$ . Therefore it certainly is sufficient to solve

$$(1.39) \quad \sum_{k=1}^m \oplus \mathcal{L}_i z_k = \sum_{k=1}^m \oplus v_i(\mathcal{B}, \mathcal{L}) z_k, \quad i = 1, 2, \dots, n.$$

Written in terms of  $\mathcal{T}_Z(\cdot)$  defined by (1.34) this system is the same as the equation

$$(1.40) \quad \mathcal{T}_Z(\mathcal{L}) = \sum_{i=1}^n \oplus \sum_{k=1}^m \oplus v_i(\mathcal{B}, \mathcal{L}) z_k = \sum_{i=1}^n \oplus v_i(\mathcal{B}, \mathcal{L}) \sum_{k=1}^m \oplus z_k \\ = v(\mathcal{B}, \mathcal{L}) Z.$$

From (1.34) we observe that the linear transformation  $\mathcal{T}_Z(\cdot)$  is automatically

bounded with  $\|\mathcal{T}_Z\| \leq |Z|$ . Since  $\mathcal{T}_Z(\cdot)$  is assumed to be a homeomorphism,  $\mathcal{T}_Z^{-1}$  is likewise a bounded linear transformation on its domain  $\mathcal{T}_Z(L)$ . Hence to prove the theorem we may deal with the equation

$$(1.41) \quad \mathcal{L} = \mathcal{T}_Z^{-1}v(\mathcal{B}, \mathcal{L})Z \stackrel{\text{def}}{=} \tilde{v}(\mathcal{B}, \mathcal{L}).$$

With Lemma 1.4 we estimate

$$(1.42) \quad \begin{aligned} \|\tilde{v}(\mathcal{B}, \mathcal{L}^1) - \tilde{v}(\mathcal{B}, \mathcal{L}^2)\| &\leq \|\mathcal{T}_Z^{-1}\| |Z| \|v(\mathcal{B}, \mathcal{L}^1) - v(\mathcal{B}, \mathcal{L}^2)\| \\ &\leq \|\mathcal{T}_Z^{-1}\| |Z| \beta \left[ \frac{3\lambda + 2\|\mathcal{W}\|\beta}{(1 - \beta\lambda)^3} \right] \|\mathcal{L}^1 - \mathcal{L}^2\| \\ &= k \|\mathcal{L}^1 - \mathcal{L}^2\|, \end{aligned}$$

with  $k = \|\mathcal{T}_Z^{-1}\| |Z| \beta [(3\lambda + 2\|\mathcal{W}\|\beta)/(1 - \beta\lambda)^3]$  holding for  $\|\mathcal{B}\| < \beta$ ,  $\|\mathcal{L}^1\| < \lambda$ ,  $\|\mathcal{L}^2\| < \lambda$  and  $\beta\lambda < 1$ . Inequalities (1.36)–(1.37) imply that  $k < 1$ . Using (1.41) and (1.35) we have

$$(1.43) \quad \begin{aligned} \|\tilde{v}(\mathcal{B}, 0) - 0\| &= \|\mathcal{T}_Z^{-1}v(\mathcal{B}, 0)Z\| \\ &\leq \|\mathcal{T}_Z^{-1}\| |Z| \sum_1^n \|v_i(\mathcal{B}, 0)\| \leq \|\mathcal{T}_Z^{-1}\| |Z| \sum_1^n \|\mathcal{B}_i^* \mathcal{W}_i\| \\ &\leq \|\mathcal{T}_Z^{-1}\| |Z| \|\mathcal{W}\| \beta \leq \lambda(1 - k), \end{aligned}$$

where the last inequality follows from (1.36) and division of (1.37) by  $(1 - \beta\lambda)^3$ . The remainder of the proof is a consequence of Theorem 1.2, which, in view of (1.42)–(1.43), can be applied to  $\tilde{v}$ .

*Remark.* Note that both (1.36) and (1.37) could be satisfied if  $\beta$  could independently be made small. This fact will be exploited in establishing the “local playability” of dynamic games.

**2. Application to dynamic games.** We are now interested in applying the results of the previous section to more concrete problems. As an example we choose a particular differential game. In the one player case it reduces to the optimal regulator problem which has been studied extensively in control theory [6], [7], [8] and has had useful applications in control engineering.

**2.1. A linear-quadratic differential game.** The data of the problem is given in the form of real matrix-valued functions  $A(t)$ ,  $B_i(t)$ ,  $W_i(t) \geq 0$ ,  $\hat{W}_i \geq 0$ ,  $U_i(t) > 0$ , and real, finite-dimensional, vector-valued functions  $\bar{y}_i(t)$ ,  $\hat{y}_i$ ,  $i = 1, 2, \dots, n$ , all defined and continuous on a finite interval  $t_0 \leq t \leq t_1$ .

The dynamics of the game are given by the linear differential equation in  $R^m \times [t_0, t_1]$ ,

$$(2.1) \quad \frac{dy}{dt} = A(t)y + \sum_1^n B_f(t)u_j,$$

where  $u_i$  is the  $m_i$ -dimensional control variable under the command of the  $i$ th player  $P_i$ .  $P_i$  is allowed to play any real, Borel measurable control command  $u_i = u_i(t)$  satisfying

$$(2.2) \quad \int_{t_0}^{t_1} |u_i(t)|^2 dt < \infty.$$

With the initial state of the game  $y(t_0) = y_0$  fixed, any such  $n$ -tuple of control commands  $u(t) = (u_1, u_2, \dots, u_n)(t)$  is called an *open loop strategy*. An open loop strategy together with the unique absolutely continuous solution of (2.1) that it generates will be called a *play* of the game.

Player  $P_i$  is also allowed to generate control commands by means of *feedback controls*. In the particular game that we are considering we allow  $u_i = \mathcal{L}_i(t)y$ , where  $\mathcal{L}_i(t)$  is any real, continuous  $m_i \times m$  matrix function. Hence  $P_i$  automatically generates control commands  $u_i(t) = \mathcal{L}_i(t)y(t)$  in response to the action of the other players via  $y(t)$ , the solution to

$$(2.3) \quad \frac{dy}{dt} = [A(t) + B_i(t)\mathcal{L}_i(t)]y + \sum_{j \neq i}^n B_j(t)u_j(t).$$

In a similar manner any number of players can choose to play feedback control simultaneously.

$P_i$  determines his cost incurred in a play of the game by a functional of the form

$$(2.4) \quad \mathcal{C}_i = \int_{t_0}^{t_1} [|y(t) - \bar{y}_i(t)|_{\bar{w}_i(t)}^2 + |u_i(t)|_{\bar{v}_i(t)}^2] dt + |y(t_1) - \hat{y}_i|_{\bar{w}_i}^2$$

which he would like to minimize.

**2.2. The differential game as a game in Hilbert space.** In order to cast this game into the framework of § 1 we consider the Hilbert spaces

$$(2.5) \quad L_j^2 = \left\{ f(\cdot): [t_0, t_1] \rightarrow R^j, \int_{t_0}^{t_1} |f(t)|^2 dt < \infty \right\}$$

and rewrite (2.1) using the variations of parameters formula

$$(2.6) \quad y(t) = \phi_0(t, t_0)y_0 + \sum_1^n \int_{t_0}^t \phi_0(t, \tau)B_j(\tau)u_j(\tau) d\tau,$$

where the fundamental matrix  $\phi_0(t, \tau)$  is defined as the solution to

$$(2.7) \quad \frac{\partial \phi_0(t, \tau)}{\partial t} = A(t)\phi_0(t, \tau), \quad \phi_0(\tau, \tau) = I_m.$$

We select

$$(2.8) \quad H = L_m^2 \oplus R^m, \quad H_i = L_{m_i}^2, \quad i = 1, 2, \dots, n,$$

and define vectors  $x, x_0, \bar{x}_i, u_i$  by

$$(2.9) \quad x = \begin{pmatrix} y(t) \\ y(t_1) \end{pmatrix} \in H, \quad x_0 = \begin{pmatrix} \phi_0(t, t_0)y_0 \\ \phi_0(t_1, t_0)y_0 \end{pmatrix} \in H,$$

$$\bar{x}_i = \begin{pmatrix} \bar{y}_i(t) \\ \hat{y}_i \end{pmatrix} \in H, \quad u_i = u_i(t) \in H_i,$$

and operators  $\mathcal{B}_i: H_i \rightarrow H$  by

$$(2.10) \quad (\mathcal{B}_i u_i)(t) = \begin{pmatrix} \int_{t_0}^t \phi_0(t, \tau) B_i(\tau) u_i(\tau) d\tau \\ \int_{t_0}^{t_1} \phi_0(t_1, \tau) B_i(\tau) u_i(\tau) d\tau \end{pmatrix},$$

$\mathcal{W}_i: H \rightarrow H$  by

$$(2.11) \quad (\mathcal{W}_i x)(t) = \begin{pmatrix} W_i(t)y(t) \\ \hat{W}_i y(t_1) \end{pmatrix},$$

and finally  $\mathcal{U}_i: H_i \rightarrow H_i$  by

$$(2.12) \quad (\mathcal{U}_i u_i)(t) = U_i(t)u_i(t), \quad i = 1, 2, \dots, n.$$

To complete the identification of the differential game with the abstract model of § 1 we must single out the  $L_i$ , the Banach spaces of operator feedbacks. Since the game allows feedback controls of the type  $u_i = \mathcal{L}_i(t)y$  it is clear we must select the linear space of matrix operators corresponding to multiplication by  $m_i \times 2m$  matrices of the type  $[\mathcal{L}_i(t) \ ; \ \mathbf{0}]$  with  $\mathcal{L}_i(t)$  of size  $m_i \times m$  and continuous and  $\mathbf{0}$  a zero matrix. For efficiency in notation we simply write  $\mathcal{L}_i(t)$  and for later reference

$$(2.13) \quad L_i = \{ \mathcal{L}_i(t) | m_i \times m, \text{ continuous} \},$$

$i = 1, 2, \dots, n$ . For the purpose of turning  $L_i$  into a complete normed linear space we consider the norms

$$\| \mathcal{L}_i \| = \sup_{|x|=1} | \mathcal{L}_i x | = \sup_{|x|=1} \left( \int_{t_0}^{t_1} | \mathcal{L}_i(t)x(t) |^2 dt \right)^{1/2}$$

called the *uniform operator norm* and

$$\| \mathcal{L}_i \|' = \sup_{[t_0, t_1]} \| \mathcal{L}_i^*(t) \mathcal{L}_i(t) \|^{1/2}$$

called the *uniform norm*.

LEMMA 2.1. *The uniform norm and uniform operator norm agree on the space of matrix operators (2.13) and make  $L_i$  into a Banach space.*

*Proof.* It will be sufficient to compute the norms of the operators in  $L_i$ . Let  $0 \neq \mathcal{L}_i \in L_i$ ,  $x \in L_m^2$ , and being careful not to confuse the various vector and matrix norms involved we compute

$$(2.14) \quad \begin{aligned} | \mathcal{L}_i x |^2 &= \int_{t_0}^{t_1} | \mathcal{L}_i(t)x(t) |^2 dt \leq \int_{t_0}^{t_1} \| \mathcal{L}_i^*(t) \mathcal{L}_i(t) \| |x(t)|^2 dt \\ &\leq \sup_{[t_0, t_1]} \| \mathcal{L}_i^*(t) \mathcal{L}_i(t) \| \int_{t_0}^{t_1} |x(t)|^2 dt \\ &= \sup_{[t_0, t_1]} \| \mathcal{L}_i^*(t) \mathcal{L}_i(t) \| |x|^2 \end{aligned}$$

which shows that  $\| \mathcal{L}_i \| \leq \sup_{[t_0, t_1]} \| \mathcal{L}_i^*(t) \mathcal{L}_i(t) \|^{1/2}$ .



Since  $\|\mathcal{L}_i^*(t)\mathcal{L}_i(t)\|$  is a continuous function of  $t$  on  $[t_0, t_1]$  there exists a  $t_* \in [t_0, t_1]$  at which it takes on its maximum. Since  $\mathcal{L}_i^*(t_*)\mathcal{L}_i(t_*)$  is a positive semidefinite symmetric matrix, there exists an eigenvector  $x_* \in R^m$  with  $|x_*| = 1$  such that

$$(2.15) \quad \|\mathcal{L}_i^*(t_*)\mathcal{L}_i(t_*)\| = x_* \cdot \mathcal{L}_i^*(t_*)\mathcal{L}_i(t_*)x_*$$

We first examine the case where  $t_* \in (t_0, t_1)$ . Consider the sequence in  $L_m^2$ ,

$$(2.16) \quad x_n(t) = \sqrt{n}\chi_n(t)x_*, \quad n = 1, 2, \dots,$$

in which  $\chi_n(t)$  is the characteristic function of the interval  $[t_*, t_* + 1/n] \cap [t_0, t_1]$ . For this sequence we compute

$$(2.17) \quad |x_n|^2 = \int_{t_0}^{t_1} |x_n(t)|^2 dt = n \int_{t_0}^{t_1} \chi_n(t) dt \leq 1$$

for  $n = 1, 2, \dots$  (with equality holding for all large  $n$ ) and

$$(2.18) \quad \begin{aligned} |\mathcal{L}_i x_n|^2 &= \int_{t_0}^{t_1} |\mathcal{L}_i(t)x_n(t)|^2 dt = n \int_{t_0}^{t_1} \chi_n(t) |\mathcal{L}_i(t)x_*|^2 dt \\ &= \int_{t_*}^{t_* + 1/n} |\mathcal{L}_i(t)x_*|^2 dt / (1/n) \end{aligned}$$

for all large  $n$ . Using (2.18) and the continuity of  $|\mathcal{L}_i(t)x_*|$  we conclude

$$(2.19) \quad \lim_{n \rightarrow \infty} |\mathcal{L}_i x_n|^2 = |\mathcal{L}_i(t_*)x_*|^2.$$

From (2.17), (2.19) and (2.15) we get

$$(2.20) \quad \|\mathcal{L}_i\|^2 = \sup_{|x| \leq 1} |\mathcal{L}_i x|^2 \geq \lim_{n \rightarrow \infty} |\mathcal{L}_i x_n|^2 = \|\mathcal{L}_i^*(t_*)\mathcal{L}_i(t_*)\|$$

which shows

$$(2.21) \quad \|\mathcal{L}_i\| \geq \sup_{[t_0, t_1]} \|\mathcal{L}_i^*(t)\mathcal{L}_i(t)\|^{1/2}$$

and hence that

$$(2.22) \quad \|\mathcal{L}_i\| = \sup_{[t_0, t_1]} \|\mathcal{L}_i^*(t)\mathcal{L}_i(t)\|^{1/2}.$$

A similar argument leads to (2.22) in the cases  $t_* = t_0$  and  $t_* = t_1$ . But (2.22) shows that the norm generating the uniform operator topology on  $L_i$  agrees with the norm generating the uniform metric topology on  $L_i$ . The completeness of the latter space is well known. This concludes the proof.

With definitions (2.5)–(2.13) we may write (2.4) as

$$(2.23) \quad \mathcal{C}_i = |x - \bar{x}_i|_{W_i}^2 + |u_i|_{\mathcal{U}_i}^2$$

and have shown that the differential game given by (2.1), (2.4) and (2.13) is an example of the abstract model (1.1), (1.2') presented in §1. Therefore all the results developed there apply in particular to the differential game. Our attention is now focused upon the application of those results.

**2.3. Local playability of the game.** Let  $[t_0, t_1]$  denote any compact subinterval of some underlying interval  $[t_\alpha, t_\beta]$  on which  $A(t), B_i(t), W_i(t) \geq 0, U_i(t) > 0$  and  $\bar{y}_i(t)$  are defined and continuous. The differential game is called *playable* on  $[t_0, t_1]$  if there exists a system of equilibrium feedback matrices in the corresponding  $L = \sum_1^n \oplus L_i$  for that interval. (Recall definitions (1.9) and (2.13).) If there exists a number  $\Delta > 0$  such that the game is playable on every subinterval  $[t_0, t_1]$  of length  $t_1 - t_0 \leq \Delta$ , then we say the game is *locally playable* on  $[t_\alpha, t_\beta]$ . The next lemma will play an important role in proving the differential game is always locally playable.

LEMMA 2.2. For  $\mathcal{B}_i$  as defined by (2.10),

$$(2.24) \quad \|\mathcal{B}\| \leq \gamma(t_1 - t_0)^{1/2},$$

where  $\gamma$  is a number giving the bound

$$(2.25) \quad \|\phi_0(t, \tau)B_i(\tau)\| \leq \frac{\gamma}{(t_\beta - t_\alpha)^{1/2} + 1}$$

for  $t_\alpha \leq t, \tau \leq t_\beta, i = 1, 2, \dots, n$ .

*Proof.* Being careful not to confuse the norms involved we estimate

$$(2.26) \quad \begin{aligned} |\mathcal{B}_i u_i| &= \left[ \int_{t_0}^{t_1} \left| \int_{t_0}^t \phi_0(t, \tau)B_i(\tau)u_i(\tau) d\tau \right|^2 dt \right]^{1/2} \\ &\quad + \left| \int_{t_0}^{t_1} \phi_0(t_1, \tau)B_i(\tau)u_i(\tau) d\tau \right| \\ &\leq \left[ \int_{t_0}^{t_1} \left( \int_{t_0}^t \|\phi_0(t, \tau)B_i(\tau)\| |u_i(\tau)| d\tau \right)^2 dt \right]^{1/2} \\ &\quad + \int_{t_0}^{t_1} \|\phi_0(t_1, \tau)B_i(\tau)\| |u_i(\tau)| d\tau \\ &\leq \frac{\gamma}{(t_\beta - t_\alpha)^{1/2} + 1} [(t_1 - t_0)^{1/2} + 1] \int_{t_0}^{t_1} |u_i(\tau)| d\tau \\ &\leq \gamma \int_{t_0}^{t_1} |u_i(\tau)| d\tau \leq \gamma(t_1 - t_0)^{1/2} |u_i|. \end{aligned}$$

Summing we get  $|\mathcal{B}u| \leq \gamma(t_1 - t_0)^{1/2}|u|$  and hence

$$(2.27) \quad \|\mathcal{B}\| = \sup_{|u|=1} |\mathcal{B}u| \leq \gamma(t_1 - t_0)^{1/2}$$

which completes the proof.

It will be useful now to compute some of the operator functions which arose in the general theory of § 1. To accomplish this end let  $(\mathcal{L}_1, \mathcal{L}_2, \dots, \mathcal{L}_n)(t)$  be a system of feedback matrices from  $L = \sum_1^n \oplus L_i$ . Substituting  $u_j = \mathcal{L}_j(t)y$  ( $j \neq i$ ) into (2.1) produces the equation

$$(2.28) \quad \dot{y}_i = \left[ A(t) + \sum_{j \neq i}^n B_j(t)\mathcal{L}_j(t) \right] y_i + B_i(t)u_i$$

whose solution can be written as

$$(2.29) \quad y_i(t) = \phi_i(t, t_0)y_0 + \int_{t_0}^t \phi_i(t, \tau)B_i(\tau)u_i(\tau) d\tau,$$

where the fundamental matrix  $\phi_i(t, \tau)$  is defined as the solution to

$$(2.30) \quad \frac{\partial \phi_i(t, \tau)}{\partial t} = \left[ A(t) + \sum_{j \neq i}^n B_j(t)\mathcal{L}_j(t) \right] \phi_i(t, \tau)$$

with  $\phi_i(t, t) = I_m$  for all  $t, \tau \in [t_\alpha, t_\beta]$ ,  $i = 1, 2, \dots, n$ . Equation (2.29) shows that

$$(2.31) \quad \left[ \left( \mathcal{I} - \sum_{j \neq i}^n \mathcal{B}_j \mathcal{L}_j \right)^{-1} B_i u_i \right](t) = \begin{pmatrix} \int_{t_0}^t \phi_i(t, \tau) B_i(\tau) u_i(\tau) d\tau \\ \int_{t_0}^{t_1} \phi_i(t_1, \tau) B_i(\tau) u_i(\tau) d\tau \end{pmatrix},$$

$i = 1, 2, \dots, n$ , and consequently

$$(2.32) \quad \left[ \mathcal{B}_i^* \left( \mathcal{I} - \sum_{j \neq i}^n \mathcal{B}_j \mathcal{L}_j \right)^{-1} x \right](t) = \int_t^{t_1} B_i^*(t) \phi_i^*(\tau, t) y(\tau) d\tau + B_i^*(t) \phi_i^*(t, t_1) y(t_1)$$

for all

$$x(t) = \begin{pmatrix} y(t) \\ y(t_1) \end{pmatrix} \in H = L_m^2(t_0, t_1) \oplus R^m.$$

Recall that this operator function, as well as the following one, plays an important role in (1.10).

When all players play feedback,  $u_j = \mathcal{L}_j(t)y$ , in the integral form of

$$(2.33) \quad \dot{y} = A(t)y + \sum_1^n B_j(t)u_j,$$

$$(2.34) \quad y(t) = \phi_0(t, t_0)y_0 + \sum_1^n \int_{t_0}^t \phi_0(t, \tau) B_j(\tau) u_j(\tau) d\tau,$$

we have the solution  $y(t) = \phi(t, t_0)y_0$  in which  $\phi(t, \tau)$  is the fundamental matrix solution to

$$(2.35) \quad \frac{\partial \phi(t, \tau)}{\partial t} = \left[ A(t) + \sum_1^n B_j(t)\mathcal{L}_j(t) \right] \phi(t, \tau)$$

with  $\phi(t, t) = I_m$  for all  $t, \tau \in [t_\alpha, t_\beta]$ . This shows that

$$(2.36) \quad \left( \mathcal{I} - \sum_1^n \mathcal{B}_j \mathcal{L}_j \right)^{-1} x_0 = \begin{pmatrix} \phi(t, t_0)y_0 \\ \phi(t_1, t_0)y_0 \end{pmatrix},$$

where

$$x_0 = \begin{pmatrix} \phi_0(t, t_0)y_0 \\ \phi_0(t_1, t_0)y_0 \end{pmatrix} \in H.$$

Recall that we made the identification  $\mathcal{L}_i \leftrightarrow$  multiplication by  $[\mathcal{L}_i(t) \ ; \ \mathbf{0}]$ . In terms of the concrete operator functions given by (2.32) and (2.36), equation (1.10) for the system of equilibrium feedback matrices is

$$(2.37) \quad \mathcal{L}_i(t)\phi(t, t_0)y_0 + \int_t^{t_1} B_i^*(t)\phi_i^*(\tau, t)W_i(\tau)\phi(\tau, t_0)y_0 \, d\tau + B_i^*(t)\phi_i^*(t, t_1)\widehat{W}_i\phi(t_1, t_0)y_0 = 0,$$

$i = 1, 2, \dots, n$ . Note that (2.37) is a complicated functional equation in  $\mathcal{L}$  since  $\phi_i$  and  $\phi$  depend upon  $\mathcal{L}$  according to their defining equation (2.30) and (2.35). Nevertheless, as stated in the next theorem, a solution does exist—at least on short time intervals.

**THEOREM 2.3.** *The differential game with dynamics (2.1), costs (2.4) and feedback spaces (2.13) is locally playable on  $[t_\alpha, t_\beta]$  with equilibrium feedback matrices which are independent of the initial state of the game. Moreover, such state-independent matrices are unique.*

*Proof.* We assume that the appropriate preliminary change of control variables has been made so that we can take  $U_i(t) = I, i = 1, 2, \dots, n$ , and for simplicity treat the case in which  $\bar{y}_i(t) = \hat{y}_i = 0$ . The proof will follow from Theorem 1.5 and Lemma 2.2.

Recall that we are dealing with the Hilbert spaces

$$H = L_m^2(t_0, t_1) \oplus R^m, \quad H_i = L_{m_i}^2(t_0, t_1)$$

and feedback operators identified as matrices  $[\mathcal{L}_i(t) \ ; \ \mathbf{0}] \in L_i$ . The “abstract initial state” is given as

$$(2.38) \quad x_0 = \begin{pmatrix} \phi_0(t, t_0)y_0 \\ \phi_0(t_1, t_0)y_0 \end{pmatrix} = \begin{pmatrix} \phi_0(t, t_0) & \mathbf{0} \\ \mathbf{0} & \phi_0(t_1, t_0) \end{pmatrix} \begin{pmatrix} y_0 \\ y_0 \end{pmatrix}.$$

In anticipation of applying Theorem 1.5 we consider  $Z = \sum_1^{2m} \oplus z_k$ , where  $z_k$  is the  $k$ th column of the  $2m \times 2m$  matrix in (2.38). Suppose that we identify  $Z$  with that matrix. Then define  $\mathcal{T}_Z$  in (1.32) according to the right matrix multiplication

$$(2.39) \quad \mathcal{T}_Z^i(\mathcal{L})(t) = [\mathcal{L}_i(t) \ ; \ \mathbf{0}]Z(t), \quad i = 1, 2, \dots, n.$$

Since  $Z(t)$  is continuous and nonsingular it follows that  $\mathcal{T}_Z(\cdot)$  is a homeomorphism and  $\mathcal{T}_Z(L) = L$ . In fact,  $\mathcal{T}_Z^{-1} = \mathcal{T}_{Z^{-1}}$ . Furthermore, from (2.32) and (2.36) we see that

$$(2.40) \quad v(\mathcal{B}, L)Z \subset L = \mathcal{T}_Z(L).$$

Remembering that the subinterval  $[t_0, t_1]$  is arbitrary in  $[t_\alpha, t_\beta]$  we conclude from continuity that  $|Z|$  and  $\|\mathcal{T}_Z^{-1}\|$  are bounded when regarded as functions of  $t_0, t_1 \in [t_\alpha, t_\beta]$ . Since it can easily be shown that

$$(2.41) \quad \|\mathcal{W}_i\| \leq \sup_{[t_0, t_1]} \|W_i(t)\| + \|\widehat{W}_i\|,$$

it follows that  $\|\mathcal{W}\|$  is also bounded as a function of  $t_0, t_1 \in [t_\alpha, t_\beta]$ . Hence for any  $\lambda > 0$  the hypothesis of Theorem 1.5 can be satisfied by first selecting  $\beta > 0$

sufficiently small so that (1.36)–(1.37) hold for all  $t_0, t_1 \in [t_\alpha, t_\beta]$ . In view of Lemma 2.2 we can then maintain  $\|\mathcal{B}\| < \beta$  by restricting  $\Delta = t_1 - t_0 < \beta^2/\gamma^2$ . The state-independence is clear since  $x_0 \in \text{span}\{z_1, z_2, \dots, z_{2m}\}$  for all  $y_0 \in R^m$ . The uniqueness argument will rely partly upon the calculations of the next section and is treated in the concluding comments.

*Remark.* By carrying out the last steps in the proof of Theorem 2.3 numerically one can compute the length of intervals  $\Delta = t_1 - t_0$  over which the game with constraint  $\|\mathcal{L}\| < \lambda$  is playable. In doing this the inequalities corresponding to (1.36)–(1.37) which must be solved can be somewhat simplified by, for example, replacing the first occurrence of the term  $(1 - \beta\lambda)^3$  by 1, etc.

**2.4. Calculation of the equilibrium feedback.** The results up to this point show (2.37), the concrete form of (1.10), has a solution in  $L$  which moreover is a solution to

$$(2.42) \quad \begin{aligned} \mathcal{L}_i(t)\phi(t, t_0) + \int_t^{t_1} B_i^*(t)\phi_i^*(\tau, t)W_i(\tau)\phi(\tau, t_0) d\tau \\ + B_i^*(t)\phi_i^*(t_1, t)\widehat{W}_i\phi(t_1, t_0) = 0, \end{aligned}$$

$i = 1, 2, \dots, n$ , for  $\Delta = t_1 - t_0$  sufficiently small. We rewrite this equation as

$$(2.43) \quad \mathcal{L}_i(t) = -B_i^*(t)Q_i(t),$$

where

$$\begin{aligned} Q_i(t) = \int_t^{t_1} \phi_i^*(\tau, t)W_i(\tau)\phi(\tau, t_0)\phi^{-1}(t, t_0) d\tau \\ + \phi_i^*(t_1, t)\widehat{W}_i\phi(t_1, t_0)\phi^{-1}(t, t_0), \end{aligned}$$

$i = 1, 2, \dots, n$ . Using elementary properties of fundamental matrices one can easily show that the matrix product  $\phi(\tau, t_0)\phi^{-1}(t, t_0)$  is independent of  $t_0$ , and hence, as already anticipated in our notation, upon setting  $t_0 = t$  we have

$$(2.44) \quad Q_i(t) = \int_t^{t_1} \phi_i^*(\tau, t)W_i(\tau)\phi(\tau, t) d\tau + \phi_i^*(t_1, t)\widehat{W}_i\phi(t_1, t),$$

$i = 1, 2, \dots, n$ . Since  $\phi_i$  and  $\phi$  satisfy (2.30) and (2.35), respectively, from elementary properties of fundamental matrices it follows that

$$(2.45) \quad \frac{\partial \phi_i(\tau, t)}{\partial t} = -\phi_i(\tau, t) \left[ A(t) + \sum_{j \neq i}^n B_j(t)\mathcal{L}_j(t) \right],$$

$$(2.46) \quad \frac{\partial \phi(\tau, t)}{\partial t} = -\phi(\tau, t) \left[ A(t) + \sum_1^n B_j(t)\mathcal{L}_j(t) \right],$$

$i = 1, 2, \dots, n$ . By differentiating (2.44) with the aid of (2.45)–(2.46) and using (2.43) to substitute out the  $\mathcal{L}_i(t)$ , a long but otherwise simple calculation shows

that the  $Q_i(t)$  satisfy the system of ordinary matrix differential equations

$$\begin{aligned}
 (2.47) \quad & -\dot{Q}_i = W_i(t) + A^*(t)Q_i + Q_iA(t) - Q_iB_i(t)B_i^*(t)Q_i \\
 & - Q_i \left[ \sum_{j \neq i}^n B_j(t)B_j^*(t)Q_j \right] - \left[ \sum_{j \neq i}^n Q_j^*B_j(t)B_j^*(t) \right] Q_i, \\
 & Q_i(t_1) = \hat{W}_i, \quad i = 1, 2, \dots, n.
 \end{aligned}$$

The following uniqueness argument further shows  $Q_i(t) = Q_i^*(t)$ . The differential equation obtained from (2.47) by replacing the term  $Q_j^*$  by  $Q_j$  has a unique and hence symmetric solution. But that solution also satisfies (2.47), and hence again by uniqueness the solutions are the same.

Now we show that the  $Q_i(t)$  are positive definite for  $t < t_1$ . Let  $\mathcal{L}(t)$  be the equilibrium solution to (2.42). Then (2.43) and (2.47) hold for all  $t < t_1$  and near  $t_1$ . Letting  $y(t)$  be the equilibrium response to  $\mathcal{L}(t)$ ; i.e., the solution to

$$(2.48) \quad \dot{y} = \left[ A(t) + \sum_1^n B_j(t)\mathcal{L}_j(t) \right] y$$

with  $y(t_0) = y_0$ , using (2.43) and (2.47) one may easily verify

$$\begin{aligned}
 (2.49) \quad & \frac{d}{dt}[y(t) \cdot Q_i(t)y(t)] = -[|y(t)|_{W_i(t)}^2 + |B_i^*Q_iy(t)|^2] \\
 & = -[|y(t)|_{W_i(t)}^2 + |\mathcal{L}_i(t)y(t)|^2].
 \end{aligned}$$

Integration of (2.49) gives

$$(2.50) \quad |y_0|_{Q_i(t_0)}^2 = \int_{t_0}^{t_1} [|y(t)|_{W_i(t)}^2 + |\mathcal{L}_i(t)y(t)|^2] dt + |y(t_1)|_{\hat{W}_i}^2$$

which shows that the quadratic form in  $y_0$  associated with  $Q_i(t)$  gives  $P_i$ 's cost in the game with all players playing equilibrium feedback. Equation (2.50) also shows that each  $Q_i(t)$  is positive definite for  $t < t_1$ .

**2.5. Concluding comments.** We have shown that the equilibrium feedback matrices whose existence was established by a fixed-point argument can be computed by integrating a system of differential equations, (2.47), for the cost matrices and then doing the matrix multiplication (2.43). One can easily show that the matrix feedbacks computed in this manner provide an equilibrium solution on the entire left-maximal interval of existence for (2.47).

The uniqueness of these initial-state-independent equilibrium matrices can now be argued. Of course equilibria in general are not unique. For example, if we choose  $x_0 = 0$  in (1.10), then any system of operators for which the equation is defined constitute equilibrium solutions. Moreover, the example following Theorem 1.3 shows that even initial-state-independent equilibria in general are not unique. However Theorem 1.1 guarantees that any equilibrium must be a solution to (1.10) and in particular any initial-state-independent equilibrium solution to the differential game must be a solution to (2.42). But the technique used in computing the solution to that equation shows that it can have but one solution. That is, the uniqueness of the solution of (2.42) is a consequence of the

uniqueness of solutions to the differential equations of the type (2.47). These observations complete the proof of uniqueness stated in Theorem 2.3. We remark that when one considers the complexity of the operator calculations (that one in principle could do for the concrete operators arising in the differential game) which were necessary in proving the only if part of Theorem 1.1, it is understandable that previous investigators in differential games were not successful in treating the uniqueness question. Clearly the abstract approach of this paper makes a definite contribution here.

The existence of the solution to the differential game could have been established directly by construction of the solution from the solution of (2.47). However the fixed-point approach provides a means for estimating the interval of playability. From (2.42) we see that  $\mathcal{L}_i(t_1) = -B_i^*(t_1)\hat{W}_i$ . Hence if  $\hat{W}_i = 0$ ,  $i = 1, 2, \dots, n$ , then we can obtain an estimate using (1.37) as described in the remark following Theorem 2.3. In general one would expect to get improved estimates by utilizing this prior information as pointed out in the remark following Theorem 1.3. That is, the fixed-point theory should be applied about the point  $-B_i^*(t_1)\hat{W}_i$ .

The equilibrium feedbacks computed above agree with the solution derived in [4] using a Hamilton–Jacobi approach. In the one player game; i.e.,  $n = 1$ , equation (2.47) reduces to the Kalman–Riccati equation

$$(2.51) \quad -\dot{Q} = W(t) + A^*(t)Q + QA(t) - QB(t)B^*(t)Q,$$

$Q(t_1) = \hat{W}$ , studied in [6] in connection with the optimal regulator problem. This equation has a solution on every interval  $[t_0, t_1] \subseteq [t_\alpha, t_\beta]$ .

In [1] an equilibrium was found in  $\sum_1^n \oplus H_i$ , the space of open loop controls, and then synthesized by a matrix feedback system in  $L$  generated in the same way as the equilibrium in  $L$  but with an equation lacking the last term in (2.47). In general the two feedback systems are not the same and produce different control responses.

When the differential game is playable on every subinterval of  $[t_\alpha, t_\beta]$  we call the game *globally playable*. A sufficient condition for global playability is

$$(2.52) \quad W_i(t) = \lambda_i W_0(t), \quad \hat{W}_i = \lambda_i \hat{W}_0$$

and

$$(2.53) \quad B_i(t)B_i^*(t) = \frac{1}{\lambda_i} S_0(t),$$

$i = 1, 2, \dots, n$ , for some choice of scalars  $\lambda_i > 0$  and matrices  $W_0(t)$ ,  $\hat{W}_0$ ,  $S_0(t)$ . This result can be checked by using (2.52)–(2.53) to show  $Q_i(t) = \lambda_i Q(t)$  is the required solution of (2.47) if we take  $Q(t)$  to be the solution of (2.51) corresponding to the selection  $W(t) = W_0(t)$ ,  $\hat{W} = \hat{W}_0$  and  $B(t)$  any symmetric square root of  $(2n - 1)S_0(t)$ .

As a final remark we point out that the results of this paper can be applied to the study of equilibrium solutions in other feedback spaces and bear upon games with dynamics described by other types of equations; e.g., integral equations.

## REFERENCES

- [1] D. L. LUKES AND D. L. RUSSELL, *A global theory for linear quadratic differential games*, J. Math. Anal. Appl., 33 (1971).
- [2] P. P. VARAIYA, *N-person non-zero sum differential games with linear dynamics*, this Journal, 8 (1970), pp. 441–449.
- [3] D. L. LUKES AND A. STRAUSS, *Two countable systems of differential inequalities with applications to the stability of linear quadratic differential games*, Tech. Summary Rep. 1003, Mathematics Research Center, Univ. of Wisconsin, Madison, 1969.
- [4] J. H. CASE, *Equilibrium points of N-person differential games*, Tech. Rep. 1967-1 (Thesis), Dept of Industrial Engineering, Univ. of Michigan, Ann Arbor, 1967.
- [5] J. DIEUDONNÉ, *Foundations of Modern Analysis*, Academic Press, New York, 1960.
- [6] D. L. LUKES, *Stabilizability and optimal control*, Funkcial. Ekvac., 11 (1968), pp. 39–50.
- [7] E. B. LEE AND L. MARKUS, *Foundations of Optimal Control Theory*, John Wiley, New York, 1967.
- [8] D. L. LUKES, *Optimal regulation of nonlinear dynamical systems*, this Journal, 7 (1969), pp. 75–100.



## CONTROLLABILITY OF LINEAR OSCILLATORY SYSTEMS USING POSITIVE CONTROLS\*

STEPHEN H. SAPERSTONE† AND JAMES A. YORKE‡

**Abstract.** A linear autonomous control process is considered where the null control is an extreme point of the restraint set  $\Omega$ . In the event that  $\Omega = [0, 1]$  (hence, scalar control) necessary and sufficient conditions are given so that the reachable set from the origin (in phase space) contains the origin as an interior point. For vector-valued controls with each component in  $[0, 1]$ , sufficient conditions are given so that the reachable set from the origin of a nonlinear autonomous control process contains the origin as an interior point.

**1. Introduction and example.** The results in this paper are best motivated by the following example of controllability. Can the motion of a simple pendulum be brought to rest in a finite time by the application of a unit force acting only in one direction? In terms of a differential equation, the problem can be stated as follows: Let  $u: [0, \infty) \rightarrow \{0, 1\}$  be measurable. The linearized equation of motion of the pendulum is  $\ddot{\theta} + k\theta = u$ , where  $k$  is a positive constant of the motion and  $\theta$  is the angular displacement of the pendulum from the vertical. Is there an open neighborhood  $V$  of the origin (in  $\theta, \dot{\theta}$  space) such that if  $(\theta(0), \dot{\theta}(0)) \in V$ , then there exist a controller  $u(\cdot)$  and some  $T > 0$  such that  $(\theta(T), \dot{\theta}(T)) = (0, 0)$ ? The answer is yes. In fact all values of  $\theta(0)$  and  $\dot{\theta}(0)$  can be steered to the origin in finite time.

Similar questions of controllability may be asked of more complicated oscillatory (i.e., without real eigenvalues) systems having many degrees of freedom. In § 7 we shall analyze a double pendulum. Related questions for nonlinear systems are discussed in Corollary 6.3. A more general approach will be given in [5].

**2. Definitions and statement of main result.** Consider the system of linear differential equations in real Euclidean  $d$ -space  $R^d$ ,

$$(2.1) \quad \dot{x} = Ax + bu,$$

where  $A$  is a real constant  $d \times d$  matrix and  $x, b$  are real column  $d$ -vectors. Let  $R^+$  denote the nonnegative real numbers, and  $\Omega$ , the restraint set, be a nonempty interval in  $R$ . Let  $U_\Omega$  be the set of all bounded measurable functions  $u(\cdot)$ , where  $u: R^+ \rightarrow \Omega$ . For each  $u \in U_\Omega$  let  $x(t) = x(t; u(\cdot))$  be the unique absolutely continuous function satisfying (2.1) such that  $x(0; u(\cdot)) = 0$ . Define the *reachable set* (attainable set) at time  $t \geq 0$  by

$$K_\Omega^+(t) = \{x(t; u(\cdot)) : u \in U_\Omega\}$$

and the *reachable set* by

$$K_\Omega^+ = \bigcup_{t \geq 0} K_\Omega^+(t).$$

\* Received by the editors March 23, 1970, and in revised form August 13, 1970.

† Center for Naval Analyses of the University of Rochester, Arlington, Virginia and University of Maryland, College Park, Maryland. Now at Department of Mathematics, Howard University, Washington, D.C. 20001. The work of this author was supported by a CNA Fellowship.

‡ Institute for Fluid Dynamics and Applied Mathematics, University of Maryland, College Park, Maryland 20742. The work of this author was supported in part by the National Science Foundation under Grant NSF-GP-9347.

$K_{\Omega}^+(t)$  is convex and connected, and compact provided  $\Omega$  is compact (see [4, p. 69]). We note that  $x(t; u(\cdot))$  is given by

$$x(t; u(\cdot)) = \int_0^t e^{A(t-s)}bu(s) ds.$$

Let  $x'$  denote the transpose of the column vector  $x \in R^d$ . The *controllability matrix* of (2.1) is the  $d \times d$  matrix whose columns are  $b, Ab, A^2b, \dots, A^{d-1}b$ . We say that the origin can be *steered* (or *controlled*) to the point  $x_0 \in R^d$  provided  $x_0 \in K_{\Omega}^+(t_0)$  for some  $t_0 < \infty$ . The system (2.1) is called *locally controllable at the origin* if there exists  $t_1 < \infty$  such that  $K_{\Omega}^+(t_1)$  contains a neighborhood of the origin.

Let  $\Omega = [0, 1]$ . The main result we will prove is the following theorem.

**THEOREM 2.1.** *The system (2.1) with controls  $u(\cdot)$  belonging to  $U_{[0,1]}$  is locally controllable at the origin if and only if*

- (i) *all eigenvalues of  $A$  have nonzero imaginary parts, and*
- (ii) *the controllability matrix for (2.1) has rank  $d$ .*

**Remark 2.2.** Observe in Theorem 2.1 that if  $d$  is odd, then the system (2.1) is not controllable since a matrix  $A$  of odd order must have at least one real eigenvalue. Likewise  $A$  must be nonsingular. Thus we are led to consider only systems of an ‘‘oscillatory’’ nature.

**Remark 2.3.** We note that (2.1) is locally controllable at the origin for the class  $U_{\Omega}$  if and only if there exists a neighborhood  $V$  of the origin such that every point of  $V$  can be steered to  $x = 0$  in finite time. In fact, observe that  $\lambda$  is an eigenvalue of  $A$  if and only if  $-\lambda$  is an eigenvalue of  $-A$  and the vectors  $\{b, Ab, A^2b, \dots, A^{d-1}b\}$  span the same space as  $\{-b, Ab, -A^2b, \dots, (-1)^dA^{d-1}b\}$ . It follows that (2.1) is locally controllable at the origin if and only if

$$(2.2) \quad \dot{x} = -Ax - bu$$

is. But  $u(t)$  steers the origin to  $x_0$  along a solution of (2.1) on  $[0, t_0]$  if and only if  $u(t_0 - t)$  steers  $x_0$  to the origin along a solution of (2.2) on  $[0, t_0]$ . Thus every point of  $K_{[0,1]}^+(t)$  for any  $t \geq 0$  may be steered to the origin in time  $t$  if and only if conditions (i) and (ii) of Theorem 2.1 hold.

Theorem 2.1 is related to the following well-known result in controllability. We give a new proof of this result and will make use of parts of the proof in establishing Theorem 2.1. Let  $\Omega = [-\varepsilon, \varepsilon]$ .

**THEOREM 2.4.** *For the system (2.1) with controls  $u(\cdot)$  belonging to  $U_{[-\varepsilon,\varepsilon]}$ ,  $x = 0$  is interior to  $K_{[-\varepsilon,\varepsilon]}^+$  if and only if the controllability matrix for (2.1) has rank  $d$ .*

*Proof.* Let  $L^r = \text{span} \{b, Ab, \dots, A^{r-1}b\}$  for  $r = 1, 2, \dots$ . Since dimension  $L^r \leq d$ , we may choose  $\alpha$  (where  $0 \leq \alpha \leq d$ ) to be the smallest integer such that  $A^\alpha b \in L^\alpha$ . For any  $y \in L^\alpha$  there exist  $a_0, a_1, \dots, a_{\alpha-1}$  such that  $y = a_0b + a_1Ab + \dots + a_{\alpha-1}A^{\alpha-1}b$ . It follows that  $Ay = a_0Ab + a_1A^2b + \dots + a_{\alpha-1}A^\alpha b$  is in  $L^\alpha$ . Since  $Ay \in L^\alpha$ , we have  $A^2y \in L^\alpha$ , etc. In particular,  $A^r b \in L^\alpha$  for all  $r$  and so  $L^\alpha = L^s$  if  $s \geq \alpha$ . Also, for any scalar function  $u(t)$ ,

$$(2.3) \quad Ay + bu(t) \in L^\alpha \quad \text{if } y \in L^\alpha.$$

From (2.3) we have that for any function  $u(t)$ , equation (2.1) (which is a differential equation on  $R^d$ ) may also be considered as a differential equation on the subspace

$L^\alpha$  in that  $A:L^\alpha \rightarrow L^\alpha$  and  $bu \in L^\alpha$ . Since (2.1) is a differential equation on  $L^\alpha$ ,  $x(t, u(\cdot)) \in L^\alpha$  for all  $t$ . When (2.1) is considered as a differential equation on  $R^d$ , we must still have  $x(t, u(\cdot)) \in L^\alpha$ . Hence  $K_\Omega^+(t) \subseteq L^\alpha$  for all  $t \geq 0$ .

Denote by  $x(t)$  the solution of (2.1) with controller  $u(t) \equiv \varepsilon$  and satisfying  $x(0) = 0$ . Choose  $T > 0$ . Let  $L$  be the linear subspace of  $R^d$  generated by  $x(t)$  for all  $t \in [0, T]$ ; that is,

$$L = \text{span} \{x(t) : t \in [0, T]\}.$$

Then for  $t$  and  $t + \tau$  in  $[0, T]$ , we have

$$\tau^{-1}[x(t + \tau) - x(t)] \in L.$$

It follows that  $\dot{x}(t) \in L$  for all  $t \in [0, T]$ . Similarly  $\tau^{-1}[\dot{x}(t + \tau) - \dot{x}(t)] \in L$ . Arguing inductively we see that all derivatives of  $x(\cdot)$  are in  $L$ . Setting  $t = 0$ , (2.1) implies

$$\dot{x}(0) = A0 + b\varepsilon = b\varepsilon.$$

Taking higher derivatives of both sides of (2.1) and evaluating at  $t = 0$  we get, by induction,

$$x^{(2)}(0) = A\dot{x}(0) = Ab\varepsilon,$$

...

$$x^{(r)}(0) = Ax^{(r-1)}(0) = A^{r-1}b\varepsilon.$$

Therefore  $b, Ab, \dots, A^r b \in L$  for all  $r > 0$ , so  $L^\alpha \subseteq L$ .

For  $\tau \in [0, T]$ , let  $u_\tau(t) = \varepsilon$  for  $t \in [T - \tau, T]$  and  $u_\tau(t) = 0$  elsewhere. Then  $x(T, u_\tau(\cdot)) = x(\tau)$ ; hence  $x(t) \in K_\Omega^+(T)$  for  $0 \leq t \leq T$ . We have thus far shown that  $\{x(t) : 0 \leq t \leq T\} \subseteq K_\Omega^+(T) \subseteq L^\alpha \subseteq L$ . Taking spans we get  $L \subseteq L^\alpha \subseteq L$ , and hence,  $L = L^\alpha$ .

Choose  $0 < t_1 < t_2 < \dots < t_\alpha < T$  such that  $B \stackrel{\text{def}}{=} \{x(t_i)\}_{i=1}^\alpha$  is a linearly independent set. Such a set exists because  $L^\alpha = L$  has dimension  $\alpha$  and by definition of  $L$ ,  $\{x(t) : 0 \leq t \leq T\}$  spans  $L$ . Since  $B$  is a basis for  $L^\alpha$ , there exists for each  $y \in L^\alpha$  a unique set  $\{c_1, c_2, \dots, c_\alpha\}$  such that  $y = c_1x(t_1) + \dots + c_\alpha x(t_\alpha)$ . Now define  $\|y\|_\alpha = |c_1| + \dots + |c_\alpha|$ . ( $\|y\|_\alpha$  is a norm on  $L^\alpha$  because it is the sum of the absolute values of the coordinates  $c_i$  using the basis  $B$ .) Since  $\Omega = [-\varepsilon, \varepsilon]$  is convex, it is easy to verify that  $K_\Omega^+(T)$  is convex. Furthermore if  $y \in K_\Omega^+(T)$ , then  $-y \in K_\Omega^+(T)$  since  $x(T, u(\cdot)) = -x(T, -u(\cdot))$ . Now  $x(t_i) = x(T, u_{t_i}(\cdot)) \in K_\Omega^+(T)$ . Therefore  $\pm x(t_i) \in K_\Omega^+(T)$ . So,  $\{y : \|y\|_\alpha \leq 1\} \subset K_\Omega^+(T)$ . Hence  $K_\Omega^+(T)$  contains a neighborhood of  $x = 0$  in the topology of  $L^\alpha$ . Therefore  $K_\Omega^+(T) \subset L^\alpha$  contains a neighborhood of  $x = 0$  in the topology of  $R^d$  if and only if the rank of  $\{b, Ab, \dots, A^{d-1}b\}$  is  $d$ .

*Remark 2.5.* Since  $x(T, u(\cdot))$  is linear in  $u$ , the point

$$y = \sum_{i=1}^\alpha c_i x(t_i) = \sum_{i=1}^\alpha c_i x(T, u_{t_i}(\cdot)) = x\left(T, \sum_{i=1}^\alpha c_i u_{t_i}(\cdot)\right)$$

(hence any point in  $L^\alpha$ ) can be reached by a controller  $u(\cdot) = \sum_{i=1}^\alpha c_i u_{t_i}(\cdot)$  which is piecewise constant and has only  $\alpha$  switching times  $t_1, t_2, \dots, t_\alpha$ . (Of course,  $|u| \leq \varepsilon$  does not necessarily hold.) These switching times do not depend upon the

point  $y$ . If the controllability matrix has rank  $d$ , then  $L^\alpha = R^d$  and the number of switches  $\alpha$  is  $d$ .

Though both Theorems 2.1 and 2.4 characterize controllability in terms of the range of the control functions allowed, the theorems differ significantly with respect to the time required to reach any desired point in  $K_\Omega^+$ . In particular for any arbitrary small time  $t_1 > 0$ ,  $K_{[-\epsilon, \epsilon]}^+(t_1)$  contains a neighborhood of  $x = 0$  (see [4, p. 83]). Such a claim is false for the case in Theorem 2.1 when  $\Omega = [0, \delta]$  for any  $\delta > 0$ . We can see from the simple pendulum in §1 that  $0 \in \text{int } K_{\{0,1\}}^+(t_1)$  only if  $t_1 > \pi k^{-1/2}$ . Physical considerations require that the time to bring the pendulum to rest from each initial position in a neighborhood of the origin be as much as half a period. This can be easily verified by using the Pontryagin maximum principle. That this is also true for  $\Omega = [0, \delta]$ ,  $\delta > 0$ , follows from the ‘‘bang-bang’’ principle [3] which we now state. Let  $E\Omega$  be the set of extreme points of  $\Omega$ .

**THEOREM 2.6.** *Suppose the range  $\Omega$  of the control function is compact. Then for any  $t \geq 0$ ,  $K_\Omega^+(t) = K_{E\Omega}^+(t)$  for the system (2.1). Furthermore, if  $t_0$  is the minimal time required to steer the origin to  $x_0$  for  $u \in U_\Omega$ , then  $t_0$  is also minimal for  $u \in U_{E\Omega}$ .*

Thus the reachable set at any time  $t$  for the motion of the simple pendulum is unaffected if  $\Omega = [0, 1]$  is replaced by  $\Omega = \{0, 1\}$ . By virtue of Theorem 2.6 the oscillatory systems we describe can now be controlled by ‘‘off-on’’ controls.

In the example of the simple pendulum, the eigenvalues are  $\pm ik^{1/2}$ . Furthermore we can show that  $K_{\{0,1\}}^+ = R^2$ . That the controllability matrix had rank 2 just depended on the fact that  $b \neq 0$ . The crucial aspect of generalizing the example to Theorem 2.1 lies in the oscillatory nature of the flow. But we even have controllability when the free system (with  $u \equiv 0$ ) is completely unstable. That is, the eigenvalues of  $A$  can have positive real parts. The following result is an immediate corollary.

**COROLLARY 2.7.** *Consider (2.1) for a 2-dimensional system with  $\Omega = \{0, 1\}$ . Suppose  $b \neq 0$ . Then the system is locally controllable at the origin if and only if no eigenvalue of  $A$  is real.*

**3. Preliminaries and lemmas.**

**LEMMA 3.1.** *If the matrix  $A$  of the system (2.1) has a real eigenvalue  $\lambda$  (possibly zero), then  $x = 0$  belongs to  $\partial K_{\{0,1\}}^+$ .*

*Proof.* Let  $\lambda$  be a real eigenvalue of  $A$ . Then  $\lambda$  is also an eigenvalue of  $A'$  (transpose of  $A$ ), and there is some  $v \in R^d$  ( $v \neq 0$ ) such that  $A'v = \lambda v$ . We may assume  $v \cdot b \geq 0$  (replacing  $v$  by  $-v$  if necessary). Let  $x(t)$  be a solution of (2.1) for  $t \geq 0$ , with  $x(0) = 0$  for some  $u \in U_{\{0,1\}}$ . Write  $\rho(t) = v \cdot x(t)$ . Then

$$\begin{aligned} \frac{d}{dt}\rho(t) &= v \cdot \dot{x}(t) = v \cdot Ax(t) + v \cdot bu(t) \\ &= \lambda\rho(t) + (v \cdot b)u(t). \end{aligned}$$

Since  $u(t) \geq 0$  for all  $t$ , and  $\rho(0) = 0$ , it follows that  $\rho(t) \geq 0$ . Furthermore, if  $v \cdot b = 0$ , then  $\rho(t) \equiv 0$ . That is, for all  $y \in K_{\{0,1\}}^+$ ,  $y$  is in the half-space  $\{z \in R^d : v \cdot z \geq 0\}$  when  $v \cdot b > 0$  (and  $y$  is in the hyperplane  $v \cdot z = 0$  if  $v \cdot b = 0$ ). Considering the control function  $u(t) \equiv 0$ , we see that  $0 \in K_{\{0,1\}}^+$ . Hence  $0 \in \partial K_{\{0,1\}}^+$ .

From Lemma 3.1 for local controllability with  $\Omega = [0, 1]$ , we see that  $A$  must be nonsingular and of even dimension since matrices of odd dimension have

at least one real eigenvalue. Thus we are led to consider systems of even order  $d = 2n$ , such that the eigenvalues of  $A$  have nonzero imaginary parts. In particular let the distinct eigenvalues of  $A$  be

$$\lambda_k = \alpha_k + i\beta_k, \quad \lambda_k^* = \alpha_k - i\beta_k, \quad \alpha_k, \beta_k \in R^1, \quad \beta_k \neq 0 \quad \text{for } k = 1, 2, \dots, s,$$

with multiplicities  $m_1, m_2, \dots, m_s$ . (\* denotes complex conjugate.) Let  $\eta$  be any nonzero row vector in  $R^{2n}$ . We find a convenient representation and important property for the product  $\eta e^{A\tau} b$ ,  $\tau \in R^+$ . (We assume the controllability matrix has rank  $2n$ .)

LEMMA 3.2.  $\eta e^{A\tau} b$  is not identically zero on  $R^+$  and can be written

$$(3.1) \quad \eta e^{A\tau} b = \tau^m e^{\alpha\tau} \left\{ \mu(\tau) + \sum_{k=1}^s h_k \sin(\beta_k \tau + \delta_k) \right\}$$

for constants  $m, \alpha, \delta_k, h_k$  (where not all  $h_k$  are zero) and where  $\mu(\tau) \rightarrow 0$  as  $\tau \rightarrow \infty$ , and the summation is not identically zero.

*Proof.* On the contrary, assume  $\eta e^{A\tau} b$  is identically zero. Then  $e^{A\tau} b$  lies in some proper subspace  $V$  of  $R^{2n}$  for all  $\tau \geq 0$ . Now  $e^{A\tau} b$  is the unique solution to the differential equation  $\dot{x} = Ax$ ,  $x(0) = b$ . We saw in the proof of Theorem 2.4 that  $\dot{x}(0), x^{(2)}(0), \dots, x^{(2n-1)}(0)$  lie in  $V$ ; hence  $b, Ab, \dots, A^{2n-1}b$  belong to  $V$ . Since the controllability matrix has rank  $2n$ , it follows that  $V = R^{2n}$  and therefore  $\eta e^{A\tau} b \neq 0$ .

There exists an invertible matrix  $P$  such that  $J = PAP^{-1}$  has Jordan normal form (cf. [2, p. 76]). Then  $\eta e^{A\tau} b$  can be written

$$\eta P^{-1} e^{J\tau} P b.$$

Consider the complex vector  $e^{J\tau} P b$ . It is the solution to the complex system  $\dot{y} = Jy$ ,  $y(0) = P b$  and is of the form

$$\sum_{\lambda} e^{\lambda\tau} \sum_{j=1}^{m_{\lambda}} P_{\lambda,j}(\tau) v_{\lambda,j},$$

where the first sum is taken over all eigenvalues  $\lambda$  of  $A$ ,  $P_{\lambda,j}(\tau)$  is a polynomial in  $\tau$  which depends on  $\lambda$  and has order  $m_{\lambda} - 1$  (where  $m_{\lambda}$  is the multiplicity of the eigenvalue  $\lambda$ ), and  $v_{\lambda,j}$  is a (complex) vector depending on  $\lambda$ . Replacing  $e^{\lambda\tau}$  by  $e^{\alpha\tau}(\cos \beta\tau + i \sin \beta\tau)$ , we observe that since  $\eta e^{A\tau} b$  is real, it has the form

$$(3.2) \quad \sum_{k=1}^s e^{\alpha_k\tau} [R_k(\tau) \cos \beta_k\tau + S_k(\tau) \sin \beta_k\tau],$$

where  $R_k(\tau), S_k(\tau)$  for  $k = 1, 2, \dots, s$  are real polynomials in  $\tau$  and are not all zero. If  $\alpha = \max \{ \alpha_1, \alpha_2, \dots, \alpha_s \}$ , we can factor  $e^{\alpha\tau}$  out of (3.2) along with an appropriate power of  $\tau$ , say  $\tau^m$ , and obtain

$$\eta e^{A\tau} b = \tau^m e^{\alpha\tau} \left\{ \mu(\tau) + \sum_{k=1}^s h_k \sin(\beta_k \tau + \delta_k) \right\},$$

where not all the  $\{h_k\}_{k=1}^s$  are zero and  $\mu(\tau) \rightarrow 0$  as  $\tau \rightarrow \infty$ .

**4. Almost periodic functions.** Let  $\mathbb{C}$  denote the field of complex numbers. Recall a continuous function  $f: R \rightarrow \mathbb{C}$  is *almost periodic* (a.p.) if for any  $\varepsilon > 0$  there exists a number  $L(\varepsilon) > 0$  with the property that any interval of length  $L(\varepsilon)$  of  $R$

contains at least one point  $\xi$  such that  $|f(t + \xi) - f(t)| < \varepsilon$  for all  $t \in R$ . The number  $\xi$  is called an  $\varepsilon$ -translation number of  $f$ . We note the following properties about a.p. functions [1]. If  $f$  and  $g$  are a.p. functions, then  $a_1f + a_2g$  ( $a_1, a_2 \in \mathbb{C}$ ) and  $f^*$  (the complex conjugate of  $f$ ) are a.p. Furthermore  $f$  is bounded and uniformly continuous on  $R$ . Also, the integral

$$M\{f\} \stackrel{\text{def}}{=} \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T f(t) dt$$

exists and is called the mean value of  $f$  (see [1, p. 12]).

LEMMA 4.1. Suppose  $f$  is an a.p. function such that  $M\{f\} = 0$  and

$$(4.1) \quad f(t) \leq |\mu(t)| \quad \text{for all } t \in R^+,$$

where  $|\mu(t)| \rightarrow 0$  as  $t \rightarrow \infty$ . Then  $f \equiv 0$ .

*Proof.* First we show that  $f(t) \leq 0$  for  $t \in R$ . Suppose the contrary, that there exists  $t_0$  such that  $f(t_0) = m > 0$ . Then there exists  $\delta > 0$  such that  $f(t) > 2m/3$  for all  $t \in (t_0 - \delta, t_0 + \delta)$ . There exists a number  $L(m/3)$  such that any interval of length  $L(m/3)$  contains at least one number of the form  $t_0 + \xi$ , where  $|f(t + \xi) - f(t)| < m/3$ . Choose  $L = L(m/3) > 2\delta$ . Any interval of length  $L$  contains at least one of the intervals  $(t_0 + \xi - \delta, t_0 + \xi)$  or  $(t_0 + \xi, t_0 + \xi + \delta)$ . In each of these intervals, which are the intervals  $(t_0 - \delta, t_0)$ ,  $(t_0, t_0 + \delta)$  translated by  $\xi$ , we have  $f(t) > m/3$ . But this contradicts the assumption (4.1). Hence we must have  $f(t) \leq 0$  for all  $t \in R^+$ , which reduces our lemma to a result in [1, p. 36]. This completes the proof.

Functions of the form  $f(t) = \sum_{k=1}^N a_k e^{i w_k t}$ , where  $a_k \in \mathbb{C}$ ,  $w_k \in R$ ,  $k = 1, 2, \dots, N$ , are a.p. In particular  $\text{Im}(f)$  is a.p.

**5. Proof of Theorem 2.1.** We have already proved that a necessary condition for  $x = 0$  to be interior to  $K_{[0,1]}^+$  is that no eigenvalue of  $A$  is real (Lemma 3.1). That the controllability matrix have rank  $d = 2n$  is also necessary for  $x = 0$  to be interior to  $K_{[0,1]}^+$  follows directly from Theorem 2.4. It only remains to prove the sufficiency of the conditions (i) and (ii).

So assume the controllability matrix has rank  $2n$  and no eigenvalue of  $A$  is real. We claim that for large enough  $t$ ,  $x = 0$  is interior to  $K_{[0,1]}^+(t)$ . Suppose not; i.e., for all  $t \geq 0$ ,  $x = 0$  belongs to  $\partial K_{[0,1]}(t)$ . Since  $K_{[0,1]}^+(t)$  is convex for any  $t \geq 0$ , we can find supporting hyperplanes to  $K_{[0,1]}^+(t)$  at  $x = 0$ . Let  $N(t)$  denote the set of all unit outer normals of such supporting hyperplanes to  $K_{[0,1]}^+(t)$  at  $x = 0$ . Then for each  $t \geq 0$ ,  $N(t) \neq \emptyset$ , since we have assumed  $0 \in K_{[0,1]}^+(t)$ .  $N(t)$  clearly is bounded. Moreover, it is closed. In fact, let  $\{\eta_i\}$  be a sequence in  $N(t)$  converging to  $\eta \in R^{2n}$ . Then for any  $x(t) \in K_{[0,1]}^+(t)$ ,  $\eta_i \cdot x(t) \leq 0$ . Passing to the limit we have  $\eta \cdot x(t) \leq 0$  and  $\|\eta\| = 1$ . Thus  $\eta \in N(t)$ , and  $N(t)$  is a nonempty compact subset of  $R^{2n}$ . Furthermore, if  $t_1 < t_2$ , then  $N(t_2) \subset N(t_1)$  since  $K_{[0,1]}^+(t)$  increases with increasing  $t$ . Consequently, for any finite set of times  $0 \leq t_1 < t_2 < \dots < t_k < \infty$ ,  $\bigcap_{i=1}^k N(t_i) = N(t_k) \neq \emptyset$ . Hence  $\bigcap_{t \geq 0} N(t) \neq \emptyset$  from the finite intersection property of compact sets. Then there exists a unit outer normal vector  $\eta_0$ , and hence a supporting hyperplane to  $K_{[0,1]}^+(t)$  good for all  $t \geq 0$ . In particular, for any

controller  $u(\cdot) \in U_{[0,1]}$  and corresponding solution  $x(t; u(\cdot))$ , we have

$$\eta_0 \cdot x(t; u(\cdot)) = \int_0^t \eta_0 e^{A(t-\tau)} b u(\tau) d\tau \leq 0$$

for all  $t \geq 0$ . It follows by continuity and special choice of  $u(\cdot)$  that for all  $\tau \geq 0$ ,

$$\varphi(\tau) \stackrel{\text{def}}{=} \eta_0 e^{A\tau} b \leq 0.$$

According to Lemma 3.2,  $\varphi(\tau)$  is not identically zero and

$$\varphi(\tau) = \tau^m e^{a\tau} \left\{ \mu(\tau) + \sum_{k=1}^s h_k \sin(\beta_k \tau + \delta_k) \right\},$$

where  $|\mu(\tau)| \rightarrow 0$  as  $\tau \rightarrow \infty$ , and the summation is not identically zero. Denote this sum by  $v(\tau)$ . Observe that  $v(\tau)$  is an a.p. function with mean  $M\{v\} = 0$ . We now claim that there exists a  $\tau^* \in R^+$  such that  $v(\tau^*) > |\mu(\tau^*)|$ . If, on the contrary, no such  $\tau^*$  exists,  $v(t) \leq |\mu(t)|$  for all  $t \in R^+$ . From Lemma 4.1 we see that  $v(\tau) \equiv 0$ . But we know this to be false. Therefore we conclude there exists  $\tau^* \in R^+$  such that  $v(\tau^*) > |\mu(\tau^*)|$ . Thus  $\varphi(\tau^*) > 0$ , a contradiction. So we conclude that  $x = 0$  is interior to  $K_{[0,1]}^+(t)$  for some  $t < \infty$ .

**6. Extensions and generalizations.** Theorem 2.1 has a generalization in the event that the vector  $b$  is replaced by a matrix  $B$  of order  $d \times q$ . Let  $\Omega_0$  denote the unit cube in  $R^q$ ,  $\Omega_0 = \{(u_1, u_2, \dots, u_q) \in R^q : 0 \leq u_i \leq 1, i = 1, 2, \dots, q\}$  and define

$$U_{\Omega_0} = \{u : R^+ \rightarrow \Omega_0, u \text{ measurable}\}.$$

Consider

$$(6.1) \quad \dot{x} = Ax + Bu,$$

$u \in U_{\Omega_0}$ . Then  $K_{\Omega_0}^+(t)$  is defined accordingly. Denote by  $C$  the controllability matrix of (6.1), where  $C$  is now the matrix of order  $d \times dq$  whose columns are the columns of  $B, AB, A^2B, \dots, A^{d-1}B$ . We note that for any nonsingular matrix  $T$  of order  $d \times d$ ,  $\text{rank}(TC) = \text{rank}(C)$ . We can now give sufficient conditions for controllability of (6.1) at the origin.

**THEOREM 6.1.** *Suppose all the eigenvalues of  $A$  have nonzero imaginary parts and that the controllability matrix for (6.1) has rank  $d$ . Then  $x = 0$  is interior to  $K_{\Omega_0}^+(t_1)$  for some  $t_1 < \infty$ .*

*Proof.* We proceed as in the proof of Theorem 2.1 to establish a unit normal vector,  $\eta_0$ . Denote  $\eta_0 e^{A\tau} B$  by  $\varphi(\tau)$ . If the columns of  $B$  are denoted by  $b_1, b_2, \dots, b_q$ , then the matrix  $\varphi$  satisfies

$$\varphi(\tau) = [\eta_0 e^{A\tau} b_1, \eta_0 e^{A\tau} b_2, \dots, \eta_0 e^{A\tau} b_q].$$

We claim that not all of the terms  $\eta_0 e^{A\tau} b_j$  for  $j = 1, 2, \dots, q$  are identically zero. In fact, suppose they were all zero. Then for all  $\tau \geq 0$  and every  $j = 1, 2, \dots, q$ ,  $e^{A\tau} b_j$  lies in some proper subspace of  $R^d$ . Denote this subspace by  $V$ . Now each  $e^{A\tau} b_j$  is the unique solution to the differential equation

$$\dot{x} = Ax, \quad x(0) = b_j.$$

As in the proof of Theorem 2.1 we show  $b_j, Ab_j, \dots, A^{d-1}b_j$  lie in  $V$  for each  $j = 1, 2, \dots, q$ . It follows that  $\eta_0 B = \eta_0 AB = \dots = \eta_0 A^{d-1} B = 0$ , thus contradicting the fact that the controllability matrix  $C$  has rank  $d$ . As in the case of Theorem 2.1 we can show that if  $x = 0$  is not interior to  $K_{\Omega_0}^+(t)$  for any  $t \geq 0$ , then each nonzero component  $\varphi_j(\tau) \stackrel{\text{def}}{=} \eta_0 e^{A\tau} b_j$  of  $\varphi(\tau)$  is never positive on  $R^+$ . But each  $\varphi_j(\tau)$  has the form (3.1), and eventually is positive. From this contradiction we conclude that  $x = 0$  is interior to  $K_{\Omega_0}^+(t)$  for some  $t_1 < \infty$ , thus completing the proof.

It is clear from the above proof that the conditions for local controllability at the origin with positive controls are too strong. In fact, the following simple example illustrates how we can have  $x = 0$  interior to  $K_{\Omega_0}^+$  in the event (6.1) is a scalar equation.

*Example 6.2.* Consider the scalar equation

$$\dot{x} = \lambda x + b_1 u_1 + b_2 u_2$$

with  $0 \leq u_1, u_2 \leq 1$ . It may be seen that  $0 \in \text{int } K_{\Omega_0}^+$  if and only if  $b_1 b_2 < 0$ .

Theorem 6.1 can be extended to nonlinear systems which admit a linearization of the form (6.1).

*COROLLARY 6.3.* Consider the following system:

$$(6.2) \quad \dot{x} = f(x, u),$$

where  $f: R^d \times R^q \rightarrow R^d$ ,  $f$  is  $C^1$  in  $R^d \times R^q$  and  $u \in U_{\Omega_0}$ . Define

$$A = \left. \frac{\partial f}{\partial x} \right|_{(0,0)}, \quad B = \left. \frac{\partial f}{\partial u} \right|_{(0,0)}.$$

Suppose:

- (i)  $f(0, 0) = 0$ ,
- (ii) all eigenvalues of  $A$  have nonzero imaginary parts,
- (iii)  $\text{rank } [B \ : \ AB \ : \ A^2 B \ : \ \dots \ : \ A^{d-1} B] = d$ .

Then  $x = 0$  is interior to  $K_{\Omega_0}^+$ .

*Proof.* The proof is exactly the same as used by Markus and Lee in [4, p. 366]. In fact, nowhere in their proof do they require controllers which assume both positive and negative values. They only require that the solution to the linearized form of (6.2) be controllable to any point in some neighborhood of  $x = 0$  (starting from  $x = 0$ ). But conditions (ii) and (iii) guarantee this.

**7. The double pendulum.** As an example of controllability with positive controls, we consider the linearized equations of motion of a double pendulum (Fig. 7.1). Each mass  $m_k$  is fixed at the end of a rigid weightless rod of length  $L_k$ ,  $k = 1, 2$ . Both masses are nonzero. The system swings without friction about the pivot points  $P_0$  and  $P_1$ , where  $P_0$  is fixed to a rigid structure. A time-varying force,  $u(t)$ ,  $0 \leq u(t) \leq 1$ ,  $t \in R^+$ , is applied simultaneously at masses  $m_1$  and  $m_2$  as indicated in Fig. 7.1.

The linearized equations of motion of the system are given by (with gravitational constant  $g = 1$ )

$$\begin{aligned} m_1 L_1 \ddot{\theta}_1 + (m_1 + m_2) \theta_1 - m_2 \theta_2 &= u, \\ m_2 L_2 \ddot{\theta}_2 + (m_1 + m_2) \theta_2 - (m_1 + m_2) \theta_1 &= (m_1/m_2 - 1)u. \end{aligned}$$



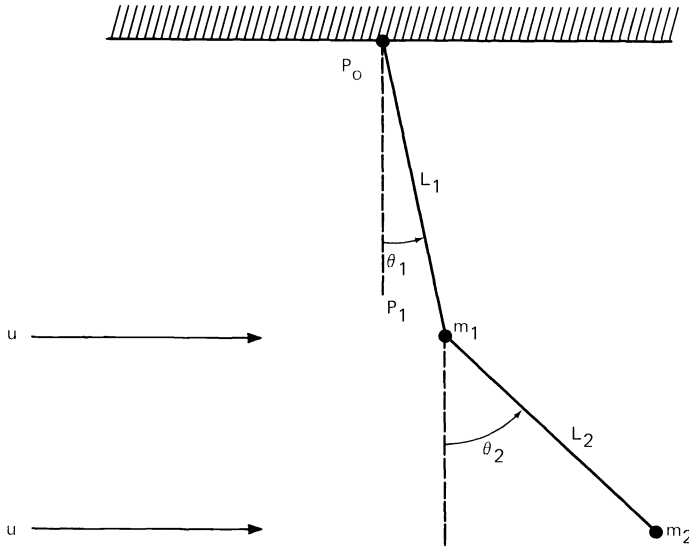


FIG. 7.1

Set  $x_1 = \theta_1, x_2 = \dot{\theta}_1, x_3 = \theta_2, x_4 = \dot{\theta}_2$ . In vector form we get  $\dot{x} = Ax + bu$ , where

$$A = \begin{bmatrix} 0 & 1 & 0 & 0 \\ -k/L_1 & 0 & (k-1)/L_1 & 0 \\ 0 & 0 & 0 & 1 \\ k/L_2 & 0 & -k/L_2 & 0 \end{bmatrix}, \quad b = \begin{bmatrix} 0 \\ (m_1 L_1)^{-1} \\ 0 \\ L_2^{-1}(m_2^{-1} - m_1^{-1}) \end{bmatrix}$$

and  $k = 1 + m_2/m_1 > 1$ .

The eigenvalues of  $A$  are pure imaginary. Denote them by  $\pm i\omega_1, \pm i\omega_2$ . We have (letting  $\pm$  be  $+$  for  $\omega_1^2$  and  $-$  for  $\omega_2^2$ )

$$\omega_1^2, \omega_2^2 = (2L_1 L_2)^{-1} \{ (L_1 + L_2)k \pm [(L_1 + L_2)^2 k^2 - 4L_1 L_2 k]^{1/2} \}.$$

We take  $\omega_1, \omega_2 > 0$ . It is easy to show that  $\omega_1^2 \neq \omega_2^2$ . Define the matrix  $Q$ , whose  $j$ th row is given by

$$[(-1)^{j+1}\omega, 1, (-1)^{j+1}(1 - L_1\omega^2)/(L_1\omega), (1 - L_1\omega^2)/(L_1\omega^2)],$$

where  $\omega = \omega_1$  for  $j = 1, 2$  and  $\omega = \omega_2$  for  $j = 3, 4$ . Under the coordinate transformation defined by  $Q$ ,  $A$  becomes

$$\begin{bmatrix} 0 & \omega_1 & 0 & 0 \\ -\omega_1 & 0 & 0 & 0 \\ 0 & 0 & 0 & \omega_2 \\ 0 & 0 & -\omega_2 & 0 \end{bmatrix},$$

and  $b$  becomes  $[b_1, b_1, b_2, b_2]$ , where

$$b_j = (m_1 L_1)^{-1} + (1 - L_1 \omega_j^2)(m_2^{-1} - m_1^{-1}) / (L_1 L_2 \omega_j^2), \quad j = 1, 2.$$

An easy computation shows that the determinant  $\Delta$  of the controllability matrix is  $4b_1^2 b_2^2 \omega_1 \omega_2 (\omega_1^2 - \omega_2^2)^2$ . In general,  $\Delta \neq 0$  except for certain critical values of the parameters  $L_1, L_2$  and  $k$ . Under the assumption that  $\Delta \neq 0$  we conclude that the reachable set for the double pendulum contains a neighborhood of the origin. To put it another way, for sufficiently small initial values of  $(\theta_1(0), \dot{\theta}_1(0), \theta_2(0), \dot{\theta}_2(0))$ , there exists an appropriate controller and a finite time  $T \geq 0$  such that  $(\theta_1(T), \dot{\theta}_1(T), \theta_2(T), \dot{\theta}_2(T)) = (0, 0, 0, 0)$ . Furthermore, one can show that if  $T > \pi / \min(\omega_1, \omega_2)$ , then  $K_{[0,1]}^+(T)$  contains the origin as an interior point.

**Acknowledgment.** We would like to thank the referee who made many useful comments throughout the paper. In particular, his suggestions allowed us to simplify and clarify § 5 which was at best confusing.

#### REFERENCES

- [1] A. S. BESICOVITCH, *Almost Periodic Functions*, Dover, New York, 1954.
- [2] E. A. CODDINGTON AND N. LEVINSON, *Theory of Ordinary Differential Equations*, McGraw-Hill, New York, 1955.
- [3] R. CONTI, *Contribution to linear control theory*, J. Differential Equations, 1 (1965), pp. 427–445.
- [4] E. B. LEE AND L. MARKUS, *Foundations of Optimal Control Theory*, John Wiley, New York, 1967.
- [5] J. A. YORKE, *The maximum principle and controllability of nonlinear equations*, to appear.

## ON THE CONVERGENCE OF APPROXIMATING SOLUTIONS FOR LINEAR DISTRIBUTED PARAMETER OPTIMAL CONTROL PROBLEMS\*

HITOSHI SASAI† AND ETSUJURO SHIMEMURA‡

**Abstract.** Optimal control problems for distributed parameter systems, particularly systems described by partial differential equations, are often treated using mathematical function space techniques. As a result, the equations which define the optimal control are frequently obtained in abstract terms. Numerical solutions are obtained by approximating the abstract operations in a computationally feasible manner. In obtaining approximating solutions, the finite difference method is widely used.

After an approximate optimal control has been found, the question arises whether a sequence of these approximating optimal controls converges to an optimal control of the original system.

A condition of convergence of a sequence of approximating solutions of initial value problems by the finite difference method was given by H. F. Trotter. This theorem is concerned with a homogeneous system and does not give the condition of convergence of approximating optimal controls for a distributed parameter system.

In this paper the condition of convergence of a sequence of approximating solutions is given for the time-optimal control problem and the final value control problem for a class of linear systems with distributed parameters.

**1. Introduction.** Optimal control problems for distributed parameter systems, particularly systems described by partial differential equations, are often treated using mathematical function space techniques. As a result, the equations which define the optimal control are frequently obtained in abstract terms. Numerical solutions are obtained by approximating the abstract operations in a computationally feasible manner. In obtaining approximating solutions, the method of expanding in eigenfunctions [1] and the finite difference method are widely used.

After having found an approximate optimal control, the question arises whether a sequence of these approximating optimal controls converges to an optimal control of the original system. A condition of convergence of a sequence of approximating solutions of initial value problems by the finite difference method in which a differential operator is replaced by a difference operator was given by H. F. Trotter [2], [3]. This theorem is concerned with a homogeneous system and does not give the condition of convergence of approximating optimal controls for a distributed parameter system, but its result becomes very useful for our discussion as shown below.

In this paper the condition of convergence of a sequence of approximating solutions is given for the time-optimal control problem and the final value control problem for a class of linear systems with distributed parameters.

To discuss the convergence of approximating solutions, it is important to define a concept of metric. In the following, the space is considered to be a Hilbert space. Practically it would be convenient to consider the space to be  $L^2(S)$ ,  $S$  being a region of Euclidean space with Lebesgue measure.

---

\* Received by the editors December 16, 1969, and in final revised form November 13, 1970.

† Department of Aeronautical Engineering, Nagoya University, Nagoya, Japan.

‡ Department of Electrical Engineering, Waseda University, Tokyo, Japan.

**2. Definitions and basic theorem.** For convenience, the following notations, definitions and basic theorem are introduced. Let  $X$  be a Hilbert space and let  $X_n$  be a subspace of  $X$ . We consider a linear mapping  $P_n: X \rightarrow X_n$  which satisfies the following conditions:

$$P_n^2 = P_n, \quad \|P_n\| = 1, \\ \lim \|P_n f - f\| = 0 \quad \text{for every } f \in X.$$

**DEFINITION.** The *limit of a sequence of operators*  $\{A_n\}$ , where  $A_n$  is an operator on  $X_n$ , is the operator on  $X$  whose domain consists of  $f \in X$ , for which  $\{A_n P_n f\}$  converges, and whose value is  $\lim A_n P_n f$ .

Let  $\{S_n(t)\}$  be a sequence of strongly continuous semigroups of operators on  $X_n$  into  $X_n$ , and  $\{A_n\}$  be the sequence of associated infinitesimal generators. Then we have the following theorem, which is due to H. F. Trotter.

**THEOREM (Trotter [2]).** *If the range of  $\lambda I - A$ , denoted by  $R(\lambda I - A)$ , is dense in  $X$  for some  $\lambda > K$  and the following conditions are satisfied, where  $M$  and  $K$  are positive constants:*

(C)  $A = \lim A_n$  and  $D(A)$  is dense in  $X$ ,

(S)  $\|S_n(t)\| \leq M e^{Kt}$ ,

then a closed extension of  $A$  is the infinitesimal generator of  $S(t)$ , where  $S(t) = \lim S_n(t)$  is a strongly continuous semigroup of operators on  $X$ .

In this theorem it is demonstrated that if conditions (C) and (S) are satisfied, approximating solutions of approximating equations starting from given approximating initial states converge in the sense of norm to a solution of the original homogeneous partial differential equation, which satisfies a given initial condition.

**3. Statement of the problem.** Let us consider the following distributed parameter system:

(1)  $\partial x(t, s) / \partial t = Ax(t, s) + u(t, s)$ ,

where for each  $t, 0 \leq t < \infty$ ,  $x(t, s)$  is an element of  $X$ , which is an arbitrary Hilbert space  $X(\Omega)$ ,  $s \in \Omega$ , consisting of functions on a bounded domain  $\Omega$ ;  $A$  is an unbounded linear operator from  $X$  to  $X$ , for instance,  $A$  being a partial differential operator on  $\Omega$ ; and  $u(t, s)$  is an element of  $X$  for each  $t$  and is bounded measurable in  $t$ . We denote  $x(t, s)$  and  $u(t, s)$  as  $x(t)$  and  $u(t)$ , respectively. Here  $A$  is assumed to be an infinitesimal generator of the strongly continuous semigroup  $S(t)$ . Then (1) becomes

(1')  $\dot{x}(t) = Ax(t) + u(t)$ .

Along with the system (1'), the following equation is considered:

(2)  $\dot{x}_n(t) = A_n x_n(t) + u_n(t)$ ,

where  $x_n(t)$  is an element of  $X_n$  for each  $t, 0 \leq t < \infty$ ,  $u_n(t)$  is an element of  $X_n$  for each  $t$  and is bounded measurable in  $t$ , and  $A_n$  is the operator from  $X_n$  into  $X_n$  which is an infinitesimal generator of a strongly continuous semigroup  $S_n(t)$ .

In the case of  $A_n$  being the difference operator on  $X_n$  into  $X_n$ , (1') is the original system which we consider and (2) is the approximating system given by the finite difference method. In what follows  $A_n$  is assumed to be a bounded linear operator.

As an example of the mapping  $P_n$  consider the following :

$$X = L_2(0, 1), \quad x(s) \in X, \quad s \in [0, 1],$$

$$P_n x(s) = \sum_0^{n-1} \alpha_k c_k(s),$$

where

$$\alpha_k = n \int_{k/n}^{k/n + 1/n} x(s) ds,$$

$$C_k(s) = \begin{cases} 1, & k/n \leq S \leq k/n + 1/n, \\ 0, & k/n \geq S \quad \text{or} \quad S \geq k/n + 1/n. \end{cases}$$

In this case  $x_n$  is a space of finitely-valued functions.

The solution of (1') is formally written as

$$(3) \quad x(t) = S(t)x(0) + \int_0^t S(t - \sigma)u(\sigma) d\sigma.$$

Since  $u(t)$  is bounded and measurable, the integral on the right-hand side of (3) exists in the sense of Bochner. Sufficient conditions under which (3) represents the strong solution to (1') (differentiable in the strong sense) have been given by Bala-krishnan [4]. But since the right-hand side of (3) exists in the sense of Bochner, we define (3) as the solution of (1') for any bounded measurable function  $u(t)$  and a given initial condition  $x(0)$ ; i.e., we deal with the mild solution [5] in this paper.

For simplicity's sake, one can assume without loss of generality that  $x(0) = 0$  in the sequel. Throughout the following discussion, we assume that the conditions (C) and (S) are satisfied.

Now the mappings  $L(T)$  and  $L_n(T)$  are defined as follows :

$$(4) \quad L(T)u = \int_0^T S(T - \sigma)u(\sigma) d\sigma, \quad u(t) \in X \quad \text{for each } t,$$

and

$$(5) \quad L_n(T)u_n = \int_0^T S_n(T - \sigma)u_n(\sigma) d\sigma, \quad u_n(t) \in X_n \quad \text{for each } t.$$

Under the hypothesis of Trotter's theorem, it is shown that  $\lim \|S(t)x - S_n(t)P_n x\| = 0$  for any  $x \in X$ , and again by the condition (S), the following equation is obtained for any  $u \in X$  :

$$\lim \|L(T)u - L_n(T)P_n u\| = 0.$$

The constraints  $C, C_n$  and the sets  $U, U_n$  of admissible controls are defined respectively as follows :

$$C = \{x \in X : \|x\| \leq 1\},$$

$$C_n = \{x_n \in X_n : \|x_n\| \leq 1\},$$

$$U = \{ \text{bounded measurable function } u(t) : u(t) \in C \text{ for almost all } t \},$$

$$U_n = \{ \text{bounded measurable function } u_n(t) : u_n(t) \in C_n \text{ for almost all } t \}.$$

From the property of  $P_n$  it can be shown that  $P_n C = C_n$ .

In this article we discuss the convergence of approximating solutions of the following Problems A and B, which are commonly seen in optimal control problems.

**PROBLEM A. Final value problem.**

(A1) The original final value problem in system (1') is to find the admissible control minimizing  $\|x(T) - y\| = \|L(T)u + S(T)x(0) - y\|$  for given  $y \in X$  at a fixed time  $T$ , starting from a given initial state  $x(0)$ .

(A2) The approximating final value problem in system (2) is to find the admissible control minimizing  $\|x_n(T) - P_n y\| = \|L_n(T)u_n + S_n(T)x_n(0) - P_n y\|$  for given  $P_n y \in X_n$  at a fixed time  $T$ , starting from a given approximating initial state  $P_n x(0)$ , with  $n \rightarrow \infty$ .

Then our problem is reduced to, "Does a sequence of approximating optimal controls of the problem (A2) converge to an optimal control of the problem (A1)?"

**PROBLEM B. Time-optimal problem.**

(B1) The original time-optimal problem in system (1') is to find the admissible control which transfers a system from a given initial state  $x(0)$  to a given target state  $x \in X$  in minimum time.

(B2) The approximating time-optimal problem is to find the admissible control which transfers a system from a given initial state  $P_n x(0)$  to a given approximating target state  $P_n x \in X_n$  in minimum time.

Then our problem is, "Does a sequence of these approximating time-optimal controls of the problem (B2) converge to a time-optimal control of the problem (B1)?"

**4. Condition of convergence.**

**PROBLEM A. Final value problem.** Let  $u^0(t)$  be the optimal control of the original final value problem and  $u_n^0(t)$  be the optimal control of the approximating final value problem at the  $n$ th degree of approximation.

Now we make the assumption that the approximating optimal solutions  $u_n^0(t)$  exist and the conditions (C), (S) and

$$(C^*) \quad \lim \|S_n^*(t)P_n f - S^*(t)f\| = 0, \quad f \in X,$$

are satisfied, where  $S^*$  and  $S_n^*$  are the adjoint operators of  $S$  and  $S_n$  considered on the space  $X$  and  $X_n$  respectively.

Moreover, let  $B_2[X, T]$  be the space of strongly measurable functions  $x(t)$  with range in  $X$  such that  $\int_0^T \|x(t)\|^2 dt < \infty$ .  $B_2[C, T]$  is the set of all  $u(t)$  such that  $\text{ess sup } \|u(t)\| \leq 1$  and  $B_2[C_n, T]$  is the set of all  $u_n(t)$  such that  $\text{ess sup } \|u_n(t)\| \leq 1$ .  $B_2[C, T]$  is the admissible control set for the exact problem, and  $B_2[C_n, T]$  is the admissible control set for the approximate problem.

**LEMMA A1.** *Let the conditions (C), (S) and (C\*) be satisfied. Then there exists a weakly convergent subsequence  $\{u_n^0(t)\}$  of  $\{u_n^0(t)\}$  in  $B_2[X, T]$ , and its weak limit  $u(t)$  belongs to  $B_2[C, T]$ . Moreover  $u(t)$  is an optimal control of the original final value problem.*

*Proof.* The first part of the lemma is easily shown from the fact that  $u_n^0(t) \in B_2[C_n, T] \subseteq B_2[C, T], B_2[X, T]$  is a reflexive Banach space, and the set  $B_2[C, T]$  is bounded, closed and convex in  $B_2[X, T]$ .

Next we shall show that  $L_n(T)u_n^0$  converges weakly to  $L(T)u$ . In the following we write  $u_n^0(t)$  for  $u_n^0(t)$ .

Since  $\{X - X_n\} \perp X_n$  and  $P_n$  is a projection operator, for  $f \in X$ , we have

$$[L_n(T)u_n^0, f]_X = [L_n(T)u_n^0, P_n f]_X = [L_n(T)u_n^0, P_n f]_{X_n},$$

$$[u_n^0, S_n^*(\sigma)P_n f]_{X_n} = [u_n^0, S_n^*(\sigma)P_n f]_X,$$

where  $[\cdot, \cdot]_X$  and  $[\cdot, \cdot]_{X_n}$  denote inner product on  $X$  and  $X_n$ . Hence, we have

$$\begin{aligned}
 [L_n(T)u_n^0 - L(T)u, f]_X &= [L_n(T)u_n^0, f]_X - [L(T)u_n^0, f]_X + [L(T)u_n^0 - L(T)u, f]_X \\
 &= [L_n(t)u_n^0, P_n f]_{X_n} - [L(T)u_n^0, f]_X + [L(T)u_n^0 - L(T)u, f]_X \\
 &= \int_0^T [S_n(T - \sigma)u_n^0(\sigma), P_n f]_{X_n} d\sigma \\
 &\quad - \int_0^T [S(T - \sigma)u_n^0(\sigma), f]_X d\sigma \\
 &\quad + \int_0^T [S(T - \sigma)(u_n^0(\sigma) - u(\sigma)), f]_X d\sigma \\
 &= \int_0^T [u_n^0(\sigma), S_n^*(T - \sigma)P_n f]_{X_n} d\sigma \\
 &\quad - \int_0^T [u_n^0(\sigma), S^*(T - \sigma)f]_X d\sigma \\
 &\quad + \int_0^T [u_n^0(\sigma) - u(\sigma), S^*(T - \sigma)f]_X d\sigma \\
 &= \int_0^T [u_n^0(\sigma), S_n^*(T - \sigma)P_n f - S^*(T - \sigma)f]_X d\sigma \\
 &\quad + \int_0^T [u_n^0(\sigma) - u(\sigma), S^*(T - \sigma)f]_X d\sigma.
 \end{aligned}$$

The first term of the right-hand side goes to zero as  $n \rightarrow \infty$  by conditions (C\*), (S), and the second term goes to zero as  $n \rightarrow \infty$  since the integrand is bounded and converges to zero for almost all (a.a.)  $\sigma$ .

Next we shall show that  $u(t)$  is optimal. It is clear that the following inequality is satisfied for any  $v(t) \in U$ :

$$\|L_n(T)u_n^0 - P_n v\| \leq \|L_n(T)P_n v - P_n v\|.$$

By the resonance theorem [6] we have

$$\begin{aligned} \|L(T)u - y\| &\leq \liminf \|L_n(T)u_n^0 - P_n y\| \\ &\leq \liminf \|L_n(T)P_n v - P_n y\| \\ &= \|L(T)v - y\|. \end{aligned}$$

Hence  $u(t)$  is optimal.

**THEOREM A1.** *Let the conditions (C), (S) and (C\*) be satisfied. Suppose that  $P_n y \notin L_n(T)u_n$ ,  $y \notin L(T)u$ , where  $u_n \in U_n$ ,  $u \in U$  and  $u^0(t)$  is the unique optimal control of the original final value problem. Uniqueness has been assumed in this theorem. Then,  $u_n^0(t)$  converges to  $u^0(t)$  strongly for a.a.  $t$ .*

*Proof.* By Lemma A1, there exists a subsequence  $\{u_{n_j}^0(t)\}$  of  $\{u_n^0(t)\}$  which converges weakly to an optimal control  $u(t)$  in  $B_2[X, T]$ .  $u^0(t)$  being unique, it follows that  $u(t) = u^0(t)$  for a.a.  $t$ . Since  $P_n y \notin L_n(T)u_n$ ,  $y \notin L(T)u$  and  $u^0(t)$ ,  $u_n^0(t)$  are optimal, we see that [5]

$$\|u^0(t)\| = 1 \text{ for a.a. } t, \quad \|u_{n_j}^0(t)\| = 1 \text{ for a.a. } t.$$

If we define the norm  $\|x\|$  of the space  $B_2[X, T]$  by  $\|x\|^2 = \int_0^T \|x(t)\|^2 dt$ , it follows that  $\|u_{n_j}^0\| = \|u^0\|$ .

Therefore it can be shown that  $\{u_{n_j}^0(t)\}$  converges strongly to  $u^0(t)$  in  $B_2[X, T]$ , because  $B_2[X, T]$  is a Hilbert space,  $\lim_{n_j \rightarrow \infty} \|u_{n_j}^0\| = \|u^0\| = T$  and  $\{u_{n_j}(t)\}$  converges weakly to  $u^0(t)$ . Any subsequence of  $\{u_n^0(t)\}$  has its subsequence converging weakly to  $u^0(t)$ , and therefore it has the subsequence converging strongly to  $u^0(t)$ . Consequently  $\{u_n^0(t)\}$  converges strongly to  $u^0(t)$  in  $B_2[X, T]$ . By the definition of the norm of  $B_2[X, T]$ , it can be concluded that  $u_n^0(t)$  converges to  $u^0(t)$  for a.a.  $t$ .

**PROBLEM B.** *Time-optimal problem.*

**LEMMA B1.** *Let the conditions (C), (S) and (C\*) be satisfied. Further, if the following conditions (L1) and (L2) are satisfied :*

- (L1) *there exist some finite time  $T_n$  and a control  $u_n \in U_n$  such that  $L_n(T_n)u_n = P_n x$  for each  $n$ ,*
- (L2)  *$T_n$  converges to  $T$ ;*

*then there exists  $u(t) \in U$  such that  $L(T)u = x$ .*

*Proof.* We consider the space  $B_2[X, T_0]$  and the set  $B_2[C, T_0]$ , where  $T_0$  is a sufficiently large value. For  $t \geq T_n$ , we put  $u_n(t) = 0$ . Clearly  $u_n(t) \in B_2[C_n, T_0] \subseteq B_2[C, T_0]$ . Since  $B_2[X, T_0]$  is a reflexive space and  $B_2[C, T_0]$  is closed, convex and bounded in  $B_2[X, T_0]$ , we can find a subsequence  $\{u_{n_j}(t)\}$  which converges weakly to  $u(t) \in B_2[C, T_0]$  (i.e.,  $u(t) \in U$ ).

Next it will be shown that  $L(T)u = x$ . We prove that  $[P_{n_j} x, f]_X$  converges to  $[x, f]_X$  and  $[L_{n_j}(T_{n_j})u_{n_j}, f]_X$  converges to  $[L(T)u, f]_X$ , where  $f \in X$ :

$$\begin{aligned} &[L_{n_j}(T_{n_j})u_{n_j} - L(T)u, f]_X \\ &= [L_{n_j}(T_{n_j})u_{n_j}, P_{n_j} f]_{X_{n_j}} - [L(T)u_{n_j}, f]_X + [L(T)u_{n_j} - L(T)u, f]_X \\ &= \int_0^{T_{n_j}} [S_{n_j}(T_{n_j} - \sigma)u_{n_j}(\sigma), P_{n_j} f]_{X_{n_j}} d\sigma + \int_0^T [S_{n_j}(T_{n_j} - \sigma)u_{n_j}(\sigma), P_{n_j} f]_{X_{n_j}} d\sigma \\ &\quad - \int_0^T [S(T - \sigma)u_{n_j}(\sigma), f]_X d\sigma + \int_0^T [S(T - \sigma)(u_{n_j}(\sigma) - u(\sigma)), f]_X d\sigma \end{aligned}$$



$$\begin{aligned}
 &= \int_T^{T_{n_j}} [S_{n_j}(T_{n_j} - \sigma)u_{n_j}(\sigma), P_{n_j}f]_{X_{n_j}} d\sigma + \int_0^T [u_{n_j}(\sigma), S_{n_j}^*(T_{n_j} - \sigma)P_{n_j}f]_{X_{n_j}} d\sigma \\
 &\quad - \int_0^T [u_{n_j}(\sigma), S^*(T - \sigma)f]_X d\sigma + \int_0^T [u_{n_j}(\sigma) - u(\sigma), S^*(T - \sigma)f]_X d\sigma \\
 &= \int_T^{T_{n_j}} [S_{n_j}(T_{n_j} - \sigma)u_{n_j}(\sigma), f]_X d\sigma \\
 &\quad + \int_0^T [u_{n_j}(\sigma), S_{n_j}^*(T_{n_j} - \sigma)P_{n_j}f - S^*(T - \sigma)f]_X d\sigma \\
 &\quad + \int_0^T [u_{n_j}(\sigma) - u(\sigma), S^*(T - \sigma)f]_X d\sigma.
 \end{aligned}$$

The first term from the last equality sign goes to zero as  $n \rightarrow \infty$  by condition (S) and the third term goes to zero since the integrand is bounded and converges to zero for a.a.  $\sigma$ . The second term goes to zero by condition (S), condition (C\*) and the continuity of  $S^*(t)$  in  $t$ . Hence  $L(T)u = x$ .

*Remark.* In order to obtain the condition (C\*), we again make use of Trotter's theorem. Since condition (S) is always satisfied for  $S_n^*(t)$ , it is sufficient to check the validity of condition (C) for  $A_n^*$ .

**LEMMA B2.** *Let  $S(t)$  be a strongly continuous group of operators. If  $x$  is attainable at time  $T$  with  $u(t) \in U$ , then for sufficiently large  $n$ , there exist  $u_n(t) \in U_n$ , such that  $L_n(T + \varepsilon(n))u_n = P_n x$ , where  $\varepsilon(n) \rightarrow 0$  as  $n \rightarrow \infty$ .*

*Proof.* By  $L(T)u = x$ , the following inequality is obtained:

$$\|L_n(T)P_n u - P_n x\| \leq \|L_n(T)P_n u - P_n L(T)u\| + \|P_n L(T)u - P_n x\| \leq \varepsilon^2(n).$$

Thus we can set

$$P_n x = L_n(T)P_n u + \varepsilon^2(n)Z_n,$$

where  $\|Z_n\| = 1$ .

Since  $S(t)$  is a group,

$$\begin{aligned}
 P_n x &= \int_0^T S_n(T - \sigma)P_n u(\sigma) d\sigma + \varepsilon(n) \int_0^{\varepsilon(n)} S_n(T + \varepsilon(n) - \sigma)S_n(\sigma - \varepsilon(n) - T)Z_n d\sigma \\
 &= \int_{\varepsilon(n)}^{T + \varepsilon(n)} S_n(T + \varepsilon(n) - \sigma)P_n u(\sigma - \varepsilon(n)) d\sigma + \int_0^{\varepsilon(n)} S_n(T + \varepsilon(n) - \sigma)w_n(\sigma) d\sigma,
 \end{aligned}$$

where

$$w_n(\sigma) = \varepsilon(n)S_n(\sigma - \varepsilon(n) - T)Z_n.$$

By condition (S), we have  $w_n(\sigma) \in U_n$  for sufficiently large  $n$ .

Finally we may only define a new admissible control  $u_n(t)$  as follows:

$$u_n(t) \begin{cases} w_n(t), & 0 \leq t \leq \varepsilon(n), \\ P_n u(t - \varepsilon(n)), & \varepsilon(n) \leq t \leq T + \varepsilon(n). \end{cases}$$

**THEOREM B1.** *Let  $S(t)$  be a strongly continuous group of operators, and let the conditions (C), (S), (C\*) and the following two conditions be satisfied:*

(T1)  $T_n^0$  is an optimal time and  $u_n^0$  is an associated optimal control such that  $L_n(T_n^0)u_n^0 = P_n x$ .

(T2) There exists  $w(t) \in U$  such that  $L(T)w = x$  at some finite time  $T$ .

Then  $T_n^0$  converges to an optimal time  $T^0$  as  $n \rightarrow \infty$ .

*Proof.* By condition (T2) and Lemma B2, there exists a subsequence  $\{T_{n_j}^0\}$  which converges to  $T'$ . Hence we see that an optimal time exists, by condition (C\*) and Lemma B1; i.e., there exists an admissible control  $u \in U$  such that  $L(T')u = x$ .

If we denote the optimal time by  $T^0$ , the following inequality can be obtained by Lemma B2 and the optimality of  $T_{n_j}^0$ :

$$T_{n_j}^0 \leq T^0 + \varepsilon(n).$$

Since  $T_{n_j}$  converges to  $T'$ , we have  $T' \leq T^0$  and hence  $T' = T^0$  by the optimality of  $T^0$ .

The above proof is true for any convergent subsequence and hence  $T_n^0$  converges to  $T^0$ .

**LEMMA B3.** *Let  $S(t)$  be a holomorphic semigroup of operators, and let conditions (C) and (S) be satisfied. Suppose that  $u(t)$  is continuously differentiable in  $[0, T]$  and  $\|A_n P_n x\| < 1$ . If  $x$  is attainable at time  $T$  with  $u(t) \in U$ , where  $x \in D(A)$ , then for sufficiently large  $n$ , there exists  $u_n(t) \in U_n$  such that  $L_n(T + \gamma(n))u_n = P_n x$ , where  $\gamma(n) \rightarrow 0$  as  $n \rightarrow \infty$ .*

*Proof.* Let us define  $y_n = L_n(T)P_n u$ .

Since the function

$$x_n(t) = \frac{K-t}{K} y_n + \frac{t}{K} P_n x, \quad 0 \leq t \leq K,$$

satisfies (2) in which we put

$$v_n(t) = \frac{1}{K}(P_n x - y_n) - \frac{K-t}{K}(A_n y_n - A_n P_n x) - A_n P_n x,$$

the control  $v_n(t)$  transfers  $y_n$  to  $P_n x$  in time  $K$ . By the assumption  $x = L(T)u$ , we have

$$\|P_n x - y_n\| = \|P_n L(T)u - L_n(T)P_n u\| = \varepsilon^2(n),$$

where  $\varepsilon(n) \rightarrow 0$  as  $n \rightarrow \infty$ .

Next it will be shown that

$$\|A_n y_n - A_n P_n x\| = \eta^2(n),$$

where  $\eta(n) \rightarrow 0$  as  $n \rightarrow \infty$ .

Since  $A_n$  is a bounded operator and  $u(t)$  is continuously differentiable in  $[0, T]$ , we have  $A_n S_n(T - \tau)P_n u(\tau) \in L_1(0, T)$  and  $S_n(T - \tau)P_n \partial u(\tau)/\partial \tau \in L_1(0, T)$ .

From the above facts and the boundedness of  $A_n$  it follows that  $S_n(T - \tau)P_n u(\tau)$  is strongly differentiable and  $\partial(S_n(T - \tau)P_n u(\tau))/\partial \tau \in L_1(0, T)$ . Hence by the

continuity of  $\partial(S_n(T - \tau)P_n u(\tau))/\partial\tau$  and again by the boundedness of  $A_n$ , we have

$$\begin{aligned} A_n y_n &= A_n \int_0^T S_n(T - \tau)P_n u(\tau) \, d\tau = \int_0^T A_n S_n(T - \tau)P_n u(\tau) \, d\tau \\ &= - \int_0^T \frac{\partial}{\partial\tau} (S_n(T - \tau)P_n u(\tau)) \, d\tau + \int_0^T S_n(T - \tau)P_n \frac{\partial}{\partial\tau} u(\tau) \, d\tau \\ &= - P_n u(T) + S_n(T)P_n u(0) + \int_0^T S_n(T - \tau)P_n \frac{\partial}{\partial\tau} u(\tau) \, d\tau. \end{aligned}$$

Now we define  $Z_\delta$  as follows :

$$Z_\delta = S(\delta)x = \int_0^T S(T + \delta - \tau)u(\tau) \, d\tau.$$

Since  $S(t)$  is a holomorphic semigroup, i.e.,  $\|AS(t)\| < C/t$ , where  $C$  is some constant, we see that  $AS(T + \delta - \tau)u(\tau) \in L_1(0, T)$  and  $S(T + \delta - \tau)\partial u(\tau)/\partial\tau \in L_1(0, T)$ . From these and the holomorphic property of  $S(t)$ , it can be easily shown that  $S(T + \delta - \tau)u(\tau)$  is strongly differentiable and  $\partial(S(T + \delta - \tau)u(\tau))/\partial\tau \in L_1(0, T)$ . Hence the next equality can be obtained by noticing that  $A$  is a closed operator :

$$\begin{aligned} AZ_\delta &= A \int_0^T S(T + \delta - \tau)u(\tau) \, d\tau \\ &= \int_0^A AS(T + \delta - \tau)u(\tau) \, d\tau \\ &= - S(\delta)u(T) + S(T + \delta)u(0) + \int_0^T S(T + \delta - \tau) \frac{\partial}{\partial\tau} u(\tau) \, d\tau. \end{aligned}$$

It is clear that the following inequality is satisfied :

$$\lim_{n \rightarrow \infty} \|A_n y_n - Ax\| \leq \lim_{\delta \rightarrow 0} \lim_n \|A_n y_n - Az_\delta\| + \lim_{\delta \rightarrow 0} \|AZ_\delta - Ax\|.$$

Now by conditions (C) and (S), we have

$$\begin{aligned} \lim_{\delta \rightarrow 0} \lim_n \|A_n y_n - AZ_\delta\| &= \lim_{\delta \rightarrow 0} \lim_n \| - P_n u(T) + S_n(T)P_n u(0) \\ &\quad + \int_0^T S_n(T - \tau)P_n \frac{\partial}{\partial\tau} u(\tau) \, d\tau + S(\delta)u(T) - S(T + \delta)u(0) \\ &\quad - \int_0^T S(T + \delta - \tau) \frac{\partial}{\partial\tau} u(\tau) \, d\tau \| \\ &= \lim_{\delta \rightarrow 0} \| - u(T) + S(T)u(0) + \int_0^T S(T - \tau) \frac{\partial}{\partial\tau} u(\tau) \, d\tau \\ &\quad + S(\delta)u(T) - S(T + \delta)u(0) - \int_0^T S(T + \delta - \tau) \frac{\partial}{\partial\tau} u(\tau) \, d\tau \| = 0. \end{aligned}$$

On the other hand, by the fact that  $x \in D(A)$  and strong continuity of  $S(t)$ , we

have that

$$\lim_{\delta \rightarrow 0} \|AZ_\delta - AX\| = \lim_{\delta \rightarrow 0} \|S(\delta)AX - AX\| = 0.$$

Consequently we have  $\lim \|A_n y_n - A_n P_n x\| = \eta^2(n)$ , since  $\lim \|A_n y_n - Ax\| = 0$ , and by the triangular inequality,

$$\|A_n y_n - A_n P_n x\| \leq \|A_n y_n - Ax\| + \|Ax - A_n P_n x\|.$$

Since  $\|A_n P_n x\| < 1$ , if we choose  $K = \gamma(n) = \max(\varepsilon(n), \eta(n))$ , we have  $v_n(t) \in U_n$  for sufficiently large  $n$ .

Finally we define a new admissible control  $u_n(t)$  as follows:

$$u_n(t) = \begin{cases} P_n u(t), & 0 \leq t \leq T, \\ v_n(t - T), & T \leq t \leq T + \gamma(n). \end{cases}$$

Then we have the conclusion of Lemma B3, i.e.,  $P_n x = L_n(T + \gamma(n))u_n$ .

Now we can easily get the following theorem from Lemma B3.

**THEOREM B2.** *Let  $S(t)$  be a holomorphic semigroup of operators, and let the conditions (C), (S), (C\*) and the following conditions be satisfied:*

- (T1)  $T_n^0$  is an optimal time control and  $u_n^0$  is an associated optimal control such that  $L_n(T_n^0)u_n^0 = P_n x$ .
- (T2) There exists  $u(t) \in U$  such that  $L(T)u = x$  at some time  $T$ , where  $x \in D(A)$ .
- (T3)  $u(t)$  is continuously differentiable in  $[0, T]$ .
- (T4)  $\|A_n P_n x\| < 1$ .

Then  $T_n^0$  converges to an optimal time  $T^0$  as  $n \rightarrow \infty$ .

*Remark.* One should note that in the problem of getting to the origin,  $\|A_n P_n x\| < 1$  is automatically satisfied, since  $x = 0 \in D(A)$ .

*Remark.* In the case of the space being a Hilbert space, it is clear that an optimal control sequence  $\{u_n^0(t)\}$  converges to an optimal control for a.a.  $t$ , since the optimal control associated with the target  $x$  is unique [4], [5].

**5. Conclusion.** In this paper we have treated the convergence of approximating solutions for both the final value and time optimal control problems for linear distributed parameter systems. For the final value problem we have shown, under stated conditions, that where the sequence of approximating solutions exists, a subsequence converges weakly to the true optimal control (Lemma A1). In addition, if the desired system state is not in the range of certain system operations, and if the true optimal control is unique, then the preceding sequence converges strongly to the optimal control (Theorem A1).

For the time optimal control problem we have shown, under stated conditions, that the sequence of optimal times for the approximate problems exists and converges to the optimal time for the exact problem (Theorems B1, B2).

Theorem A1 is true even if  $S(t)$  is not a group but a strongly continuous semigroup of operators. But, Lemma B2 and hence Theorem B1, Lemma B3 and hence Theorem B2 are valid only if  $S(t)$  is a group of operators and a holomorphic semigroup of operators respectively.

In case  $S(t)$  is a semigroup of operators, our theorems are also likely to be true, if we can choose more refined projection operators  $P_n$  and difference operators  $A_n$ .

**Acknowledgment.** The authors wish to thank Associate Professor H. Ishigaki, the School of Education, Waseda University, for many helpful suggestions and discussions.

## REFERENCES

- [1] E. I. AXELBAND, *Function space methods for the optimal control of a class of distributed parameter control systems*, Proc. Joint Automatic Control Conf., 1965, pp. 374–380.
- [2] H. F. TROTTER, *Approximation of semi-groups of operators*, Pacific J. Math., 8 (1958), pp. 887–919.
- [3] P. D. LAX, *Survey of the stability of linear finite difference equations*, Comm. Pure Appl. Math., 9 (1956), pp. 267–293.
- [4] A. V. BALAKRISHNAN, *Optimal control problems in Banach spaces*, this Journal, 3 (1968), pp. 109–127.
- [5] A. FRIEDMAN, *Optimal control in Banach spaces*, J. Math. Anal. Appl., 19 (1968), pp. 35–55.
- [6] K. YOSHIDA, *Functional Analysis*, Springer-Verlag, Berlin, 1965.

## SOME EFFECTS OF MEASUREMENT UNCERTAINTY IN LINEAR MULTISTAGE GAMES\*

W. W. WILLMAN†

**Abstract.** A class of linear multistage games is examined in which the players have noisy measurements of the state. If a linear solution exists for such a game, it is shown to be related to one of the solutions of the corresponding deterministic game by the certainty-equivalence principle. This solution is generally not the one ordinarily associated with this deterministic game. Results are obtained by constructing state variables from the players' measurement sequences and applying known control theoretic results to pairs of optimal control problems associated with linear game solutions.

**1. Introduction.** A certainty-equivalence principle holds for a class of multistage games in which the players have noisy measurements of the state. As in the control theory context, the games in this class have linear dynamics, quadratic payoffs, and measurements perturbed by sequences of independent additive Gaussian random variables. The measurement noise sequences affecting the players are statistically independent. Sufficient conditions for this certainty-equivalence result, and the corresponding control laws, can be expressed in terms of a set of implicit difference equations.

Unlike its counterpart in optimal control theory, however, this result does not apply to the state feedback control laws usually computed as a solution to the corresponding deterministic game, but to a more complicated pair of optimal control laws. One implication of this certainty-equivalence result is that the mean sample path generated by the optimal control laws in the stochastic game coincides with the optimal trajectory in the deterministic game.

The type of game examined here is an extension of a deterministic game solved earlier by Ho, Bryson and Baron [4]. Several other stochastic extensions of this game have already been investigated: by Behn and Ho [1] for the case in which one of the players has perfect measurements, by Rhodes and Luenberger [9] for the case in which one player has no measurements, and by Rhodes and Luenberger [10] for a case in which both players have noisy measurements and constrained state estimators. The mean optimal sample paths of these other extensions were shown by their investigators to possess the property described above. The solution to the game considered here, however, is not obtained as explicitly as in these other cases. The formal extension of the following results to the continuous time case is discussed in Willman [12].

**2. Notation.** All matrices and vectors in the following have real components. Capital letters are used to denote matrices, lower-case letters to denote vectors. It is also convenient to introduce the notation

$$\langle A \rangle_{i,j}$$

---

\* Received by the editors April 16, 1970, and in revised form August 6, 1970.

† Operations Research Group, Naval Research Laboratory, Washington, D.C. 20390. This paper represents part of a dissertation submitted to Harvard University in partial fulfillment of the requirements for the degree of Doctor of Philosophy. The work was supported in part by the National Aeronautics and Space Administration under Contract NGR-22-007-068, by the Office of Naval Research under Contract N00014-67-A-0298-0006, and by the Division of Engineering and Applied Physics, Harvard University.

for the  $(i, j)$ th partition of a partitioned matrix  $A$ . Only one subscript is used for a partitioned vector, or a vertically or horizontally partitioned matrix.  $A^T$  is used to denote the transpose of a matrix  $A$ .  $I_r$  is used to denote the  $r$ -dimensional identity matrix when the dimension is not clear from the context.

**3. A particular class of games.** For greater clarity, the main results are derived for a restricted class of zero-sum games. In accordance with the original motivation, the two players are called the pursuer and evader. The extent to which the results generalize will be discussed later.

The generic game of interest here has dynamics defined by the following linear vector difference equation, where  $i$  is an integer index (to be called time) varying from 0 to  $N - 1$  for some fixed integer  $N$ :

$$(1) \quad x(i + 1) = x(i) + G_p(i)u(i) - G_e(i)v(i) + n(i), \quad i = 0, 1, \dots, N - 1,$$

where

- $x$ : state of dynamic system,
- $u$ : pursuer's control,
- $v$ : evader's control,
- $n$ : random variable to model uncertainty in the dynamics (process noise),
- $x(0)$ : Gaussian  $(\bar{x}, P_0)$  random variable.

The players receive the following noisy measurements of the state:

$$(2) \quad z_p(i) = H_p(i)x(i) + w_p(i) \quad (\text{pursuer}),$$

$$(3) \quad z_e(i) = H_e(i)x(i) + w_e(i) \quad (\text{evader}).$$

The vectors of the sequence

$$\left\{ \begin{bmatrix} n(i) \\ \text{---} \\ w_p(i) \\ \text{---} \\ w_e(i) \end{bmatrix} \right\}$$

are independent zero-mean Gaussian random variables, statistically independent of  $x(0)$ , and have covariance matrices

$$\begin{bmatrix} Q(i) & 0 & 0 \\ \text{---} & \text{---} & \text{---} \\ 0 & R_p(i) & 0 \\ \text{---} & \text{---} & \text{---} \\ 0 & 0 & R_e(i) \end{bmatrix}.$$

The criterion associated with this game, which the pursuer should be regarded as trying to minimize and the evader as trying to maximize, is the quadratic form

$$(4) \quad J \triangleq \mathcal{E} \left\{ x^T(N)S_N x(N) + \sum_{i=0}^{N-1} [u^T(i)B(i)u(i) - v^T(i)C(i)v(i)] \right\},$$

where  $\bar{e}$  denotes prior expected value. The matrices  $B(i)$ ,  $R_p(i)$ ,  $R_e(i)$  and  $C(i)$  are symmetric and positive definite;  $S_N$  is symmetric and positive semidefinite.

The control values in equations (1) and (4) should be regarded as being generated from the available measurements by a pair of pursuit and evasion strategies.

DEFINITION. In the context of this game, a *pursuit strategy* is a function assigning a control sequence  $\{u(i)\}$  to each measurement sequence  $\{z_p(i)\}$ . *Evasion strategies* are similarly defined.

DEFINITION. A pursuit strategy is *admissible* if and only if for  $i = 0, \dots, N - 1$ :

(i) the value of  $u(i)$  is determined by the measurement subsequence  $\{z_p(0), \dots, z_p(i)\}$ , a “nonanticipatory” requirement; and

(ii)  $u(i)$  is Borel measurable as a function of the measurements, a technical requirement for the expectation operator in (4) to be meaningful.

This holds similarly for evasion strategies.

Strategies will be denoted by capital letters to distinguish them from the control values they determine. “ $\mathcal{U}$ ” and “ $\mathcal{V}$ ” will be used to denote the sets of admissible pursuit and evasion strategies, respectively.

The value of the criterion  $J$  depends only on the pair of strategies employed by the two players. As usual, a solution to this game is defined to be a particular pair of admissible pursuit and evasion strategies  $(U^*, V^*)$  such that

$$(5) \quad J(U^*, V) \leq J(U^*, V^*) \leq J(U, V^*)$$

for all admissible  $U$  and  $V$ .

Strategy pairs which satisfy this saddle point condition are called *minimax*, and have the property that each strategy of the pair optimizes against the other. Minimax strategy pairs are generally agreed to be the only reasonable kind of “optimum” solution for zero-sum games, although other distinct solution concepts become important in nonzero-sum games, as discussed in Starr and Ho [11]. A strategy is called minimax if it is a component of a minimax strategy pair.

Following Luce and Raiffa [8], it can be shown by two applications of the saddle-point condition that, if strategy pairs  $(U^1, V^1)$  and  $(U^2, V^2)$  are both minimax, then

$$J(U^1, V^2) \leq J(U^1, V^1) \leq J(U^2, V^1) \leq J(U^2, V^2) \leq J(U^1, V^2).$$

Therefore, all four of these strategy pairs are minimax, and each pair gives the same value of  $J$ .

**4. The corresponding deterministic game.** The main objective of this paper is to establish relationships between solutions to the above stochastic game and those of the corresponding deterministic game, in which no randomness is present. More precisely, the deterministic game corresponding to the preceding stochastic game is defined by the equations

$$(6) \quad \begin{aligned} x(i + 1) &= x(i) + G_p(i)u(i) - G_e(i)v(i), & i = 0, \dots, N - 1, \\ x(0) &= \bar{x} \quad (\text{dynamics}), \end{aligned}$$



and

$$(7) \quad J_d \triangleq x^T(N)S_N x(N) + \sum_{i=0}^{N-1} [u^T(i)B(i)u(i) - v^T(i)C(i)v(i)] \quad (\text{criterion}),$$

where the players have perfect measurements of the state. It is a simple extension of the results of Ho, Bryson and Baron [4] that the states and controls generated by any solution of this game are related by the equations

$$(8) \quad u(i) = -F_p(i)x(i), \quad i = 0, \dots, N-1,$$

and

$$(9) \quad v(i) = -F_e(i)x(i), \quad i = 0, \dots, N-1,$$

where the matrix sequences  $\{F_p(i)\}$  and  $\{F_e(i)\}$  are determined explicitly from the parameters defining the game by a matrix difference equation which is too intricate to be profitably displayed here.

Since (6), (8) and (9) determine a unique trajectory (the state and control variable sequences) for this game, there exists at most one minimax trajectory associated with such a game.

If this game has a solution, (8) and (9) can be used to define a pair of linear feedback strategies by interpreting  $x$  as a state measurement instead of a time function. This particular pair of strategies has two desirable properties. Each component strategy satisfies the principle of optimality in the sense that its restriction to any terminal segment  $(k, \dots, N-1)$  of the time sequence is also minimax in the subgame that "starts" at time  $k$  with *any* value of the state  $x(k)$  which can be generated by an opponent's (possibly nonoptimal) strategy. Also, it is an easy modification of a result due to Berkovitz [2] to show that if any solution of this deterministic game exists, this strategy pair is also a solution. Because of these properties, this particular solution, referred to here as the *feedback solution*, is often the only one considered for this game.

However, the feedback solution is usually not the only one for the corresponding deterministic game, if it exists. For example, it is often true under these circumstances that mixtures of the feedback strategies and the corresponding "open loop" strategies, defined by (6), (8) and (9), are also minimax. In the stochastic game this indeterminacy of the solution often disappears in an interesting way.

It is appropriate to comment here that the above deterministic game is often called a *closed loop* game, since the players can base their control decisions on the current value of the state  $x$ . This is in distinction to the corresponding *open loop* game in which the players know the initial state exactly, and also the parameters defining the game, but have available no information on the current value of the state as the game progresses. Thus their control decisions must be a function of time only. The stochastic game under investigation here can be regarded as one in which the quality of information available to the players is intermediate between the open and closed loop extremes, although it will be seen later that it is perhaps more closely related to the closed loop type of game.

**5. A useful construction.** The basic approach used here is to obtain results about solutions to the stochastic game by examining pairs of associated stochastic optimal control problems. As it happens, it is useful to express these control problems in terms of new state variable sequences, each term of which is a vector with  $N$  partitions of equal dimension. These “enlarged state variables” are constructed, following Wonham [14], as

$$(10) \quad \langle \zeta_p(i) \rangle_j \triangleq \begin{cases} z_p(j), & j < i, \\ 0, & j \geq i, \end{cases}$$

$$(11) \quad \langle \zeta_e(i) \rangle_j \triangleq \begin{cases} z_e(j), & j < i, \\ 0, & j \geq i, \end{cases}$$

for  $i = 0, \dots, N$  and  $j = 0, \dots, N - 1$ . In terms of the information they contain, these enlarged state variables are essentially the sequences of initial segments of the players’ measurement sequences.

**6. Some preliminary results.** In view of the particular nature of the stochastic game under consideration, it is not surprising that it is worthwhile to confine the search for minimax strategies to those which specify the players’ controls as linear functions of their available measurements. It is convenient to begin this search by establishing some preliminary results concerning such linear strategies in this game.

LEMMA 1. *If the controls of the evader are determined by the strategy*

$$(12) \quad V' : v(i) = A_e(i)\bar{x} + \sum_{j=0}^i L_e(i,j)z_e(j), \quad i = 0, \dots, N - 1,$$

where  $\{A_e(i)\}$  and  $\{L_e(i,j)\}$  are arbitrary sequences of real matrices, then any admissible pursuit strategy  $U^*$  which minimizes  $J$  against  $V'$  (over  $\mathcal{U}$ ) satisfies the conditions

$$(13) \quad u(i) = K_p(i)\bar{x} + M_p(i)\hat{\sigma}_p(i), \quad i = 0, \dots, N - 1,$$

where  $\hat{\sigma}_p(i)$  is generated from the pursuer’s measurements by

$$(14) \quad \hat{\sigma}_p(i) = N_p(i)\bar{x} + T_p(i)\hat{\sigma}_p(i - 1) + E_p(i)z_p(i), \quad \hat{\sigma}_p(-1) = \begin{bmatrix} \bar{x} \\ 0 \end{bmatrix},$$

for  $i = 0, \dots, N - 1$  and almost all measurement sequences  $\{z_p(i)\}$ , where the matrices  $\{K_p(i)\}$ ,  $\{M_p(i)\}$ ,  $\{N_p(i)\}$ ,  $\{E_p(i)\}$  and  $\{T_p(i)\}$  are determined explicitly from  $\{A_e(i)\}$  and  $\{L_e(i,j)\}$  by an effective procedure (to be described in more detail in the proof).

*Proof sketch.* Define for  $i, j = 0, \dots, N - 1$  the sequence of horizontally partitioned matrices  $\{\Lambda_e(i)\}$  such that

$$(15) \quad \langle \Lambda_e(i) \rangle_j = \begin{cases} L_e(i,j), & j < i, \\ 0, & j \geq i, \end{cases}$$

and the sequence of vertically partitioned matrices  $\{\Theta_e(i)\}$  such that

$$(16) \quad \langle \Theta_e(i) \rangle_j = \begin{cases} I_r, & j = i, \\ 0, & j \neq i, \end{cases}$$

where  $r$  is the dimension of  $z_e$ .

It is tedious but straightforward to verify by substitution that, if the evader uses the strategy  $V'$ , the dynamics of (1) are equivalent to

$$(17) \quad \begin{bmatrix} x(i+1) \\ \zeta_e(i+1) \end{bmatrix} = \begin{bmatrix} I - G_e(i)L_e(i)H_e(i) & -G_e(i)\Lambda_e(i) \\ \Theta_e(i)H_e(i) & I \end{bmatrix} \begin{bmatrix} x(i) \\ \zeta_e(i) \end{bmatrix} + \begin{bmatrix} G_p(i) \\ 0 \end{bmatrix} u(i) \\ - \begin{bmatrix} G_e(i)A_e(i) \\ 0 \end{bmatrix} \bar{x} + \begin{bmatrix} I & -G_e(i)L_e(i,i) \\ 0 & \Theta_e(i) \end{bmatrix} \begin{bmatrix} n(i) \\ w_e(i) \end{bmatrix},$$

where  $\begin{bmatrix} x(0) \\ \zeta_e(0) \end{bmatrix}$  is

$$\text{Normal} \left( \begin{bmatrix} \bar{x} \\ 0 \end{bmatrix}, \begin{bmatrix} P_0 & 0 \\ 0 & 0 \end{bmatrix} \right)$$

a priori, the criterion can be expressed as

$$(18) \quad J = \mathcal{E} \left\{ \begin{bmatrix} x^T(N) & \zeta_e^T(N) \end{bmatrix} \begin{bmatrix} S_N & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x(N) \\ \zeta_e(N) \end{bmatrix} \right. \\ + \sum_{i=0}^{N-1} \left( u^T(i)B(i)u(i) - \left( \bar{x}^T A_e^T(i) + [x^T(i) \mid \zeta_e^T(i)] \begin{bmatrix} H_e^T(i)L_e^T(i,i) \\ -\Lambda_e^T(i) \end{bmatrix} \right) \right. \\ \cdot C(i) \left( \begin{bmatrix} L_e(i,i)H_e(i) & \Lambda_e(i) \end{bmatrix} \begin{bmatrix} x(i) \\ \zeta_e(i) \end{bmatrix} + A_e(i)\bar{x} \right) \left. \right) \\ \left. - \mathcal{E} \left\{ \sum_{i=0}^{N-1} w_e^T(i)L_e^T(i,i)C(i) \left( L_e(i,i)w_e(i) \right. \right. \right. \\ \left. \left. \left. + 2 \left( A_e(i)\bar{x} + [L_e(i,i)H_e(i) \mid \Lambda_e(i)] \begin{bmatrix} x(i) \\ \zeta_e(i) \end{bmatrix} \right) \right) \right\} \right\},$$

and the pursuer's measurements are

$$(19) \quad z_p(i) = [H_p(i) \mid 0] \begin{bmatrix} x(i) \\ \zeta_e(i) \end{bmatrix} + w_p(i).$$

Since  $w_e(i)$  is statistically independent of  $x(i)$ ,  $\zeta_e(i)$ , and  $\bar{x}$  for any admissible strategy, the second term in (18) is independent of the control law. Therefore, the problem of minimizing  $J(U, V')$  over  $U$  in  $\mathcal{U}$  is equivalent to the stochastic control problem defined by (17), (19) and the first term of (18), where  $J$  is to be minimized.

As reformulated in terms of the enlarged state variables, this minimization problem has been converted into a stochastic optimal control problem with linear dynamics, quadratic cost, and additive independent Gaussian process and measurement noise sequences. This means that the well-known necessary conditions developed for this type of control problem developed by Joseph and Tou [5]

and Gunckel and Franklin [3] apply here. As applied to the particular example at hand, these necessary conditions are precisely (13) and (14), where the unspecified matrix sequences are determined from the matrices  $\{A_e(i)\}$  and  $\{L_e(i, j)\}$  and the parameters defining the game by a certain set of matrix difference equations with specified boundary conditions. Since the solution to this type of optimal control problem is well known, the details of these equations are not presented here. This completes the proof.

The random sequence  $\{\hat{\sigma}_p(i)\}$  in equation system (14) is to be interpreted as the sequence of estimates of the composite vector  $\sigma_p(i) \triangleq \begin{bmatrix} x(i) \\ \zeta_e(i) \end{bmatrix}$  produced by the appropriate Kalman-Bucy filter for the dynamics and measurements of (17) and (19). This means, from well-known results due to Kalman [6], that  $\hat{\sigma}_p(i)$  is the conditional expected value of  $\sigma_p(i)$ , given the sequence of measurements  $\{z_p(0), \dots, z_p(i)\}$  for  $i = 0, \dots, N - 1$ . For this reason, (13) and (14) will be referred to as the *certainty-equivalent* form of  $U^*$ .

**COROLLARY.** Since the preceding necessary conditions are satisfied by only one strategy, there is at most one admissible pursuit strategy (up to a set of measurement sequences of measure zero) which minimizes  $J$  against  $V'$  over  $\mathcal{V}$ .

**LEMMA 2.** In the context of Lemma 1, the strategy  $U^*$  can be expressed as

$$(20) \quad U^* : u(i) = A_p(i)\bar{x} + \sum_{j=0}^i L_p(i, j)z_p(j), \quad i = 0, \dots, N - 1,$$

where the matrices  $\{A_p(i)\}$  and  $\{L_p(i, j)\}$  are explicitly determined from the matrices  $\{A_e(i)\}$  and  $\{L_e(i, j)\}$  by an effective procedure.

*Proof sketch.* This proof merely consists of showing by induction that (20) is equivalent to (13) and (14) (i.e., gives the same control sequence for every measurement sequence) when  $\{A_p(i)\}$  and  $\{L_p(i, j)\}$  are computed from the matrices  $\{K_p(i)\}$ ,  $\{M_p(i)\}$ ,  $\{N_p(i)\}$ ,  $\{T_p(i)\}$  and  $\{E_p(i)\}$  according to a certain set of matrix difference equations developed according to Kalman and Bucy [7]. Again, the details are not given here.

For future reference, the functional dependence of the matrices  $\{A_p(i)\}$  and  $\{L_p(i, j)\}$  in (20) on  $\{A_e(i)\}$  and  $\{L_e(i, j)\}$  will be expressed as

$$(21) \quad \begin{aligned} A_p(i) &= f_{p,i}(\{A_e(i)\}, \{L_e(i, j)\}), \\ L_p(i, j) &= g_{p,i,j}(\{A_e(i)\}, \{L_e(i, j)\}). \end{aligned}$$

Equation (20) will be called the *canonical form* of the pursuit strategy  $U^*$ .

Finally, it is useful to establish sufficient conditions for the pursuit strategy  $U^*$  to minimize the criterion against the evasion strategy  $V'$ . This is done by using the well-known sufficient conditions for "linear-quadratic-Gaussian" optimal control problems. As applied to the control problem specified by (17), (19), and the first term of (18), these sufficient conditions lead to the following result.

**LEMMA 3.** In the context of the preceding lemmas, if the pursuit strategy  $U^*$  is determined from the evasion strategy  $V'$  according to (20) and (21), and if the matrices

$$B(i) + G_p^T(i)S_p(i + 1)G_p(i)$$

are positive definite for  $i = 0, \dots, N - 1$ , then

$$J(U^*, V') \leq J(U, V') \quad \text{for all } U \in \mathcal{U}.$$

The matrix  $S_p(i + 1)$  is the partition of the coefficient matrix for the quadratic form in  $x$  in the cost function for this optimal control problem at time  $i$ . This sequence of matrices is determined by the usual Riccati-like difference equation associated with this control problem. The details of this equation are straightforward and are not presented here.

It is possible to prove similarly a set of lemmas analogous to Lemmas 1, 2 and 3 with the roles of the pursuer and evader interchanged. These lemmas will henceforth be designated respectively as Lemmas 4, 5 and 6. In particular, let

$$(22) \quad U' : u(i) = A_p(i)\bar{x} + \sum_{j=0}^i L_p(i, j)z_p(j), \quad i = 0, \dots, N - 1,$$

$$(23) \quad v(i) = K_e(i)\bar{x} + M_e(i)\hat{\sigma}_p(i),$$

$$(24) \quad \hat{\sigma}_e(i) = N_e(i)\bar{x} + T_e(i)\hat{\sigma}_e(i - 1) + E_e(i)z_e(i),$$

$$\langle \hat{\sigma}_e(-1) \rangle_j \triangleq \begin{bmatrix} \bar{x} \\ 0 \end{bmatrix},$$

$$(25) \quad V^* : v(i) = A_e(i)\bar{x} + \sum_{j=0}^i L_e(i, j)z_e(j), \quad i = 0, \dots, N - 1,$$

and

$$(26) \quad \begin{aligned} A_e(i) &= f_{e,i}(\{A_p(i)\}, \{L_p(i, j)\}), \\ L_e(i, j) &= g_{e,i,j}(\{A_p(i)\}, \{L_p(i, j)\}) \end{aligned}$$

be the respective analogues of (12), (13), (14), (20) and (21).

**7. Stochastic game solution.** By combining the results of Lemmas 2 and 5, it follows that any solution to this stochastic game of the form

$$(27) \quad \begin{aligned} U^* : u(i) &= A_p(i)\bar{x} + \sum_{j=0}^i L_p(i, j)z_p(j) \quad (\text{pursuer}), \\ V^* : v(i) &= A_e(i)\bar{x} + \sum_{j=0}^i L_e(i, j)z_e(j) \quad (\text{evader}), \end{aligned}$$

$i = 0, \dots, N - 1$ , must have coefficients  $\{A_p(i)\}$ ,  $\{A_e(i)\}$ ,  $\{L_p(i, j)\}$  and  $\{L_e(i, j)\}$  which satisfy both (21) and (26). Furthermore, by Lemmas 3 and 6, any strategy pair  $(U^*, V^*)$  of the form of (27) which satisfies these equations and also has the property that

$$(28) \quad B(i) + G_p^T(i)S_p(i + 1)G_p(i)$$

and

$$(29) \quad C(i) + G_e^T(i)S_e(i + 1)G_e(i)$$

are positive definite for  $i = 0, \dots, N - 1$ , where the matrices  $\{S_p(i)\}$  and  $\{S_e(i)\}$  are determined by the equations described earlier, is a game solution.

Equations (21) and (26) together can be regarded as a set of implicit difference equations to be solved for a set of matrix parameters defining a strategy pair  $(U^*, V^*)$  via (27). Not much is known at present about the existence or computation of solutions to this set of implicit equations. Numerical evidence indicates that linear minimax strategies of the form considered here do indeed exist for a variety of two-stage scalar games of this type, and that the implicit equations can be solved by iterative substitution in these cases. Also, it is shown in the next section that no such linear minimax strategy exists if there is no solution to the corresponding deterministic game.

The foregoing approach to finding minimax strategies, of course, is limited in that it eliminates from consideration all possibilities except strategy pairs of the form of (27). However, if a minimax strategy pair  $(U^*, V^*)$  of this form does exist, then it follows from the remark at the end of § 3 that any other minimax strategy pair results in the same value of the criterion  $J$ .

In this particular class of games, moreover, it is also true that if such a linear solution  $(U^*, V^*)$  exists, then any solution  $(U', V')$  is in fact the same strategy pair (except perhaps for a set of measurements of measure zero). This follows from the fact that  $(U', V')$  and  $(U^*, V')$  must also be minimax, which in turn implies that  $U'$  and  $V'$  satisfy the necessary conditions of Lemmas 1 and 4, respectively. Since  $U^*$  and  $V^*$  also satisfy these conditions, and since these conditions are satisfied by only one pair of strategies,  $U' = U^*$  and  $V' = V^*$ .

This uniqueness argument, unfortunately, depends on the equivalence and interchangeability of all solutions, and hence does not generalize to the case of nonzero-sum games. This point will be discussed further in § 9.

**8. Relation to the corresponding deterministic game.** Although there is only one minimax trajectory in the corresponding deterministic game, assuming that solutions exist, there is usually a continuum of minimax strategies for each player. Two such minimax strategies always generate the minimax trajectory when the opponent uses a minimax strategy, but do not in general generate the same trajectory otherwise. It is shown in this section that if a linear solution exists for the stochastic game, minimax strategy pairs also exist for the corresponding deterministic game, *one* of which is related to the stochastic game solution by the certainty-equivalence principle.

It is assumed throughout the rest of this section that a stochastic game solution of the form of (27) has been found, and hence also a solution to (21) and (26).

**DEFINITION.** The *noiseless sample path* of the stochastic game for strategy pair  $(U, V)$  is defined as the trajectory generated by this strategy pair according to (1)–(3) when all the quantities defined as random variables assume their mean values; that is, when  $x(0) = \bar{x}$ ,  $n(i) \equiv 0$ ,  $w_p(i) \equiv 0$ , and  $w_e(i) \equiv 0$ .

Since all of the above equations and (27) are linear in these random variables, the noiseless sample path for the strategy pair  $(U^*, V^*)$  is the mean trajectory generated by this strategy pair, and hence describes the “average behavior” of the optimally controlled stochastic game. This particular sample trajectory

is useful in relating the solutions of the stochastic and corresponding deterministic game.

**THEOREM.** *If a linear solution  $(U^*, V^*)$  of the form of (27) exists for the stochastic game, then:*

(i) *A minimax strategy pair exists for the corresponding deterministic game which is related to  $(U^*, V^*)$  by the certainty-equivalence principle.*

(ii) *The noiseless sample path of the stochastic game for  $(U^*, V^*)$  coincides with the minimax trajectory of the corresponding deterministic game.*

*Proof sketch.* Construct the following two pursuit strategies for the corresponding deterministic game:

$U_c^*$ —in which the pursuer's controls are generated by the canonical form of  $U^*$  (equation (20)) with the substitution of  $H_p(i)x(i)$  for  $z_p(i)$ .

$U_{ce}^*$ —in which the controls are generated by the certainty-equivalent form of  $U^*$  with the substitution of  $\sigma_p(i)$  for  $\hat{\sigma}_p(i)$  in (13). Equation (14) is superfluous in this strategy.  $\sigma_p(i)$  is defined as  $\begin{bmatrix} x(i) \\ \zeta_e(i) \end{bmatrix}$  in this deterministic context, where  $\zeta_e(i)$  is defined by (11) and (3) with  $w_e(i) \equiv 0$ . Note that  $\sigma_p(i)$  is directly observable to the pursuer in this context.

By Lemma 2, the strategy  $U_c^*$  can also be defined by (13) and (14) with the substitution of  $H_p(i)x(i)$  for  $z_p(i)$  in (14).

Similarly define evasion strategies  $V_c^*$  and  $V_{ce}^*$  for the corresponding deterministic game.

The argument is now arranged into four intermediate results.

**Result 1.** The trajectory generated by  $(U_c^*, V_c^*)$  in the corresponding deterministic game is the same as the noiseless sample path of the optimally controlled stochastic game.

This result is verified by showing that (6) and  $(U_c^*, V_c^*)$  define the same trajectory as do (1), (2), (3), (20) and (25) when  $x(0) = \bar{x}$ ,  $n \equiv 0$ ,  $w_p \equiv 0$  and  $w_e \equiv 0$ .

**Result 2.** In the optimally controlled stochastic game,  $\hat{\sigma}_p \equiv \sigma_p$  and  $\hat{\sigma}_e \equiv \sigma_e$  on the noiseless sample path.

If  $e_p(i)$  is defined as  $\sigma_p(i) - \hat{\sigma}_p(i)$ , it obeys a linear difference equation whose inhomogeneous term is zero if  $n \equiv 0$ ,  $w_p \equiv 0$  and  $w_e \equiv 0$ . The initial value of  $e_p$  is proportional to  $x(0) - \bar{x}$ . Similarly this holds for the evader. The result follows by a simple induction argument on  $i$ .

**Result 3.** In the corresponding deterministic game,  $U_{ce}^*$  minimizes the criterion against  $V_c^*$ , and  $V_{ce}^*$  maximizes it against  $U_c^*$ .

This result is established by showing that the minimization of  $J_d(U, V_c^*)$  over  $U$  is equivalent to solving the optimal control problem of Lemma 1 with all random variables replaced by their mean values. Since this control problem is of the classical "linear-quadratic-Gaussian" type, the certainty-equivalence principle applies to it. Therefore, by the definition of  $U_{ce}^*$ ,

$$J_d(U_{ce}^*, V_c^*) \leq J_d(U, V_c^*)$$

for all admissible  $U$ . An analogous argument shows that  $V_{ce}^*$  maximizes against  $U_c^*$ .

**Result 4.** The strategy pair  $(U_{ce}^*, V_{ce}^*)$  is minimax in the corresponding deterministic game.

By Results 1 and 2, the substitution of  $H_p(i)x(i)$  for  $z_p(i)$  in (14) generates  $\hat{\sigma}_p \equiv \sigma_p$  in the corresponding deterministic game when the strategy pair  $(U_c^*, V_c^*)$  is played. It follows from the second definition of  $U_c^*$  that the same trajectory is produced in this game by using the (observable) value of  $\sigma_p$  in (13) instead of generating  $\hat{\sigma}_p$  by this modification of (14). But this is equivalent to replacing the strategy  $U_c^*$  by  $U_{ce}^*$ . Therefore,

$$J_d(U_c^*, V_c^*) = J_d(U_{ce}^*, V_c^*)$$

since the value of  $J_d$  is determined by the trajectory. Similarly,

$$J_d(U_c^*, V_c^*) = J_d(U_c^*, V_{ce}^*).$$

These equalities and Result 3 imply that  $(U_c^*, V_c^*)$  satisfies the saddle-point condition.

Since the strategy pair  $(U_c^*, V_c^*)$  is minimax in the corresponding deterministic game, the trajectory it generates must be the (unique) minimax trajectory defined by (6), (8) and (9). Since  $(U_{ce}^*, V_{ce}^*)$  also generates the noiseless sample path, these two trajectories are identical.

By the argument for Result 4, both  $U_c^*$  and  $V_c^*$  can be replaced by  $U_{ce}^*$  and  $V_{ce}^*$  without changing the trajectory produced in the corresponding deterministic game. Therefore the strategy pair  $(U_{ce}^*, V_{ce}^*)$  is also minimax there. The strategies of this pair are related to  $U^*$  and  $V^*$  by the certainty-equivalence principle in the sense that one can be obtained from the other by replacing  $\sigma_p$  (or  $\sigma_e$ ) by the estimate  $\hat{\sigma}_p$  (or  $\hat{\sigma}_e$ ). This completes the proof.

Although a certainty-equivalence principle holds for this class of games, it does not have the same significance as the one for optimal control problems. For one thing, it does not relate the stochastic game solution to the feedback solution of the corresponding deterministic game, but to a much more complicated solution. It has been verified by numerical counterexample [13] that this other solution is not the same as the feedback solution in general. A knowledge of the feedback solution for the corresponding deterministic game is therefore of little use in computing a solution of the stochastic game.

This certainty-equivalence result is also predicated on the existence of a linear solution to the stochastic game. It can be inferred from the results established here that such a solution cannot exist unless the feedback solution of the corresponding deterministic game exists, which is easy to determine. It is not known, however, whether or not the converse is true. It is also not known whether or not a solution of some nonlinear form ever exists for the stochastic game.

The fact that the noiseless sample path of a linear stochastic game solution coincides with the minimax trajectory of the corresponding deterministic game, referred to here as the *certainty-coincidence property*, is of computational significance. It implies that if the stochastic game has such a solution, its average minimax behavior can be determined by finding the minimax trajectory of the corresponding deterministic game, which is relatively easy. If the noise magnitudes are small, moreover, its minimax behavior can be expected to approximate closely this average trajectory most of the time. Thus, this property indicates that there is no discontinuity of behavior as the noise magnitudes go from small values to zero, although conditions for which this is true have not been established.



It should be emphasized that the certainty-coincidence property is not a statement about the behavior of the optimally controlled stochastic game as the noise variances approach zero, but rather about its behavior when the noise values are zero and the variances are not zero. Intuitively, it states that the presence of the noise does not make an optimal player more or less cautious on the average, provided that a linear solution exists in its presence.

**9. Extensions.** A careful examination of the arguments in the preceding sections reveals that most of the results obtained there extend to the many-player nonzero-sum case as well. Specifically, all of these results except those depending on the equivalence of different solutions hold for the general many-person nonzero-sum game with linear dynamics, quadratic criteria (in general not related) for the players to minimize, linear measurements for each player, and additive independent Gaussian perturbations in the dynamics and measurements. The appropriate solution concept here is the *Nash equilibrium solution*, which is an  $M$ -tuple of admissible strategies  $(U_1^*, \dots, U_M^*)$  for the  $M$  players such that the inequalities

$$J_k(U_1^*, \dots, U_k^*, \dots, U_M^*) \leq J_k(U_1^*, \dots, U_k, \dots, U_M^*), \quad k = 1, \dots, M,$$

hold for all admissible strategy  $M$ -tuples  $(U_1, \dots, U_M)$ , where  $J_k$  is the criterion of the  $k$ th player.

If they exist, Nash solutions linear in the measurements can be found in this more general context by solving a system of  $M$  implicit equation schemes, and verifying the positive definiteness of some matrices. Furthermore, the noiseless sample path generated by any such solution coincides with the unique Nash trajectory in the corresponding closed loop deterministic game (which is guaranteed to have a Nash solution if such a solution exists for the stochastic game). This last point is interesting because Starr and Ho [11] have shown that the Nash trajectories are in general distinct for the various combinations of open and closed loop assumptions about the several players in nonzero-sum games of this sort. The fact that the certainty-coincidence property takes the form that it does in this nonzero-sum context indicates that noisy measurements, no matter how poor, are more closely related to perfect measurements than to no measurements, at least as far as the Nash behavior of such games is concerned. If this phenomenon seems counterintuitive, it should be recalled that, in the general nonzero-sum game, Nash solutions do not have all the properties that one would intuitively expect of a completely satisfactory solution. Nash solutions are not in general minimax or noninferior, as they are in zero-sum games (see Starr and Ho [11]).

As mentioned earlier, the general nonequivalence and noninterchangeability of Nash solutions prevents the proof of uniqueness at the end of §7 from being extended to nonzero-sum games. The lack of such a result is a more serious shortcoming here, however, precisely because of this nonequivalence. Nash strategies lose much of their significance in nonzero-sum games if they are not unique, because the rationale for playing any particular one of them becomes ambiguous.

**Acknowledgment.** The author is indebted to Professor Y. C. Ho of Harvard University for suggesting this area of research and for his help and encouragement during its development.

## REFERENCES

- [1] R. D. BEHN AND Y. C. HO, *On a class of linear stochastic differential games*, IEEE Trans. Automatic Control, AC-13 (1968), pp. 227–239.
- [2] L. D. BERKOVITZ, *A variational approach to differential games*, Advances in Game Theory, Annals of Math. Studies, No. 52, Princeton University Press, Princeton, 1964, pp. 127–174.
- [3] T. F. GUNCKEL AND G. F. FRANKLIN, *A general solution for linear sampled-data control systems*, Trans. ASME Ser. D. J. Basic Engrg., 85 (1963), pp. 197–203.
- [4] Y. C. HO, A. E. BRYSON AND S. BARON, *Differential games and optimal pursuit-evasion strategy*, IEEE Trans. Automatic Control, AC-10 (1965), pp. 385–389.
- [5] P. JOSEPH AND J. TOU, *On linear control theory*, AIEE Trans., Part II, Sept. 1961.
- [6] R. E. KALMAN, *A new approach to linear filtering and prediction problems*, Trans. ASME Ser. D. J. Basic Engrg., 82 (1960), pp. 35–45.
- [7] R. E. KALMAN AND R. BUCY, *New results in linear filtering and prediction*, Ibid., 83 (1961), pp. 95–108.
- [8] R. D. LUCE AND H. RAIFFA, *Games and Decisions*, John Wiley, New York, 1957.
- [9] I. B. RHODES AND D. G. LUENBERGER, *Differential games with imperfect state information*, IEEE Trans. Automatic Control, AC-14 (1969), pp. 29–39.
- [10] ———, *Nondeterministic differential games with constrained state estimators*, Ibid., AC-14 (1969), pp. 476–481.
- [11] A. W. STARR AND Y. C. HO, *Nonzero-sum differential games*, Tech. Rep. 564, Division Engrg. Appl. Phys., Harvard Univ., Cambridge, Mass., 1968.
- [12] W. W. WILLMAN, *Formal solutions for a class of stochastic pursuit-evasion games*, IEEE Trans. Automatic Control, AC-14 (1969), pp. 504–509.
- [13] ———, *On a class of stochastic pursuit-evasion games*, Tech. Rep. 585, Division Engrg. Appl. Phys., Harvard Univ., Cambridge, Mass., 1969.
- [14] W. M. WONHAM, *On the separation theorem of stochastic control*, this Journal, 6 (1968), pp. 312–326.

## CLOSURE, LOWER CLOSURE, AND SEMICONTINUITY THEOREMS IN OPTIMAL CONTROL\*

LAMBERTO CESARI†

**Abstract.** First we discuss in detail the concept of lower closure for Lagrange problems, and we extend in various ways previous closure and lower closure theorems. In particular, we show that the present lower closure theorems and related concepts are extensions to Lagrange problems of optimal control of well-known semicontinuity theorems for free problems of the calculus of variations and the related concept of seminormality. Finally, we prove by a new approach that the convexity condition usually requested in lower closure theorems is, in a suitable sense, both a necessary and sufficient condition for lower closure.

In previous papers [1, abcdefgh] we have given existence theorems for one- and multidimensional Lagrange and Mayer problems of optimal control based on closure and lower closure properties of the functional under consideration. In the present paper we discuss in detail some of the relevant concepts, in particular the concept of lower closure for Lagrange problems. Also, we extend in various ways previous closure and lower closure theorems. In particular, we show that the present lower closure theorems and related concepts are extensions to Lagrange problems of optimal control of well-known lower semicontinuity theorems for free problems and the concept of seminormality of Tonelli [7] and McShane [3]. Finally, we prove by a new approach that the convexity condition usually requested in lower closure theorems is—in a suitable sense—both a necessary and sufficient condition for lower closure.

We limit ourselves in this paper to one-dimensional problems. We deal, therefore, with the (Lagrange) problem of the minimum of an integral

$$I[x, u] = \int_{t_1}^{t_2} f_0(t, x(t), u(t)) dt,$$

in classes of pairs of functions  $x(t) = (x^1, \dots, x^n)$ ,  $u(t) = (u^1, \dots, u^m)$ ,  $t_1 \leq t \leq t_2$ , satisfying a system of ordinary differential equations of the form

$$dx/dt = f(t, x(t), u(t)) \quad \text{a.e. in } [t_1, t_2],$$

and constraints of the form

$$(t, x(t)) \in A, \quad u(t) \in U(t, x(t)) \quad \text{a.e. in } [t_1, t_2].$$

**1. A closure theorem.** We shall denote by  $t$  a scalar variable (time), and by  $x = (x^1, \dots, x^n)$  the state variables. Let  $A$  be a given subset of the  $tx$ -space  $E_{n+1}$ , and for any  $(t, x) \in A$  let  $Q(t, x)$  be a given subset of the  $z$ -space  $E_n$ ,  $z = (z^1, \dots, z^n)$ . A vector function  $\chi(t) = (x^1, \dots, x^n)$ ,  $t_1 \leq t \leq t_2$ , is said to be a solution of the orientor field

$$(1.1) \quad dx/dt \in Q(t, x),$$

---

\* Received by the editors February 13, 1970, and in revised form August 10, 1970.

† Department of Mathematics, University of Michigan, Ann Arbor, Michigan 48104. This research was supported in part by the United States Air Force Office of Scientific Research under Grant 69-1662.

provided :

- (a)  $x$  is absolutely continuous in  $[t_1, t_2]$ ;
- (b)  $(t, x(t)) \in A$  for all  $t \in [t_1, t_2]$ ;
- (c)  $dx/dt \in Q(t, x(t))$  a.e. in  $[t_1, t_2]$ .

For any  $(\bar{t}, \bar{x}) \in A$  and  $\delta > 0$  we shall denote by  $N_\delta(\bar{t}, \bar{x})$ , or  $\delta$ -neighborhood of  $(\bar{t}, \bar{x})$  in  $A$ , the set of all  $(t, x) \in A$  at a distance  $\leq \delta$  from  $(\bar{t}, \bar{x})$ . We say that the sets  $Q(t, x)$  satisfy property (Q) at the point  $(\bar{t}, \bar{x}) \in A$  provided

$$Q(\bar{t}, \bar{x}) = \bigcap_\delta cl\ co \cup Q(t, x),$$

where  $\cup$  is taken for all  $(t, x) \in N_\delta(\bar{t}, \bar{x})$ . We say that the sets  $Q(t, x)$  satisfy property (Q) in  $A$  if they satisfy this property at every  $(\bar{t}, \bar{x}) \in A$ . The weaker property (U) can be defined analogously by replacing the operation  $cl\ co$  above by the operation  $cl$ .

If  $x(t), a \leq t \leq b, y(t), c \leq t \leq d$ , are any two continuous vector functions with values in  $E_n$ , we denote as their  $\rho$ -distance the number  $\rho(x, y) = |a - c| + |b - d| + \max |x(t) - y(t)|$ , where the maximum is taken in  $(-\infty, +\infty)$ , and  $x(t), y(t)$  are thought of as extended in  $(-\infty, +\infty)$  by continuity and constancy outside their intervals of definition.

Note that if  $x(t), t_1 \leq t \leq t_2$ , is any continuous vector function, and  $x_k(t), t_{1k} \leq t \leq t_{2k}, k = 1, 2, \dots$ , is any sequence of continuous vector functions with  $\rho(x_k, x) \rightarrow 0$  as  $k \rightarrow \infty$ , then  $(t_{1k}, x_k(t_{1k})) \rightarrow (t_1, x(t_1))$  and  $(t_{2k}, x_k(t_{2k})) \rightarrow (t_2, x(t_2))$  as  $k \rightarrow \infty$ . Also, if  $A$  is a closed subset of the  $tx$ -space  $E_{n+1}$ , and  $(t, x_k(t)) \in A$  for all  $t \in [t_{1k}, t_{2k}]$  and all  $k$ , then we also have  $(t, x(t)) \in A$  for all  $t \in [t_1, t_2]$ . Finally, if  $B$  is a closed subset of  $E_{2n+2}$  and  $(t_{1k}, x_k(t_{1k}), t_{2k}, x_k(t_{2k})) \in B$  for all  $k$ , then  $(t_1, x(t_1), t_2, x(t_2)) \in B$ .

In [1a] we proved the following closure theorem.

**THEOREM 1.1.** *If  $A$  is closed, if the sets  $Q(t, x)$  are closed, convex, and satisfy property (Q) at every point of  $A$  (with exception perhaps of a set of points whose  $t$ -coordinate lies in a set of measure zero on the  $t$ -axis), if  $x_k(t), t_{1k} \leq t \leq t_{2k}, k = 1, 2, \dots$ , is a sequence of solutions of the orientor field (1.1), if  $x(t), t_1 \leq t \leq t_2$ , is an absolutely continuous vector function, and  $\rho(x_k, x) \rightarrow 0$  as  $k \rightarrow \infty$ , then  $x(t), t_1 \leq t \leq t_2$ , is a solution of the orientor field (1.1).*

**2. A closure theorem involving singular components.** Let  $\alpha$  be any integer,  $0 \leq \alpha \leq n$ , and if  $x = (x^1, \dots, x^n)$ , let  $y$  denote the  $\alpha$ -vector  $y = (x^1, \dots, x^\alpha)$  and  $z$  the  $(n - \alpha)$ -vector  $z = (x^{\alpha+1}, \dots, x^n)$ , so that we can write  $x = (y, z)$ . Let  $A_0$  be a subset of the  $ty$ -space  $E_{\alpha+1}$ , let  $A = A_0 \times E_{n-\alpha}$ , and for any  $(t, y) \in A_0$  let  $Q(t, y)$  be a given subset of the  $x$ -space  $E_n$ . We shall consider the orientor field

$$(2.1) \quad dx/dt \in Q(t, y).$$

We proved in [1a] the following closure theorem.

**THEOREM 2.1.** *If  $A_0$  is closed, if the sets  $Q(t, y)$  are closed, convex, and satisfy property (Q) at every point of  $A_0$  (with exception perhaps of a set of points whose  $t$ -coordinate lies in a set of measure zero on the  $t$ -axis), if  $x_k(t) = (y_k, z_k), t_{1k} \leq t \leq t_{2k}, k = 1, 2, \dots$ , is a sequence of solutions of the orientor field (2.1) and  $x(t) = (y, z)$ ,*

$t_1 \leqq t \leqq t_2$ , is any vector function with  $y(t)$  absolutely continuous in  $[t_1, t_2]$ ,  $z(t) = Z(t) + S(t)$ ,  $t_1 \leqq t \leqq t_2$ ,  $Z(t)$  absolutely continuous in  $[t_1, t_2]$ , and  $S(t)$  singular, if  $\rho(y_k, y) \rightarrow 0$  as  $k \rightarrow \infty$ , and  $z_k(t) \rightarrow z(t)$  as  $k \rightarrow \infty$  pointwise in  $(t_1, t_2)$ ; then  $X(t) = (y(t), Z(t))$ ,  $t_1 \leqq t \leqq t_2$ , is a solution of the orientor field (2.1).

**3. Lower closure of functionals in integral form.** As usual we denote by  $t$  the independent variable, by  $x = (x^1, \dots, x^n)$  the state variables, and by  $u = (u^1, \dots, u^m)$  the control variables.

As usual, let  $A$  be a closed subset of the  $tx$ -space  $E_1 \times E_n$ ; for every  $(t, x) \in A$  let  $U(t, x)$  be a given subset of the  $u$ -space  $E_m$ ; let  $M$  be the set of all  $(t, x, u)$  with  $(t, x) \in A$ ,  $u \in U(t, x)$ ; and let  $\tilde{f}(t, x, u) = (f_0, f_1, \dots, f_n) = (f_0, f)$  be a given continuous vector function on  $M$ . Let  $B$  be a closed subset of the  $t_1x_1t_2x_2$ -space  $E_{2n+2}$ . We consider now the functional

$$(3.1) \quad I[x, u] = \int_{t_1}^{t_2} f_0(t, x(t), u(t)) dt.$$

We shall say that a pair  $x(t), u(t)$ ,  $t_1 \leqq t \leqq t_2$ , is *admissible* provided  $x(t)$  is absolutely continuous in  $[t_1, t_2]$ ,  $u(t)$  is measurable in  $[t_1, t_2]$ ,  $(t, x(t)) \in A$  for  $t \in [t_1, t_2]$ ,  $u(t) \in U(t, x(t))$  a.e. in  $[t_1, t_2]$ ,  $dx/dt = f(t, x(t), u(t))$  a.e. in  $[t_1, t_2]$ ,  $f_0(t, x(t), u(t))$  is  $L$ -integrable in  $[t_1, t_2]$ , and  $(t_1, x(t_1), t_2, x(t_2)) \in B$ . Whenever we wish to disregard boundary conditions, we have only to take  $B = E_{2n+2}$ . We shall say that a vector function  $x(t)$ ,  $t_1 \leqq t \leqq t_2$ , is an *admissible trajectory* if there is at least one vector function  $u(t)$ ,  $t_1 \leqq t \leqq t_2$ , such that the pair  $x, u$  is admissible. Analogously, we shall say that a vector function  $u(t)$ ,  $t_1 \leqq t \leqq t_2$ , is an *admissible strategy* if there is at least one vector function  $x(t)$ ,  $t_1 \leqq t \leqq t_2$ , such that the pair  $x, u$  is admissible.

Let  $x(t)$ ,  $t_1 \leqq t \leqq t_2$ , be any absolutely continuous vector function (which is the limit in the metric  $\rho$  of admissible trajectories). If, for any sequence  $x_k(t)$ ,  $u_k(t)$ ,  $t_{1k} \leqq t \leqq t_{2k}$ ,  $k = 1, 2, \dots$ , of admissible pairs with  $\rho(x_k, x) \rightarrow 0$ ,  $\liminf I[x_k, u_k] < +\infty$  as  $k \rightarrow \infty$ , there is some measurable function  $u(t)$ ,  $t_1 \leqq t \leqq t_2$ , such that  $x(t), u(t)$ ,  $t_1 \leqq t \leqq t_2$ , is admissible, and

$$(3.2) \quad I[x, u] \leqq \liminf_{k \rightarrow \infty} I[x_k, u_k],$$

then we say that  $I[x, u]$  has the property of *lower closure* at the trajectory  $x(t)$ ,  $t_1 \leqq t \leqq t_2$ .

Before we prove a sufficient condition for lower closure, the following remarks are needed. First, if  $x$  is the limit in the  $\rho$ -metric of admissible trajectories as assumed, then by the remark in § 1 we know that  $(t, x(t)) \in A$  for all  $t \in [t_1, t_2]$ , and  $(t_1, x(t_1), t_2, x(t_2)) \in B$ .

Furthermore, if we know that the set  $M$  is closed, and that for every  $(t, x) \in A$  the sets  $Q(t, x) = f(t, x, U(t, x))$  are closed convex subsets of  $E_n$  satisfying property (Q) in  $A$ , then certainly  $x'(t) \in Q(t, x(t))$  a.e. in  $[t_1, t_2]$  by force of the closure Theorem 1.1; and then there is some measurable  $u(t)$ ,  $t_1 \leqq t \leqq t_2$ , such that

$$(3.3) \quad u(t) \in U(t, x(t)), \quad x'(t) = f(t, x(t), u(t)) \quad \text{a.e. in } [t_1, t_2],$$

by force of the implicit function theorem (see, e.g., [4]). As usual, we say that any such strategy  $u(t)$  generates  $x(t)$ ,  $t_1 \leqq t \leqq t_2$ . Obviously, in the concept of lower semicontinuity we require more; namely, we need a strategy  $u$  generating  $x$  for which (3.2) holds.

It may well occur that  $x$  is generated by some strategy  $\bar{u}$  for which (3.2) does not hold. The following example displays two strategies  $u$  and  $\bar{u}$ , both generating the same trajectory  $x$ , such that (3.2) holds for  $u$  but not for  $\bar{u}$ .

Indeed, take  $m = n = 1$ ,  $t_1 = 0$ ,  $t_2 = 1$ ,  $f_0 = 1 + \cos \pi u$ ,  $f = f_1 = \sin \pi u$ ,  $u \in U = [-1 \leqq u \leqq 1]$ ,  $x(t) = 0$ ,  $0 \leqq t \leqq 1$ ,  $A = E_2$ . Now take  $u_k(t) = \pm 2^{-1}$  according as  $k^{-1}i \leqq t < k^{-1}i + (2k)^{-1}$  or  $k^{-1} + (2k)^{-1} \leqq t < (i + 1)k^{-1}$ ,  $i = 0, 1, \dots, k - 1, k = 1, 2, \dots$ ; and take  $x_k(t) = t - k^{-1}i$  or  $x_k(t) = k^{-1}(i + 1) - t$ , according as  $t$  is in one or the other set of intervals above. Then  $x_k, u_k, k = 1, 2, \dots$ , is a sequence of admissible pairs,  $0 \leqq x_k(t) \leqq (2k)^{-1}$ , and  $x_k \rightarrow x$  as  $k \rightarrow \infty$  uniformly in  $[0, 1]$ . The trajectory  $x(t) = 0, 0 \leqq t \leqq 1$ , is now generated by both  $u(t) = 1, 0 \leqq t \leqq 1$ , and by  $\bar{u}(t) = 0, 0 \leqq t \leqq 1$ . On the other hand,

$$I[x, u] = 0, \quad I[x, \bar{u}] = 2, \quad I[x_k, u_k] = 1, \quad k = 1, 2, \dots,$$

and thus relation (3.2) holds for  $u$  but not for  $\bar{u}$ .

As we shall see in § 7, the concept of lower closure introduced above contains as a particular case the usual concept of lower semicontinuity, in particular, the concept of lower semicontinuity for free problems.

**4. Sufficient conditions for lower closure.** Let  $A, U(t, x), M, B, f_0(t, x, u)$  and  $f(t, x, u) = (f_1, \dots, f_n)$  be defined as in § 3. For any  $(t, x) \in A$  let  $\tilde{Q}(t, x)$  be the set of all  $\tilde{z} = (z^0, z^1, \dots, z^n) = (z^0, z)$  with  $z^0 \geqq f_0(t, x, u)$ ,  $z = f(t, x, u)$  for some  $u \in U(t, x)$ .

**THEOREM 4.1.** *If the sets  $A, M, B$  are closed, and  $f_0(t, x, u), f(t, x, u) = (f_1, \dots, f_n)$  are continuous on  $M$ , let us assume that the sets  $\tilde{Q}(t, x)$  are closed, convex, and satisfy property (Q) at every point  $(t, x) \in A$  with the exception perhaps of a set of points whose  $t$ -coordinate lies on a set of measure zero on the  $t$ -axis. Let us assume that, for some locally  $L$ -integrable scalar function  $\psi(t) \geqq 0$  we have:*

$$(\psi) \quad f_0(t, x, u) \geqq -\psi(t) \quad \text{for all } (t, x, u) \in M,$$

*with the exception perhaps of another set of points whose  $t$ -coordinate lies on a set of measure zero on the  $t$ -axis. Then the integral (3.1) has the property of lower closure at every absolutely continuous vector function  $x(t) = (x^1, \dots, x^n), t_1 \leqq t \leqq t_2$ , which is the limit in the  $\rho$ -metric of admissible trajectories. In other words, for every absolutely continuous vector  $x(t) = (x^1, \dots, x^n), t_1 \leqq t \leqq t_2$ , and sequence  $x_k(t), u_k(t), t_{1k} \leqq t \leqq t_{2k}, k = 1, 2, \dots$ , of admissible pairs with  $\rho(x_k, x) \rightarrow 0$ ,  $\liminf I[x_k, u_k] < +\infty$  as  $k \rightarrow \infty$ , there is a measurable function  $u(t), t_1 \leqq t \leqq t_2$ , such that  $x(t), u(t), t_1 \leqq t \leqq t_2$ , is admissible and  $I[x, u] \leqq \liminf I[x_k, u_k]$ .*

We proved Theorem 4.1 in [1a]. If  $i$  denotes the number  $i = \liminf I[x_k, u_k]$  as  $k \rightarrow \infty$ , then it is not restrictive in this theorem to consider only a subsequence, say still  $[k]$ , with  $I[x_k, u_k] \rightarrow i$  as  $k \rightarrow \infty$ . In the proof in [1a] of the lower closure Theorem 4.1 the functions

$$x_k^0(t) = \int_{t_{1k}}^t f_0(\tau, x_k(\tau), u_k(\tau)) d\tau, \quad t_{1k} \leqq t \leqq t_{2k}, \quad k = 1, 2, \dots,$$

are taken into consideration, Helly's selection theorem is applied, and the closure Theorem 2.1 is used. This process to prove lower closure (or lower semicontinuity) (Cesari [1a]) was later used by E. J. McShane [3f], T. Nishiura [5], C. Olech [6].

By a well-known remark by C. S. Goodman, the hypothesis in Theorem 4.1 that the functions  $f_0, f$  are continuous on  $M$  can be replaced by the weaker assumption that  $f_0, f$  are continuous in  $x, u$  for every  $t$  and are measurable in  $t$  for every  $x, u$ . The proofs are essentially the same. Also, the following further remark concerning Theorem 4.1 may be of interest. To formulate it we shall denote by  $[\bar{t}, T]$  a fixed interval containing all intervals  $[t_{1k}, t_{2k}], [t_1, t_2]$ , and we shall extend  $x_k(t), x(t)$  on the whole interval  $[\bar{t}, T]$  by continuity and constancy in each interval  $[\bar{t}, t_{1k}], [t_{2k}, T], [\bar{t}, t_1], [t_2, T]$ . If in Theorem 4.1 we make the further assumption that the derivatives  $x'_k(t)$  converge weakly in  $L_1[\bar{t}, T]$  toward  $x'(t)$  as  $k \rightarrow \infty$ , then the assumption concerning the sets  $\tilde{Q}(t, x)$  satisfying property (Q) in Theorem 4.1 can be replaced by the following weaker assumption: There is a countable decomposition of  $[\bar{t}, T]$  into disjoint measurable sets  $H_\lambda, \lambda = 1, 2, \dots$ , such that, if  $A_\lambda$  denotes the set  $A_\lambda = \{(t, x) | (t, x) \in A, t \in H_\lambda\}$ , then the sets  $\tilde{Q}(t, x)$  satisfy property (Q) in  $A_\lambda$  for almost every  $t, \lambda = 1, 2, \dots$ . The assumption concerning the weak convergence of  $x'_k$  to  $x'$  in  $L_1$  is usually satisfied in applications, but it is not requested in Theorem 4.1. A proof of the modification just mentioned is essentially contained in Cesari [1e, pp. 94–101, particularly p. 98]. The assumption of weak convergence of  $x'_k$  to  $x'$  in  $L_1$  is essential. A counterexample, essentially due to A. Lasota and C. Olech [2], is reported in the Appendix at the end of this paper. On the other hand, under the hypothesis of weak convergence of the derivatives, the contention in Theorem 4.1 can be proved also under a different set of hypotheses (see C. Olech [6b]).

As we shall see in detail in § 7, Theorem 4.1 contains a number of lower semicontinuity theorems as corollaries.

Let us note here that whenever  $f = u, n = m$  (e.g., for free problems); hence  $dx/dt = u$ , the functional can be written simply as  $I[x]$ , and the lower closure theorem reduces to a lower semicontinuity theorem.

**COROLLARY 1.** *Under the hypotheses of Theorem 4.1 with  $n = m, f = u$ ; hence  $dx/dt = u$  and*

$$\tilde{Q}(t, x) = \{(z^0, u) | z^0 \geq f_0(t, x, u), u \in U(t, x)\} \subset E_{n+1};$$

*if  $x_k \rightarrow x$  in the  $\rho$ -metric, and  $\liminf I[x_k] < \infty$ , then  $I[x] \leq \liminf I[x_k]$ .*

Note that in Theorem 4.1 (as well as in the closure theorems, Theorems 1.1 and 2.1) no topology has been chosen for the measurable functions  $u(t), t_1 \leq t \leq t_2$ , under consideration, and hence the question of what happens if the functions  $u_k$  converge toward some function  $u_0$  does not arise. Here we may well assume that all intervals  $[t_{1k}, t_{2k}], [t_1, t_2]$  are all contained in a fixed interval  $[\bar{t}, T], \bar{t}, T$  finite, and we may extend all functions  $u_k$  to the fixed interval  $[\bar{t}, T]$  by defining them to be constant and equal to a fixed  $\bar{u}$  (say,  $\bar{u} = 0$ ) outside  $[t_{1k}, t_{2k}]$ . If we assume that all functions  $u_k$  are in  $L_p$  (that is, each component  $u_k^i$  is in  $L_p[\bar{t}, T]$ ) and that the functions  $u_k$  converge weakly in  $L_p$  toward some element  $u_0$  of  $L_p$ , then we may ask whether the pair  $x, u_0$  is admissible, and whether the relation  $I[x, u_0]$

$\leq \liminf I[x_k, u_k]$  holds. The following corollaries of Theorem 4.1 give affirmative answers to the question being investigated under suitable assumptions.

To state Corollary 2 we shall need the sets

$$\begin{aligned} \tilde{Q}'(t, x) &= \{(z^0, z, \zeta) | z^0 \geq f_0(t, x, u), z = f(t, x, u), \zeta = u, u \in U(t, x)\} \\ &= \{(z^0, z, u) | z^0 \geq f_0(t, x, u), z = f(t, x, u), u \in U(t, x)\} \subset E_{n+m+1}. \end{aligned}$$

**COROLLARY 2.** *Under the conditions of Theorem 4.1 with the sets  $\tilde{Q}$  replaced by the sets  $\tilde{Q}'(t, x)$ , if the functions  $u_k$  are in  $(L_p)^m$ ,  $p > 1$ , and  $u_k \rightarrow u_0$  weakly in  $(L_p)^m$  as  $k \rightarrow \infty$ , then the pair  $x, u_0$  is admissible, and  $I[x, u_0] \leq \liminf I[x_k, u_k]$ . The same statement holds for  $p = 1$  provided we know in addition that the functions  $u_k$  are equiabsolutely integrable.*

*Proof of Corollary 2.* Let us consider the auxiliary control problem with  $n$  and  $m$  replaced by  $n + m$  and  $m$ , the same control variable  $u = (u^1, \dots, u^m)$ , the state variable  $x = (x^1, \dots, x^n)$ ,  $y = (y^1, \dots, y^m)$ , same functional and constraints, and the  $n + m$  differential equations

$$dx/dt = f(t, x, u), \quad dy/dt = u.$$

Here  $A$  is replaced by the closed set  $A' = A \times E_m$ , and the sets  $\tilde{Q}(t, x)$  by the sets  $\tilde{Q}'(t, x)$  which are assumed to be closed, convex, and to satisfy property (Q). Now the functions

$$y_k(t) = \int_{t_{1k}}^t u_k(\tau) d\tau, \quad t_{1k} \leq t \leq t_{2k}, \quad k = 1, 2, \dots,$$

are equiabsolutely continuous and equibounded. We first extract a subsequence  $[k_s]$  such that  $I[x_{k_s}, u_{k_s}] \rightarrow i = \liminf I[x_k, u_k]$ , and we further refine it in such a way that also  $y_{k_s} \rightarrow y$  in the  $\rho$ -metric, where  $y$  is now absolutely continuous in  $[t_1, t_2]$ . By Theorem 4.1, there is a measurable  $u(t)$ ,  $t_1 \leq t \leq t_2$ , such that  $x, u$  is admissible for the new problem and  $I[x, u] \leq i$ ; besides,  $dx/dt = f(t, x(t), u(t))$ ,  $dy/dt = u(t)$  a.e. in  $[t_1, t_2]$ . On the other hand,  $u_{k_s} \rightarrow u_0$  weakly in  $L_p$ ,

$$\int_i^t u_{k_s}(\tau) d\tau \rightarrow \int_i^t u_0(\tau) d\tau \quad \text{as } s \rightarrow \infty$$

for every  $t$ ; hence,  $y_{k_s} \rightarrow y(t)$  for every  $t \in (t_1, t_2)$ , or  $y(t) = \int_i^t u_0(\tau) d\tau$ , and  $u(t) = dy/dt = u_0(t)$  a.e. in  $[t_1, t_2]$ . This proves the corollary.

Note that in Corollary 2 the sequence  $\|u_k\|_p$  is certainly bounded ( $p \geq 1$ ), and for  $p > 1$ , the functions  $u_k$  are certainly equiabsolutely integrable. This is the case even for  $p = 1$  under suitable ‘‘growth’’ conditions. A suitable growth condition is the following one:  $(\varepsilon_0)$  For each  $\varepsilon > 0$  there is some integrable function  $\psi_\varepsilon(t) \geq 0$ ,  $i \leq t \leq T$ , such that

$$|u| \leq \psi_\varepsilon(t) + \varepsilon f_0(t, x, u)$$

for all  $(t, x, u) \in M$ . An analogous growth condition has been considered in [1b] for a different purpose.

It is easy to see that in Corollary 2 the weak convergence of  $u_k$  in  $L_p$  for  $p \geq 1$  and the additional assumption that the  $u_k$  are equiabsolutely integrable if  $p = 1$



can all be replaced by the single simple assumption :  $y_k \rightarrow y$  in the  $\rho$ -metric, where

$$y_k(t) = \int_{t_{1k}}^t u_k(\tau) d\tau, \quad t_{1k} \leqq t \leqq t_{2k}, \quad k = 1, 2, \dots,$$

and

$$y(t) = \int_{t_1}^t u_0(\tau) d\tau, \quad t_1 \leqq t \leqq t_2.$$

The proof is the same with obvious simplifications, and is still based on Theorem 4.1. For instance, for  $m = 1$ , the sequence  $u_k(t) = k^{1/2} \sin kt, 0 \leqq t \leqq 2\pi, k = 1, 2, \dots$ , converges in this sense to  $u_0 = 0$ , but the sequence  $u_k$  is not weakly convergent, say in  $L_1$ , since  $\|u_k\|_1 = 4k^{1/2}$  is not bounded.

In Corollary 2 the hypothesis that the sets  $\tilde{Q}(t, x)$  be convex cannot be disregarded, as the following example shows. Take  $n = 1, m = 1, f_0 = u, f = u^2, U = \{\pm 1\}$ ; hence  $\tilde{Q}$  is the fixed closed nonconvex set made up of the two half-straight-lines  $(z = u = 1, z^0 \geqq 1)$  and  $(z = 1, u = -1, z^0 \geqq -1)$ . Also, take  $t_{1k} = t_1 = 0, t_{2k} = t_2 = 1, u_k(t) = 1$  for  $i/k \leqq t \leqq i/k + 1/2k, i = 0, 1, \dots, k - 1, u_k(t) = -1$  in the complementary intervals,  $x_k(t) = x(t) = t, 0 \leqq t \leqq 1, k = 1, 2, \dots$ . Then  $u_k \rightarrow u_0$  weakly in any  $L_p (p \geqq 1)$  with  $u_0(t) = 0, 0 \leqq t \leqq 1$ , and now the pair  $x, u_0$  is not admissible since  $0 \notin U$ . On the other hand,

$$\tilde{Q} = \{(z^0, z) | z^0 \geqq -1, z = 1\}$$

is convex and all conditions of Theorem 4.1 are satisfied. If we take  $u(t) = -1, 0 \leqq t \leqq 1$ , then the pair  $x, u$  is admissible with  $I[x, u] = -1$ , while  $I[x_k, u_k] = 0, k = 1, 2, \dots$ .

An analogous example where both  $\tilde{Q}$  and  $U$  are convex, where  $x, u_0$  is admissible, and yet  $I[x, u_0] > \liminf I[x_k, u_k]$ , is as follows. Take  $m = n = 1, U = E_1, f_0(-u) = f_0(u), f(-u) = -f(u), f_0 = 2 - u$  for  $0 \leqq u \leqq 2, f_0 = 0$  for  $u \geqq 2, f(u) = u$  for  $0 \leqq u \leqq 1, f(u) = 2 - u$  for  $1 \leqq u \leqq 2, f(u) = 0$  for  $u \geqq 2$ . Then  $\tilde{Q}$  is the fixed closed convex set

$$\{(z^0, z) | z^0 \geqq |z| \text{ for } -1 \leqq z \leqq 1\} \subset E_2,$$

while  $\tilde{Q}'$  is a fixed closed nonconvex set in  $E_3$ . If we take  $t_{1k} = t_1 = 0, t_{2k} = t_2 = 1, k = 1, 2, \dots$ , and  $u_k(t) = 1$  for  $i/k \leqq t \leqq i/k + 1/2k, i = 0, 1, \dots, k - 1, u_k(t) = -1$  in the complementary intervals, then  $u_k \rightarrow u_0$  weakly in any  $L_p, p \geqq 1$ , with  $u_0(t) = 0, 0 \leqq t \leqq 1$ . If we take  $x_k(t) = t - i/k$ , or  $x_k(t) = (i + 1)/k - t$  according as  $t$  belongs to one or the other set of intervals above, then  $x'_k = \bar{u}_k, x_k \rightarrow x$  uniformly and  $x'_k \rightarrow x'$  weakly in any  $L_p$  as  $k \rightarrow \infty$  with  $x(t) = 0, 0 \leqq t \leqq 1$ . The pair  $x, u_0$  is admissible with  $I[x, u_0] = 2$ , while  $I[x_k, u_k] = 1$  for all  $k$ . On the other hand, the conditions of Theorem 4.1 are all satisfied, and there is therefore some measurable  $u$  such that  $x, u$  is admissible and  $I[x, u] \leqq \liminf I[x_k, u_k]$ . One such  $u$  is  $u(t) = 2, 0 \leqq t \leqq 1$ . Another example of the same type, where  $\tilde{Q}$  is a fixed closed convex set,  $f_0$  is convex in  $u, U$  is convex, but  $\tilde{Q}'$  is fixed and closed but not convex, is as follows. Take  $m = n = 1, f_0 = 0, f = u^2, U = [-1 \leqq u \leqq 1]$ , so that  $\tilde{Q} = \{(z^0, z) | z^0 \geqq 0, 0 \leqq z \leqq 1\}$ . Take  $u_k(t) = \pm 1$  as in the previous example, so that  $u_k \rightarrow u_0$  weakly in any  $L_p$  with  $u_0(t) = 0, 0 \leqq t \leqq 1$ . On the

other hand,  $x_k(t) = x(t) = t, 0 \leq t \leq 1, I[x_k, u_k] = 0, k = 1, 2, \dots$ , while  $x, u_0$  is not admissible. For  $u(t) = 1, 0 \leq t \leq 1, x, u$  is admissible and  $I[x, u] = 0$ .

Other examples of the same kind can be found in [1i, particularly p. 13] and [1j, particularly p. 53].

Note that the requirement in Corollary 2 that the sets  $\tilde{Q}(t, x)$  be convex implies that, for every  $(t, x) \in A$ , the sets  $U(t, x)$  are convex,  $f(t, x, u)$  is linear in  $u$  (hence of the form  $f = B(t, x)u + C(t, x)$ ,  $B$  an  $n \times m$  matrix,  $C$  an  $n \times 1$  matrix), and  $f_0(t, x, u)$  is convex in  $u$ .

The following Corollary 3 for general Lagrange problems (with  $f$  linear in  $u$ ) extends both Corollaries 1 and 2. To state Corollary 3 we shall need the sets

$$\begin{aligned} \tilde{Q}^*(t, x) &= \{(z^0, \zeta) | z^0 \geq f_0(t, x, u), \zeta = u, u \in U(t, x)\} \\ &= \{(z^0, u) | z^0 \geq f_0(t, x, u), u \in U(t, x)\} \subset E_{m+1}. \end{aligned}$$

**COROLLARY 3.** *Under the conditions of Theorem 4.1 with the sets  $\tilde{Q}^*(t, x) \subset E_{m+1}$  replacing the sets  $\tilde{Q}$ , if  $f$  is linear in  $u$  (or  $f = B(t, x)u + C(t, x)$ ,  $(t, x) \in A, u \in U(t, x)$ , the matrices  $B$  and  $C$  with entries continuous in  $A$ ), if the functions  $u_k$  are in  $(L_p)^m, p > 1$ , and  $u_k \rightarrow u_0$  weakly in  $(L_p)^m$  as  $k \rightarrow \infty$ , then the pair  $x, u_0$  is admissible, and  $I[x, u_0] \leq \liminf I[x_k, u_k]$ . The same statement holds for  $p = 1$  provided we know in addition that the functions  $u_k$  are equiabsolutely integrable.*

*Proof of Corollary 3.* Let us consider the auxiliary control problem with  $n + m$  control variables  $u = (u^1, \dots, u^m), v = (v^1, \dots, v^n)$ , with  $n + m$  state variables  $x = (x^1, \dots, x^n), y = (y^1, \dots, y^m)$ , the same functional  $I = \int_{t_1}^{t_2} f_0(t, x, u) dt$ , with constraints  $(t, x) \in A, u \in U(t, x), v \in V = E_n$ , and  $n + m$  differential equations

$$dx/dt = v, \quad dy/dt = u.$$

Here  $A$  is replaced by the closed set  $A' = A \times E_m$ , and the sets  $Q(t, x)$  of Theorem 4.1 are replaced by the sets  $Q^*(t, x) \times E_n \subset E_{n+m+1}$  which certainly are closed convex and satisfy property (Q). We take

$$\begin{aligned} v_k(t) &= x'_k(t), \\ y_k(t) &= \int_{t_{1k}}^t u_k(\tau) d\tau, \quad t_{1k} \leq t \leq t_{2k}, \quad k = 1, 2, \dots, \end{aligned}$$

and then the functions  $y_k(t)$  are equiabsolutely continuous and equibounded. We first extract a subsequence  $[k_s]$  such that  $I_{k_s} \rightarrow i = \liminf I_k$ , and we further refine it in such a way that also  $y_{k_s} \rightarrow y$  in the  $\rho$ -metric, where  $y$  is now absolutely continuous in  $[t_1, t_2]$ . By Theorem 4.1 there are measurable  $u(t), v(t), t_1 \leq t \leq t_2$ , such that  $x, y, u, v$ , is admissible for the new problem, and  $I \leq i$ , besides

$$x'(t) = v(t), \quad y'(t) = u(t) \in U(t, x(t)) \text{ a.e.},$$

$$I = \int_{t_1}^{t_2} f_0(t, x(t), u(t)) dt.$$

On the other hand, for every  $k$ , we have  $(t, x_k(t)) \in A$ , and

$$x'_k(t) = v_k(t), \quad x'_k(t) = B(t, x_k(t))u_k(t) + C(t, x_k(t)), \quad y'_k(t) = u_k(t)$$

a.e. in  $[t_{1k}, t_{2k}]$ ; and hence in a fixed interval  $[\bar{t}, T]$  after suitable extensions,

$[t_{1k}, t_{2k}] \subset [\bar{t}, T], k = 1, 2, \dots$ . Then  $B(t, x_k(t)) \rightarrow B(t, x(t)), C(t, x_k(t)) \rightarrow C(t, x(t))$  as  $k \rightarrow \infty$  in the  $\rho$ -metric, as well as uniformly in  $[\bar{t}, T]$ , while  $u_k \rightarrow u_0$  weakly in  $L_p$ . Then  $x'_k \rightarrow x'$  weakly in  $L_p$  and hence

$$x'(t) = B(t, x(t))u_0(t) + C(t, x(t)) \quad \text{a.e. in } [t_1, t_2].$$

If we denote by  $W$ -lim the usual weak limit in  $L_p$ , we have also

$$u(t) = y'(t) = \underset{s \rightarrow \infty}{W\text{-lim}} y_{k_s}(t) = \underset{k \rightarrow \infty}{W\text{-lim}} u_k(t) = u_0(t)$$

a.e. in  $[t_1, t_2]$ , or  $u = u_0$  a.e., and

$$x'(t) = f(t, x(t), u_0(t)), \quad (t, x(t)) \in A, \quad u_0(t) \in U(t, x(t))$$

a.e. in  $[t_1, t_2]$ . Then

$$I[x, u_0] = I[x, u] \leq i = \liminf_{k \rightarrow \infty} I[x_k, u_k].$$

This proves Corollary 3. Note that the hypothesis of convexity of  $Q^*(t, x)$  certainly implies the convexity of  $U(t, x)$ .

In Corollary 3 for  $p = 1$  the functions  $u_k$  are certainly equiabsolutely integrable under the growth condition  $(\epsilon_0)$  mentioned below Corollary 2.

Note that if both  $f$  and  $f_0$  are linear in  $u$ , or

$$f = B(t, x)u + C(t, x), \quad f_0 = b(t, x)u + c(t, x),$$

where  $B, C, b, c$  are  $n \times m, 1 \times n, 1 \times m, 1 \times 1$  matrices with entries continuous in  $A$ , then in Corollary 3 the conditions on the sets  $Q^*$  can be dropped. In other words we have: If  $x_k(t), u_k(t), t_{1k} \leq t \leq t_{2k}, k = 1, 2, \dots$ , is a sequence of admissible pairs with  $x_k \rightarrow x$  in the  $\rho$ -metric and  $u_k \rightarrow u_0$  weakly in  $L_p, p > 1$ , then  $x, u_0$  is admissible and  $I[x, u_0] = \lim I[x_k, u_k]$ . The same holds for  $p = 1$  if the functions  $u_k$  are known to be equiabsolutely integrable. The proof is the same as above with obvious simplifications.

*Remark 1.* Condition  $(\psi)$  in Theorem 4.1 and in the corollaries will be drastically reduced in Theorem 4.2 below. Simple conditions under which the sets  $\tilde{Q}$  above, if convex, are closed and satisfy property (Q) as requested will be given in § 8.

*Remark 2.* In the sufficient condition for lower closure, Theorem 4.1, and in the corollaries, it is enough to request that the sets  $\tilde{Q}(t, x)$  be closed, convex, and have property (Q) at the points  $(t, x(t)) \in A$  for almost all  $t \in [t_1, t_2]$ . In this form, and under suitable regularity hypotheses, the convexity assumption of the sufficient condition, Theorem 4.1, for lower closure will be shown to be necessary also.

*Remark 3.* Condition  $(\psi)$  in Theorem 4.1 and in the corollaries is satisfied if, for instance,  $f_0(t, x, u) \geq 0$  for all  $(t, x, u) \in M$ , or  $f_0(t, x, u) \geq v$  for all  $(t, x, u) \in M$ , where  $v$  is some real constant. Nevertheless, condition  $(\psi)$  in Theorem 4.1 can be reduced. For instance, we may replace it by the following weaker assumption:

$(\psi')$  For every compact subset  $A_0$  of  $A$  there is a locally integrable function  $\psi_0(t)$  (which may depend on  $A_0$ ) such that  $f_0(t, x, u) \geq \psi_0(t)$  for all  $(t, x) \in A_0, u \in U(t, x)$ .

The proof is the same since we can include all trajectories  $x$  and  $x_k$  in a unique compact subset  $A_0$  of  $A$ .

A more drastic generalization of Theorem 4.1 will be given below (Theorem 4.2), where we shall use the following much weaker form of the condition.

( $\psi^*$ ) For every  $(\bar{t}, \bar{x}) \in A$  there are a neighborhood  $N(\bar{t}, \bar{x})$  of  $(\bar{t}, \bar{x})$  in  $A$ , a locally integrable function  $\psi(t)$ , and real numbers  $b_1, \dots, b_n$  (all  $b_1, \dots, b_n$  and  $\psi$  may depend on  $\bar{t}, \bar{x}, N$ ) such that

$$(4.1) \quad \bar{f}(t, x, u) = f_0(t, x, u) - \sum_{j=1}^n b_j f_j(t, x, u) \geq \psi(t)$$

for all  $(t, x) \in N(\bar{t}, \bar{x}), u \in U(t, x)$ , with the exception perhaps of a set of points  $(t, x)$  whose  $t$ -coordinate lies on a set of measure zero on the  $t$ -axis.

*Remark 4.* We shall note here that, under condition ( $\psi^*$ ), it is natural to consider the sets

$$\tilde{Q}(t, x) = \{(z^0, z) | z^0 \geq f_0(t, x, u), z = f(t, x, u), u \in U(t, x)\},$$

or the analogous sets

$$\tilde{Q}^*(t, x) = \{(Z^0, Z) | Z^0 \geq \bar{f}_0(t, x, u), Z = f(t, x, u), u \in U(t, x)\}.$$

It is easy to see that the sets  $\tilde{Q}$  are closed, or convex, or satisfy property (Q) if and only if the same occurs for the sets  $\tilde{Q}^*$ . Indeed, the sets above are transformed into one another by the fixed affine transformation  $Z^0 = z^0 - bz, Z = z$ .

**THEOREM 4.2.** *Let  $A, B, U(t, x), M, f(t, x, u), f_0(t, x, u)$  be as in Theorem 4.1, and let us assume that condition ( $\psi^*$ ) holds. With  $N(\bar{t}, \bar{x})$  as in condition ( $\psi^*$ ), and for every  $(t, x) \in N(\bar{t}, \bar{x})$ , let  $\tilde{Q}(t, x)$  denote the set of all  $\bar{z} = (z^0, z^1, \dots, z^n) = (z^0, z)$  with  $z^0 \geq \bar{f}_0(t, x, u), z = f(t, x, u)$  for  $u \in U(t, x)$ , and assume that the sets  $\tilde{Q}(t, x)$  are closed, convex, and satisfy property (Q) at all points  $(t, x) \in N(\bar{t}, \bar{x})$ , with the exception perhaps of a set of points whose  $t$ -coordinate lies on a set of measure zero on the  $t$ -axis. Then the integral (3.1) has the property of lower closure at every absolutely continuous vector function  $x(t), t_1 \leq t \leq t_2$ , which is the limit in the  $\rho$ -metric of admissible trajectories.*

Theorem 4.2 holds even if the continuity of  $f_0, f$  is replaced by the weaker assumption that  $f_0, f$  are continuous in  $x, u$  for every  $t$ , and measurable in  $t$  for every  $x, u$ . The same argument mentioned under Theorem 4.1 applies.

*Proof.* Let  $x(t), t_1 \leq t \leq t_2$ , be any absolutely continuous function as in the text, and  $x_k(t), u_k(t), t_{1k} \leq t \leq t_{2k}, k = 1, 2, \dots$ , be a sequence of admissible pairs with  $\rho(x_k, x) \rightarrow 0, \liminf I[x_k, u_k] < +\infty$ . Let  $A_0$  be a compact neighborhood (containing the graph of  $x$  and all  $x_k$ ). By hypothesis, for every  $(\bar{t}, \bar{x}) \in A_0$  there are numbers  $\delta > 0, b_1, \dots, b_n$  real, and a locally integrable function  $\psi(t), -\infty < t < +\infty$ , such that  $\bar{f}_0(t, x, u) \geq \psi(t)$  for all  $(t, x) \in N_{2\delta}(\bar{t}, \bar{x})$  and  $u \in U(t, x)$ . We consider the smaller neighborhoods  $N_\delta(\bar{t}, \bar{x})$  which we think of as open (in  $A$ ). These too form a cover of the compact set  $A_0$ . Thus, finitely many of these  $N_\delta$  cover  $A_0$ , say  $N_{\delta_\gamma}(t_\gamma, x_\gamma), \gamma = 1, \dots, s$ . Let  $\delta_\gamma > 0, b_{\gamma 1}, \dots, b_{\gamma n}, \psi_\gamma$  be the corresponding elements, so that

$$\bar{f}_{0_\gamma}(t, x, u) = f_0(t, x, u) - \sum_{j=1}^n b_{\gamma j} f_j(t, x, u) \geq \psi_\gamma(t)$$

for all  $(t, x) \in N_{2\delta_\gamma}(t_\gamma, x_\gamma), u \in U(t, x), \gamma = 1, \dots, s$ , and  $\bigcup_{\gamma=1}^s N_{\delta_\gamma}(t_\gamma, x_\gamma) \supset A_0$ . Let  $b = \max [ |b_{\gamma j}|, j = 1, \dots, n, \gamma = 1, \dots, s ], \delta_0 = \min [ \delta_\gamma, \gamma = 1, \dots, s ]$ .

Since  $(t, x(t)) \in A_0$  for all  $t_1 \leqq t \leqq t_2$ , we can divide the arc  $C_0: x = x(t)$ ,  $t_1 \leqq t \leqq t_2$ , into finitely many subarcs, say  $C_\sigma$ ,  $\sigma = 1, \dots, N$ , each  $C_\sigma$  completely contained in some neighborhood  $N_{\delta_\gamma}(t_\gamma, x_\gamma)$ . Thus, we have for the arcs  $C_\sigma$  the representations  $C_\sigma: x = x(t)$ ,  $\tau_{\sigma-1} \leqq t \leqq \tau_\sigma$ , with  $t_1 = \tau_0 < \tau_1 < \dots < \tau_N = t_2$ , and each  $C_\sigma$  lies in a certain  $N_{b_\gamma}(t_\gamma, x_\gamma)$  which now remains associated with  $C_\sigma$ . Since  $\rho(x_k, x) \rightarrow 0$  as  $k \rightarrow \infty$ , and hence  $t_{1k} \rightarrow t_1$ ,  $t_{2k} \rightarrow t_2$ , we see that for all  $k$  sufficiently large we have  $t_{1k} < \tau_1 < \dots < \tau_{N-1} < t_{2k}$ . Thus, for all  $k$  sufficiently large, the arc  $C_k: x = x_k(t)$ ,  $t_{1k} \leqq t \leqq t_{2k}$ , is divided into the same number  $N$  of subarcs, say  $C_{k\sigma}: x = x_k(t)$ ,  $\tau_{\sigma-1} \leqq t \leqq \tau_\sigma$ ,  $\sigma = 1, \dots, N$ , where now  $\tau_0 = t_1$  must be replaced by  $t_{1k}$  and  $\tau_N = t_2$  must be replaced by  $t_{2k}$ . Also, for all  $k$  sufficiently large, say for  $k \geqq k_0$ , the arc  $C_{k\sigma}$  is completely contained in  $N_{2\delta_\gamma}(t_\gamma, x_\gamma)$  for the same  $\gamma$  we have already associated with  $C_\sigma$ . Thus, for  $k \geqq k_0$ ,  $C_\sigma$  lies in some  $N_{\delta_\gamma}(t_\gamma, x_\gamma)$  and  $C_{k\sigma}$  in  $N_{2\delta_\gamma}(t_\gamma, x_\gamma)$ . Also,  $C_{k\sigma} \rightarrow C_\sigma$  as  $k \rightarrow \infty$  in the sense that the  $\rho$ -distance approaches zero as  $k \rightarrow \infty$ . We shall now consider for each  $\sigma = 1, \dots, N$ , the auxiliary functional

$$J = \int_{t'}^{t''} \bar{f}_0(t, x(t), u(t)) dt$$

for all admissible pairs  $x, u$  with the graph of  $x$  lying in  $N_{2\delta_\gamma}(t_\gamma, x_\gamma)$ . Here by admissible we mean that the conditions listed in § 3 are satisfied with  $A$  replaced by  $N_{2\delta_\gamma}(t_\gamma, x_\gamma)$ , and of course  $\bar{f}_0(t, x(t), u(t))$   $L$ -integrable as usual.

For each  $\sigma$  we may now apply Theorem 4.1 to the arc  $C_\sigma$ , the sequence  $C_{k\sigma}$ ,  $k = 1, 2, \dots$ , and the functional  $J$ . We conclude that each  $C_\sigma$  is admissible and that

$$(4.2) \quad J[C_\sigma] \leqq \liminf_{k \rightarrow \infty} J[C_{k\sigma}], \quad \sigma = 1, \dots, N.$$

More precisely, for each  $\sigma$ , there is a measurable  $u(t)$ ,  $\tau_{\sigma-1} \leqq t \leqq \tau_\sigma$ , such that the pair  $x(t)$ ,  $u(t)$ ,  $\tau_{\sigma-1} \leqq t \leqq \tau_\sigma$ , is admissible for the functional  $J$ ; in particular,  $u(t) \in U(t, x(t))$ ,  $dx/dt = f(t, x(t), u(t))$ ,  $\tau_{\sigma-1} \leqq t \leqq \tau_\sigma$  (a.e.),  $\sigma = 1, \dots, N$ , and the expression

$$\bar{f}_0(t, x(t), u(t)) = f_0(t, x(t), u(t)) - \sum_{j=1}^n b_{j\sigma} f_j(t, x(t), u(t))$$

is  $L$ -integrable in  $[\tau_{\sigma-1}, \tau_\sigma]$ . Since the functions  $f_j$  here are certainly  $L$ -integrable in the same interval (as derivatives of the absolutely continuous functions  $x^i(t)$  in  $[\tau_{\sigma-1}, \tau_\sigma]$ ), we conclude that  $f_0(t, x(t), u(t))$  itself is  $L$ -integrable in each  $[\tau_{\sigma-1}, \tau_\sigma]$  and hence in the whole of  $[t_1, t_2]$ . We have proved that the pair  $x(t)$ ,  $u(t)$ ,  $t_1 \leqq t \leqq t_2$ , is admissible for the original integral  $I$ .

Now, given  $\varepsilon > 0$ , we deduce from (4.2) that there is some  $\bar{k} \geqq k_0$  such that, for  $k \geqq \bar{k}$ , we have

$$(4.3) \quad J[C_{k\sigma}] - J[C_\sigma] > -\varepsilon/N, \quad \rho(C_{k\sigma}, C_\sigma) < \varepsilon/(Nnb), \quad \sigma = 1, \dots, N.$$

Now we have

$$\begin{aligned} I[x, u] &= \int_{t_1}^{t_2} f_0(t, x(t), u(t)) dt = \sum_{\sigma=1}^N \int_{\tau_{\sigma-1}}^{\tau_\sigma} f_0(t, x(t), u(t)) dt \\ &= \sum_{\sigma=1}^N \left\{ \int_{\tau_{\sigma-1}}^{\tau_\sigma} \bar{f}_0(t, x(t), u(t)) dt + \sum_{j=1}^n b_{\sigma j} [x^j(\tau_\sigma) - x^j(\tau_{\sigma-1})] \right\}, \end{aligned}$$

$$\begin{aligned}
 I[x_k, u_k] &= \int_{t_{1k}}^{t_{2k}} f_0(t, x_k(t), u_k(t)) dt = \sum_{\sigma=1}^N \int_{\tau_{\sigma-1}}^{\tau_{\sigma}} f_0(t, x_k(t), u_k(t)) dt \\
 &= \sum_{\sigma=1}^N \left\{ \int_{\tau_{\sigma-1}}^{\tau_{\sigma}} \bar{f}_0(t, x_k(t), u_k(t)) dt + \sum_{j=1}^n b_{\sigma j} [x_k^j(\tau_{\sigma}) - x^j(\tau_{\sigma-1})] \right\},
 \end{aligned}$$

where we have written  $b_{\sigma j}$  instead of  $b_{\gamma j}$  where  $\gamma$  is the index we have associated to  $\sigma$ , and where we have written  $\tau_1 = t_1$  and  $\tau_N = t_2$  instead of  $t_{1k}$  and  $t_{2k}$  in the expression for  $I[x_k, u_k]$ . By taking the difference we now have, remembering relations (4.3),

$$\begin{aligned}
 I[x_k, u_k] - I[x, u] &> -N(\varepsilon/N) \\
 &+ \sum_{\sigma=1}^N \sum_{j=1}^n b_{\sigma j} \{ [x_k^j(\tau_{\sigma}) - x^j(\tau_{\sigma})] + [x_k^j(\tau_{\sigma-1}) - x^j(\tau_{\sigma-1})] \},
 \end{aligned}$$

where each bracket is now in absolute value less than or equal to  $(Nnb)^{-1}\varepsilon$ . We conclude that for all  $k \geq \bar{k}$  we have

$$I[x_k, u_k] - I[x, u] > -\varepsilon - Nnb[2(Nnb)^{-1}\varepsilon] = -3\varepsilon.$$

Because  $\varepsilon$  is arbitrary, we have proved the lower closure of  $I$  at  $x$ .

*Remark 5.* As shown in [1ab], lower closure theorems easily yield existence theorems for the minimum. Indeed, for  $A$  compact, we have only to consider any closed nonempty class  $\Omega$  of admissible pairs  $x, u$  with  $I[x, u] \leq M$ , and guarantee that the corresponding trajectories  $x$  are equiabsolutely continuous. To this purpose, growth conditions on  $f$  and  $f_0$  such as  $(\gamma)$  or  $(\varepsilon)$  below suffice, or alternatively, a condition such as  $(\mu)$  below (see [1ab] for details, for the case of  $A$  not compact, and for the simplifications occurring when the control space  $U$  is bounded).

*Remark 6.* We report below for the sake of comparison a lower closure and existence theorem proved by C. Olech [6a]. Let  $A, U(t, x), M, B, f_0(t, x, u) = (f_1, \dots, f_n), Q(t, x) \subset E_{n+1}$ , be defined as in § 3. Let  $\tilde{z} = (z^0, z) = (z^0, z^1, \dots, z^n)$  be the usual variable in  $E_{n+1}$ . We shall denote by  $\tilde{c} = (c^0, c) = (c^0, c^1, \dots, c^n)$ ,  $\tilde{d} = (d^0, d) = (d^0, d^1, \dots, d^n)$  points in  $E_{n+1}$ , and by  $\tilde{c} \cdot \tilde{z} = c^0 z^0 + c^1 z^1 + \dots + c^n z^n$  the usual inner product in  $E_{n+1}$ . For any point  $(t, x) \in A$  let  $C(t, x)$  denote the set

$$C(t, x) = \{ \tilde{c} \in E_{n+1} | \tilde{z} + \lambda \tilde{c} \in \tilde{Q}(t, x) \text{ for all } \lambda \geq 0 \text{ and all } \tilde{z} \in \tilde{Q}(t, x) \}.$$

Obviously,  $C(t, x)$  is a cone with vertex the origin in  $E_{n+1}$ . The set  $C(t, x)$  is said to be the *asymptotic cone* of the set  $\tilde{Q}(t, x)$ . It is easy to see that if  $\tilde{Q}(t, x)$  is convex, so is  $C(t, x)$ ; if  $\tilde{Q}(t, x)$  is closed, so is  $C(t, x)$ . Also, we shall denote by  $\Gamma$  the cone in  $E_{n+1}$  made up of only the positive half-straight-line  $\Gamma = \{ (c^0, c) | c^2 \geq 0, c = 0 \}$ . Obviously,  $\Gamma \subset \tilde{Q}(t, x)$  for all  $(t, x) \in A$ .

By the polar cone  $C^0(t, x)$  of  $C(t, x)$  is meant as usual the set

$$C^0(t, x) = \{ \tilde{d} = (d^0, d) \in E_{n+1} | \tilde{d} \cdot \tilde{c} \leq 0 \text{ for all } \tilde{c} \in C(t, x) \}.$$

Then, for  $C(t, x) = \Gamma$ , we have

$$C^0(t, x) = \Gamma^0 = \{ \tilde{d} = (d^0, d) | d^0 \leq 0, d \in E_n \},$$

and hence  $\text{int } \Gamma^0$  is the set of all  $\tilde{d} = (d^0, d)$  with  $d^0 < 0, d \in E_n$ .

**THEOREM 4.3** (A lower closure theorem). *Let the sets  $M, B$  be closed and  $A$  compact, let  $f_0(t, x, u), f(t, x, u) = (f_1, \dots, f_n)$  be continuous on  $M$ , and let us assume that the sets  $\tilde{Q}(t, x)$  have the following properties:*

- (a) *the sets  $\tilde{Q}(t, x)$  are closed and convex;*
- (b) *the sets  $\tilde{Q}(t, x)$  satisfy property (U) in  $x$  for every  $t$ ;*
- (c)  *$C(t, x) = \Gamma$  for all  $(t, x) \in A$ ;*
- (d) *for every  $r > 0$  and  $\tilde{c} = (c^0, c)$  with  $c^0 < 0$  there is a locally integrable function  $\psi(t; r, c), t \in E_1$ , which may depend on  $r$  and  $c$ , such that  $\tilde{c} \cdot \tilde{z} \leq \psi(t; r, c)$  for all  $(t, x, u) \in M$  with  $|x| \leq r$  and all  $\tilde{z} = (z^0, z) \in \tilde{Q}(t, x)$ .*

*Then the integral (3.1) has the property of lower closure at every absolutely continuous vector function  $x(t) = (x^1, \dots, x^n), t_1 \leq t \leq t_2$ , which is the limit in the  $\rho$ -metric of admissible trajectories, and the minimum of the functional exists in every closed class  $\Omega$  of admissible trajectories.*

As for Theorem 4.1 and Theorem 4.2 the hypothesis of continuity of  $f_0$  and  $f$  can be replaced by the weaker assumption that  $f_0$  and  $f$  are continuous in  $x, u$  for every  $t$ , and measurable in  $t$  for every  $x, u$  [see [6a)]. Theorem 4.3 is proved by C. Olech in [6a] (reported here in the setting of pp. 175–176, in our notations, and for the problem in which we are interested). Condition (c) of Theorem 4.3 is not required in our Theorems 4.1 and 4.2, and condition (d) in Theorem 4.3 is more demanding than condition  $(\psi^*)$  required in Theorem 4.2. Indeed, if we take

$$b = (b_1, \dots, b_n) = (-c^0)^{-1}c, \quad \bar{\psi}(t) = (-c^0)^{-1}\psi,$$

then (d) can be written in the form  $f_0 + b \cdot f \geq \bar{\psi}$ . In other words, in Theorem 4.3 it is required that for every  $b \in E_n$  and  $r > 0$  there is some  $\bar{\psi}$  satisfying  $f_0 + b \cdot f \geq \bar{\psi}$ . Condition  $(\psi^*)$  of Theorem 4.2 only requires that for every  $(t, x)$  there is some  $b \in E_n$  such that  $f_0 + b \cdot f \geq \psi$  holds in a neighborhood of  $(t, \bar{x})$ . The following examples illustrate these points. For instance, for  $f_0 = (1 + u^2)^{1/2}, f = u, n = 1$ , condition (c) is not satisfied; the same occurs for  $f_0 = (1 + x^2u^2)^{1/2}, f = u, n = 1$ . For  $f_0 = 1 + |u|, f = |u|, n = 1$ , condition (d) is not satisfied, while  $(\psi^*)$  is certainly satisfied.

It may be pointed out that condition (c) actually means that  $f(t, x, u)$  is “of slower growth” than  $f_0(t, x, u)$  in the following sense.

( $\gamma$ ) Given  $\varepsilon > 0$  there is an  $N \geq 0$  such that  $|f| \leq \varepsilon f_0$  whenever  $|f_0| \geq N$ .

A condition of this kind has been used in existence theorems, but lower closure Theorems 4.1 and 4.2 do not need a similar condition.

Another condition which has been used in existence theorems is the following one.

( $\varepsilon$ ) Given  $\varepsilon > 0$  there is some function  $\psi_\varepsilon(t) \geq 0$  which is locally integrable and such that

$$|f(t, x, u)| \leq \psi_\varepsilon(t) + \varepsilon f_0(t, x, u).$$

Under this condition then (d) holds. Indeed, given  $b \in E_n$  take  $\varepsilon > 0$  such that  $\varepsilon|b| < 1$  and note that

$$f_0 + b \cdot f \geq f_0 - |b||f| \geq \varepsilon^{-1}(\varepsilon f_0 - |f|) \geq \varepsilon^{-1}\psi(t).$$

Again, lower closure Theorems 4.1 and 4.2 require much less than (d) or ( $\varepsilon$ ).

An example offered by C. Olech [5, p. 179] concerns the existence of the minimum and is similar to the one considered in [9, particularly § 2 and the first part of Example 4.1, p. 300], namely  $\int_0^1 t^\alpha x'^2 dt$ ,  $0 < \alpha < 1$ , where our existence theorems apply and use is made of condition  $(\varepsilon)$ . Nevertheless, growth conditions  $(\gamma)$ , or  $(\varepsilon)$ , are not needed for existence in classes  $\Omega$  satisfying constraints of the following type:

$$(\mu) \quad \int_{t_1}^{t_2} |x'|^p dt \leq M \quad \text{for some constants } p > 1, M > 0.$$

Condition  $(\mu)$  and the less demanding conditions of a lower closure theorem such as Theorem 4.1 or 4.2 suffice.

Note that all conditions  $(\gamma)$ ,  $(\varepsilon)$ , or  $(\mu)$ , and similar ones, only guarantee the equiabsolute integrability of the derivatives  $x'_k$  of the element  $x_k$  of a minimizing sequence  $[x_k]$ , so that we can extract a subsequence  $x_{k_s}$  whose elements  $x_{k_s}$  converge uniformly (since  $A$  is compact), and we could even request that the derivatives  $x'_{k_s}$  converge weakly in  $L_1$  or  $L_p$ . In this sense we could state in Theorem 4.1 that it is enough in applications to know that the sets  $\tilde{Q}(t, x)$  satisfy property (Q) in the subsets  $A_\lambda$  of a partition  $[A_\lambda]$  of  $A$  into measurable sets as mentioned above. In Theorem 4.3 the combined hypotheses (a) to (d) guarantee that the sets  $\tilde{Q}(t, x)$  satisfy property (Q) in  $x$  for almost every  $t$ , as proved by C. Olech [6a, pp. 168–169].

**5. A variant of the lower closure property.** Theorem 4.1 holds in a slightly stronger form. To formulate it we need, besides the sets  $\tilde{Q}(t, x) \subset E_{n+1}$  of § 4, also the sets  $Q(t, x) = f(t, x, U(t, x)) \subset E_n$ . These sets  $Q(t, x)$  are the projections on the  $z$ -space  $E_n$  of the sets  $\tilde{Q}(t, x)$  of the  $z^0 z$ -space  $E_{n+1}$ . Thus, if the sets  $\tilde{Q}(t, x)$  are convex, so are the sets  $Q(t, x)$ . On the other hand, the sets  $\tilde{Q}(t, x)$  may be closed, without the sets  $Q(t, x)$  being so. This is shown by the example  $n = 1, m = 1, U = \{u | -\infty < u < +\infty\}, f_0 = (1 + u^2)^{1/2}, f = \tan^{-1} u, -\pi/2 < f < \pi/2$ . Then,  $Q$  and  $\tilde{Q}$  are the fixed sets

$$Q = \{z | -\pi/2 < z < \pi/2\} \subset E_1,$$

$$\tilde{Q} = \{(z^0, z) | z^0 \geq \sec z, -\pi/2 < z < \pi/2\} \subset E_2,$$

and  $\tilde{Q}$  is closed, but  $Q$  is not. This example shows also that property (Q) for the sets  $\tilde{Q}$  does not imply the same property for the sets  $Q$ . In the statement below we shall require that both the sets  $\tilde{Q}(t, x)$  and the sets  $Q(t, x)$  have property (Q).

**THEOREM 5.1.** *If we assume, in addition to the hypotheses of Theorem 4.1, that both the sets  $\tilde{Q}(t, x) \subset E_{n+1}$  and the sets  $Q(t, x) \subset E_n$  are closed, convex, and satisfy property (Q) at all points of  $A$  with the exception perhaps of a set of points whose  $t$ -coordinate lies on a set of measure zero on the  $t$ -axis, then for every sequence  $x_k(t), u_k(t), t_{1k} \leq t \leq t_{2k}, k = 1, 2, \dots$ , of admissible pairs, and any absolutely continuous vector function  $x(t) = (x^1, \dots, x^n), t_1 \leq t \leq t_2$ , with  $\rho(x_k, x) \rightarrow 0$  as  $k \rightarrow \infty$ , there is a measurable function  $u(t), t_1 \leq t \leq t_2$ , such that  $u(t) \in U(t, x(t)), x'(t) = f(t, x(t), u(t))$  a.e. in  $[t_1, t_2]$ , and*

$$(5.1) \quad I[x, u] \leq \liminf_{k \rightarrow \infty} I[x_k, u_k].$$

If  $\liminf I[x_k, u_k] < +\infty$ , then certainly the pair  $x, u$  is admissible, and (5.1) holds.



An analogous variant of Theorem 4.2 also holds.

*Proof of Theorem 5.1.* First note that  $f_0(t, x(t), u(t))$  is measurable in  $[t_1, t_2]$  and is  $\geq -\psi(t)$ ; hence,  $I[x, u]$  is finite or  $+\infty$ . If the second member of (5.1) is finite, then  $I[x, u]$  must be finite, hence  $f_0(t, x(t), u(t))$  must be  $L$ -integrable, and the conclusion of Theorem 5.1 reduces to the conclusion of Theorem 4.1 in the case under consideration. If the second member of (5.1) is  $+\infty$ , then (5.1) in itself is trivial, but we still have to prove that a measurable  $u(t)$ ,  $t_1 \leq t \leq t_2$ , exists with  $u(t) \in U(t, x(t))$ ,  $x'(t) = f(t, x(t), u(t))$  a.e. in  $[t_1, t_2]$ . This, however, is a consequence of closure Theorem 1.1 applied to the absolutely continuous  $n$ -vector function  $x$ , the  $n$ -vector function  $f$ , and the sets  $Q(t, x) \subset E_n$ . Theorem 5.1 is thereby proved.

Finally, let us show by an example that an integral  $I[x, u]$  may possess the properties of Theorem 4.1, and thus the property of lower closure as defined at the beginning of § 3, and yet not possess the stronger property of the present § 5.

Indeed, take  $m = n = 1$ ,  $U = \{u | -\infty < u < +\infty\}$ ,  $f = \exp(u)$ ,  $f_0 = \exp(u^2)$ ,  $A = E_2$ , and take  $x(t) = 0$ ,  $0 \leq t \leq 1$ ,  $x_k(t) = k^{-1}t$ ,  $0 \leq t \leq 1$ ,  $k = 1, 2, \dots$ . Here  $x_k \rightarrow x$  uniformly in  $[0, 1]$  as  $k \rightarrow \infty$  and  $I[x_k, u_k] = \exp(\log k)^2 \rightarrow +\infty$  as  $k \rightarrow +\infty$ . Obviously, there is no measurable  $u(t)$ ,  $0 \leq t \leq 1$ , with  $-\infty < u(t) < +\infty$ , such that  $0 = x'(t) = \exp(u(t))$  a.e. in  $[0, 1]$ . The integral  $I$  does not have the strong property represented by the conclusion of Theorem 5.1. Yet the integral  $I$  has the property of lower closure as defined in § 3 as a consequence of Theorem 4.1. Indeed, here

$$\tilde{Q} = \{(z^0, z) | z^0 \geq \exp(u^2), z = \exp(u), u \in E_1\},$$

or

$$\tilde{Q} = \{(z^0, z) | z^0 \geq \exp(\log z)^2, 0 < z < +\infty\},$$

is a fixed closed convex subset of  $E_2$ , and all conditions of Theorem 4.1 are satisfied. Instead,  $Q = \{z | z = \exp(u), u \in E_1\}$  is the set  $Q = \{z | 0 < z < +\infty\}$ , a fixed convex set, and  $Q$  is not closed.

**6. Criteria for property (Q) of the sets  $\tilde{Q}(t, x)$ .** We assume here that the sets  $A, U(t, x), M, Q(t, x), \tilde{Q}(t, x)$  are defined as usual, that the sets  $A$  and  $M$  are closed, and that the functions  $f_0(t, x, u), f(t, x, u) = (f_1, \dots, f_n)$  are continuous on  $M$ .

(a) We say that a function  $g(t, x, u)$  on  $M$  is of *slower growth* than  $f_0(t, x, u)$  as  $|u| \rightarrow \infty$  in a subset  $A_0$  of  $A$  if, for every  $\varepsilon > 0$ , there is some number  $H$ , which may depend on  $\varepsilon, f_0$  and  $A_0$ , such that  $(t, x) \in A_0, |u| \geq H, u \in U(t, x)$  implies  $|g| \leq \varepsilon f_0$ .

**THEOREM 6.1.** *If 1 and  $f$  are of slower growth than  $f_0$  as  $|u| \rightarrow \infty$  in a neighborhood  $N_\delta(\bar{t}, \bar{x})$  of  $(\bar{t}, \bar{x})$  in  $A$ , and  $\tilde{Q}(\bar{t}, \bar{x})$  is convex, then the sets  $\tilde{Q}(t, x)$  satisfy property (Q) at  $(\bar{t}, \bar{x})$  (in particular,  $\tilde{Q}(\bar{t}, \bar{x})$  is closed).*

This statement is proved in Cesari [1b, (2.2.ii)]. Note that if 1 and  $f$  are of slower growth than  $f_0$  as  $|u| \rightarrow \infty$  in  $A$ , then not only the sets  $\tilde{Q}(t, x)$  of Theorem 4.1 satisfy property (Q) in  $A$ , but also condition  $(\psi)$  of Theorem 4.1 is trivially satisfied with  $\psi = \text{const}$ .

(b) In Theorem 6.2 below we shall use a different set of hypotheses. At the beginning of § 5 we noticed that the sets  $Q(t, x) \subset E_n$  are the projections of the sets  $\tilde{Q}(t, x) \subset E_{n+1}$  on the  $z$ -space  $E_n$ ; hence, the convexity of any set  $\tilde{Q}(t, x)$  in  $E_{n+1}$  implies the convexity of the corresponding set  $Q(t, x)$  in  $E_n$ . Nevertheless,

as we proved by an example at the beginning of § 5 the sets  $\tilde{Q}(t, x)$  may be closed and even satisfy property (Q) at any given point  $(\bar{t}, \bar{x})$  without this being the case for the sets  $Q(t, x)$ .

However, the following holds: If the sets  $\tilde{Q}(t, x)$  satisfy property (Q) at  $(\bar{t}, \bar{x})$ , then

$$(*) \quad (z^0, z) \in \bigcap_{\delta} \text{cl co } \cup \tilde{Q}(t, x)$$

implies  $z \in Q(\bar{t}, \bar{x})$ . Indeed, (\*) yields  $(z^0, z) \in \tilde{Q}(\bar{t}, \bar{x})$  by property (Q) at  $(\bar{t}, \bar{x})$ , and then  $z \in Q(\bar{t}, \bar{x})$ .

We shall say that condition  $(\alpha)$  holds at the point  $(\bar{t}, \bar{x}) \in A$  provided:

$$(\alpha) \quad (z^0, z) \in \bigcap_{\delta} \text{cl co } \cup \tilde{Q}(t, x) \text{ implies } z \in Q(\bar{t}, \bar{x}).$$

As mentioned, this condition is necessary for property (Q) of the sets  $\tilde{Q}(t, x)$  at  $(\bar{t}, \bar{x})$ . This same condition  $(\alpha)$  alone is not sufficient for property (Q) as the following example shows: Take  $m = n = 1, U = E_1, f_0 = \exp(tu), f = 0, 0 \leq t \leq 1$ . Then

$$\tilde{Q}(0) = \{(z^0, z) | z^0 \geq 1, z = 0\}, \quad \tilde{Q}(t) = \{(z^0, z) | z^0 > 0, z = 0\}$$

if  $t > 0$ , the sets  $\tilde{Q}$  do not satisfy condition (Q) at  $t = 0$ , but condition  $(\alpha)$  certainly holds at the same point. Note that condition  $(\alpha)$  is trivially satisfied for free problems ( $m = n, f = u, U = E_n$ ) since  $Q = U = E_n$ , and all points  $z \in E_n$  are in  $Q$ .

Now we shall say that condition (X) holds at the point  $(\bar{t}, \bar{x}) \in A$  provided the following holds.

(X) For every  $\bar{z} \in Q(\bar{t}, \bar{x})$  there is at least one point  $\bar{u} \in U(\bar{t}, \bar{x})$  with  $\bar{z} = f(\bar{t}, \bar{x}, \bar{u})$  and the following property: Given  $\varepsilon > 0$  there are numbers  $\delta > 0$ , and  $r, b = (b_1, \dots, b_n)$  real such that

$$(X') \quad f_0(t, x, u) \geq r + \sum_j b_j f_j(t, x, u)$$

for all  $(t, x) \in N_{\delta}(\bar{t}, \bar{x})$  and  $u \in U(t, x)$ ,

$$(X'') \quad f_0(\bar{t}, \bar{x}, \bar{u}) \leq r + \sum_j b_j f_j(\bar{t}, \bar{x}, \bar{u}) + \varepsilon.$$

As we have shown in [1c] this is a very weak requirement. For free problems, for instance, this condition reduces to a weak form of the well-known "seminormality convexity condition" (see below).

**THEOREM 6.2.** *If conditions  $(\alpha)$  and (X) hold at a point  $(\bar{t}, \bar{x}) \in A$ , then the sets  $\tilde{Q}(t, x)$  are convex, closed, and satisfy property (Q) at the point  $(\bar{t}, \bar{x})$ .*

This statement was proved in [1c].

(c) *A partial converse of Theorem 6.2.* We have already seen that the convexity of  $\tilde{Q}(t, x)$  in  $E_{n+1}$  implies the convexity of  $Q(t, x)$  in  $E_n$ . We shall denote by  $R$  the linear manifold in  $E_n$  of minimum dimension  $r$  containing  $Q(t, x)$ . Thus  $Q(t, x) \subset R \subset E_n, 0 \leq r \leq n$ . As usual we denote by  $\text{int } Q(t, x)$  the set of all  $z \in Q(t, x)$  which are in the interior of the convex set  $Q(t, x)$  with respect to  $E_n$ . Analogously, we denote by  $R\text{-int } Q(t, x)$  the set of all  $z \in Q(t, x)$  which are in the interior of  $Q(t, x)$  with respect to  $R$ . Then

$$\text{int } Q(t, x) \subset R\text{-int } Q(t, x) \subset Q(t, x) \subset R \subset E_n.$$

If  $Q(t, x)$  is reduced to a single point, then  $\text{int } Q(t, x) = \emptyset, R\text{-int } Q(t, x) = \emptyset,$

$r = 0$ ,  $R = Q(t, x)$ . If  $Q(t, x)$  contains at least two points, then  $1 \leq r \leq n$ , and  $R\text{-int } Q(t, x) \neq \emptyset$ .

For any point  $(\bar{t}, \bar{x}) \in A$  and  $z \in Q(\bar{t}, \bar{x}) \subset E_n$  let us denote by  $T(z; \bar{t}, \bar{x})$  the number,  $-\infty \leq T < +\infty$ , defined by

$$(6.1) \quad \begin{aligned} T(z; \bar{t}, \bar{x}) &= \inf \{z^0 | (z^0, z) \in \tilde{Q}(\bar{t}, \bar{x}) \text{ for some } z^0 \in E_1\} \\ &= \inf \{z^0 | z^0 = f_0(\bar{t}, \bar{x}, u) \text{ for all } u \in U(\bar{t}, \bar{x}) \text{ with } f(\bar{t}, \bar{x}, u) = z\}. \end{aligned}$$

We shall consider  $T(z; \bar{t}, \bar{x})$  as a function of  $z$  in  $Q(t, x) \subset E_n$ . We proved in [1c] the following properties of the function  $T(z; \bar{t}, \bar{x})$ .

( $\pi_1$ ) If  $\tilde{Q}(t, x)$  is convex, then either  $T(z; \bar{t}, \bar{x}) = -\infty$  for all  $z \in R\text{-int } Q(\bar{t}, \bar{x})$ , or  $T(z; \bar{t}, \bar{x})$  is a convex real-valued function on the convex set  $Q(\bar{t}, \bar{x}) \subset E_n$ , is finite everywhere in  $Q(\bar{t}, \bar{x})$ , and continuous at every point  $z \in R\text{-int } Q(t, x)$ .

( $\pi_2$ ) If  $\tilde{Q}(\bar{t}, \bar{x})$  is convex and closed, and  $T(z; \bar{t}, \bar{x})$  is finite on  $Q(\bar{t}, \bar{x})$ , then  $\min$  holds instead of  $\inf$  in both relations (6.1), and  $T(z; \bar{t}, \bar{x})$  is (continuous at each  $z \in R\text{-int } Q(\bar{t}, \bar{x})$  and) lower semicontinuous at each  $z \in Q(\bar{t}, \bar{x}) - R\text{-int } Q(\bar{t}, \bar{x})$ .

In [1c] we proved by examples that the first alternative in ( $\pi_1$ ) may well occur, and that, under the conditions of ( $\pi_2$ ),  $T(z; \bar{t}, \bar{x})$  may well be discontinuous at the points  $z \in Q(\bar{t}, \bar{x}) - R\text{-int } Q(\bar{t}, \bar{x})$ .

( $\pi_3$ ) Under the conditions of ( $\pi_2$ ),  $T(z; \bar{t}, \bar{x})$  has a supporting plane at every point  $\bar{z} \in R\text{-int } Q(\bar{t}, \bar{x})$ ; that is, there are real numbers  $r, b = (b_1, \dots, b_n)$  such that  $T(\bar{z}; \bar{t}, \bar{x}) \geq r + b \cdot \bar{z}$  for all  $z \in Q(\bar{t}, \bar{x})$  and  $T(z; \bar{t}, \bar{x}) = r + b \cdot \bar{z}$ .

This last property can be reinterpreted by saying that for every  $\bar{z} \in R\text{-int } Q(\bar{t}, \bar{x})$  there is at least one point  $\bar{u} \in U(\bar{t}, \bar{x})$  and real numbers  $r, b = (b_1, \dots, b_n)$  such that

$$(6.2) \quad \begin{aligned} T(z; \bar{t}, \bar{x}) &\geq r + b \cdot z \quad \text{with } z = f(\bar{t}, \bar{x}, u) \quad \text{for all } u \in U(\bar{t}, \bar{x}), \\ T(\bar{z}; \bar{t}, \bar{x}) &= r + b \cdot \bar{z} \quad \text{with } \bar{z} = f(\bar{t}, \bar{x}, \bar{u}). \end{aligned}$$

Finally, the same property ( $\pi_3$ ) can be reinterpreted in turn by saying that for every  $\bar{z} \in R\text{-int } Q(\bar{t}, \bar{x})$  there is at least one point  $\bar{u} \in U(\bar{t}, \bar{x})$  and real numbers  $r, b = (b_1, \dots, b_n)$  such that

$$(6.3) \quad \begin{aligned} f_0(\bar{t}, \bar{x}, u) &\geq r + b \cdot f(\bar{t}, \bar{x}, u) \quad \text{for all } u \in U(\bar{t}, \bar{x}), \\ f_0(\bar{t}, \bar{x}, \bar{u}) &= r + b \cdot f(\bar{t}, \bar{x}, \bar{u}). \end{aligned}$$

In [1c] we proved by examples that under the conditions of ( $\pi_2$ ) the supporting plane for the convex set  $\tilde{Q}(t, x)$  at points  $(z^0, \bar{z})$  with  $\bar{z} \in Q(\bar{t}, \bar{x}) - R\text{-int } Q(\bar{t}, \bar{x})$ ,  $z^0 = T(z; \bar{t}, \bar{x})$ , may well be vertical. Thus there may be no supporting plane (of the form  $z^0 = r + b \cdot z$ ) for  $T(z; \bar{t}, \bar{x})$  at the points  $\bar{z} \in Q(\bar{t}, \bar{x}) - R\text{-int } Q(\bar{t}, \bar{x})$ . At each of these points  $\bar{z}$  we can only say that for every  $\varepsilon > 0$  there are real numbers  $r, b = (b_1, \dots, b_n)$  depending on  $\varepsilon, \bar{t}, \bar{x}, \bar{z}$  such that  $T(z; \bar{t}, \bar{x}) \geq r + b \cdot z$  for all  $z \in Q(\bar{t}, \bar{x})$ , and  $T(\bar{z}; \bar{t}, \bar{x}) < r + b \cdot \bar{z} + \varepsilon$ . Properties analogous to (6.2) and (6.3) then also hold.

Finally, in [1c] we proved the following partial converse of Theorem 6.2.

**THEOREM 6.3.** *If  $\tilde{Q}(t, x)$  is closed and convex for all  $(t, x) \in A$ , if  $\tilde{Q}(t, x)$  satisfies property (Q) at  $(\bar{t}, \bar{x}) \in A$ , and  $T(z; \bar{t}, \bar{x})$  is finite on  $Q(\bar{t}, \bar{x})$ , then for every  $\bar{z} \in Q(\bar{t}, \bar{x})$  and  $\varepsilon > 0$  there are real numbers  $r, b = (b_1, \dots, b_n)$  and  $\delta > 0$  such that*

$$\begin{aligned} T(z; \bar{t}, \bar{x}) &\geq r + b \cdot z - \varepsilon \quad \text{for all } z \in Q(t, x) \quad \text{and } (t, x) \in N_\delta(\bar{t}, \bar{x}), \\ T(\bar{z}; \bar{t}, \bar{x}) &= r + b \cdot \bar{z}. \end{aligned}$$

Hence, there is at least one point  $\bar{u} \in U(\bar{t}, \bar{x})$  such that

$$T(z; t, x) \geq r + b \cdot z \quad \text{for all } (t, x) \in N_\delta(\bar{t}, \bar{x}) \quad \text{and} \quad u \in U(t, x) \\ \text{such that } z = f(t, x, u),$$

$$T(\bar{z}; \bar{t}, \bar{x}) \leq r + b \cdot \bar{z} + \varepsilon.$$

Also, we have

$$f_0(t, x, u) \geq r + b \cdot f(t, x, u) \quad \text{for all } u \in U(t, x) \quad \text{and} \quad (t, x) \in N_\delta(\bar{t}, \bar{x}), \\ f_0(\bar{t}, \bar{x}, \bar{u}) \leq r + b \cdot f(\bar{t}, \bar{x}, \bar{u}) + \varepsilon.$$

Thus, under the mild restriction of assuming  $T(z; t, x)$  finite, we see that condition (Q) for the sets  $\tilde{Q}(t, x)$  and the set of conditions (α) and (X) are equivalent. Under the mild restriction above, we have obtained a characterization of property (Q) for the sets  $\tilde{Q}(t, x)$ .

(d) *The case of  $f$  linear in  $u$ .* We shall assume here that  $A$  is a given closed subset of the  $tx$ -space  $E_{1+n}$ , that  $U = E_m$ , that  $f_0(t, x, u)$  and  $f(t, x, u) = (f_1, \dots, f_n)$  are continuous on  $M = A \times E_m$ , and that  $f$  is linear in  $u$ ; that is,

$$f_i(t, x, u) = \sum_{j=1}^m b_{ij}(t, x)u^j + c_i(t, x), \quad i = 1, \dots, n,$$

or

$$f(t, x, u) = B(t, x)u + C(t, x),$$

where  $B, C$  are  $n \times m$  and  $n \times 1$  matrices with entries continuous in  $A$ . For every compact subset  $A_0$  of  $A$ , the functions  $b_{ij}, c_i$  are continuous and bounded on  $A_0$ ; hence, there are constants  $G_0, F_0$  such that  $|f(t, x, u)| \leq G_0|u| + F_0$  for all  $(t, x) \in A_0$  and  $u \in E_n$ .

**THEOREM 6.4.** *If  $f_0(t, x, u)$  is convex in  $u$ , and  $f$  is linear in  $u$  with  $U = E_m$ , then the sets  $\tilde{Q}(t, x)$  are convex.*

**THEOREM 6.5.** *If  $A$  is closed,  $U = E_m, M = A \times E_m$ , if  $f_0(t, x, u)$  is continuous on  $M$ , convex in  $u$ , and “seminormal” in  $u$  at a point  $\bar{x} \in A$  (see definition (SN) in (e) below), if  $f(t, x, u) = B(t, x)u + C(t, x)$ , where the matrices  $B, C$  have entries continuous in  $A$ , then the sets*

$$\tilde{Q}(t, x) = \{(z^0, z) | z^0 \geq f_0(t, x, u), z = f(t, x, u), u \in E_m\}$$

satisfy property (Q) at  $(\bar{t}, \bar{x})$ .

A proof was given in [1c]. This statement for  $f$  linear in  $u$ , or  $f = B(t, x)u + C(t, x)$ , is much stronger than the analogous statement, Theorem 6.1. Indeed, we would deduce from Theorem 6.1 an analogous statement as Theorem 6.5 under a growth condition  $f_0(t, x, u) \geq \Phi(|u|)$  with  $\Phi(\zeta)/\zeta \rightarrow +\infty$  as  $\zeta \rightarrow +\infty$ .

(e) *The free problem  $m = n, f = u, U = E_n$ .* Here the sets  $Q$  reduce to the fixed, closed, and convex set  $Q = U = E_n$ . The sets  $\tilde{Q}(t, x)$  reduce here to

$$\tilde{Q}(t, x) = \{(z^0, u) | z^0 \geq f_0(t, x, u), u \in E_n\}.$$

These sets are closed whenever  $f_0$  is continuous, and convex whenever  $f_0(t, x, u)$  is convex in  $u$ . As mentioned, condition (α) is trivially satisfied. Condition (X) at a

point  $(\bar{t}, \bar{x}) \in A$  reduces to the following simple (and well-known) requirement :

$(X_f)$  (weak seminormality condition). We say that the real-valued function  $f_0(t, x, u), (t, x) \in A, u \in E_n$ , is weakly seminormal at the point  $(\bar{t}, \bar{x}) \in A$  provided, for every  $\bar{u} \in E_n$  and  $\varepsilon > 0$ , there are numbers  $\delta > 0$ , and  $r, b = (b_1, \dots, b_n)$  real such that

$$(X_f') \quad f_0(t, x, u) \geq r + b \cdot u \quad \text{for all } (t, x) \in N_\delta(\bar{t}, \bar{x}) \quad \text{and } u \in E_n,$$

$$(X_f'') \quad f_0(\bar{t}, \bar{x}, \bar{u}) \leq r + b \cdot \bar{u} + \varepsilon.$$

For this concept and a number of variants, see L. Tonelli [7a, b] and E. J. McShane [3a, b]. Theorem 6.2 now yields the following theorem.

**THEOREM 6.6.** *For free problems ( $m = n, f = u, U = E_n$ ), if  $A$  is closed, if  $f_0(t, x, u)$  is continuous on  $M = A \times E_n$  and convex in  $u$ , and if  $f_0$  is weakly seminormal at a point  $(\bar{t}, \bar{x}) \in A$ , then the sets  $\tilde{Q}(t, x)$  satisfy property (Q) at  $(\bar{t}, \bar{x})$ .*

Convexity of  $f_0$  alone does not imply that the sets  $\tilde{Q}(t, x)$  have property (Q) in  $A$ . This is shown by the following simple example. Take  $n = 1, f_0(t, u) = tu, 0 \leq t \leq 1, u \in U = E_1$ . Then  $f_0$  is continuous and convex in  $u$  for every  $t$ , but at every  $t, 0 \leq t \leq 1$ , we have

$$\tilde{Q}(t) = \{(z^0, u) | z^0 \geq tu, u \in E_1\},$$

a half-plane in  $E_2$ , while  $\bigcap_\delta \text{cl co } \tilde{Q}(t; \delta)$  is the entire plane  $E_2$ . Thus the sets  $\tilde{Q}$  do not satisfy property (Q) at any  $t, 0 \leq t \leq 1$ . Clearly, the function  $f_0$  is not weakly seminormal.

Let us consider now the seminormality condition, which is a somewhat stronger requirement than the weak seminormality condition. The seminormality condition too was used by Tonelli [7] and McShane [3].

$(SN)$  (Seminormality condition). We say that the real-valued function  $f_0(t, x, u), (t, x) \in A, u \in E_n$ , is seminormal at the point  $(\bar{t}, \bar{x}) \in A$  provided, for every  $\bar{u} \in E_n$ , there are numbers  $\delta > 0, v > 0$ , and  $r, b = (b_1, \dots, b_n)$  real such that

$$(SN') \quad f_0(t, x, u) \geq r + b \cdot u + v|u - \bar{u}| \quad \text{for all } (t, x) \in N_\delta(\bar{t}, \bar{x}) \quad \text{and } u \in E_n,$$

$$(SN'') \quad f_0(\bar{t}, \bar{x}, \bar{u}) \leq r + b \cdot \bar{u} + \varepsilon.$$

The seminormality condition has a very simple elegant characterization.

**THEOREM 6.7.** *For free problems ( $m = n, f = u, U = E_n$ ), if the real-valued function  $f_0(t, x, u), (t, x) \in A, u \in E_n$ , is continuous on  $A \subset E_n$ , and convex in  $u$  at some  $(t, x) \in A$ , then  $f_0$  is seminormal at  $(t, x)$  if and only if for no  $u, u_1 \in E_n, u_1 \neq 0$ , it occurs that*

$$f_0(t, x, u) = 2^{-1}[f_0(t, x, u + \lambda u_1) + f_0(t, x, u - \lambda u_1)] \quad \text{for all } \lambda \geq 0.$$

A proof of this statement under smoothness hypotheses was given by L. Tonelli [7a, b]. A proof under the sole hypotheses of continuity and convexity stated in Theorem 6.7 can be found in L. Turner [8] and is reported in [1c]. Note that, if we denote by  $\tilde{Q}(t, x)$  the set  $\{(z^0, u) | z^0 = f_0(t, x, u), u \in E_n\}$ , then  $\tilde{Q}(t, x)$  is often denoted as the "figurative" of  $f_0$  (at the point  $(t, x) \in A$ ). Theorem 6.7 then states that  $f_0$  is seminormal at  $(\bar{t}, \bar{x})$  if and only if the figurative contains no straight line. In particular, if say  $f_0(\bar{t}, \bar{x}, u) \rightarrow +\infty$  as  $|u| \rightarrow +\infty$ , and  $f_0(\bar{t}, \bar{x}, u)$  is convex in  $u$ , then the figurative  $\tilde{Q}(\bar{t}, \bar{x})$  cannot contain any straight line,  $f_0$  is seminormal at  $(\bar{t}, \bar{x})$ ,  $f_0$  is weakly seminormal, and certainly the sets  $\tilde{Q}(t, x)$  satisfy property (Q) at  $(\bar{t}, \bar{x})$ .

**7. Lower semicontinuity.** The rather general concept of lower closure defined in § 3 is a natural extension of the usual concept of lower semicontinuity. Indeed, the definition of lower closure in § 3 reduces to the usual concept of lower semicontinuity whenever the strategy  $u$  is “determined” by the (admissible) trajectory  $x$ , and then the functional (3.1) can be thought of as depending on the (admissible) trajectory  $x$  only :

$$(7.1) \quad I[x] = \int_{t_1}^{t_2} f_0(t, x(t), u(t)) dt.$$

This occurs, for instance, for free problems where  $u(t) = x'(t)$  (a.e.). The purpose of the present § 7 and the next one, § 8, is to clarify the concepts, to deduce theorems of lower semicontinuity from our previous theorems of lower closure in § 4 and § 5, and to show that the usual theorems of lower semicontinuity for free problems are corollaries of our theorems of lower closure. We already presented a statement of this sort in Corollary 1 of § 4.

(a) It may happen that the data  $A, U(t, x), B, f_0(t, x, u), f(t, x, u) = (f_1, \dots, f_n)$  are so arranged that, for any admissible pair  $x(t), u(t), t_1 \leq t \leq t_2$ , the trajectory  $x$  uniquely determines the strategy  $u$  (a.e. in  $[t_1, t_2]$ ). Then the functional (7.1) can be thought of as being defined for every admissible trajectory  $x$ , and we may denote it as  $I[x]$ . In this section we shall refer to these systems, for the sake of brevity, as TDS systems, or systems in which any admissible trajectory uniquely determines the corresponding strategy.

For all these systems the concept of lower closure (§ 3) reduces to the one of lower semicontinuity. Let  $x(t), t_1 \leq t \leq t_2$ , be any absolutely continuous vector function which is the limit in the  $\rho$ -metric of a sequence of admissible trajectories  $x_k(t), t_{1k} \leq t \leq t_{2k}, k = 1, 2, \dots$ , with  $\rho(x_k, x) \rightarrow 0$  and  $\liminf I[x_k] < +\infty$  as  $k \rightarrow \infty$  (thus, of course  $t_{1k} \rightarrow t_1, t_{2k} \rightarrow t_2$ ). The functional (7.1) is said to be *lower semicontinuous* at  $x$  provided, from any such sequence we can conclude that  $x$  is admissible, and that  $I[x] \leq \liminf I[x_k]$  as  $k \rightarrow \infty$ .

(b) For general TDS systems Theorems 4.1 and 5.1 reduce to the following ones.

**THEOREM 7.1.** *For TDS systems, and under the same conditions of Theorem 4.1, let  $x(t), t_1 \leq t \leq t_2$ , be any absolutely continuous function with  $(t, x(t)) \in A$  for all  $t \in [t_1, t_2]$ , and let  $x$  be the uniform limit of admissible trajectories  $x_k(t), t_{1k} \leq t \leq t_{2k}, k = 1, 2, \dots$ , with  $\rho(x_k, x) \rightarrow 0, \liminf I[x_k] < +\infty$  as  $k \rightarrow \infty$ . Then,  $x$  is admissible and  $I[x] \leq \liminf I[x_k]$ .*

**THEOREM 7.2.** *For TDS systems, and under the conditions of Theorem 5.1, let  $x(t), t_1 \leq t \leq t_2$ , be any absolutely continuous function with  $(t, x(t)) \in A$  for all  $t \in [t_1, t_2]$  and let  $x$  be the uniform limit of admissible trajectories  $x_k(t), t_{1k} \leq t \leq t_{2k}, k = 1, 2, \dots$ , that is,  $\rho(x_k, x) \rightarrow 0$  as  $k \rightarrow \infty$ . Then,  $x$  is admissible and  $I[x] \leq \liminf I[x_k]$ .*

In the last statement we understand that there is a measurable function  $u(t), t_1 \leq t \leq t_2$ , such that  $x(t), u(t), t_1 \leq t \leq t_2$ , satisfy all conditions for admissibility but perhaps the  $L$ -integrability of  $f_0(t, x(t), u(t))$  in  $[t_1, t_2]$ , and that this condition also is satisfied whenever  $\liminf I[x_k] \leq +\infty$ .

Theorem 4.2 also has its counterpart here, but we leave its formulation to the reader.

**8. Theorems of lower semicontinuity for free problems.** Let us consider here free problems, that is, systems with  $m = n$ ,  $f = u$ ,  $U = E_n$ ; hence, the strategy  $u(t) = x'(t)$  (a.e.) is determined by the trajectory (a.e.). If  $A$  and  $B$  are closed sets as usual, then  $M = A \times E_n$  is also closed, and  $f_0(t, x, u)$  is a given continuous scalar function on  $M$ . Here a function  $x(t) = (x^1, \dots, x^n)$ ,  $t_1 \leq t \leq t_2$ , is an *admissible trajectory* provided  $x$  is absolutely continuous in  $[t_1, t_2]$ ,  $(t, x(t)) \in A$  for all  $t \in [t_1, t_2]$ ,  $(t_1, x(t_1), t_2, x(t_2)) \in B$ , and  $f_0(t, x(t), x'(t))$  is  $L$ -integrable in  $[t_1, t_2]$ . Then, the cost functional is

$$(8.1) \quad I[x] = \int_{t_1}^{t_2} f_0(t, x(t), x'(t)) dt.$$

The corresponding sets  $Q(t, x)$  and  $\tilde{Q}(t, x)$  have already been discussed in § 6, part (d), and the concept of weak seminormality has been introduced there.

Our general statement, Theorem 4.2, in conjunction with Theorem 6.5 yields the following theorem.

**THEOREM 8.1** (A theorem of lower semicontinuity for free problems). *For free problems ( $m = n$ ,  $f = u$ ,  $U = E_n$ ), if  $A$  is closed, if  $f_0(t, x, u)$  is continuous on  $M = A \times E_n$ , convex in  $u$ , and weakly seminormal with respect to  $u$  in  $A$ , then the functional (8.1) has the property of lower semicontinuity; that is, if  $x(t) = (x^1, \dots, x^n)$ ,  $t_1 \leq t \leq t_2$ , is an absolutely continuous function which is the limit in the  $\rho$ -metric of admissible trajectories  $x_k(t)$ ,  $t_{1k} \leq t \leq t_{2k}$ ,  $k = 1, 2, \dots$ , with  $\rho(x_k, x) \rightarrow 0$  and  $\liminf I[x_k] < +\infty$  as  $k \rightarrow +\infty$ , then  $x$  is admissible and  $I[x] \leq \liminf I[x_k]$ .*

The condition of weak seminormality is certainly satisfied if  $f_0(t, x, u)$  is continuous in  $(t, x, u)$ , convex in  $u$  for every  $(t, x)$ , and  $f_0(t, x, u) \rightarrow +\infty$  as  $|u| \rightarrow +\infty$  for every  $(t, x) \in A$ .

Theorem 8.1 is due to L. Tonelli [7a] who proved it for  $f_0$  of class  $C'$  in  $u$ . A proof under the present sole continuity hypotheses was given by L. Turner [8]. The lower semicontinuity Theorem 8.1 is here a corollary of Theorem 4.2 for lower closure of general Lagrange problems.

Theorem 8.1 without the hypothesis of weak seminormality is not true, as the following simple example shows. Take  $n = 2$ ,  $A = E_3$ ,  $f_0 = yx' - xy'$ ,  $x, y$  state variables. Then  $f_0$  is certainly convex in  $(x', y')$ , namely linear. Nevertheless,  $I = \int_{t_1}^{t_2} (yx' - xy') dt$  is not lower semicontinuous. Indeed, if we take  $C: x = 0, y = 0, 0 \leq t \leq 2\pi$ , and  $C_k: x = k^{-1} \cos k^2 t, y = k^{-1} \sin k^2 t, 0 \leq t \leq 2\pi, k = 1, 2, \dots$ , then  $C_k \rightarrow C, I[C_k] = -2\pi, k = 1, 2, \dots$ , and  $I[C] = 0$ . An analogous example for  $n = 1$  has been given by Tonelli [7b, vol. 2, pp. 390–392]. Nevertheless, Tonelli proved that, for  $n = 1$  and  $f_0$  continuous in  $(t, x, x')$  with continuous first order partial derivatives  $f_{0x}$  and  $f_{0x'x'}$ , Theorem 8.1 holds without the weak seminormality requirement [7a, pp. 205–206]. Again, the example above shows that this is not the case for  $n \geq 2$ . (See, for analogous examples, McShane [3b].)

**9. Convexity as a necessary and sufficient condition for lower closure.** We are now in a position to prove the statement we mentioned in § 4, Remark 2, that the convexity of the sets  $\tilde{Q}$ , that is, the convexity part of property (Q), is essentially a necessary and sufficient condition for lower closure.

The sufficiency part is covered by Theorems 4.1, 4.2 and 5.1. For the necessity part we shall prove that the convexity of the sets  $\tilde{Q}$  is essentially necessary for lower closure.

We shall need a few more definitions. Here again, as in §§ 1–7,  $A$  is a closed subset of the  $tx$ -space  $E_{1+n}$ ,  $U(t, x) \subset E_m$ ,  $M$  defined as usual is a closed set of the  $txu$ -space  $E_{1+n+m}$ ,  $f_0(t, x, u)$ ,  $f(t, x, u) = (f_1, \dots, f_n)$  are continuous on  $M$ , and  $I$  denotes the functional for Lagrange problems

$$(9.1) \quad I[x, u] = \int_{t_1}^{t_2} f_0(t, x(t), u(t)) dt.$$

In this section we shall denote as admissible any pair  $x(t), u(t)$ ,  $t_1 \leq t \leq t_2$ , satisfying all requirements listed in § 3 (disregarding boundary conditions, or equivalently, taking  $B = E_{2n+2}$ ). The functional  $I[x, u]$  is defined for every admissible pair  $x, u$ .

Also, we shall assume below that  $f(t, x, u)$  is locally Lipschitzian with respect to  $x$  in  $M$ .

Given a point  $(\bar{t}, \bar{x}) \in A$ , we shall say that a convex combination of points in  $E_n$ ,

$$(9.2) \quad \begin{aligned} \bar{z} &= \sum_{j=1}^{n+1} \lambda_j z_j, & \sum_{j=1}^{n+1} \lambda_j &= 1, \\ \bar{z} &= f(\bar{t}, \bar{x}, \bar{u}), & z_j &= f(\bar{t}, \bar{x}, u_j), \quad j = 1, \dots, n+1, \end{aligned}$$

is generic at  $(\bar{t}, \bar{x})$  provided:

- (a) there is some  $\delta > 0$  such that  $\bar{u}, u_j \in U(t, x)$ ,  $j = 1, \dots, n+1$ , for all  $(t, x) \in N_\delta(\bar{t}, \bar{x})$ ;
- (b)  $\Delta = \det(F_{ij}, i, j = 1, \dots, n+1) \neq 0$ , where  $F_{ij} = f_i(\bar{t}, \bar{x}, u_j)$  for  $i = 1, \dots, n$ ,  $j = 1, \dots, n+1$ ;  $F_{ij} = 1$  for  $i = n+1, j = 1, \dots, n+1$ ;
- (c)  $0 < \lambda_j < 1, j = 1, \dots, n+1$ .

For instance, for  $m = n = 1$ ,  $f = u$ ,  $U = E_1$ ,  $\lambda_1 = \lambda_2 = 1/2$ ,  $\bar{u} = 0$ ,  $u_1 = 1$ ,  $u_2 = -1$ , we have  $\Delta = \det(1, -1; 1, 1) = 2 \neq 0$ , and (a), (b), (c) are satisfied. For free problems ( $m = n$ ,  $f_i = u^i$ ,  $i = 1, \dots, n$ ,  $U = E_n$ ,  $z = u$ ), (a) is always trivially satisfied, and near any convex combination there are as many generic convex combinations as we want.

The statement below, Theorem 9.1, depicts a situation where there is no “lower semicontinuity” at suitable trajectories issued from a point  $(\bar{t}, \bar{x}) \in A$  when the convexity requirement is not satisfied.

**THEOREM 9.1.** *Let  $A, U(t, x), M$  be defined as usual and closed, and let  $f_0(t, x, u)$ ,  $f(t, x, u) = (f_1, \dots, f_n)$  be continuous on  $M$  and locally Lipschitzian with respect to  $x$ . If  $(\bar{t}, \bar{x})$  is any point interior to  $A$ , and if  $\bar{z} = \sum_{j=1}^{n+1} \lambda_j z_j$  is a convex combination of points of  $E_n$ , with  $\bar{z} = f(\bar{t}, \bar{x}, \bar{u})$ ,  $z_j = f(\bar{t}, \bar{x}, u_j), j = 1, \dots, n+1$ , which is generic at  $(\bar{t}, \bar{x})$ , and*

$$(9.3) \quad f_0(\bar{t}, \bar{x}, \bar{u}) > \sum_{j=1}^{n+1} \lambda_j f_0(\bar{t}, \bar{x}, u_j),$$

*then there are admissible pairs  $x(t), u(t)$ ,  $\bar{t} \leq t \leq \bar{t} + \delta$ , and  $x_k(t), u_k(t)$ ,  $\bar{t} \leq t \leq \bar{t} + \delta$ ,*



$k = 1, 2, \dots$ , such that  $x_k \rightarrow x$  as  $k \rightarrow \infty$  uniformly in  $[\bar{t}, \bar{t} + \delta]$ , and  $I[x, u] > \liminf I[x_k, u_k]$ .

*Proof.* By force of (a) there is a  $\delta_0 > 0$  such that  $\bar{u}, u_j \in U(t, x)$  for all  $(t, x) \in N_\delta(\bar{t}, \bar{x})$ ,  $j = 1, \dots, v$ . Take any number  $\delta'$  with  $0 < \delta' < \delta_0(n^2 + 1)^{-1/2}$ , and note that all points  $(t, x)$  with  $\bar{t} \leq t \leq \bar{t} + \delta'$ ,  $|x - \bar{x}| \leq \delta'$  are at a distance from  $(\bar{t}, \bar{x})$  which is  $\leq (\delta'^2 + n^2\delta'^2)^{1/2} = \delta'(n^2 + 1)^{1/2} < \delta_0$ . We can take  $\delta' > 0$  sufficiently small so that all mentioned points are in  $A$ .

Let us consider the differential equation

$$(9.4) \quad dx/dt = f(t, x, \bar{u}),$$

with initial condition  $x(\bar{t}) = \bar{x}$ . Having assumed  $f$  continuous on  $M$  and locally Lipschitzian with respect to  $x$  in  $M$ , we know that  $f(t, x, \bar{u})$  is continuous and uniformly Lipschitzian with respect to  $x$  in the set  $\{(t, x) | \bar{t} \leq t \leq \bar{t} + \delta', |x - \bar{x}| \leq \delta'\}$ . By differential equation theory we know that for some  $\delta$  sufficiently small,  $0 < \delta \leq \delta'$ , there is one and only one absolutely continuous solution  $x(t)$ ,  $\bar{t} \leq t \leq \bar{t} + \delta$ , to (9.4) with  $x(\bar{t}) = \bar{x}$  and  $|x(t) - \bar{x}| \leq \delta'$  for all  $\bar{t} \leq t \leq \bar{t} + \delta$ .

Now let us consider the system of  $n + 1$  linear equations in the  $n + 1$  unknowns  $p_1(t), \dots, p_{n+1}(t)$ :

$$(9.5) \quad \begin{aligned} &\sum_{j=1}^{n+1} p_j(t) f(t, x(t), u_j) - f(t, x(t), \bar{u}) = 0, \\ &\sum_{j=1}^{n+1} p_j(t) - 1 = 0, \quad \bar{t} \leq t \leq \bar{t} + \delta. \end{aligned}$$

Note that for  $t = \bar{t}$  these equations have the trivial solution  $p_j(\bar{t}) = \lambda_j, j = 1, \dots, n + 1$ . The functional determinant of these equations is

$$\Delta(t) = \det (F_{ij}(t), i, j = 1, \dots, n + 1),$$

with  $F_{ij}(t) = f_i(t, x(t), u_j), i = 1, \dots, n, j = 1, \dots, n + 1$ , and  $F_{ij} = 1, i = n + 1, j = 1, \dots, n + 1$ . Hence,  $\Delta(\bar{t}) \neq 0$  by force of (b).

Thus, we can take  $\delta > 0$  sufficiently small so that equations (9.5) have a continuous solution  $p(t) = (p_1, \dots, p_{n+1})$  in  $[\bar{t}, \bar{t} + \delta]$ . Since  $p_j(\bar{t}) = \lambda_j$ , and  $0 < \lambda_j < 1$ , we can take  $\delta > 0$  so small that we have also  $0 \leq p_j(t) \leq 1, \bar{t} \leq t \leq \bar{t} + \delta, j = 1, \dots, n + 1$ .

Note that  $x(t), \bar{t} \leq t \leq \bar{t} + \delta$ , is an absolutely continuous solution of the differential system

$$(9.6) \quad \frac{dx}{dt} = \sum_{j=1}^n p_j(t) f(t, x(t), u_j), \quad \bar{t} \leq t \leq \bar{t} + \delta,$$

with initial value  $x(\bar{t}) = \bar{x}$ , since the second member coincides with  $f(t, x(t), \bar{u})$  in  $[\bar{t}, \bar{t} + \delta]$ .

Note that, by hypothesis,

$$(9.7) \quad f_0(\bar{t}, \bar{x}, \bar{u}) > \sum_{j=1}^{n+1} \lambda_j f_0(\bar{t}, \bar{x}, u_j).$$

If we denote by  $2\sigma > 0$  the difference between the first and second member in (9.7), we see that by simple continuity argument, we can take  $\delta > 0$  sufficiently

small so that

$$f_0(t, x(t), \bar{u}) > \sum_{j=1}^n p_j(t) f_0(t, x(t), u_j) + \sigma$$

for all  $\bar{t} \leqq t \leqq \bar{t} + \delta$ .

Note that we have here a generalized system  $x(t), p(t), v(t), \bar{t} \leqq t \leqq \bar{t} + \delta$ , with  $v(t) = \{u^{(j)}(t) = u_j, j = 1, \dots, n + 1\}, \bar{t} \leqq t \leqq \bar{t} + \delta$ , and

$$\begin{aligned} I[x, p, v] &= \int_{\bar{t}}^{\bar{t}+\delta} \sum_{j=1}^{n+1} p_j(t) f_0(t, x(t), u_j) dt \\ &< \int_{\bar{t}}^{\bar{t}+\delta} f_0(t, x(t), \bar{u}) dt - \sigma\delta = I[x, u] - \sigma\delta, \end{aligned}$$

where  $u$  denotes here the constant strategy  $u(t) = \bar{u}, \bar{t} \leqq t \leqq \bar{t} + \delta$ .

Now we shall apply the general theorem of approximation of generalized solutions by means of usual solutions proved in [1a, § 14, (i), p. 416]. To do this we need to reduce drastically the control space since in that theorem it is assumed that the control space has to depend on  $t$  only. In the present situation, however, we can take as auxiliary control space  $U^*$  the fixed space made up of the  $n + 2$  points  $\bar{u}, u_j, j = 1, \dots, n + 1$ , of  $E_m$ ; these are the only points we need in the control space. Then, by force of [1a, § 14, (i)] we know that there exists a sequence of admissible pairs  $x_k(t), u_k(t), \bar{t} \leqq t \leqq \bar{t} + \delta, k = 1, 2, \dots$ , with  $x_k \rightarrow x$  uniformly in  $[\bar{t}, \bar{t} + \delta]$  as  $k \rightarrow \infty$ , and  $I[x_k, u_k] \rightarrow I[x, p, v]$ . Thus, for  $k$  sufficiently large we certainly have  $|x_k(t) - x(t)| < \delta'$  and thus  $(t, x_k(t)) \in A$  for all  $\bar{t} \leqq t \leqq \bar{t} + \delta$ , and  $I[x_k, u_k] < I[x, p, v] + \sigma\delta/2$ ; hence,

$$I[x_k, u_k] < I[x, u] - \sigma\delta/2$$

for all  $k$  sufficiently large. Theorem 9.1 is thereby proved.

The statement below, Theorem 9.2, depicts a situation where there is no "lower semicontinuity" at a given trajectory  $x$  in  $A$ , that is, at a given admissible pair  $x(t), u(t), t_1 \leqq t \leqq t_2$ , where  $u(t)$  may be bounded or unbounded. The statement is similar to Theorem 9.1. Since we shall use the same theorem [1a, § 14, (i)] concerning the approximation of generalized solutions by means of usual solutions, we have to make sure that the conditions of that theorem are satisfied. For the case in which the strategy  $u(t)$  is bounded, the local Lipschitz condition suffices; for the case in which  $u(t)$  is unbounded we shall need the following assumption (S) which is of the type considered by McShane in [3f] for the same purpose of extending to unbounded strategies statements proved for bounded strategies.

(S) There is a number  $\delta > 0$  and an  $L$ -integrable function  $S(t), t_1 \leqq t \leqq t_2$ , such that  $t \in [t_1, t_2], |x' - x(t)| \leqq \delta, |x'' - x(t)| \leqq \delta$  implies

$$|f(t, x', u(t)) - f(t, x'', u(t))|, |f_0(t, x', u(t)) - f_0(t, x'', u(t))| \leqq |x' - x''|S(t).$$

We are now in a position to state and prove the following theorem.

**THEOREM 9.2.** *Let  $A, U(t, x), M$  be defined as usual and closed, and let  $f_0(t, x, u), f(t, x, u) = (f_1, \dots, f_n)$  be continuous on  $M$  and locally Lipschitzian with respect to  $x$ . Let  $x(t), u(t), t_1 \leqq t \leqq t_2$ , be any admissible pair, whose trajectory  $x$  is interior to  $A$ . Let us assume that there is a subset  $E$  of positive measure in  $[t_1, t_2]$  such that for*

every  $\bar{i} \in E$  there is a convex combination of points of  $E_n$ :

$$(9.8) \quad \begin{aligned} \bar{z} &= \sum_{j=1}^{n+1} \lambda_j z_j, & \sum_{j=1}^{n+1} \lambda_j &= 1, \\ \bar{z} &= f(\bar{i}, \bar{x}, \bar{u}), & z_j &= f(\bar{i}, \bar{x}, u_j), & \bar{x} &= x(\bar{i}), & \bar{u} &= u(\bar{i}), \end{aligned}$$

which is generic at  $(\bar{i}, \bar{x})$ , and such that

$$(9.9) \quad f_0(\bar{i}, \bar{x}, \bar{u}) > \sum_{j=1}^{n+1} \lambda_j f_0(\bar{i}, \bar{x}, u_j).$$

Finally, if  $u(t)$ ,  $t_1 \leq t \leq t_2$ , is unbounded, we assume that condition (S) holds. Then, there is a sequence of admissible pairs  $x_k(t), u_k(t)$ ,  $t_1 \leq t \leq t_2$ ,  $k = 1, 2, \dots$ , such that  $x_k \rightarrow x$  as  $k \rightarrow \infty$  uniformly in  $[t_1, t_2]$  and  $I[x, u] > \liminf I[x_k, u_k]$ .

*Proof.* By a suitable reduction we may well assume  $E$  to be a closed, hence compact, subset of  $[t_1, t_2]$ . Also, since  $u(t)$  is measurable in  $[t_1, t_2]$ , and hence continuous on compact subsets of measure as close to  $t_2 - t_1$  as we want, we may well assume that  $u(t)$  is continuous on  $E$ . If a point  $\bar{i} \in E$  and hence relations (9.8), (9.9) hold for certain points  $u_j$  and numbers  $\lambda_j, j = 1, \dots, n + 1$ , let us prove that there is some closed interval  $I = [\bar{i} - \delta, \bar{i} + \delta]$  and functions  $p_j(t), t \in E \cap I$ , such that  $p_j(\bar{i}) = \lambda_j, j = 1, \dots, n + 1$ , and

$$(9.10) \quad \begin{aligned} f(t, x(t), u(t)) &= \sum_{j=1}^{n+1} p_j(t) f(t, x(t), u_j), & t \in E \cap I, \\ 1 &= \sum_{j=1}^{n+1} p_j(t), & t \in E \cap I, \\ f_0(t, x(t), u(t)) &> \sum_{j=1}^{n+1} p_j(t) f_0(t, x(t), u_j), & t \in E \cap I, \\ 0 < p_j(t) < 1, & j = 1, \dots, n + 1, & t \in E \cap I, \\ \bar{u} \in U(t, x(t)), & u_j \in U(t, x(t)), & j = 1, \dots, n + 1, & t \in I. \end{aligned}$$

Indeed, the first two equations of (9.10), thought of as equations in  $p_1, \dots, p_{n+1}$ , are satisfied at  $t = \bar{i}$  with  $p_j(\bar{i}) = \lambda_j, j = 1, \dots, n + 1$ , and the functional determinant  $\Delta(t)$  is  $\neq 0$  at  $t = \bar{i}$ . To see that we are applying the usual implicit function theorem properly, we should first choose a fixed interval  $I_0$  as  $I$  above, then modify  $u(t)$  on  $I_0 - E$  so as to make  $u$  continuous in the whole of  $I_0$ , and then determine the interval  $I \subset I_0$  by the usual implicit function theorem of calculus on the system so modified. We then disregard the values chosen for  $u(t)$  in  $I_0 - E$ . Finally, we can further reduce  $I$  if necessary so as to satisfy also the third, fourth, and fifth of relations (9.10). Note that  $\bar{u} = \bar{u}(\bar{i}) \in U(t, x(t))$  together with  $u(t) \in U(t, x(t))$  hold for all  $t \in I$ .

By the process indicated we associate with each point  $\bar{i} \in E$  a closed interval  $I = [\bar{i} - \delta, \bar{i} + \delta]$  with the properties above. Hence, finitely many of these intervals cover  $E$ , say  $I_1, \dots, I_{N_0}$ . The endpoints of these intervals can be used to define a subdivision of  $[t, t_2]$  into finitely many parts  $J_1, \dots, J_N$ . Let  $\{J\}'$  be the collection of those  $J_s$  which are parts of at least one interval  $I$  above, and then let us make

a choice of the corresponding points  $u_j$  and multipliers  $p_j(t), t \in J_s, j = 1, \dots, n + 1$ , for which all relations (9.10) hold. Let  $\{J\}''$  be the collection of the remaining intervals  $J_s$ .

Note that now relations (9.10) hold in each closed interval  $J_s \in \{J\}'$  and for the points  $u_j$  and multipliers  $p_j(t)$  chosen for the closed interval  $J_s$ . By a continuity argument and the closedness of the intervals  $J_s$  we conclude that there are some constants  $\sigma > 0, \sigma_0 > 0$  such that

$$\begin{aligned}
 f(t, x(t), u(t)) &= \sum_{j=1}^{n+1} p_j(t) f(t, x(t), u_j), & t \in E \cap J_s, \\
 1 &= \sum_{j=1}^{n+1} p_j(t), & t \in E \cap J_s, \\
 f_0(t, x(t), u(t)) &\cong \sum_{j=1}^{n+1} p_j(t) f_0(t, x(t), u_j) + \sigma, & t \in E \cap J_s, \\
 0 < \sigma_0 &\leq p_j(t) \leq 1 - \sigma_0 < 1, & t \in E \cap J_s, \\
 \bar{u} \in U(t, x(t)), & \quad u_j \in U(t, x(t)), & t \in J_s,
 \end{aligned}$$

for every  $J_s \in \{J\}'$ .

Let  $\mu > 0$  denote the measure of  $E$ .

Now let us define  $p_j(t), u_j(t)$  on the whole of  $[t_1, t_2]$ . For  $t \in E$ , hence  $t \in E \cap J_s$  for some  $s, J_s \in \{J\}'$ , let us take for  $p_j(t), u_j(t) = u_j$ , the values already chosen in  $J_s$ ; for  $t \in H = [t_1, t_2] - E$  let us take  $p_j(t) = (n + 1)^{-1}, u_j(t) = u(t), j = 1, \dots, n + 1$ . It is immaterial what choice we make at the endpoints of the intervals  $J_s$ . If we consider the differential system

$$\frac{dy}{dt} = \sum_{j=1}^{n+1} p_j(t) f(t, y, u_j(t)), \quad t_1 \leq t \leq t_2,$$

we see that  $x(t), t_1 \leq t \leq t_2$ , is a solution since the second member coincides with  $f(t, x(t), u(t))$  for all  $t \in [t_1, t_2]$ . In other words,  $x(t), p(t) = (p_1(t), \dots, p_{n+1}(t)), u(t) = (u_1(t), \dots, u_{n+1}(t))$  is a generalized system.

Now we shall apply the general theorem of approximation of generalized solutions by means of usual solutions proved in [1a, § 14, (i)]. We need to reduce drastically the control space since in [1a, § 14, (i)] it is assumed that such a control space has to depend on  $t$  only. For every  $\bar{t} \in E$ , hence  $\bar{t} \in E \cap J_s, J_s \in \{J\}'$ , we take as auxiliary control space  $U^*(\bar{t})$  the space made up of the  $n + 1$  points  $u_j, j = 1, \dots, n + 1$ . These are the only points we need, since  $\bar{z} = f(\bar{t}, \bar{x}, \bar{u}), \bar{x} = x(\bar{t}), \bar{u} = u(\bar{t})$  is used here only as the convex combination  $\sum_{j=1}^{n+1} p_j(\bar{t}) f(\bar{t}, \bar{x}, u_j)$ . For every  $\bar{t} \in [t_1, t_2] - E$ , we take  $U^*(\bar{t}) = \{u(t)\}$ , made up, that is, of the sole point  $u(t)$ . No other points  $u$  are needed.

Let  $A_0$  be a compact neighborhood of the graph  $G$  of the trajectory  $x$ , so that  $G \subset \text{int } A_0 \subset A_0 \subset A$ . Let  $M^*$  be the set of all  $(t, x, u)$  with  $(t, x) \in A_0, u \in U^*(t)$ . If the original strategy  $u(t)$  was bounded, then  $M^*$  is a closed bounded set, and  $f_0$  and  $f$  are Lipschitzian on  $M^*$ . If the original strategy  $u(t)$  was unbounded, then  $M^*$  is only a closed set, and  $f_0$  and  $f$  are locally Lipschitzian on  $M^*$ , and property (S) is satisfied. In either case the theorem in [1a, § 14, (i)] holds, and hence there is a

sequence

$$x_k(t) = (x_k^1, \dots, x_k^n), \quad u_k(t) = (u_k^1, \dots, u_k^n), \quad t_1 \leq t \leq t_2, \quad k = 1, 2, \dots,$$

of (usual) admissible pairs with  $x_k \rightarrow x$  as  $k \rightarrow \infty$  uniformly in  $[t_1, t_2]$  and such that  $\lim I[x_k, u_k] = I[x, p, v]$ , or

$$\lim \int_{t_1}^{t_2} f_0(t, x_k(t), u_k(t)) dt = \int_{t_1}^{t_2} \sum_{j=1}^{n+1} p_j(t) f_0(t, x(t), u_j(t)) dt.$$

On the other hand,

$$\begin{aligned} I[x, u] &= \int_{t_1}^{t_2} f_0(t, x(t), u(t)) dt = \left( \int_H + \int_E \right) f_0 dt \\ &\geq \int_H f_0(t, x(t), u(t)) dt + \int_E \left[ \sum_{j=1}^{n+1} p_j(t) f_0(t, x(t), u_j) + \sigma \right] dt \\ &= I[x, p, v] = \sigma \mu = \lim I[x_k, u_k] + \sigma \mu. \end{aligned}$$

This proves Theorem 9.2.

*Remark.* Theorems 9.1 and 9.2 correspond to analogous statements proved by Tonelli concerning the necessity of the convexity hypothesis for lower semicontinuity in free problems (Tonelli [7]) (see also L. Turner [8]). The present treatment of the question of convexity as a necessary condition for lower closure differs from these essentially because of the use of the theorem of approximation of generalized solutions by means of usual solutions which simplifies an otherwise rather difficult argument in the Lagrange problems under consideration. The present treatment differs also from the recent one by P. Brunovsky [10] because of the use of the same theorem and because of a completely different set of underlying hypotheses.

**Appendix.** We give the counterexample mentioned in § 4, which is essentially due to A. Lasota and C. Olech [2]. Let  $n = 1$  and  $A = [0, 1] \times E_1$ , let  $C$  be a closed Cantor subset of  $[0, 1]$  whose measure  $m(C)$  is positive, and let  $C' = [0, 1] - C$ . Then  $C'$  is the countable union of disjoint subintervals of  $[0, 1]$ . Let  $s(t)$  be a continuous function on  $C'$ , which is positive, integrable on  $C'$ , and which tends to  $+\infty$  whenever  $t$  tends to an end of any interval component of  $C'$ . Let  $m = 1$  and define

$$U(t, x) = U(t) = \begin{cases} \{-1\} & \text{if } t \in C, \\ \{u | u \geq s(t)\} & \text{if } t \in C'. \end{cases}$$

Let  $f(t, x, u) = u$ . Then  $Q(t, x) = Q(t) = U(t)$  and we note that the sets  $Q(t)$  so defined satisfy property (Q) in either  $A_1 = C \times E_1$  or  $A_2 = C' \times E_1$ , a decomposition of  $A$  into disjoint measurable subsets as described in § 4. Let us extend the function  $s$  by taking  $s(t) = 0$  when  $t \in C$ , and consider the decomposition of  $[0, 1]$  into  $k$  intervals of equal lengths:

$$J_{ks} = [t_{k,s-1}, t_{ks}], \quad s = 1, \dots, k.$$

Then  $t_{ks} = s/k$ . Define  $u_k$  by taking

$$u_k(t) = s(t) + v_k(t),$$

where  $v_k(t) = -1$  if  $t \in C$  and

$$v_k(t) = m(C \cap J_{ks})/m(C' \cap J_{ks})$$

if  $t \in C' \cap J_{ks}$ . Then  $u_k$  is integrable in  $[0, 1]$  and  $u_k(t) \in U(t)$  for every  $t \in [0, 1]$  and  $k$ . Let  $x_k$  denote the trajectory corresponding to  $u_k$  with  $x_k(0) = 0$ . Thus

$$x_k(t) = \int_0^t v_k(\tau) d\tau + \int_0^t s(\tau) d\tau = x(t) + y_k(t).$$

It is easy to see that  $y_k(t_{ks}) = 0$  and that  $|y_k(t)| \leq 1/k$ . Hence  $x_k \rightarrow x$  uniformly in  $[0, 1]$ , where  $x'(t) = s(t)$  and  $x(0) = 0$ . Now  $x$  is not an admissible solution of the orientor field under consideration. If we take  $f_0 = 0$ , we obtain a situation where the closure property of Theorem 4.1 is not true, with the property (Q) satisfied only in  $A_1$  and  $A_2$  separately. As mentioned in § 4, whenever  $x'_k \rightarrow x'$  weakly, then it is enough to know that property (Q) is satisfied at each set  $A_k$  of a decomposition  $[A_1, A_2, \dots]$  of  $A$  into countably many disjoint measurable subsets as described in § 4. The counterexample above was suggested by the referee.

#### REFERENCES

- [1a] L. CESARI, *Existence theorems for weak and usual optimal solutions in Lagrange problems with unilateral constraints. I and II*, Trans. Amer. Math. Soc., 124 (1966), pp. 369–412, 413–429.
- [1b] ———, *Existence theorems for optimal controls of the Mayer type*, this Journal, 6 (1968), pp. 517–552.
- [1c] ———, *Seminormality and upper semicontinuity in optimal control*, J. Optimization Theory and Applications, 6 (1970), pp. 114–137.
- [1d] ———, *Existence theorems for multidimensional problems of optimal control*, Differential Equations and Dynamic Systems, Academic Press, New York, 1967, pp. 115–132.
- [1e] ———, *Existence theorems for multidimensional Lagrange problems*, J. Optimization Theory and Applications, 1 (1967), pp. 87–112.
- [1f] ———, *Sobolev spaces and multidimensional Lagrange problems of optimization*, Ann. Scuola Norm. Sup. Pisa, 22 (1960), pp. 193–227.
- [1g] ———, *Multidimensional Lagrange problems of optimization in a fixed domain and an application to a problem of magnetohydrodynamics*, Arch. Rational Mech. Anal., 29 (1968), pp. 81–104.
- [1h] ———, *Existence theorems for abstract multidimensional control problems*, J. Optimization Theory and Applications, 6 (1970), pp. 210–236.
- [1i] ———, *An existence theorem in problems of optimal control*, this Journal, 3 (1965), pp. 7–22.
- [1j] ———, *Un teorema di esistenza in problemi di controlli ottimi*, Ann. Scuola Norm. Sup. Pisa (3), 19 (1965), pp. 35–78.
- [2] A. LASOTA AND C. OLECH, *On Cesari's semicontinuity condition for set valued mappings*, Bull. Acad. Polon. Sci. Sér. Sci. Math. Astronom. Phys., 16 (1968), pp. 711–716.
- [3a] E. J. MCSHANE, *On the semicontinuity of integrals in the calculus of variations*, Ann. of Math., 33 (1932), pp. 460–484.
- [3b] ———, *Existence theorems for ordinary problems of the calculus of variations*, Ann. Scuola Norm. Sup. Pisa (2), 3 (1934), pp. 181–211, pp. 287–315.
- [3c] ———, *Semicontinuity of integrals in the calculus of variations*, Duke Math. J., 2 (1936), pp. 597–616.
- [3d] ———, *Some existence theorems for problems in the calculus of variations*, Ibid., 4 (1938), pp. 132–156.
- [3e] ———, *A navigation problem in the calculus of variations*, Amer. Math. J., 59 (1937), pp. 327–334.
- [3f] ———, *Relaxed controls and variational problems*, this Journal, 5 (1967), pp. 438–485.
- [4] E. J. MCSHANE AND R. B. WARFIELD, *On Filippov's implicit functions lemma*, Proc. Amer. Math. Soc., 18 (1967), pp. 41–47.

- [5] T. NISHIURA, *On an existence theorem for optimal control*, this Journal, 5 (1967), pp. 532–544.
- [6a] C. OLECH, *Existence theorems for optimal problems with vector valued cost functions*, Trans. Amer. Math. Soc., 136 (1969), pp. 157–180.
- [6b] ———, *Existence theorems for optimal control problems involving multiple integrals*, J. Differential Equations, 6 (1969), pp. 512–526.
- [7a] L. TONELLI, *Sugli integrali del calcolo delle variazioni in forma ordinaria*, Ann. Scuola Norm. Sup. Pisa (2), 3 (1934), pp. 401–450 = Opere Scelte, vol. 3, Edizioni Cremonese, Roma, 1962, pp. 192–254.
- [7b] ———, *Fondamenti di Calcolo delle Variazioni*, 2 vols., Zanichelli, Bologna, 1921–23.
- [8] L. TURNER, *The direct method in calculus of variations*. Doctoral thesis, Purdue University, Lafayette, Ind., 1957.
- [9] L. CESARI, J. R. LA PALM AND T. NISHIURA, *Remarks on some existence theorems for optimal control*, J. Optimization Theory Appl., 3 (1969), pp. 296–305.
- [10] P. BRUNOVSKY, *On the necessity of a certain convexity condition for lower closure of control problems*, this Journal, 6 (1968), pp. 174–185.

## TRACKING VIA FEEDBACK FOR SYSTEMS WITH IRRATIONAL TRANSFER FUNCTION\*

M. I. FREEDMAN† AND R. GLASSEY‡

**Abstract.** Given stable convolution filters  $G$  and  $H$ , possibly with irrational transfer functions, find a feedback filter  $C$  such that for all inputs (discrete stationary processes  $(x_n)_0^\infty$  in an ensemble with fixed covariance  $(\Gamma_k)_{-\infty}^\infty$ ) the closed-loop system  $G(I + CG)^{-1}$  will eventually track the open system  $H$  in optimal fashion. The filter  $C$  may be stable or unstable but must satisfy suitable realizability conditions. Focus will be on the case where  $G$  is not of minimum phase. In this case the system must satisfy an additional restriction, called  $M$ -stability, which will exclude the occurrence of "pole-zero" cancellation in the product  $CG$ .

**1. Introduction.** After the fundamental work of Wiener [1], Wiener and Masani [2] and Kolmogorov [3] on the frequency domain approach to prediction and filtering of random processes the viewpoint of many control oriented researchers, following the lead of Kalman–Bucy [4], turned to the time domain. This viewpoint has led to great success in dealing with models described by a system consisting of a finite number of ordinary differential equations influenced by white noise. Application of the Wiener–Hopf technique does not require this "finiteness" condition on the model description, i.e., the transfer function of the operators involved may be irrational.

From one point of view Kalman and Bucy's notion of state estimator, when coupled with the optimal state variable feedback technique (see [5]), amounts to the use of (possibly-unstable) feedback about an open-loop operator in order to cause the overall system to behave in accordance with an a priori prescribed plan.

The applications the authors of this present paper are mainly interested in arise from situations where partial differential equations are involved. As such the transfer functions involved are generally not rational. Thus, a frequency domain approach seems appropriate.

The main problem considered in this paper is heuristically as follows (see Fig. 1):

Given stable convolution filters  $G$  and  $H$ , possibly with irrational transfer functions, find a feedback filter  $C$  such that for all inputs (discrete stationary processes  $(x_n)_0^\infty$  in an ensemble with fixed covariance  $(\Gamma_k)_{-\infty}^\infty$ ) the closed-loop system  $G(I + CG)^{-1}$  will eventually track the open system  $H$  in optimal fashion. The filter  $C$  may be stable or unstable but must satisfy suitable realizability conditions. Focus will be on the case where  $G$  is not of minimum phase. In this case the system must satisfy an additional restriction, called  $M$ -stability, which will exclude the occurrence of "pole-zero" cancellation in the product  $CG$ . See § 2 for details.

\* Received by the editors May 11, 1970, and in revised form October 28, 1970.

† NASA Electronics Research Center, Cambridge, Massachusetts, and Department of Mathematics, Boston University, Boston, Massachusetts 02215.

‡ Division of Applied Mathematics, Brown University, Providence, Rhode Island 02912.



The results found are explicit and computational. They involve the spectral density function  $\Gamma(\theta)$  for the processes  $(x_n)_0^\infty$  and a Hermite interpolation procedure. Conditions are derived for the existence of an optimal  $C$  and a formula for this  $C$  is given. Finally an expression for the “tracking error” is determined.

As a final remark we mention that discrete rather than continuous time systems are considered. This allows the convenience of dealing with functions analytic on the interior of the unit disc rather than in a half-plane. As such, certain technical and notational points are simplified, particularly in § 6 and § 7. The translation of the results to the continuous time situation offers no real obstacle.

**2. General preliminaries.** This short preliminary section begins by briefly reviewing facts about the  $l_p$ -spaces and related power series. A few definitions specific to this paper, such as the notion of  $M$ -stability, are then presented. The section ends with some very basic remarks about discrete time series.

**2.1. Remarks about  $l_p$ -spaces and related series.** For  $1 < p < \infty$ ,  $l_p$  will denote the space of infinite sequences of complex numbers  $\{a_k\}_{k=-\infty}^\infty$  such that  $\sum_{k=-\infty}^\infty |a_k|^p < \infty$  and  $l_p^+$  will denote the subspace of  $l_p$  consisting of those sequences  $\{a_k\}_{k=-\infty}^\infty$  with  $a_k = 0$  for  $k < 0$ .

$l_p$  and  $l_p^+$  are Banach spaces (with the norm  $\|\{a_k\}\|_p = (\sum_{k=-\infty}^\infty |a_k|^p)^{1/p}$ ) and their properties are both standard and well known. In this paper our dealings will be with the cases  $p = 1$  and  $p = 2$ .  $l_1$  and  $l_1^+$  are Banach algebras under convolution “\*”, where  $\{a_k\}_{k=-\infty}^\infty * \{b_k\}_{k=-\infty}^\infty = \{c_k\}_{k=-\infty}^\infty$  with  $c_k = \sum_{m=-\infty}^\infty a_{k-m} b_m$  and with the identity  $I = \{a_k\}_{k=-\infty}^\infty$ , where  $a_0 = 1$ ,  $a_k = 0$  for  $k \neq 0$ . Also  $l_2$  and  $l_2^+$  are Hilbert spaces with the inner product  $\langle \{a_k\}_{k=-\infty}^\infty, \{b_k\}_{k=-\infty}^\infty \rangle = \sum_{k=-\infty}^\infty a_k \bar{b}_k$ . We remark that  $l_1 \subset l_2$  and  $l_1 \subset l_2^+$ .

As a shorthand we shall alternatively denote a sequence  $\{a_k\}_{k=-\infty}^\infty$  by the capital letter  $A$ ,  $\{b_k\}_{k=-\infty}^\infty$  by  $B$ , etc.  $AB$  will refer to  $\{a_k\}_{k=-\infty}^\infty * \{b_k\}_{k=-\infty}^\infty$  where defined. We note that any element  $A \in l_1$  induces a bounded linear map  $A : l_p \rightarrow l_p$  (for any  $p$ ,  $1 \leq p$ ) given by  $A(B) = AB$  for  $B \in l_p$  and it is by virtue of this association that we often call elements of  $l_1$ , operators.

**DEFINITION 2.2.** A sequence  $A = \{a_k\}_{k=-\infty}^\infty$  will be called *causal* if  $a_k = 0$  for  $k < 0$  and *stable* (alternatively *bounded*) if  $A \in l_1$ .

For  $1 \leq p < \infty$ ,  $L_p(d\theta)$  will denote the Banach space of complex-valued measurable functions  $f(e^{i\theta})$  with  $|f|^p$  summable on the unit circle, while for  $p = \infty$ ,  $L_\infty(d\theta)$  will denote the space of the essentially bounded complex-valued unit circle. Recall that since the unit circle represents a finite measure space,  $L_\infty(d\theta) \subset L_2(d\theta) \subset L_1(d\theta)$ .

Given  $A = \{a_k\}_{k=-\infty}^\infty$  in  $l_1$  we shall denote by  $A(e^{i\theta})$  the continuous function on the unit circle with  $k$ th Fourier coefficient equal to  $a_k$ , i.e.,  $A(e^{i\theta}) = \sum_{k=-\infty}^\infty a_k e^{ik\theta}$ , where this series converges absolutely. Further if  $A \in l_1^+$  we define  $A(z) = \sum_{k=0}^\infty a_k z^k$  for  $z$  complex,  $|z| < 1$ . This series converges uniformly absolutely on  $|z| < 1$ .

If  $A = \{a_k\}_{k=-\infty}^\infty$  is in  $l_2$ , the series  $A(e^{i\theta}) = \sum_{k=-\infty}^\infty a_k e^{ik\theta}$  converges in  $L_2(d\theta)$ -norms and so represents a square integrable function on the unit circle. In this case, if  $A \in l_2^+$ , then  $A(z) = \sum_{k=0}^\infty a_k z^k$  converges uniformly absolutely on any compact subset of  $|z| < 1$  and represents an analytic function on  $|z| < 1$ . The

Hilbert subspace of  $L_2(d\theta)$  corresponding to  $A$  in  $l_2^+$  is often denoted by  $H_2$  (Hardy space).

DEFINITION 2.3. For  $A$  an element of either  $l_1^+$  or  $l_2^+$  we shall call  $A(z)$  the  $1/z$ -transform of  $A$ .<sup>1</sup>

For this paper we require a few more facts about the  $l_p$ -spaces ( $p = 1, 2$ ) and their transforms:

(a)  $(AB)(e^{i\theta}) = A(e^{i\theta})B(e^{i\theta})$  a.e. for  $A$  in  $l_1$  and for  $B$  in either  $l_1$  or  $l_2$ .

(b) If  $A, B$  are in  $l_2$ , then  $\langle A, B \rangle = (1/(2\pi)) \int_0^{2\pi} A(e^{i\theta})\overline{B(e^{i\theta})} d\theta$  (the Plancherel theorem).

(c) Given  $A \in l_1$ , then a necessary and sufficient condition for  $A$  to have an inverse in  $l_1^+$  is that  $A(z) \neq 0$  on  $|z| \leq 1$ , and if this is so,  $A^{-1}(z) = 1/A(z)$  (Wiener's theorem, see [6, pp. 72]).

This completes the short summary of standard facts about  $l_p, l_p^+, p = 1, 2$ , which will be needed.

#### 2.4. The notion of model-stability.

DEFINITION 2.5. Given  $G \in l_1^+, C \in l_1$  and  $Y \in l_1^+$  the 3-tuple  $\{G, C; Y\}$  will be termed a *feedback triplet* if  $I + CG$  is invertible in  $l_1$  with  $Y = (I + CG)^{-1}$ .

DEFINITION 2.6. A *feedback triplet*  $\{G, C; Y\}$  will be called *model-stable* (or simply *M-stable*) if given any  $G_1 \in l_1^+$  then for all real  $\varepsilon$  sufficiently small in magnitude,  $Y_\varepsilon \stackrel{\text{def}}{=} (I + C(G + \varepsilon G_1))^{-1}$  remains in  $l_1^+$ .

This definition is of importance in the sequel. In § 8, given a fixed  $G \in l_1^+$ , we desire to minimize a certain integro-quadratic form over an appropriate class of feedback operators  $\mathcal{C}$ . This class  $\mathcal{C}$  must be such as to assure that the "closed-loop system" (i.e., the system defined via feedback, with impulse response  $G(I + CG)^{-1}$ ) corresponds to a legitimate "physical situation." What we require is that not only must  $G(I + CG)^{-1}$  be causal and stable (lie in  $l_1^+$ ) but that under sufficiently small perturbations in  $G$  the new closed-loop system must remain causal and stable. The exact definition of the class  $\mathcal{C}$  must wait until § 4. For now we conclude with the following easily checked remark which motivated Definition 2.6.

Remark 2.7. In the situation  $G \in l_1^+$  and  $C \in l_1$  with  $G(e^{i\theta})$  and  $C(e^{i\theta})$  both representable as rational functions (of the variable  $e^{i\theta}$ ), then *M-stability* is equivalent to the assumption that no "pole-zero" cancellation occurs (in the interior of the unit disc) in the definition of  $Y$  as  $(I + CG)^{-1}$ .

2.8. Remarks on discrete time series. Let  $\Omega$  be a space having a Borel field of subsets over which a probability measure  $P$  is defined. Let  $L_2(\Omega)$  be the set of complex-valued  $P$ -measurable functions  $f$  on  $\Omega$  for which  $\int_\Omega |f(\omega)|^2 dP(\omega) < \infty$ ; then  $L_2(\Omega)$  is a Hilbert space with inner product given by

$$E(f, g) = \int_\Omega f(\omega)\overline{g(\omega)} dP(\omega).$$

<sup>1</sup> In engineering literature the  $z$ -transform of a causal sequence  $\{a_k\}_{k=-\infty}^\infty = A$  is usually defined to be  $\sum_{k=0}^\infty a_k z^{-k}$  making  $A(z)$  analytic in the exterior of the unit disc. For our analysis it was decided that the convenience of dealing with functions analytic on the interior of the disc was great enough to justify the definition of  $A(z)$  as given below. To emphasize the relationship between  $A(z)$  here and the usual  $z$ -transform we call our transform the " $1/z$ -transform."

DEFINITION 2.9. (a) An element  $f$  of  $L_2(\Omega)$  is called a *random variable*.

(b) A sequence  $(x_n)_0^\infty$  with  $x_n$  a random variable for each  $n$  is called a *stochastic process* (S.P.).

(c) A stochastic process  $(x_n)_0^\infty$  is called *stationary* (wide-sense) if for each integer  $k$  with  $k + n \geq 0$ ,  $E(x_{n+k}, x_n) = \Gamma_k$ , where  $\Gamma_k$  is a complex number independent of  $n$ .

(d) A *spectral density function*  $\Gamma(\theta)$  for a stationary S.P. is an a.e. nonnegative function  $\Gamma(\theta)$  defined and integrable on the unit circle which has  $\Gamma_k$  as its  $k$ th Fourier coefficient., i.e.,

$$\Gamma_k = \frac{1}{2\pi} \int_0^{2\pi} e^{-ik\theta} \Gamma(\theta) d\theta \quad \text{for } -\infty < k < \infty.$$

Not every stationary S.P. has a spectral density function, but throughout this paper we hypothesize their existence for the stationary S.P.'s we consider. (In general,  $\Gamma_k = (1/(2\pi)) \int_0^{2\pi} e^{-ik\theta} dF(\theta)$ ,  $-\infty < k < \infty$ , where  $F(\theta)$  is a bounded nondecreasing right-continuous function on  $[0, 2\pi]$  with  $F(0) = 0$ . This  $F(\theta)$  is called the spectral distribution function of the stationary stochastic process and its existence follows from Bochner's theorem (see [6, p. 42]).)

**3. M-stability.** In this section an alternative characterization of  $M$ -stability for a feedback triplet  $\{G, C; Y\}$  is presented (assuming  $G^{-1} \in l_1$  exists). This new description involves explicitly only the causal operators  $G$  and  $Y$  and proves useful in later sections.

A lemma based on the variational nature of the  $M$ -stability definition is considered first.

LEMMA 3.1. *The feedback triplet  $\{G, C; Y\}$  is  $M$ -stable if and only if  $Y^{n+1}C^n \in l_1^+$  for all integers  $n \geq 1$  (true for  $n = 0$  by definition of feedback triplet).*

*Proof.* Let  $G_1$  be in  $l_1^+$  and for  $\varepsilon$  a sufficiently small real number let  $Y_\varepsilon \triangleq (I + C(G + \varepsilon G_1))^{-1}$  be an element of  $l_1$ . For  $\varepsilon$  sufficiently small,  $Y_\varepsilon$  has an  $l_1$ -convergent Taylor expansion, namely:

$$\begin{aligned} (3.2) \quad Y_\varepsilon &= (I + CG)^{-1} - \varepsilon CG_1(I + CG)^{-2} + \varepsilon^2 C^2 G_1^2 (I + CG)^{-3} - \dots \\ &= Y - \varepsilon CY^2 G_1 + \varepsilon^2 C^2 G_1^2 Y^3 - \dots \end{aligned}$$

Now  $\{G, C; Y\}$  will be  $M$ -stable provided that for all sufficiently small  $\varepsilon$ ,  $Y_\varepsilon$  lies in  $l_1^+$ . Clearly this will be the case if  $Y^{n+1}C^n$  is in  $l_1^+$  for  $n \geq 1$ . Conversely let  $G_1$  be the identity  $I$  and assume that for all  $\varepsilon$  sufficiently small the corresponding  $Y_\varepsilon$  lies in  $l_1^+$ . An easy argument shows that  $Y^{n+1}C^n$  lies in  $l_1^+$  for  $n \geq 1$ , as follows:

Consider  $(1/\varepsilon)(Y_\varepsilon - Y)$  which is for each sufficiently small nonzero  $\varepsilon$  an element of  $l_1^+$ . From (3.2) it is clear that as  $\varepsilon \rightarrow 0$  this expression has an  $l_1$ -norm limit equal to  $Y^2C$  (as  $G_1$  is  $I$  here) and so  $Y^2C$  must therefore lie in the subspace  $l_1^+$ . The induction scheme showing  $Y^{n+1}C^n \in l_1^+$  for general  $n \geq 1$  follows these same simple lines.

THEOREM 3.3. *Let  $G^{-1}$  exist in  $l_1$ . Then a necessary and sufficient condition for the feedback triplet  $\{G, C; Y\}$  to be  $M$ -stable is that  $G^{-1}(Y - I)$  lie in  $l_1^+$ .*

The proof to be presented is of a computational nature and depends to an extent on the following.

*Remark 3.4.* Given  $A, B$  in  $l_1^+$  with  $A^{-1} \in l_1$ , then a necessary and sufficient condition for  $A^{-1}B$  to lie in  $l_1^+$  is that  $B(z)/A(z)$  be analytic in  $|z| < 1$ ; i.e., any zero of  $A(z)$  in  $|z| < 1$  of multiplicity  $k \geq 1$  is likewise a zero of  $B(z)$  of at least multiplicity  $k$ .

*Proof of Remark 3.4.* For any integer  $k$ ,

$$(3.5) \quad \int_0^{2\pi} e^{ik\theta} \frac{B(e^{i\theta})}{A(e^{i\theta})} d\theta = \frac{1}{2\pi i} \oint_{|z|=1} z^{k-1} \frac{B(z)}{A(z)} dz.$$

Now  $A^{-1}B \in l_1^+$  if and only if the integral on the left in (3.5) vanishes for each positive  $k$ . But the right-hand integral of (3.5) shows that this is exactly the condition that  $B(z)/A(z)$  be analytic in  $|z| < 1$ .

*Proof of Theorem 3.3. Sufficiency.* Suppose the feedback triplet  $\{G, C; Y\}$  is  $M$ -stable. By Remark 3.4 we must show for each root  $a$  of  $G(z)$ ,  $|a| < 1$ , of multiplicity  $k \geq 1$  that  $(Y-1)^{(j)}(a) = 0$ ,  $j = 0, \dots, k-1$ . Assume at first that  $G(a) = 0$  for a complex,  $|a| < 1$ . We shall show that  $Y(a) = 1$ .

To begin consider the expression  $Y(I + CG) = I$ . It follows that  $Y^2 + Y^2CG = Y$  and hence for  $|z| < 1$ ,

$$Y^2(z) + (Y^2CG)(z) = Y(z)$$

or

$$(3.6) \quad (Y(z))^2 + (Y^2C)(z)G(z) = Y(z)$$

since  $Y^2C \in l_1^+$  and as such has a "1/z-transform."

Since we are assuming  $G(a) = 0$  it follows that  $(Y(a))^2 = Y(a)$  and therefore  $Y(a) = 0$  or  $Y(a) = 1$ . We proceed to show that the possibility  $Y(a) = 0$  is unstable.

Assume therefore that  $Y(a) = 0$ . For  $k \geq 1$ ,

$$\begin{aligned} Y^{k+2}C^{k+1}G &= Y^{k+2}C^k(I + CG) - Y^{k+2}C^k \\ &= Y^{k+1}C^k - Y^{k+2}C^k. \end{aligned}$$

So

$$(3.7) \quad Y^{k+2}C^{k+1}G = (Y^{k+1}C^k)(I - Y).$$

Since  $Y^{k+2}C^{k+1}$  and  $Y^{k+1}C^k$  are both elements of  $l_1^+$  it follows on taking transforms that

$$(3.8) \quad (Y^{k+2}C^{k+1})(z)G(z) = (Y^{k+1}C^k)(z)(1 - Y(z)).$$

Therefore since  $G(a) = 0$  and we are assuming  $Y(a) = 0$ , (3.8) yields  $(Y^{k+1}C^k)(a) = 0$  for  $k \geq 1$ .

Still under the assumption  $Y(a) = 0$  and proceeding by induction we assume

$$Y(a) = Y'(a) = \dots = Y^{(n-1)}(a) = 0$$

and

$$(Y^{k+1}C^k)(a) = (Y^{k+1}C^k)'(a) = \dots = (Y^{k+1}C^k)^{(n-1)}(a) = 0$$

for some  $n \geq 0$  and all  $k \geq 0$ .

Differentiating (3.8)  $n$ -times with respect to  $z$  and using the above induction assumptions yields  $(Y^{k+1}C^k)^{(n)}(a) = 0$  for all  $k \geq 0$ .

Returning to (3.6) we likewise differentiate that expression  $n$  times and evaluate at  $z = a$ . Using the induction assumptions once again yields  $Y^{(n)}(a) = 0$  and completes the induction procedure.

Therefore  $Y^{(n)}(a) = 0$  for  $n = 0, 1, 2, \dots$ . Since  $Y(z)$  is analytic in the unit disc  $|z| < 1$ , it follows that  $Y(z) = 0$  in a neighborhood of the origin and so by continuation in the unit disc. This, of course, contradicts the fact that  $Y$  is the  $l_1$ -inverse of  $I + CG$ .

It follows, therefore, that  $Y(a) = 1$  is the correct conclusion to draw from (3.6).

Assume next that  $G$  has a zero at  $z = a$  of order  $k \geq 2$ , i.e.,  $G(a) = G'(a) = \dots = G^{(k-1)}(a) = 0$ . We know that  $Y(a) = 1$ . We additionally now show that  $Y'(a) = Y''(a) = \dots = Y^{(k-1)}(a) = 0$ .

To proceed, for each  $l, 0 < l \leq k - 1$ , differentiate (3.6)  $l$  times with respect to  $z$  and evaluate at  $z = a$ . Then

$$(3.9) \quad \left. \frac{d^l}{dz^l} (Y^2(z)) \right|_{z=a} = Y^{(l)}(a)$$

since  $G$  has a  $k$ th order zero at  $z = a$ . Now  $2Y(a)Y'(a) = Y'(a)$  shows that  $Y'(a) = 0$  since  $Y(a) = 1$ ; proceeding inductively one concludes easily that  $Y''(a) = \dots = Y^{(k-1)}(a) = 0$  completing this half of the proof.

*Necessity.* Assuming that  $G^{-1}(Y - I) \in l_1^+$  we must show that the feedback triplet  $\{G, C; Y\}$  is  $M$ -stable. The characterization of  $M$ -stability given in Lemma 3.1 proves useful here, namely:  $\{G, C; Y\}$  will be  $M$ -stable if and only if  $Y^{n+1}C^n \in l_1^+$  for all integers  $n \geq 1$ . Now, by definition,  $Y = (I + CG)^{-1}$  so  $YCG = I - Y$  and  $G^{-1} = -YC$ . Therefore  $YC \in l_1^+$  and in general for any  $n \geq 1$ ,  $Y(YC)^n = Y^{n+1}C^n \in l_1^+$ , completing the proof.

*Remark 3.10.* In engineering literature the term “minimum phase” is often applied to the situation where  $G(z) \neq 0$  in  $|z| \leq 1$ . It follows from Theorem 3.3 that given  $G \in l_1^+$  of minimum phase, then for any  $C \in l_1$  and  $Y \in l_1^+$  with  $Y = (I + CG)^{-1}$  the feedback triplet  $\{G, C; Y\}$  is automatically  $M$ -stable, as  $(Y(z) - 1)/G(z)$  is analytic on  $|z| < 1$  in this case. Our interest in this paper will focus on  $G \in l_1^+$  which is *not* of minimum phase.

DEFINITION 3.11. Let  $C \in l_1^+$  with  $G^{-1} \in l_1$  and define

$$\mathcal{W} = \left\{ C \in l_1 \left| \begin{array}{l} Y = (I + CG)^{-1} \text{ exists in } l_1^+ \text{ and the} \\ \text{feedback triplet } \{G, C; Y\} \text{ is } M\text{-stable} \end{array} \right. \right\}.$$

The subset  $\mathcal{W}$  will find use in § 8 on minimization.

**4. Quasi-extendability.** In certain situations we shall wish to associate causal but possibly unstable sequences with stable but possibly noncausal ones (see Definition 2.2). The notion of quasi-extendability is pertinent in this regard.

DEFINITION 4.1. A causal sequence  $C = \{c_k\}_{k=-\infty}^{\infty}$  (with  $c_k = 0, k < 0$ ) is called *quasi-extendable* if the power series  $C(z) = \sum_{k=0}^{\infty} c_k z^k$  is analytic on  $|z| < \varepsilon, \varepsilon$  sufficiently small, and  $C(z)$  has a meromorphic extension denoted  $C^E(z)$  to  $|z| \leq 1$  such that  $C^E(e^{i\theta})$  is continuous for  $\theta \in [0, 2\pi]$  and is the Fourier series of some  $C^E$  in  $l_1$ .  $C^E$  will be called the *quasi-extension* of  $C$ .

Example 4.2. Let  $c$  satisfy  $|c| < 1$ . Define  $B = \{b_k\}_{k=-\infty}^{\infty}$  by

$$b_k = \begin{cases} -c^{-(k+1)}, & k \geq 0, \\ 0, & k \leq -1. \end{cases}$$

Then

$$B(z) = \sum_{k=0}^{\infty} b_k z^k = -\frac{1}{c} - \frac{1}{c^2 z} - \frac{1}{c^2 z^2} - \dots,$$

so  $B(z) = (z - c)^{-1}$  for  $|z| < |c|$ . Let  $A = \{a_k\}_{k=-\infty}^{\infty}$  in  $l_1$  be defined by

$$a_k = \begin{cases} 0, & k \geq 0, \\ c^{-(k+1)}, & k \leq -1. \end{cases}$$

The  $(z - c)^{-1}$  has a meromorphic extension to  $|z| \leq 1$  and it is clear that  $B^E = A$ .

In § 8 we shall consider a minimization problem taken over a class of operators occurring in feedback. These operators will necessarily be causal due to physical constraints but need not necessarily be stable (bounded). More specifically, we shall have to deal with operators that are quasi-extendable. To help explain the sort of feedback equations involved we proceed with the following definition and theorem.

DEFINITION 4.3. Let  $G \in l_1^+$  with  $G(e^{i\theta})$  nonzero for  $\theta \in [0, 2\pi]$ . Define

$$\mathcal{C} = \{C | C \text{ is causal and quasi-extendable to } C^E \text{ and } C^E \in \mathcal{W}^{\wedge}\}.$$

Now let  $\{x_n\}_0^{\infty}$  lie in  $l_2^+$ . Consider the pair of feedback equations for  $\{e_n\}_0^{\infty}$ ,  $\{f_n\}_0^{\infty}$  in  $l_2^+$  given, for  $n \geq 0$ , by

$$(4.4) \quad \begin{aligned} e_n &= x_n - \sum_{j=0}^n c_{n-j} f_j, \\ f_n &= \sum_{l=0}^n g_{n-l} e_l. \end{aligned}$$

Symbolically we write (see Fig. 1)

$$\mathbf{e} = \mathbf{x} - \mathbf{Cf}, \quad \mathbf{f} = \mathbf{G}\mathbf{e},$$

where  $\mathbf{e} = \{e_n\}_{n=0}^{\infty}$ ,  $\mathbf{f} = \{f_n\}_{n=0}^{\infty}$  and  $\mathbf{x} = \{x_n\}_{n=0}^{\infty}$ .

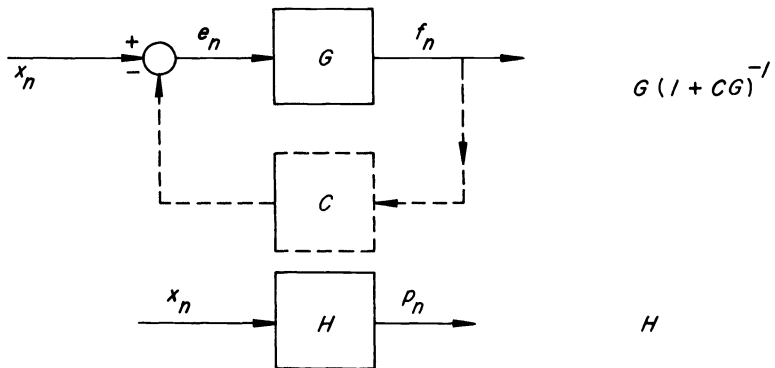


FIG. 1

**THEOREM 4.5.** Assume feedback equations (4.6) relating  $\mathbf{e}, \mathbf{f}, \mathbf{x} \in l_2^+$ . Then if  $C = \{c_n\}_{n=-\infty}^\infty$  lies in  $\mathcal{C}$  it follows that  $\mathbf{e}$  and  $\mathbf{f}$  are uniquely specified by  $\mathbf{x}$  and, in fact, there exists  $D = \{d_k\}_{k=-\infty}^\infty \in l_1^+$  such that  $\mathbf{e} = D\mathbf{x}$  or  $e_n = \sum_{k=0}^n d_{n-k}x_k$  (i.e., the closed-loop system  $(I + CG)^{-1}$  is causal and stable). Further if  $G \in l_1^+$  is replaced by  $G^* \in l_1^+$  sufficiently close to  $G$  in  $l_1$ -norm, then there will likewise exist  $D^* \in l_1$  with  $\mathbf{e}^* = D^*\mathbf{x}$ , where  $\mathbf{e}^*$  satisfies the modified (4.4) ( $\mathbf{x}$  fixed and  $G$  replaced by  $G^*$ ).

*Proof.* Let  $D = \{d_k\}_{k=-\infty}^\infty = (I + C^E G)^{-1}$ . The rest of the proof follows from the definition of  $\mathcal{C}$ .

**Remark 4.6.** Theorem 4.5 suggests that  $\mathcal{C}$  is an appropriate class of causal but possibly unstable operators for use in feedback loops. The operators in  $\mathcal{C}$  produce closed-loop systems and do so in a physically acceptable way, in that the overall system remains insensitive to slight inaccuracies in the open-loop model.

In the remainder of this section we prove a result concerning quasi-realizability.

**THEOREM 4.7.** Let  $\{G, C; Y\}$  be an  $M$ -stable feedback triplet. Then if  $Y(0) \neq 0$ ,  $C$  is the quasi-extension of some causal sequence  $A = \{b_k\}_{k=-\infty}^\infty$ , i.e.,  $C = A^E$ .

*Proof.*  $Y(z)$  being analytic on  $|z| < 1$ , continuous and nonzero on  $|z| = 1$  can have only a finite number of zeros within the unit disc. As such,  $Y(z)$  may be factored into a product  $B(z)H(z)$  valid for  $|z| \leq 1$ . Here  $H(z)$  is a nonzero analytic function on the disc, arising from an element  $H \in l_1^+$ , while  $B(z)$  is given by the finite product  $B(z) = z^p \prod_{i=1}^n (z - \alpha_i)^{k_i}$  corresponding to the roots of  $Y(z)$  in  $|z| < 1$  with appropriate multiplicity. However in this case  $p = 0$ , as  $Y(0) \neq 0$ , i.e.,  $B(z) = \prod_{i=1}^n (z - \alpha_i)^{k_i}$ . Now let  $B$  be that element of  $l_1^+$  corresponding to  $B(e^{i\theta})$ . It is clear that  $B^{-1}$  exists in  $l_1$  and from Example 4.2 one sees that  $B^{-1}$  is the quasi-extension of a causal operator (being the product of causal operators). But  $H$  is invertible in  $l_1^+$  since  $H(z) \neq 0$  for  $|z| \leq 1$ . So  $H^{-1} \in l_1^+$  is its own quasi-extension. Similarly,  $G^{-1}(I - Y) \in l_1^+$  is its own quasi-extension. The result now follows as  $C = G^{-1}(I - Y)B^{-1}H^{-1}$ .

In keeping with the above we make the following definition.

**DEFINITION 4.8.** Assume  $G \in l_1^+$  with  $G^{-1} \in l_1$ . Then let  $\mathcal{W}$  be as in § 3 and define  $\mathcal{W}_0 = \{C \in \mathcal{W} \mid Y(0) \neq 0\}$ . (Note that if  $G(0) = 0$ , then  $Y(0) = 1$  and so  $\mathcal{W}_0 = \mathcal{W}$ .)

**5. A density theorem.**  $\mathcal{W}$  and  $\mathcal{W}_0$  have previously been defined. The following additional definitions are also needed in this section.

**DEFINITION 5.1.** Assume  $G \in l_1^+$  with  $G^{-1} \in l_1$  and define  $T = \{Y \in l_2^+ \mid G^{-1}(Y - I) \in l_2^+\}$  and  $T_0 = \{Y \in T \mid Y(0) \neq 0\}$ .

Note that  $T_0 = T$  if  $G(0) = 0$ .

**ASSUMPTION 5.2.** From this point on throughout this paper we shall assume, without always explicitly stating so, that the symbol  $G$  represents an operator in  $l_1^+$  with  $G(z)$  nonzero on  $|z| = 1$ , i.e., with  $G^{-1}$  lying in  $l_1$ .

**Remark 5.3.**  $T_0$  is dense in  $T$  (in the  $l_2$ -induced topology).

**Remark and Notation 5.4.** For  $A = \{a_k\}_{k=-\infty}^\infty, A \in l_p^+, 1 \leq p \leq \infty$ , and for any  $r, 0 < r < 1$ , define  $A_r$  in  $l_1^+ \subset l_2^+$  by

$$A_r = \begin{cases} a_k r^k, & k \geq 0, \\ 0, & k < 0. \end{cases}$$

Then the Fourier series for  $A_r, A_r(e^{i\theta})$  equals  $A(re^{i\theta})$ , where  $A(z), |z| \leq 1$ , represents

as usual the  $1/z$ -transform of  $A$ . Now if  $p = 2$ , then as  $r \rightarrow 1$ ,  $A_r \rightarrow A$  in  $l_2$ -norm (as follows from Parseval's theorem, since  $A_r(e^{i\theta}) \rightarrow A(e^{i\theta})$  in  $L_2(d\theta)$ -norm; see [7, p. 32]). Note also that if  $A \in l_1^+$ , then as  $r \rightarrow 1$ ,  $A_r \rightarrow A$  in  $l_1$ -norm (by the Lebesgue bounded convergence theorem).

**THEOREM 5.5.**  $(I + \mathcal{W}G)^{-1}$  is dense in  $T$  in the  $l_2$ -induced topology.

*Proof.* Let  $Y \in T$ . Then there exists a sequence  $\{r_k\}$  of numbers,  $0 < r_k < 1$ , with  $r_k \rightarrow 1$  such that  $Y_{r_k}$  has an  $l_1$ -inverse for each  $r_k$ . This is so because the condition under which a given  $Y_r$  will fail to have such an inverse is that  $Y(z)$  has a zero on  $|z| = r$ .  $Y(z)$ , being analytic, may only have a zero on countably many such circles.

Fixing on this sequence  $\{r_k\}$ , define  $C_{r_k} = G_{r_k}^{-1}(I - Y_{r_k})Y_{r_k}^{-1}$ . The  $C_{r_k}$  all lie in  $l_1$ . The feedback triplets  $\{G_{r_k}, C_{r_k}; Y_{r_k}\}$  are all  $M$ -stable since  $G_{r_k}^{-1}(I - Y_{r_k}) = (G^{-1}(I - Y))_{r_k} \in l_1^+$  as seen from the assumption  $Y \in T$  and Remark 5.4.

To complete the proof one must check that the  $C_{r_k}$  lie in  $\mathcal{W}$  for  $r_k$  sufficiently close to 1 and that as  $r_k \rightarrow 1$ ,  $(I + C_{r_k}G)^{-1}$  tends to  $Y$  with respect to the  $l_2$ -norm. We omit the lengthy but straightforward computations.

**COROLLARY 5.6.** (a)  $(I + \mathcal{W}_0G)^{-1}$  is dense in  $T_0$  with respect to the  $l_2$ -induced topology.

(b)  $(I + \mathcal{W}G)^{-1}$  is dense in  $T \cap l_1^+$  with respect to the  $l_1$ -induced topology.

(c)  $(I + \mathcal{W}_0G)^{-1}$  is dense in  $T_0 \cap l_1^+$  (which is in turn dense in  $T \cap l_1^+$ ) with respect to the  $l_1$ -induced topology.

Essentially the same computations as in Theorem 5.5 are needed.

**6. Factorization of the spectral density function  $\theta$ .** Let  $\Gamma(\theta)$  be any nonnegative function Lebesgue integrable on the unit circle. In § 8 such a  $\Gamma(\theta)$  will represent the spectral density function of a stationary (invariant) stochastic process. We define

$$L_2(\Gamma) = \left\{ Y(e^{i\theta}) \left| \begin{array}{l} Y(e^{i\theta}) \text{ is measurable on the unit circle} \\ \text{and is square integrable with respect} \\ \text{to the measure } \Gamma(\theta) d\theta \end{array} \right. \right\}$$

and also  $L_2^+(\Gamma)$  = the closure in  $L_2(\Gamma)$  of trigonometric polynomials of the form  $P(e^{i\theta}) = \sum_{k=0}^n a_k e^{ik\theta}$ .

A classical theorem of prediction theory (attributable to Szegő, Kolmogorov and Krein in various forms; see [7]) states that for  $\Gamma(\theta)$  as given,

$$(6.1) \quad \exp \left\{ \frac{1}{2\pi} \int_0^{2\pi} \log \Gamma(\theta) d\theta \right\} = \min_{Y(e^{i\theta}) \in L_2^+(\Gamma)} \left\{ \frac{1}{2\pi} \int_0^{2\pi} |e^{i\theta} Y(e^{i\theta}) - 1|^2 \Gamma(\theta) d\theta \right\}.$$

Furthermore, the theorem also states that if the right-hand side of (6.1) is zero (perfect mean square prediction for lag 1), then  $(1/(2\pi)) \int_0^{2\pi} \log \Gamma(\theta) d\theta = -\infty$  and conversely, i.e., the equality (6.1) remains true in this case.

It is also known (see Doob [8, p. 577]) that  $(1/(2\pi)) \int_0^{2\pi} \log \Gamma(\theta) d\theta > -\infty$  is a necessary and sufficient condition for  $\Gamma(\theta)$  to be expressible as the absolute value square of an  $H^2$ -function, nonzero for  $|z| < 1$ , i.e., as the product of an  $H_2$ -function and its conjugate.

Our concern will be with the somewhat unorthodox minimization problem described below.



First let  $T' = \{Y(e^{i\theta}) \in L_2^+(\Gamma) | [G(e^{i\theta})]^{-1}[Y(e^{i\theta}) - 1] \in L_2^+(\Gamma)\}$ . Here  $G(e^{i\theta})$  is as in Assumption 5.2. We shall consider the problem

$$\min \left\{ \frac{1}{2\pi} \int_0^{2\pi} |G(e^{i\theta})Y(e^{i\theta}) - V(e^{i\theta})|^2 \Gamma(\theta) d\theta | Y(e^{i\theta}) \in T' \right\},$$

where  $V(e^{i\theta})$  lies in  $L_2^+(\Gamma)$ .

This minimization problem is somewhat different from the one of item (6.1) and arises from an attempt to utilize feedback in accomplishing a certain other minimization procedure described in § 8. The restriction  $Y \in T'$  will amount to assuring (see Theorem 3.3) that minimization takes place only over  $M$ -stable feedback triplets.

In this present section we restrict ourselves to proving the following analogue of the Szegő–Kolmogorov–Krein theorem quoted above.

**THEOREM 6.2.** *Let  $G \in l_1^+$  with  $G(z) \neq 0$  on  $|z| = 1$  and let  $V(e^{i\theta}) \in L_2^+(\Gamma)$ . If the problem*

$$(6.3) \quad \min \left\{ \frac{1}{2\pi} \int_0^{2\pi} |G(e^{i\theta})Y(e^{i\theta}) - V(e^{i\theta})|^2 \Gamma(\theta) d\theta | Y(e^{i\theta}) \in T' \right\}$$

has a nonzero minimum, then  $(1/(2\pi)) \int_0^{2\pi} \log \Gamma(\theta) d\theta > -\infty$ , i.e.,  $\Gamma(\theta)$  factors into the product of an  $H^2$ -function, nonzero for  $|z| < 1$ , and its conjugate.

*Remark 6.4.* Note that the integrability of  $\log \Gamma(\theta)$  does not in our case imply that (6.3) is nonzero. In fact, the latter will certainly be zero whenever  $G(e^{i\theta})Y(e^{i\theta}) = V(e^{i\theta})$  has a solution  $Y(e^{i\theta}) \in T'$  and this depends in no essential way on  $\log \Gamma(\theta)$ . Note also that (6.3) actually does have a minimum and not simply an infimum as  $G^{-1} \in l_1$ .

The proof of Theorem 6.2 will be preceded by the following.

**LEMMA 6.5.** *Let  $\{K_j\}_{j=0}^m$  and  $\{L_j\}_{j=1}^m$  be complex numbers and assume that*

$$K(\theta) = \sum_{j=0}^m K_j \cos j\theta + \sum_{j=1}^m L_j \sin j\theta \geq 0 \quad \text{for all } \theta \in [0, 2\pi].$$

Then

$$(6.6) \quad \int_0^{2\pi} \log K(\theta) d\theta > -\infty.$$

*Proof.* The equation  $K(\theta) = 0$  has only a finite number of isolated roots in  $[0, 2\pi]$ . Clearly if this number is zero, then (6.6) holds. Otherwise the only way that (6.6) might possibly fail to hold true would be that for at least one of these roots  $\theta_0$  and any arbitrary small neighborhood  $N$  of  $\theta_0$ ,  $\int_N |\log K(\theta)| d\theta = \infty$ . That this possibility actually fails to occur can be seen as follows:

Given  $\theta_0$ , a root of  $K(\theta) = 0$ , and any sufficiently small neighborhood  $N$  of  $\theta_0$  there exists a positive integer  $n$  and a function  $Q(\theta)$  continuous on  $[0, 2\pi]$  such that  $K(\theta) = |\theta - \theta_0|^n Q(\theta)$  and  $Q(\theta) \neq 0$  for  $\theta \in N$ . Let  $\delta$  be the diameter of  $N$  and let  $M$  be an upper bound on  $|\log Q(\theta)|$  for  $\theta \in N$ . Then

$$\begin{aligned} \int_N |\log K(\theta)| d\theta &\leq n \int_{\theta_0 - \delta}^{\theta_0 + \delta} \log |\theta - \theta_0| d\theta + 2M\delta \\ &\leq 2n \int_0^\delta \log \xi d\xi + 2M\delta < \infty. \end{aligned}$$

*Proof of Theorem 6.2.* To begin we note that the map  $f : T' \rightarrow L_2^+(\Gamma)$  given by  $f(Y(e^{i\theta})) = [G(e^{i\theta})]^{-1}(Y(e^{i\theta}) - 1)$  is one-to-one and onto and so the minimum of (6.3) must be equal to

$$(6.7) \quad \min \left\{ \frac{1}{2\pi} \int_0^{2\pi} |G(e^{i\theta})(1 + G(e^{i\theta})u(e^{i\theta}) - V(e^{i\theta}))|^2 \Gamma(\theta) d\theta \mid u(e^{i\theta}) \in L_2^+(\Gamma) \right\}.$$

Now (6.7) corresponds to the problem of finding the unique element  $w^*(e^{i\theta})$  in the closed submanifold  $S = \{w(e^{i\theta}) = G^2(e^{i\theta})u(e^{i\theta}) \mid u(e^{i\theta}) \in L_2^+(\Gamma)\}$  of  $L_2^+(\Gamma)$  which is closest (in  $L_2(\Gamma)$ -norm) to the element  $v(e^{i\theta}) - G(e^{i\theta})$  of  $L_2^+(\Gamma)$ . Clearly this element  $w^*(e^{i\theta})$  determines uniquely  $u^*(e^{i\theta})$  in  $L_2^+(\Gamma)$  with  $G^2(e^{i\theta})u^*(e^{i\theta}) = w^*(e^{i\theta})$ .

From the projection theorem for Hilbert space it follows that a necessary and sufficient condition for  $u^*(e^{i\theta})$  to be the element of  $L_2^+(\Gamma)$  minimizing (5.7) is that

$$(6.8) \quad \int_0^{2\pi} [G^2(e^{i\theta})u^*(e^{i\theta}) - (v(e^{i\theta}) - G(e^{i\theta}))]\overline{G^2(e^{i\theta})\eta(e^{i\theta})}\Gamma(\theta) d\theta = 0$$

for all  $\eta(e^{i\theta}) \in L_2^+(\Gamma)$ , i.e., that  $G^2(e^{i\theta})u^*(e^{i\theta}) - (v(e^{i\theta}) - G(e^{i\theta}))$  is orthogonal to  $S$ .

Translating back to the minimization problem (6.3) over  $T'$  we see that (6.3) must have a unique minimum for some  $Y(e^{i\theta}) \in T'$  and that a necessary and sufficient condition for  $Y^*(e^{i\theta})$  to be this minimal element is that

$$(6.9) \quad \int_0^{2\pi} [G(e^{i\theta})Y^*(e^{i\theta}) - V(e^{i\theta})]\overline{G^2(e^{i\theta})\eta(e^{i\theta})}\Gamma(\theta) d\theta = 0$$

for all  $\eta(e^{i\theta}) \in L_2^+(\Gamma)$ .

To proceed with the proof we note that since  $G \in l_1^+$  with  $G(z) \neq 0$  on  $|z| = 1$  it follows that  $G(z)$  has only a finite number  $p$  (counting multiplicity) of roots in  $|z| < 1$ . We may factor  $G(z)$  into a product  $B(z)H(z)$  (as was done to  $Y(z)$  in the proof of Theorem 4.7), where  $B(z)$  is a polynomial of degree  $p$  with roots and multiplicities corresponding to those of  $G(z)$ , while  $H(z)$  is continuous and nonzero on the closed disc and analytic on the open disc. Let

$$\eta(e^{i\theta}) = [G(e^{i\theta})Y^*(e^{i\theta}) - V(e^{i\theta})][H(e^{i\theta})]^{-2}B^2(e^{i\theta})e^{2pi\theta}r(e^{i\theta}),$$

where  $r(e^{i\theta})$  is an arbitrary element of  $L_2^+(\Gamma)$ . Then  $\eta(e^{i\theta})$  is in  $L_2^+(\Gamma)$  since the negative Fourier coefficients of  $[H(e^{i\theta})]^{-2}$  vanish as well as the negative Fourier coefficients of  $B^2(e^{i\theta})e^{2pi\theta}$ .

Using this  $\eta(e^{i\theta})$  in (6.9) gives

$$(6.10) \quad \int_0^{2\pi} |G(e^{i\theta})Y^*(e^{i\theta}) - V(e^{i\theta})|^2 |B(e^{i\theta})|^4 \Gamma(\theta) e^{2pi\theta} r(e^{i\theta}) d\theta = 0$$

for  $r(e^{i\theta}) \in L_2^+(\Gamma)$ . Making the choices  $e^{ik\theta}$ ,  $k \geq 0$ , for  $r(e^{i\theta})$  in (6.10) yields

$$\int_0^{2\pi} |G(e^{i\theta})Y^*(e^{i\theta}) - V(e^{i\theta})|^2 |B(e^{i\theta})|^4 \Gamma(\theta) e^{-il\theta} d\theta = 0$$

for  $l \geq 2p$  and by conjugation the same is true for  $l \leq -2p$ . From these orthogonality conditions it follows that

$$(6.11) \quad |G(e^{i\theta})Y^*(e^{i\theta}) - V(e^{i\theta})|^2 |B(e^{i\theta})|^4 \Gamma(\theta) = K(\theta)$$

for some  $K(\theta) = \sum_{j=0}^p k_j \cos j\theta + \sum_{j=1}^p L_j \sin j\theta$  with  $K(\theta)$  nonnegative (the  $K_j$  and  $L_j$  are real). Define

$$S(e^{i\theta}) = \frac{(G(e^{i\theta})Y^*(e^{i\theta}) - V(e^{i\theta}))^{-1}K(\theta)e^{-i\theta}}{B^2(e^{i\theta})}.$$

We claim that  $S(e^{i\theta})$  lies in  $H_2$ . First  $|S(e^{i\theta})|^2 = K(\theta)\Gamma(\theta)$  as follows from (6.11) and therefore  $S(e^{i\theta}) \in L_2(d\theta)$  since  $\Gamma(\theta) \in L_1(d\theta)$  (while of course  $K(\theta) \in L_\infty(d\theta)$ ). Now  $S(e^{i\theta})$  will lie in  $H_2$  if

$$\int_0^{2\pi} S(e^{i\theta})e^{ik\theta} d\theta = 0 \quad \text{for } k = 1, 2, \dots.$$

To see this we conjugate (6.9) to find that for any  $\eta(e^{i\theta})$  in  $L_2^+(\Gamma)$ ,

$$(6.12) \quad \int_0^{2\pi} \overline{[G(e^{i\theta})Y^*(e^{i\theta}) - V(e^{i\theta})]}G^2(e^{i\theta})\Gamma(\theta)\eta(e^{i\theta}) d\theta = 0.$$

Now for  $k \geq 1$ ,

$$(6.13) \quad \int_0^{2\pi} S(e^{i\theta})e^{ik\theta} d\theta = \int_0^{2\pi} \frac{[G(e^{i\theta})Y^*(e^{i\theta}) - V(e^{i\theta})]K(\theta)}{B^2(e^{i\theta})} e^{i(k-1)\theta} d\theta$$

as seen from using (6.11). But (6.13) is zero for  $k \geq 1$  as it corresponds to (6.12) with  $\eta(e^{i\theta}) = [H(e^{i\theta})]^{-2}e^{i(k-1)\theta}$  which is in  $L_2^+(\Gamma)$  for each  $k \geq 1$ . Thus  $S(e^{i\theta})$  is a nonzero function in  $H_2$ . Now the well-known Jensen inequality of classical complex variable theory (see [7, p. 68]) assures us that

$$\int_0^{2\pi} \log |S(e^{i\theta})| d\theta > -\infty.$$

But since  $|S(e^{i\theta})|^2 = K(\theta)\Gamma(\theta)$  it now follows that  $\int_0^{2\pi} \log K(\theta) d\theta + \int_0^{2\pi} \log \Gamma(\theta) d\theta > -\infty$ . By Lemma 6.5,  $\log K(\theta)$  is integrable and so  $\int_0^{2\pi} \log \Gamma(\theta) d\theta > -\infty$ . Hence also  $\Gamma(\theta)$  factors.

**7. A frequency domain minimization problem.** On the basis of the results in § 6 we shall assume in this section that the nonnegative function  $\Gamma(\theta)$ , Lebesgue integrable on the unit circle, also satisfies the condition  $\int_0^{2\pi} \log \Gamma(\theta) d\theta > -\infty$ . Let  $R(e^{i\theta})$  denote the  $H_2$ -function, with  $R(z)$  nonzero for  $|z| < 1$  such that  $\Gamma(\theta) = |R(e^{i\theta})|^2$  a.e.

We note that the integrability of  $\log \Gamma(\theta)$  assures that  $R(e^{i\theta}) \neq 0$  a.e. in  $[0, 2\pi]$ . Also the Hilbert spaces  $H_2$  and  $L_2^+(\Gamma)$  are isometric under the linear map  $Y(e^{i\theta}) \rightarrow Y(e^{i\theta})[R(e^{i\theta})]^{-1}$ .

In this section we shall investigate properties of the minimization problem

$$(7.1) \quad \min \left\{ \frac{1}{2\pi} \int_0^{2\pi} |G(e^{i\theta})Y(e^{i\theta}) - V(e^{i\theta})|^2 \Gamma(\theta) d\theta \mid Y(e^{i\theta}) \in T' \right\},$$

where  $V(e^{i\theta})$  in  $L_2^+(\Gamma)$  is fixed. (That this problem actually has a minimum was seen in § 6.)

Here  $G$  once again represents an operator in  $l_1^+$  with  $G(z)$  nonzero on  $|z| = 1$ . As such,  $G(z)$  can have only a finite number  $P$  of roots in  $|z| < 1$ . As in § 6 we may

write  $G(z) = B(z)H(z)$ , where  $H(z)$  is a continuous nonvanishing function on the closed unit disc analytic on the open disc (which is the  $1/z$ -transform of an element  $H$  in  $l_1^+$ ) while  $B(z) = \prod_{j=1}^n (z - a_j)^{k_j}$ , where  $\{a_j\}_{j=1}^n$ ,  $a_j$  complex, and  $\{k_j\}_{j=1}^n$ ,  $k_j$  nonnegative integers, represent the roots of  $G(z)$  in  $|z| < 1$  and their multiplicities respectively (including possible roots at  $z = 0$ ).

Letting  $k(e^{i\theta}) = V(e^{i\theta})R(e^{i\theta})$  we have  $k(e^{i\theta}) \in H_2$ . It is convenient to consider the minimization problem (7.1) in the alternate form

$$(7.2) \quad \min \left\{ \frac{1}{2\pi} \int_0^{2\pi} |G(e^{i\theta})s(e^{i\theta}) - k(e^{i\theta})|^2 d\theta \mid s(e^{i\theta}) \in H_2 \text{ with } s(e^{i\theta})[R(e^{i\theta})]^{-1} \in T' \right\}.$$

The advantage is that here we deal with underlying Hilbert spaces  $L_2(d\theta)$  and  $H_2$ . The minimum value attained by (7.2) is clearly the same as the minimum of (7.1), and if  $s^*(e^{i\theta})$  is the point where (7.2) attains its minimum then  $Y^*(e^{i\theta}) = s^*(e^{i\theta})[R(e^{i\theta})]^{-1}$  is the point in  $T'$  where (7.1) does so.

A necessary and sufficient condition for an element  $s(e^{i\theta})$  of  $H_2$  to satisfy  $s(e^{i\theta})[R(e^{i\theta})]^{-1} \in T'$  is that  $s(a_j) = R(a_j)$ ,  $s^1(a_j) = R^1(a_j)$ ,  $\dots$ ,  $s^{(k_j-1)}(a_j) = R^{(k_j-1)}(a_j)$  for  $j = 1, \dots, n$ , i.e., that  $f(z) = (s(z) - R(z))/G(z)$  is analytic in  $|z| < 1$ .

A procedure based on Hermite interpolation will be exhibited for finding the optimal  $s^*(e^{i\theta})$ . In this direction we first prove the following.

**THEOREM 7.3.** *A necessary and sufficient condition for the minimization problem (7.2) to attain its minimum at a function  $s^*(e^{i\theta})$  in  $H_2$  with  $s^*(e^{i\theta})[R(e^{i\theta})]^{-1}$  in  $T'$  is that there exist a trigonometric polynomial of order  $2P - 1$  or less,  $M(e^{i\theta}) = \sum_{k=0}^{2P-1} m_k e^{ik\theta}$ , such that*

$$(7.4) \quad [G(e^{i\theta})s^*(e^{i\theta}) - k(e^{i\theta})]\overline{B^2(e^{i\theta})}e^{2Pi\theta} = M(e^{i\theta}) \quad a.e.$$

*Proof. Sufficiency.* In § 6, (6.9) is a necessary and sufficient condition for (7.1) to have its minimum at  $Y^*(e^{i\theta})$  in  $T'$ . Translating (6.9) to the corresponding condition for (7.2) yields that a necessary and sufficient condition for (7.2) to have its minimum at  $s^*(e^{i\theta})$  in  $H_2$  with  $s^*(e^{i\theta})[R(e^{i\theta})]^{-1}$  in  $T'$  is that

$$(7.5) \quad \int_0^{2\pi} [G(e^{i\theta})s^*(e^{i\theta}) - k(e^{i\theta})]\overline{G^2(e^{i\theta})\eta(e^{i\theta})} d\theta = 0$$

for all  $\eta(e^{i\theta})$  in  $H_2$ .

Letting  $\eta(e^{i\theta}) = [H(e^{i\theta})]^{-2}r(e^{i\theta})$  for  $r(e^{i\theta})$  an arbitrary element of  $H_2$ , we may rewrite (7.5) as

$$(7.6) \quad \int_0^{2\pi} [G(e^{i\theta})s^*(e^{i\theta}) - k(e^{i\theta})]\overline{B^2(e^{i\theta})}e^{2Pi\theta}\overline{e^{2Pi\theta}r(e^{i\theta})} d\theta = 0$$

for all  $r(e^{i\theta})$  in  $H_2$ . Now the term  $M(e^{i\theta}) = [G(e^{i\theta})s^*(e^{i\theta}) - k(e^{i\theta})]\overline{B^2(e^{i\theta})}e^{2Pi\theta}$  must lie in  $H_2$ , as  $B^2(e^{i\theta})e^{2Pi\theta}$  has all negative Fourier coefficients zero. From (7.6) it follows that

$$(7.7) \quad \int_0^{2\pi} M(e^{i\theta})e^{-ik\theta} d\theta = 0 \quad \text{for all } k \geq 2P$$

as seen by making the appropriate choices for  $r(e^{i\theta})$ . But (7.7) together with the fact that  $M(e^{i\theta})$  is in  $H_2$  implies that  $M(e^{i\theta})$  is a trigonometric polynomial of degree  $\leq 2P - 1$ , as said in the statement of the theorem.

Conversely, if  $[G(e^{i\theta})s^*(e^{i\theta}) - k(e^{i\theta})]B^2(e^{i\theta})e^{2Pi\theta} = M(e^{i\theta})$  for some trigonometric polynomial  $M(e^{i\theta})$  of degree  $\leq 2P - 1$ , then (7.7) must hold, and since the  $\{e^{ik\theta}\}_{k=0}^\infty$  span  $H_2$ , (7.6) also must hold. The choice  $r(e^{i\theta}) = [H(e^{i\theta})]^2\eta(e^{i\theta})$  then yields (7.5), the previously established necessary and sufficient condition for  $s^*(e^{i\theta})$  to be the minimizing function for (7.2).

*Remark 7.8.* (a) For a given  $s^*(e^{i\theta})$  in  $H_2$  the validity of (7.4) for some  $M(e^{i\theta})$  a trigonometric polynomial of degree  $\leq 2P - 1$  combined with the conditions  $s^*(a_j) = R(a_j), \dots, s^{*(k_j-1)}(a_j) = R^{(k_j-1)}(a_j)$  for  $j = 1, \dots, n$  suffices to determine  $s^*(e^{i\theta})$  uniquely. This is true because (7.2) has a unique minimum. The procedure by which these conditions determine  $s^*(e^{i\theta})$  will be spelled out in Theorem 7.11 below.

(b) For notational convenience let  $B_1(e^{i\theta}) = \overline{B(e^{i\theta})}e^{Pi\theta} = \prod_{i=1}^n (1 - \bar{a}_i e^{i\theta})^{k_i}$  for  $\theta \in [0, 2\pi]$  and similarly  $B_1(z) = \prod_{i=1}^n (1 - \bar{a}_i z)^{k_i}$  for  $|z| \leq 1$ . Now  $B_1(z)$  so given and likewise  $[B_1(z)]^{-1}$  are continuous nonvanishing functions on  $|z| \leq 1$  which are analytic on  $|z| < 1$ . The condition  $[G(e^{i\theta})s^*(e^{i\theta}) - k(e^{i\theta})]B_1^2(e^{i\theta}) = M(e^{i\theta})$  a.e. implies that for  $|z| < 1$ ,

$$(7.9) \quad [G(z)s^*(z) - k(z)]B_1^2(z) = M(z).$$

*Definition and Remark 7.10.* Let  $F \in L_1^+$  with  $F(e^{i\theta}) \neq 0, \theta \in [0, 2\pi]$ . Then  $F(z)$  has a finite number  $q$  (counting multiplicity) of roots in  $|z| < 1$ . Let  $\{b_j\}_{j=1}^n, \{l_j\}_{j=1}^n$  be the roots and their multiplicities respectively. For any function  $W(e^{i\theta})$  in  $H_2, (W(z))_F$  (read  $W(z)$  mod  $F$ ) will denote that unique polynomial in  $z$  of order  $\leq q - 1$  such that  $(W(z) - (W(z))_F)/F(z)$  is analytic in  $|z| < 1$ .

Clearly  $(W(z))_F \equiv 0$  if  $F(z)$  has no zeros in  $|z| < 1$ . In general,  $(W(z))_F$  represents a *Hermite interpolation* with respect to conditions specified by the roots of  $F(z)$  and as such  $(W(z))_F$  depends only on the values  $W(b_j), W^1(b_j), \dots, W^{(l_j-1)}(b_j), j = 1, \dots, n$ . An explicit construction of  $(W(z))_F$  is as follows:

Let  $v(z) = \prod_{i=1}^n (z - b_i)^{l_i}$  and define

$$r_{ik}(z) = v(z) \frac{(z - z_i)^{k-l_i}}{k!} \frac{d^{(l_i-1-k)}}{dz^{(l_i-1-k)}} \left[ \frac{(z - z_i)^{l_i}}{v(z)} \right]_{z=z_i}$$

for  $1 \leq i \leq n, 0 \leq k \leq l_i - 1$ . Then

$$(W(z))_F = \sum_{i=1}^n [W(b_i)r_{i0}(z) + W^1(b_i)r_{i1}(z) + \dots + W^{(l_i-1)}(b_i)r_{i,l_i-1}(z)].$$

(This formula appears in [9].)

**THEOREM 7.11.** *The function  $s^*(e^{i\theta})$  for which the variational problem (7.2) attains its minimum is given by*

$$s^*(e^{i\theta}) = \frac{L(e^{i\theta}) + B_1^2(e^{i\theta})k(e^{i\theta})}{G(e^{i\theta})B_1^2(e^{i\theta})} \quad \text{for } \theta \in [0, 2\pi],$$

where  $L(z)$  is that polynomial of degree  $\leq 2P - 1$  given by  $L(z) = ([G(z)R(z) - k(z)]B_1^2(z))_{G^2}$ .

*Proof.*  $s^*(e^{i\theta})$  as defined clearly lies in  $L_2(d\theta)$ . Actually  $s^*(e^{i\theta})$  also lies in  $H_2$ . This can be checked by showing that  $[L(z) + B_1^2(z)k(z)]/G(z)B_1^2(z)$  is analytic in  $|z| < 1$ . But by the definition of  $L(z)$ ,

$$(7.12) \quad \frac{[G(z)R(z) - k(z)]B_1^2(z) - L(z)}{G^2(z)}$$

is analytic in  $|z| < 1$ .

Multiplying (7.12) by the analytic function  $G(z)[B_1(z)]^{-2}$  shows that

$$R(z) - \frac{L(z) + B_1^2(z)k(z)}{G(z)B_1^2(z)}$$

is analytic in  $|z| < 1$ , and since  $R(z)$  is analytic in  $|z| < 1$  it follows that  $(L(z) + B_1^2(z)k(z))/G(z)B_1^2(z)$  is, i.e.,  $s^*(e^{i\theta})$  is in  $H_2$ .

Clearly then  $[G(z)s^*(z) - k(z)]B_1^2(z) = L(z)$  for  $|z| \leq 1$  and substituting for  $L(z)$  in (7.12) yields that  $[R(z) - s^*(z)]B_1^2(z)/G(z)$  is analytic in  $|z| < 1$ . Multiplying this last expression by the analytic function  $[B_1(z)]^{-2}$  shows that  $(R(z) - s^*(z))/G(z)$  is analytic in  $|z| < 1$ . Hence  $s^*(e^{i\theta})[R(e^{i\theta})]^{-1}$  is in  $T'$ . The hypothesis (7.4) now holds with  $M(e^{i\theta}) = L(e^{i\theta})$  and so Theorem 7.3 assures that  $s^*(e^{i\theta})$  so defined is the optimal for variational problem (7.2).

**THEOREM 7.13.** *The minimizing element  $Y^*(e^{i\theta})$  in  $T'$  for variational problem (7.1) is given by*

$$Y^*(e^{i\theta}) = \frac{L(e^{i\theta}) + B_1^2(e^{i\theta})V(e^{i\theta})R(e^{i\theta})}{R(e^{i\theta})G(e^{i\theta})B_1^2(e^{i\theta})}, \quad \theta \in [0, 2\pi],$$

where  $L(z)$  is that polynomial of degree  $\leq 2P - 1$  given by  $L(z) = ((G(z) - V(z))R(z)B_1^2(z))_{G^2}$ .

*Proof.*  $R(z)Y^*(z) = s^*(z)$  and  $R(z)V(z) = k(z)$ .

**Remark 7.14.** (a) Suppose  $R(e^{i\theta})$  derives from  $R \in l_1^+$  with  $R^{-1}$  also in  $l_1^+$  (i.e., with  $R(z) \neq 0$  on  $|z| \leq 1$ ). Then if  $V(e^{i\theta})$  likewise comes from  $V \in l_1^+$  it follows that the minimizing function  $Y^*(e^{i\theta})$  for problem (7.1) must also arise from some  $Y^* \in l_1^+$ .

(b) Alternatively, suppose there exist positive constants  $K_1$  and  $K_2$  with  $0 < K_1 \leq \Gamma(\theta) \leq K_2$  a.e. Then  $R(e^{i\theta})$  and  $[R(e^{i\theta})]^{-1}$  both lie in  $L_\infty(d\theta) \cap H_2$  and if  $V(e^{i\theta}) \in L_\infty(d\theta) \cap H_2$  it follows that  $Y^*(e^{i\theta}) \in L_\infty(d\theta) \cap H_2$ .

(c) Alternatively, if  $R(z)$  is rational in  $z$  and  $V(z)$  is rational in  $z$ , then  $Y^*(z)$  will be a rational function in  $z$ .

The minimum value of the variational problems (7.1) or (7.2) which we shall henceforth denote by  $J_{\min}$  can now be computed from  $Y^*(e^{i\theta})$  or  $s^*(e^{i\theta})$ , for

$$\begin{aligned} J_{\min} &= \frac{1}{2\pi} \int_0^{2\pi} |G(e^{i\theta})Y^*(e^{i\theta}) - V(e^{i\theta})|^2 \Gamma(\theta) d\theta \\ &= \frac{1}{2\pi} \int_0^{2\pi} |G(e^{i\theta})s^*(e^{i\theta}) - k(e^{i\theta})|^2 d\theta. \end{aligned}$$

There is an alternative often more attractive computational method based on contour integration and residue theory. This method does not involve preliminary

computation of either  $Y^*(e^{i\theta})$  or  $s^*(e^{i\theta})$  (or for that matter computation of  $L(z)$ ).

To begin note that from (7.4) one has

$$J_{\min} = \frac{1}{2\pi} \int_0^{2\pi} \left| \frac{L(e^{i\theta})}{B_1^2(e^{i\theta})} \right|^2 d\theta,$$

where  $L(e^{i\theta})$  is the trigonometric polynomial of degree  $\leq 2P - 1$  defined in Theorem 7.11. Now  $B_1(e^{i\theta}) = e^{Pi\theta} B(e^{i\theta})$ . So  $|B_1(e^{i\theta})| = |B(e^{i\theta})|$  for  $\theta \in [0, 2\pi]$ . It follows that

$$J_{\min} = \frac{1}{2\pi} \int_0^{2\pi} \left| \frac{L(e^{i\theta})}{B^2(e^{i\theta})} \right|^2 d\theta.$$

The function  $L(e^{i\theta})[B(e^{i\theta})]^2$  is a rational (trigonometric) function which has a Laurent expansion in  $e^{i\theta}$  with only negative powers, i.e., all nonnegative Fourier coefficients vanish.

Letting  $L(e^{i\theta})[B(e^{i\theta})]^{-2} = \sum_{l=1}^{\infty} s_l e^{-il\theta}$  we have

$$s_k = \frac{1}{2\pi i} \oint_{|z|=1} \frac{L(z)}{B^2(z)} dz \quad \text{for } k = 1, 2, \dots,$$

and so by Parseval's theorem,  $J_{\min} = \sum_{k=1}^{\infty} |s_k|^2$ .

The calculation of the  $s_k$  as indicated by the above contour integral seems to involve prior computation of the polynomial  $L(z)$ . It is this computation we shall be able to avoid. In this direction consider the equation  $[G(z)s^*(z) - k(z)]B_1^2(z) = L(z)$  valid for  $|z| \leq 1$ . Multiplying through by  $z^{k-1}[B(z)]^{-2}$  for any  $k \geq 1$  and integrating around the unit circle gives

$$(7.15) \quad 2\pi i s_k = \oint_{|z|=1} z^{k-1} \frac{s^*(z)G(z)B_1^2(z)}{B^2(z)} dz - \oint_{|z|=1} z^{k-1} \frac{k(z)B_1^2(z)}{B^2(z)} dz$$

for  $k = 1, 2, \dots$ . Recalling that  $(s^*(z) - R(z))/G(z) = A(z)$  for some analytic function  $A(z)$  on  $|z| < 1$  (actually  $A(e^{i\theta}) \in H_2$ ) we have that

$$(7.16) \quad \oint_{|z|=1} z^{k-1} \frac{s^*(z)G(z)B_1^2(z)}{B^2(z)} dz = \oint_{|z|=1} z^{k-1} \frac{R(z)G(z)B_1^2(z)}{B(z)} dz + \oint_{|z|=1} z^{k-1} \frac{A(z)G^2(z)B_1^2(z)}{B^2(z)} dz$$

for  $k = 1, 2, \dots$ . Now the second integral on the right in (7.16) vanishes, as the function  $A(z)G^2(z)B_1^2(z)/B(z)$  is analytic in  $|z| < 1$ .

Therefore (7.15) may be rewritten as

$$(7.17) \quad 2\pi i s_k = \oint_{|z|=1} z^{k-1} \frac{R(z)G(z)B_1^2(z)}{B^2(z)} dz - \oint_{|z|=1} z^{k-1} \frac{k(z)B_1^2(z)}{B^2(z)} dz.$$

For convenience we now let  $T(z)$  be the rational function given by

$$T(z) = \frac{B(z)}{B_1(z)} = \prod_{i=1}^n \left( \frac{z - a_i}{1 - \bar{a}_i z} \right)^{k_i}.$$

Then making a substitution in (7.17) we have proved the following theorem.

**THEOREM 7.18.** *The minimum value taken by the variational problems (7.1) or (7.2) is given by*

$$J_{\min} = \sum_{k=1}^{\infty} |s_k|^2,$$

where

$$\begin{aligned} s_k &= \frac{1}{2\pi i} \oint_{|z|=1} z^{k-1} \frac{[R(z)G(z) - k(z)]}{[T(z)]^2} dz \\ &= \frac{1}{2\pi i} \oint_{|z|=1} z^{k-1} \frac{[R(z)(G(z) - V(z))]}{[T(z)]} dz. \end{aligned}$$

Here

$$T(z) = \prod_{i=1}^n \left( \frac{z - a_i}{1 - \bar{a}_i z} \right)^{k_i},$$

the  $\{a_i\}_{i=1}^n, \{k_i\}_{i=1}^n$  being the roots of  $G(z)$  in  $|z| < 1$  and their multiplicities respectively.

*Remark 7.19.* The rational analytic function  $T(z)$  (on  $|z| < 1$ ) is a (finite) Blaschke product for  $G(z)$  in the sense of [7, p. 66]; i.e.,  $|T(e^{i\theta})| = 1$  for  $\theta \in [0, 2\pi]$ ,  $G(z)/T(z)$  is analytic in  $|z| < 1$  and  $T(z)$  is a finite product of bilinear functions in  $z$ .

**8. Minimization via feedback.** Let  $(x_n)_0^\infty$  be a stationary (wide sense) stochastic process with a given known covariance sequence  $(\Gamma_n)_{-\infty}^\infty$ . Let  $G = \{g_n\}_{-\infty}^\infty$  be in  $l_1^+$ . Assume  $G(e^{i\theta}) = \sum_{n=0}^\infty g_n e^{in\theta}$  is nonzero for  $\theta \in [0, 2\pi]$ . As in § 4 define

$$\mathcal{C} = \left\{ C = \{c_k\}_{k=-\infty}^\infty \left| \begin{array}{l} \text{(i) } C \text{ is causal and quasi-extendable,} \\ \text{(ii) } Y^E = (1 + C^E G)^{-1} \text{ exists in } l_1^+, \\ \text{(iii) } \{G, C; Y\} \text{ is an } M\text{-stable feedback triplet} \end{array} \right. \right\}.$$

Consider the feedback system defining (nonstationary) stochastic processes  $(e_n)_0^\infty$  and  $(f_n)_0^\infty$  given by the equations

$$(8.1) \quad e_n = x_n - \sum_{j=0}^n c_{n-j} f_j \quad \text{and} \quad f_n = \sum_{l=0}^n g_{n-l} e_l,$$

where  $C = \{c_k\}_{k=-\infty}^\infty$  is an element of  $\mathcal{C}$ . Here we regard (8.1) as a *closed-loop* input-output system with  $(x_n)_0^\infty$  to be viewed as the *input* and  $(f_n)_0^\infty$  as the *output*. (See Fig. 1.)

Also let  $H = \{h_k\}_{k=-\infty}^\infty$  with  $h_k$  real for each  $h$  be a fixed element of  $l_1^+$ . We define the (nonstationary) stochastic process  $(p_n)_0^\infty$  by the equation

$$(8.2) \quad p_n = \sum_{j=0}^n h_{n-j} x_j.$$

Here  $(x_n)_0^\infty$  is viewed as the *input* and  $(p_n)_0^\infty$  as the *output* of the system (*open loop*) described by (8.2) (see again Fig. 1).

The problem treated in this section is that of selecting a *feedback operator*  $C$  in  $\mathcal{C}$  which most nearly (in a statistical sense) causes the system of (8.1) to *track* the system of (8.2), i.e., to find that  $C$  in  $\mathcal{C}$ , when it exists, which minimizes  $\lim_{n \rightarrow \infty} \mathcal{E}([f_n - p_n]^2)$ . This limit will be shown to exist.



Let us in what follows emphasize the dependence of the  $f_n$  on  $C$  by writing  $f_n(C)$ .

The problem just briefly described may be rephrased in a frequency domain setting. We denote by  $\Gamma(\theta)$  the spectral density function (which we assume to exist) for the S.P.  $(x_n)_0^\infty$ . As in § 2,  $\Gamma(\theta)$  is an a.e. defined nonnegative function integrable on the unit circle.

LEMMA 8.3. *Let  $C \in \mathcal{C}$  be given. Assume the S.P.'s  $(C_n)_0^\infty$  and  $(f_n(C))_0^\infty$  satisfy (8.1) while the S.P.  $(p_n)_0^\infty$  satisfies (8.2). Then  $\lim_{n \rightarrow \infty} \mathcal{E}([f_n(C) - p_n]^2)$  exists, and denoting this limit by  $\mathcal{E}^*(C)$ , we have*

$$\mathcal{E}^*(C) = \frac{1}{2\pi} \int_0^{2\pi} \left| \frac{G(e^{i\theta})}{1 + C^E(e^{i\theta})G(e^{i\theta})} - H(e^{i\theta}) \right|^2 \Gamma(\theta) d\theta,$$

where  $C^E \in l_1$  is the quasi-extension of  $C$ . (See Definition 4.1.)

*Proof.* We begin by showing that if  $\{k_n\}_{n=-\infty}^\infty$  is any element of  $l_1^+$  and  $(m_n)_0^\infty$ , the S.P., is given by  $m_n = \sum_{j=0}^n k_{n-j}x_j$ , then

$$(8.4) \quad \lim_{n \rightarrow \infty} \mathcal{E}([m_n]^2) = \frac{1}{2\pi} \int_0^{2\pi} |K(e^{i\theta})|^2 \Gamma(\theta) d\theta,$$

where  $K(e^{i\theta}) = \sum_{n=0}^\infty k_n e^{in\theta}$ .

To see this note that

$$\mathcal{E}([m_n]^2) = \sum_{j=0}^n \sum_{l=0}^n k_{n-j}k_{n-l}E(x_j, x_l).$$

But

$$\mathcal{E}(x_j, x_l) = \Gamma_{j-l} = \frac{1}{2\pi} \int_0^{2\pi} e^{-i(j-l)\theta} \Gamma(\theta) d\theta \quad \text{for any } j, l \geq 0.$$

Therefore,

$$\begin{aligned} \mathcal{E}([m_n]^2) &= \sum_{j=0}^n \sum_{l=0}^n \left\{ \frac{1}{2\pi} \int_0^{2\pi} k_{n-j}e^{-ij\theta}k_{n-l}e^{il\theta} \Gamma(\theta) d\theta \right\} \\ &= \sum_{j=0}^n \sum_{l=0}^n \left\{ \frac{1}{2\pi} \int_0^{2\pi} k_{n-j}e^{i(n-j)\theta}k_{n-l}e^{-i(n-l)\theta} \Gamma(\theta) d\theta \right\} \\ &= \frac{1}{2\pi} \int_0^{2\pi} |K_n(e^{i\theta})|^2 \Gamma(\theta) d\theta, \end{aligned}$$

where  $k_n(e^{i\theta}) = \sum_{j=0}^n k_j e^{ij\theta}$ . Therefore letting  $n \rightarrow \infty$  gives the result.

Now by a simple modification of Theorem 4.5 to handle S.P.'s as inputs and outputs it follows that  $e_n = \sum_{j=0}^n d_{n-j}x_j$ , where  $\{d_k\}_{k=-\infty}^\infty \in l_1^+$  with  $\{d_k\}_{k=-\infty}^\infty = Y^E = (I + C^E G)^{-1}$ . Therefore since  $f_n = \sum_{l=0}^n g_{n-l}e_l$  and  $p_n = \sum_{l=0}^n k_{n-l}x_l$  we have that

$$f_n - p_n = \sum_{j=0}^n \left( \sum_{l=0}^n g_{n-l}d_{l-j} - k_{n-j} \right) x_j.$$

But clearly  $G(I + C^E G)^{-1} - H$  is that element of  $l_1^+$  given by  $\{s_n\}_{n=0}^\infty$  with  $s_n = \sum_{l=0}^n g_{n-l} d_l - k_n$  for each  $n \geq 0$ . The conclusion of the lemma therefore follows by employing (8.4) with  $(m_n)_0^\infty = (f_n - p_n)_0^\infty$ .

*Remark 8.5.* For any integer  $n \geq 0$ ,

$$(8.6) \quad \inf_{C \in \mathcal{C}} \mathcal{E}^*(C) = \inf_{\substack{C(e^{i\theta}), \\ C \in \mathcal{W}_0}} \left\{ \frac{1}{2\pi} \int_0^{2\pi} \left| \frac{G(e^{i\theta})}{1 + C(e^{i\theta})G(e^{i\theta})} - H(e^{i\theta}) \right|^2 \Gamma(\theta) d\theta \right\}.$$

*Proof.* The proof is immediate.

*Remark 8.7.* In the case where an actual minimum is attained by the left-hand side of (8.6) at some  $A \in \mathcal{C}$ , then the right-hand side of (8.6) will also have a minimum and will realize it at the quasi-extension  $A^E$  of  $A$ ,  $A^E$  lying in  $\mathcal{W}_0$ . Conversely, if a minimum of the right-hand side of (8.6) is attained at some  $C \in \mathcal{W}_0$ , then the left-hand side of (8.6) will attain a minimum and it will be at that unique  $C_*$  in  $\mathcal{C}$  with  $C_*^E = C$ .

In fact,  $C_* = \{c_k^*\}_{k=0}^\infty$ , where

$$c_k^* = \frac{1}{2\pi i} \oint_{|z|=\varepsilon} z^{-(k+1)} C(z) dz$$

provided  $\varepsilon > 0$  is sufficiently small (so that  $|z| \leq \varepsilon$  lies within the domain of analyticity of  $C(z)$ ).

Utilizing the density theorem of § 4 we may now state the following theorem.

**THEOREM 8.8.** *For any  $C \in \mathcal{C}$  assume that S.P.'s  $(e_n)_0^\infty$  and  $(f_n(C))_0^\infty$  satisfy (8.1) while S.P.  $(p_n)_0^\infty$  satisfies (8.2). Then the following quantities are equal:*

- (i)  $\inf_{C \in \mathcal{C}} \mathcal{E}^*(C)$ ,
- (ii)  $\inf_{\substack{C(e^{i\theta}), \\ C \in \mathcal{W}_0}} \frac{1}{2\pi} \int_0^{2\pi} \left| \frac{G(e^{i\theta})}{1 + C(e^{i\theta})G(e^{i\theta})} - H(e^{i\theta}) \right|^2 \Gamma(\theta) d\theta$ ,
- (iii)  $\inf_{\substack{C(e^{i\theta}), \\ C \in \mathcal{W}'}} \frac{1}{2\pi} \int_0^{2\pi} \left| \frac{G(e^{i\theta})}{1 + C(e^{i\theta})G(e^{i\theta})} - H(e^{i\theta}) \right|^2 \Gamma(\theta) d\theta$ ,
- (iv)  $\inf_{\substack{Y(e^{i\theta}), \\ Y \in T \wedge l_1^+}} \frac{1}{2\pi} \int_0^{2\pi} |G(e^{i\theta})Y(e^{i\theta}) - H(e^{i\theta})|^2 \Gamma(\theta) d\theta$ ,
- (v)  $\min_{Y(e^{i\theta}) \in T'} \frac{1}{2\pi} \int_0^{2\pi} |G(e^{i\theta})Y(e^{i\theta}) - H(e^{i\theta})|^2 \Gamma(\theta) d\theta$ ,

where, as in § 6,

$$T' = \{Y(e^{i\theta}) \in L_2^+(\Gamma) [G(e^{i\theta})]^{-1} (Y(e^{i\theta}) - 1) \in L_2^+(\Gamma)\}.$$

*Proof.* By Remark 8.5, (i) = (ii). Corollary 5.8 states that  $(I + \mathcal{W}_0 G)^{-1}$  and  $(I + \mathcal{W}' G)^{-1}$  are both dense in  $T \cap l_1^+$  with respect to the  $l_1$ -induced topology. This is enough to ensure that (ii) = (iii) = (iv). Finally  $\{Y(e^{i\theta}) | Y \in T \cap l_1^+\}$  is dense in  $T'$  with respect to the  $L_2(\Gamma)$ -induced topology as seen from the definition of  $L_2(\Gamma)$ . This shows that (iv) = (v), completing the proof.

We now easily show that if  $\inf_{C \in \mathcal{C}} \mathcal{E}^*(C)$  is strictly positive (*imperfect tracking*), the function by  $\Gamma(\theta)$  must necessarily be integrable on the unit circle.

**THEOREM 8.9.** *For any  $C \in \mathcal{C}$  assume that S.P.'s  $(e_n)_0^\infty$  and  $(f_n(C))_0^\infty$  satisfy (8.1) while S.P.  $(p_n)_0^\infty$  satisfies (8.2). Then if  $\inf_{C \in \mathcal{C}} \mathcal{E}^*(C) > 0$  (this expression is independent of  $n, n \geq 0$ ) it follows that  $\int_0^{2\pi} \log \Gamma(\theta) d\theta > -\infty$ .*

*Proof.* By using Theorem 8.8 it follows under the hypothesis of this theorem that

$$\min_{Y(e^{i\theta}) \in T'} \left\{ \frac{1}{2\pi} \int_0^{2\pi} |G(e^{i\theta})Y(e^{i\theta}) - H(e^{i\theta})|^2 \Gamma(\theta) d\theta \right\} > 0.$$

But this is exactly the condition needed in order to call Theorem 6.2 into play. By that theorem,  $\int_0^{2\pi} \log \Gamma(\theta) d\theta > -\infty$ .

We are mainly interested in the case  $\inf \mathcal{E}^*(C) > 0$  (*imperfect tracking*). Without actually making this assumption we shall still assume integrability of  $\log \Gamma(\theta)$ .

*Assumption and Remark 8.10.* Following the indication of Theorem 8.9 we assume from this point on that  $\int_0^{2\pi} \log \Gamma(\theta) d\theta > -\infty$ . Equivalently,  $\Gamma(\theta) = |R(e^{i\theta})|^2$  a.e., where  $R(e^{i\theta})$  is an  $H_2$ -function with  $R(z)$  nonzero on  $|z| < 1$ . We shall however additionally assume in the remainder of this section that  $R(e^{i\theta})$  is the Fourier series for some  $R \in l_1^+$  with  $R(e^{i\theta}) \neq 0$  for  $\theta \in [0, 2\pi]$ . This assumption yields the following immediate consequences:

- (i)  $L_2(\Gamma) = L_2(d\theta)$ ,  $L_2^+(\Gamma) = H_2$  and  $T' = \{Y(e^{i\theta}) | Y \in T\}$ .
- (ii)  $L_2(\Gamma)$  is isometric with  $L_2(d\theta)$  and  $L_2^+(\Gamma)$  is isometric with  $H_2$  under the map  $Y(e^{i\theta}) \rightarrow Y(e^{i\theta})R(e^{i\theta})$ .
- (iii) The set

$$(8.11) \quad \left\{ \frac{1}{2\pi} \int_0^{2\pi} |G(e^{i\theta})Y(e^{i\theta}) - H(e^{i\theta})|^2 \Gamma(\theta) d\theta \mid Y(e^{i\theta}) \in T' \right\}$$

attains a minimum at some  $Y^*(e^{i\theta})$  corresponding to  $Y \in l_1^+$ , i.e., the inf in (iv) of Theorem 8.8 may be changed to a min.

Note that (iii) follows directly from Remark 7.14(a) since we are assuming that  $H(e^{i\theta})$  is the Fourier series for  $H \in l_1^+$ .

*Remark 8.12.* The assumptions in (8.10) are made in order to simplify certain details in assembling our results. The type of problem we are dealing with has one inherent untidiness as follows:

Minimization problems involving integrated quadratic forms are almost always taken over spaces of square integrable functions (Hilbert spaces). However, feedback relations involving convolution operators are most readily expressible for “ $l_1$ -operator algebras” (Banach algebras). More specifically, without the assumptions on  $R(e^{i\theta})$  of (8.10) the element  $Y_*(e^{i\theta})$  in  $T'$  minimizing (8.11) need not be the Fourier series of any  $Y_*$  in  $l_1^+$ . However, if this  $Y_*(e^{i\theta})$  is to arise as  $(1 + C(e^{i\theta})G(e^{i\theta}))^{-1}$  for some  $C(e^{i\theta})$  related to  $C \in \mathcal{W}$ , then  $Y_*(e^{i\theta})$  would have to lie in  $l_1^+$ .

To conclude this remark we point out that another way around this difficulty exists based on part (b) of Remark 7.14. To take this route would have involved

modification in our definition of  $M$ -stability, quasi-realizability, etc., so as to be able to handle  $L_2$ -filters.

We are now ready to state our main theorem.

**THEOREM 8.13.** *Let  $(x_n)_0^\infty$  be a stationary stochastic process with a given known covariance sequence  $(\Gamma_n)_{-\infty}^\infty$  and such that the corresponding spectral density function  $\Gamma(\theta)$  satisfies the conditions of Assumption 8.10. For any  $C \in \mathcal{C}$  let  $(e_n(C))_0^\infty$  and  $(f_n(C))_0^\infty$  denote S.P.'s satisfying feedback equations (8.1) and let  $(p_n)_0^\infty$  be the S.P. satisfying (8.2). Let the polynomial  $L(z)$  be given by the interpolation procedure:*

$$L(z) = ((G(z) - H(z))R(z)B_1^2(z))_{G^2},$$

where  $B_1(z)$  and  $(\cdot)_{G^2}$  are as defined in § 7. Denoting  $L(z) + R(z)H(z)B_1^2(z)$  by  $D(z)$  for  $|z| \leq 1$  we find the following statements are true:

(a) *If  $D(e^{i\theta}) \neq 0$  for any  $\theta \in [0, 2\pi]$  and if  $D(0) \neq 0$ , define*

$$(8.14) \quad C_\#(e^{i\theta}) = \frac{R(e^{i\theta})G(e^{i\theta})B_1^2(e^{i\theta}) - L(e^{i\theta}) - R(e^{i\theta})H(e^{i\theta})F_1^2(e^{i\theta})}{G(e^{i\theta})(L(e^{i\theta}) + R(e^{i\theta})H(e^{i\theta})B_1^2(e^{i\theta}))}.$$

Then  $C_\#(e^{i\theta})$  so defined is the Fourier series of  $C_\#$  in  $\mathcal{W}_0$  and

$$\begin{aligned} & \frac{1}{2\pi} \int_0^{2\pi} |G(e^{i\theta})[1 + C_\#(e^{i\theta})G(e^{i\theta})]^{-1} - H(e^{i\theta})|^2 \Gamma(\theta) d\theta \\ &= \min_{\substack{C(e^{i\theta}), \\ C \in \mathcal{W}_0}} \left\{ \frac{1}{2\pi} \int_0^{2\pi} |G(e^{i\theta})[1 + C(e^{i\theta})G(e^{i\theta})]^{-1} - H(e^{i\theta})|^2 \Gamma(\theta) d\theta \right\}, \end{aligned}$$

i.e., (ii) (and hence (iii)) of Theorem 8.8 attains an actual minimum at this  $C_\#$ .

Further let  $C_* \in \mathcal{C}$  be given by  $C_* = \{c_n^*\}_{n=0}^\infty$ , where

$$c_n^* = \frac{1}{2\pi i} \int_{|z|=\varepsilon} z^{-(n+1)} C(z) dz$$

provided  $\varepsilon > 0$  is sufficiently small (see Remark 8.7). Then  $\mathcal{E}^*(C_*) = \min_{C \in \mathcal{C}} \mathcal{E}^*(C)$ , i.e., (i) of Theorem 8.8 attains an actual minimum at  $C_* \in \mathcal{C}$ .

(b) Under the hypotheses of the theorem it follows that

$$\inf_{C \in \mathcal{C}} \mathcal{E}^*(C) = \sum_{k=1}^\infty |s_k|^2,$$

where

$$s_k = \frac{1}{2\pi i} \oint_{|z|=1} z^{k-1} \frac{[R(z)(G(z) - H(z))]}{[T(z)]^2} dz.$$

Here as in Theorem 7.18,

$$T(z) = \prod_{i=1}^n \left( \frac{z - a_i}{1 - \bar{a}_i z} \right)^{k_i},$$

the  $\{a_i\}_{i=1}^n$ ,  $\{k_i\}_{i=1}^n$  being the roots of  $G(z)$  in  $|z| < 1$  and their multiplicities, respectively. ( $T(z)$  is a finite Blaschke product for  $G(z)$ ; see [7, p. 66].)

*Proof.* (a) Let

$$Y_{\#}(e^{i\theta}) = \frac{L(e^{i\theta}) + B_1^2(e^{i\theta})H(e^{i\theta})R(e^{i\theta})}{R(e^{i\theta})G(e^{i\theta})B_1^2(e^{i\theta})} \quad \text{for } \theta \in [0, 2\pi].$$

Then from Theorem 7.13,  $Y_{\#}(e^{i\theta})$  lies in  $T'$  and  $Y_{\#}(e^{i\theta})$  minimizes (v) of Theorem 8.8. Also by Remark 7.14(a),  $Y_{\#}(e^{i\theta})$  is the Fourier series for some  $Y_{\#} \in l_1^+$  and so  $Y_{\#} \in T \cap l_1^+$  and minimizes (iv) of Theorem 8.8.

Now assuming that  $D(e^{i\theta}) \neq 0$  for  $\theta \in [0, 2\pi]$  is equivalent to stating that  $Y_{\#}(e^{i\theta}) \neq 0$  for  $\theta \in [0, 2\pi]$  so that  $Y_{\#}^{-1}$  exists in  $l_1^+$ . Then  $C_{\#}(e^{i\theta})$  as given in (8.14) equals  $[G(e^{i\theta})]^{-1}[1 - Y_{\#}(e^{i\theta})][Y_{\#}(e^{i\theta})]^{-1}$  so that  $C_{\#} \in l_1$ . Now  $Y_{\#} = (I + C_{\#}G)^{-1}$  and the feedback triplet  $\{G, C_{\#}; Y_{\#}\}$  is  $M$ -stable as  $Y_{\#} \in T$ . Therefore  $C_{\#} \in \mathcal{W}$  and  $C_{\#}$  clearly minimizes (iii) of Theorem 8.8.

If, additionally,  $D(0) \neq 0$  or equivalently  $Y_{\#}(0) \neq 0$ , then Theorem 4.7 assures that  $C_{\#} \in \mathcal{W}_0$  and part (a) of this theorem follows directly.

(b) The proof of (b) is immediate from Theorem 7.18 and the fact that (i) and (v) of Theorem 8.8 are equal.

*Remark 8.15.* In the cases where  $D(e^{i\theta}) = 0$ , some  $\theta \in [0, 2\pi]$  or  $D(0) = 0$ , there is no  $C \in \mathcal{C}$  minimizing  $\mathcal{E}^*(C)$  but minimizing sequences can be constructed. These are mainly of theoretical interest and we omit the details.

*Remark 8.16.* Part (a) of Theorem 8.13 represents an explicit formula for the solution of our minimization via feedback problem (under the conditions  $D(e^{i\theta}) \neq 0$  for  $\theta \in [0, 2\pi]$  and  $D(0) \neq 0$ ) and this together with the minimum mean square error formula (see Theorem 8.13 (b)) constitute our main results.

#### REFERENCES

- [1] N. WIENER, *The Extrapolation, Interpolation and Smoothing of Stationary Time Series*, MIT Press, Cambridge, Mass., 1950.
- [2] N. WIENER AND P. MASANI, *The prediction theory of multivariate stochastic processes, Parts I, II*, Acta Math., 98 (1957), pp. 111–150; 99 (1958), pp. 93–137.
- [3] A. KOLMOGOROV, *Stationary sequences in Hilbert space*, Bull. Math. Univ. Moscow, 2 (1941), no. 6, 44 pp.
- [4] R. E. KALMAN AND R. S. BUCY, *New results in linear filtering and prediction theory*, Trans. ASME Ser. D J. Basic Engrg., 83 (1961), pp. 95–108.
- [5] W. M. WONHAM, *Random differential equations in control theory*, Probabilistic Methods in Applied Math, vol. II, A. T. Bharucha-Reid, ed., Academic Press, New York, 1969.
- [6] L. H. LOOMIS, *An Introduction to Abstract Harmonic Analysis*, Van Nostrand, New York, 1953.
- [7] K. HOFFMAN, *Banach Spaces of Analytic Functions*, Prentice-Hall, Englewood Cliffs, N.J., 1962.
- [8] J. L. DOOB, *Stochastic Processes*, John Wiley, New York, 1953.
- [9] P. J. DAVIS, *Interpolation and Extrapolation*, Blaisdell, New York, 1963.

## GEOMETRIC THEORY OF TIME-OPTIMAL CONTROL\*

OTOMAR HÁJEK†

**Abstract.** Properties of the minimal-time function are considered for the problem of reaching the origin within (linear autonomous finite-dimensional) control systems. The function's continuity is first established, leading to the construction of optimal feedback controls for normal systems. For systems where the control matrix has maximal rank (= dimension of state space), the minimal-time function is locally Lipschitzian, and hence differentiable almost everywhere.

**1. Introduction.** In this paper we will consider the control system

$$(S) \quad \dot{x} = Ax + Bu,$$

where  $A, B$  are constant  $n \times n$  and  $n \times m$  matrices respectively, and the controls  $u$  are restricted by  $|u_i| \leq 1$  (where  $u' = (u_1, \dots, u_m)$ ; primes denote transposition). We will be concerned with the problem of reaching 0 from  $x \in R^n$  in minimal time, i.e., of finding (time-)optimal solutions.

This section merely reviews terminology and notation. In §2 we prove a fundamental theorem on continuity of the minimal-time function. Most of our results are more or less distant outgrowths of this; among the latter belongs a new characterization of normality, Corollary 7. In §3 the principal result is Theorem 10, existence of optimal feedback controls for normal systems. A number of intriguing questions connected with this result (see the remarks following Theorem 10) are still unresolved.

The last two sections are really appendices. Section 4 emphasizes the connection between control theory and dynamical system theory. In §5 we treat a rather special class of control systems, for which quite direct methods yield much information concerning the minimal-time function.

Geometric aspects of such problems, for the far more difficult nonlinear case, were treated in [5]; among other things, there is an apparent relationship between our Proposition 14 and Theorem 6.3 in [5] (or rather its proof: note  $\varepsilon < \eta$ , p. 321).

We will use the terminology and notation of [2] where available, and of [3] where not. In particular, the reachable set at time  $t \geq 0$  is

$$\mathcal{R}(t) = \left\{ \int_0^t e^{-As} Bu(s) ds : |u_i(s)| \leq 1, u \in \mathcal{L}_1[0, t] \right\}$$

(a not particularly happy choice of term: the set of points reached from the origin in time  $t$  is actually  $e^{At}\mathcal{R}(t)$ ); and the reachable set is

$$\mathcal{R} = \bigcup \{ \mathcal{R}(t) : t \geq 0 \}.$$

It is known that  $\mathcal{R}(t)$  is compact, convex, symmetric about 0 and has  $\mathcal{R}(s) \subset \mathcal{R}(t)$  whenever  $0 \leq s \leq t$  [2, p. 46]. It is easily shown that  $\mathcal{R}(t)$  satisfies the

---

\* Received by the editors August 4, 1970, and in final revised form, December 15, 1970.

† Department of Mathematics and Statistics, Case Western Reserve University, Cleveland, Ohio 44106. This work was supported in part by the National Science Foundation under Grant GP-22689.

stricter monotonicity condition

$$(1) \quad \mathcal{R}(s) \subset \text{Int } \mathcal{R}(t) \quad \text{if } 0 \leq s < t$$

if and only if  $\mathcal{R}$  is open, and also if and only if (S) is controllable [2, p. 72 ff.].

Without the restriction  $|u_i| \leq 1$  we obtain the controllability space

$$\mathcal{C}(t) = \left\{ \int_0^t e^{-As} Bu(s) ds : u \in \mathcal{L}[0, t] \right\},$$

$$\mathcal{C} = \bigcup \{ \mathcal{C}(t) : t \geq 0 \};$$

$\mathcal{C}(t)$  and  $\mathcal{C}$  are linear subspaces of  $R^n$ , whereupon, easily,  $\mathcal{C}(t) = \mathcal{C}$  for  $t > 0$  [3, p. 97ff.].

One (possibly) new concept will be useful, that of the  $(n + 1)$ -dimensional reachability set

$$\mathcal{R}^* = \{(x, t) \in R^{n+1} : x \in \mathcal{R}(t)\}.$$

Its connection with the sets  $\mathcal{R}(t)$  is obvious from

$$\mathcal{R}^* \cap (R^n \times \{t\}) = \mathcal{R}(t) \times \{t\};$$

in particular,  $\mathcal{R}^*$  is closed (with compact "slabs"), symmetric about the  $t$ -axis and contains it; usually  $\mathcal{R}^*$  is not convex.

It is immediate that  $x \in \mathcal{R}(t)$  (equivalently,  $-x \in \mathcal{R}(t)$ ) if and only if  $x$  can be steered to 0 in time  $t$  by some admissible controls. The minimal-time function  $T$  is defined [3, p. 145] by

$$(2) \quad T(x) = \inf \{t \geq 0 : x \in \mathcal{R}(t)\};$$

thus  $0 \leq T \leq +\infty$  with  $T(x) < +\infty$  if and only if  $x \in \mathcal{R}$ .

If  $\text{grad } T(x)$  exists, then

$$\max \{(-\text{grad } T(x), Ax + Bu) : u \in R^m, |u_i| \leq 1\} = 1,$$

and, if  $T$  has a total differential in an open set, then there

$$f(x) = -\text{sgn}(B' \text{grad } T(x))$$

is an optimal feedback control (see [3, p. 146], modulo notational changes). It is precisely this unwarranted *ex post* assumption of differentiability that we wish to avoid in the present paper.

## 2. The minimal-time function.

**THEOREM 1.** *The minimal-time function  $T: \mathcal{R} \rightarrow R^1$  is continuous, with  $\mathcal{R}$  open in the linear space  $\mathcal{C} \subset R^n$ .*

*Proof.* Directly from (2) and the closedness of the sets  $\mathcal{R}(t)$ ,  $T$  is lower semicontinuous. For controllable systems, upper semicontinuity follows from (1). For the remaining case, use Kálmán's decomposition [3, p. 99]: up to a linear equivalence, (S) may be written in the form

$$\begin{aligned} \dot{x}_1 &= A_{11}x_1 + A_{12}x_2 + B_1u, \\ \dot{x}_2 &= \quad \quad \quad A_{22}x_2 \quad \quad ; \end{aligned}$$

furthermore, the system

$$(S_1) \quad \dot{y} = A_{11}y + B_1u$$

(corresponding to  $x_2 \equiv 0$ ) is controllable. Obviously  $\mathcal{R}(t) = \mathcal{R}_1(t) \times \{0\}$ , so that

$$T \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{cases} T_1(x_1) & \text{if } x_2 = 0, \\ +\infty & \text{if } x_2 \neq 0. \end{cases}$$

Since  $(S_1)$  is controllable,  $T_1$  is continuous, and hence so is  $T$ ;  $T_1$  is finite precisely on  $\mathcal{R}_1$ , an open subset of  $\mathcal{C}_1 = \mathcal{C}$  [2, p. 79].

**COROLLARY 2.** *If  $x \in \partial\mathcal{R}$ , the boundary relative to  $\mathcal{C}$ , then*

$$\lim_{y \rightarrow x} T(y) = +\infty.$$

Thus  $T: \mathcal{C} \rightarrow R^1 \cup \{+\infty\}$  is also continuous.

*Proof.* Assume  $x_p \rightarrow x$ ,  $T(x_p) \leq t$ . Then  $x_p \in \mathcal{R}(t)$ , and hence  $x \in \mathcal{R}(t) \subset \mathcal{R}$  from closedness.

**COROLLARY 3.** *For all  $t \geq 0$  we have*

$$\{x: T(x) \leq t\} = \mathcal{R}(t), \quad \{x: T(x) = t\} = \partial\mathcal{R}(t)$$

(boundary relative to  $\mathcal{C}$ ).

*Proof.* As concerns the first formula, obviously

$$\mathcal{R}(t) \subset \{x: T(x) \leq t\} \subset \bigcap_{p=1}^{\infty} \mathcal{R} \left( t + \frac{1}{p} \right);$$

thus it is sufficient to show that the last set is contained in the first. Let  $x$  be in the intersection; thus there are admissible controls  $u_p: [0, t + 1/p] \rightarrow R^m$  such that

$$x = \int_0^{t+1/p} e^{-As} B u_p(s) ds = \int_0^t + \int_t^{t+1/p} = x_p + y_p.$$

Here  $x_p \in \mathcal{R}(t)$ , and the remainder terms  $y_p \rightarrow 0$  since all coordinates of the  $u_p$ 's are bounded by 1. Thus  $x_p \rightarrow x$ , so that  $x$  is in the closed set  $\mathcal{R}(t)$ .

For the second formula note first that

$$\{x: T(x) < t\} \subset \text{Int } \mathcal{R}(t)$$

from the continuity of  $T$ ; and we need only verify the converse inclusion for  $t > 0$ . If  $x \in \text{Int } \mathcal{R}(t)$ , then  $x \in \text{Int } \mathcal{R}(t - \varepsilon)$  for small  $\varepsilon > 0$  [2, p. 47], and hence  $T(x) \leq t - \varepsilon < t$ .

*Remark 1.* Consider the boundary  $\partial\mathcal{R}^*$  of the  $(n + 1)$ -dimensional reachable set, relative to  $\mathcal{C} \times R^1$ . Then Theorem 1 shows that any line in  $R^{n+1}$  parallel to the  $t$ -axis intersects  $\partial\mathcal{R}^*$  at most once. Thus the set  $\mathcal{R}^*$ , though often not convex, nevertheless exhibits a form of convexity in the  $t$ -direction.

*Remark 2.* The assertion of Theorem 1 may be invalid for time-dependent systems, or for target sets other than the origin.

*Remark 3.* The assertion of Corollary 2 is false if  $\mathcal{C}$  is replaced by  $R^n$ . In point of fact,  $T: R^n \rightarrow R^1 \cup \{+\infty\}$  is continuous if and only if  $(S)$  is controllable (i.e.,



$\mathcal{C} = \mathbb{R}^n$ ). For if (S) is not controllable, one may take  $x_p \rightarrow x$  with  $x \in \mathcal{R}$ ,  $x_p \notin \mathcal{C} \supset \mathcal{R}$ , whereupon  $T(x_p) = +\infty > T(x)$ .

*Remark 4.* From Corollary 3 it follows that the minimal-time function completely determines the performance of system (S) (interpreting "performance" as the correspondence to  $t \mapsto \mathcal{R}(t)$ ). In particular, the sets  $\{x : T(x) \leq t\}$  are strictly convex if and only if (S) is normal [2, p. 65].

In the next theorem we will need the following approximation result (this is probably well-known, easily proved by standard methods, and amenable to generalization).

LEMMA 4. For every  $\varepsilon > 0$  and  $t \geq 0$  there exists  $\delta > 0$  such that

$$\text{dist}(\mathcal{R}_1(t), \mathcal{R}(t)) < \varepsilon,$$

where  $\mathcal{R}_1(t)$  is the reachability set of any system

$$(S_1) \quad \dot{x} = A_1 x + B_1 u$$

with  $|A - A_1| < \delta, |B - B_1| < \delta$ .

THEOREM 5. Given a sequence of systems

$$(S_p) \quad \dot{x} = A_p x + B_p u$$

with  $A_p \rightarrow A, B_p \rightarrow B$ , then their minimum-time functions  $T_p$  converge to  $T$  in the following sense:

$$T_p(x_p) \rightarrow T(x) \quad \text{whenever } x_p \rightarrow x, \quad x_p \in \mathcal{R} \cup \mathcal{R}_p, \quad x \in \mathcal{R}.$$

In particular,  $T_p \rightarrow T$  uniformly on compact subsets of  $\mathcal{R}$ .

*Proof.* Assume that, on the contrary,  $T_p(x_p)$  is bounded away from  $T(x_0)$  for some subsequence of the  $p$ 's. Thus there are two cases: either

$$(3) \quad T_p(x_p) < t_0 < T(x_0)$$

for some  $t_0$  and subsequence, or similarly with opposite inequalities. In the first case we have  $x_p \in \mathcal{R}_p(t_0)$ ; from Lemma 4,  $y_p \in \mathcal{R}(t_0)$  for some  $y_p$  with  $y_p - x_p \rightarrow 0$ . But then  $T(y_p) \leq t_0$  and  $y_p \rightarrow x_0$ ; since  $t_0 < T(x_0)$ , this contradicts Theorem 1.

In the second case we have

$$(4) \quad T(x_0) < t_0 < T_p(x_p).$$

Then  $x_0 \in \text{Int } \mathcal{R}(t_0)$ , the interior relative to  $\mathcal{C}$ . From Lemma 4, there exist  $y_p \in \mathcal{R}_p(t_0)$  with  $y_p \rightarrow x_0$ . Since  $x_p \notin \mathcal{R}(t_0)$  by assumption (4), on the segment  $x_p y_p$  (or  $x_0 y_p$ , if  $y_p \notin \mathcal{R}(t_0)$ ) there is a point  $z_p \in \partial \mathcal{R}(t_0)$ ; necessarily,  $z_p \rightarrow x_0$ . Let  $c_p$  be a unit exterior normal to a supporting hyperplane at  $z_p$  to  $\mathcal{R}_p(t_0)$ . Then

$$(5) \quad (c_p, x - z_p) \leq 0 \quad \text{for all } x \in \mathcal{R}_p(t_0).$$

Choose a convergent subsequence  $c_p \rightarrow c \neq 0$ . For any  $w \in \mathcal{R}(t_0)$  there exist  $w_p \in \mathcal{R}_p(t_0)$  with  $w_p \rightarrow w$ ; taking limits in (5) with  $x = w_p$ , we conclude that

$$(c, w - x_0) \leq 0 \quad \text{for all } w \in \mathcal{R}(t_0).$$

However, this contradicts  $x_0 \in \text{Int } \mathcal{R}(t_0)$ , and completes the proof.

*Remark 5.* In the situation of this theorem we cannot assert that

$$T_p(x_p) \rightarrow T(x) \quad \text{whenever } x_p \rightarrow x.$$

Actually this holds if and only if (S) is controllable. Indeed, if (S) is not controllable, take  $(S_p) \equiv (S)$ , and proceed as in Remark 2 above. On the other hand, if (S) is controllable, then  $x_p \rightarrow x \in \mathcal{R}$  implies  $x_p \in \mathcal{R}$  ultimately.

*Remark 6.* Theorem 1 is, of course, a special case of Theorem 5; however, the former was needed in the proof of the latter.

*Remark 7.* Theorem 5 admits this pessimistic interpretation: by changing the system slightly one cannot improve performance too much.

**PROPOSITION 6.** *In the situation of Theorem 5, let the  $y_p$  be optimal solutions of  $(S_p)$ , and let  $y_p(0) \rightarrow x_0 \in \mathcal{R}$ . Then some subsequence converges to an optimal solution  $y$  of (S), uniformly on compact subsets of  $[0, T(x_0))$ .*

*Proof.* We have

$$y_p(0) \rightarrow x_0, \quad y_p(\theta_p) = 0 \quad \text{for } \theta_p = T_p(y_p(0)).$$

From Theorem 5,  $\theta_p \rightarrow \theta = T(x_0)$ . Using weak sequential compactness of the space of admissible controls, we conclude that  $\{y_p\}$  has a subsequence  $\{y_q\}$  converging, uniformly on compact subsets of  $[0, \theta)$ , to a solution  $y$  of (S) through  $x_0$ . Since the  $y_q$  were optimal,

$$T_q(y_q(t)) = T_q(y_q(0)) - t \quad \text{for } 0 \leq t \leq \theta_q;$$

thus, from Theorem 5 again,

$$(6) \quad T(y(t)) = T(x_0) - t \quad \text{for } 0 \leq t < \theta.$$

In particular,  $y$  is optimal for (S), at least on  $[0, \theta)$ . Finally,  $\lim_{t \rightarrow \theta^-} y(t) = 0$  follows from (6):  $y(t) \in \mathcal{R}(\theta - t)$ .

**COROLLARY 7.** *(S) is normal if and only if it satisfies the following condition. Whenever the  $(S_p)$  are as in Theorem 5, and the  $y_p$  are optimal solutions of  $(S_p)$  through  $y_p(0) \rightarrow x_0 \in \mathcal{R}$ , necessarily the sequence  $\{y_p\}$  converges (whereupon the limit is an optimal solution of (S) and convergence is uniform on compact subsets of  $[0, T(x_0))$ ).*

**3. Optimal feedback controls.** We will study an auxiliary concept associated with a given system (S). For every  $x \in \mathcal{R}$  let  $E(x)$  denote the set of all unit exterior normals of supporting hyperplanes at  $x$  to  $\mathcal{R}(\theta)$ , where  $\theta = T(x)$  (so that  $x \in \partial\mathcal{R}(\theta)$ , see Corollary 3).

**LEMMA 8.**  *$E(x)$  is a nonvoid compact subset of the  $(n - 1)$ -sphere  $S^{n-1}$ . Furthermore,*

$$(7) \quad \begin{aligned} E(-x) &= -E(x), \\ \limsup E(y) &\subset E(x) \quad \text{as } y \rightarrow x \text{ in } \mathcal{R}. \end{aligned}$$

*If (S) is controllable, then  $E(x)$  contains antipodal vectors only for  $x = 0$ .*

*Proof.* The first assertion actually concerns existence of supporting hyperplanes. For the second use symmetry of  $\mathcal{R}(\theta)$  about 0. The third follows easily from the continuity of  $T$  and a compactness argument. For the last observe that (S) is proper, so that  $\text{Int } \mathcal{R}(\theta) \neq \emptyset$  if  $\theta > 0$ , i.e., if  $0 \neq x \in \mathcal{R}$ .

It follows that the set of all the exterior normals,

$$\{tc : t > 0, c \in E(x)\},$$

is a wedge in  $R^n$ , and actually a proper cone if (S) is controllable and  $x \neq 0$ .

Since  $E(x)$  is nonvoid, for each  $x$  one can choose  $c(x) \in E(x)$ . Obviously a continuous selection need not be possible (unless, e.g., the  $\mathcal{R}(t)$  have no corners). However, one can at least make a measurable selection.

LEMMA 9. *There exists a measurable function  $c: \mathcal{R} \rightarrow R^n$  such that*

$$c(x) \in E(x) \quad \text{for all } x \in \mathcal{R}.$$

*Proof.* This paraphrases part of the proof of Filippov's theorem [2, p. 31]. For each  $x \in \mathcal{R}$  choose  $c(x) \in E(x)$  with minimal first coordinate  $c_1(x)$  (apply compactness of  $E(x)$ , from Lemma 8); if there is more than one such element in  $E(x)$ , choose that whose second coordinate  $c_2(x)$  is also minimal; and so on. Directly from (7),  $c_1(\cdot)$  is upper semicontinuous on  $\mathcal{R}$ . By Lusin's theorem, for every  $\varepsilon > 0$  there exists a subset  $M_1 \subset \mathcal{R}$  on which  $c_1(\cdot)$  is continuous, and has  $\text{meas } \mathcal{R} - M_1 < \varepsilon/n$ . Again from (7),  $c_2(\cdot)$  is upper semicontinuous on  $M_1$ , so there is a subset  $M_2 \subset M_1$  on which both  $c_1(\cdot)$ ,  $c_2(\cdot)$  are continuous, and  $\text{meas } M_1 - M_2 < \varepsilon/n$ . Continuing in this fashion we find a set  $M \subset \mathcal{R}$  on which all coordinates of  $c_k(\cdot)$  are continuous, and  $\text{meas } \mathcal{R} - M < \varepsilon$ . Since  $\varepsilon > 0$  was arbitrary,  $c$  is measurable (this is the trivial converse to Lusin's theorem).

THEOREM 10. *If (S) is normal, there exists a measurable function  $f: \mathcal{R} \rightarrow R^m$  which is an optimal feedback control for (S) in the following sense: In addition to (S), consider the autonomous differential equation*

$$(F) \quad \dot{y} = Ay + Bf(y).$$

*Then each optimal solution of (S) is a solution of (F) in  $\mathcal{R}$ ; as a partial converse, each solution of (F) is a solution (possibly not optimal) of (S).*

*Proof.* We will first construct  $f$  and then verify the assertion on the solutions.

Let  $c: \mathcal{R} \rightarrow R^n$  be as described in Lemma 9. For each  $x \in \mathcal{R}$  set

$$(8) \quad f(x) = \lim_{s \rightarrow 0^+} -\text{sgn}(B'e^{-A's}c(x)).$$

Since each coordinate of  $B'e^{-A's}c(x)$  is analytic in  $s$  (for fixed  $x$ ),  $\text{sgn}(B'e^{-A's}c(x))$  is piecewise constant in  $s$ , so that the limit exists. Finally,  $f$  is measurable since  $c$  is such. (Remark: the "obvious" choice would be  $f(x) = -\text{sgn}(B'c(x))$ ; the reason for the extra factor and limit is technical, and will become apparent in the course of the proof.)

Now let  $x: [0, \theta] \rightarrow R^n$  be an optimal solution of (S) through  $x(0) = x_0 \in \mathcal{R}$ ,  $\theta = T(x_0)$ ; we wish to show that simultaneously it satisfies (F). Take arbitrarily  $s$ ,  $0 \leq s < \theta$ ; there is an optimal control through  $x(s) \in \mathcal{R}$  at time  $s$ , which may be taken in the form

$$u_s(t) = -\text{sgn}(B'e^{-A'(t-s)}c), \quad s \leq t \leq \theta,$$

for any choice of  $c$  in  $E(x(s))$  (i.e.,  $-c \in E(-x(s))$ , see [2, p. 51]); we choose  $c = c(x(s))$ . Since (S) is normal, the response to  $u_0$  ( $u_s$  for  $s = 0$ ) is precisely  $x$ ; furthermore, on  $[s, \theta]$ ,  $u_0$  is again an optimal control for  $x(s)$  at initial time  $s$  (see (6)). Again from normality, optimal controls are essentially unique; thus, for each individual  $s$  and almost all  $t \geq s$  we have

$$u_0(t) = u_s(t) = -\text{sgn}(B'e^{-A'(t-s)}c(x(s))).$$

Now take limits  $t \rightarrow s+$  over these  $t$ 's and obtain (see (8))

$$u_0(s+) = f(x(s)) \quad \text{for } s \in [0, \theta).$$

Since  $u_0$  is piecewise constant, we have  $u_0(s) = f(x(s))$  for almost all  $s$ ; therefore the response  $x$  to  $u_0$  satisfies

$$\dot{x}(s) = Ax(s) + Bu(s) = Ax(s) + Bf(x(s))$$

almost everywhere. As  $x$  is absolutely continuous, it is a solution of (F).

Finally, that solutions of (F) solve (S) is almost immediate. Indeed, any solution  $y$  of (F) is the response of (S) to  $v(t) = f(y(t))$ , and obviously this is an admissible control (see (8) for the bound, and also Lemma 9 for measurability). This completes the proof of Theorem 10.

*Remark 8.* This raises further questions.

It would be desirable to know whether there is a complete converse in Theorem 10, more precisely, whether all solutions of (F) are optimal solutions of (S), or equivalently, whether (F) has uniqueness of solutions to the initial value problem into positive time.

In the positive case the converse assertion may even hold without normality; a proof might proceed via Proposition 6 and an approximation of (S) by normal systems.

**4. Semidynamical systems.** Assume that (S) is normal. Then with every point  $x \in \mathcal{R}$  we may associate an optimal solution  $y: [0, +\infty) \rightarrow \mathcal{R}^n$  of (S) through  $y(0) = x$ , whereupon  $y$  is determined uniquely (let  $y(t) = 0$  identically for  $t \geq T(x)$ ). Thus there is a well-defined mapping  $\pi: \mathcal{R} \times [0, +\infty) \rightarrow \mathcal{R}$ ,

$$\pi(x, t) = y(t) \quad \text{for } t \geq 0.$$

It is obvious that

$$\pi(x, 0) = x, \quad \pi(\pi(x, t), s) = \pi(x, t + s);$$

and it is easily shown that  $\pi$  is continuous. Therefore [1, p. 12] we have the following theorem.

**THEOREM 11.** *If (S) is normal, then  $\pi$  is a global semidynamical system on  $\mathcal{R}$ .*

Note that even in the most reasonable cases,  $\pi$  need not have negative uniqueness (i.e., solutions are not determined uniquely by initial data in the negative time direction); see, e.g., the situation indicated in [3, p. 10].

An examination of the relation between (S) and  $\pi$  is intriguing. Thus, 0 is the only critical point of  $\pi$ , and it is asymptotically stable, with  $\mathcal{R}$  as its region of attraction. Furthermore, the minimum-time function of (S) (or rather  $-T$ ) is a Lyapunov function for  $\pi$ , strictly decreasing along trajectories outside 0, and with  $T(x) = 0$  if and only if  $x = 0$ ,  $T(x) \rightarrow +\infty$  as  $x$  approaches the boundary of  $\mathcal{R}$  (Corollary 2). Finally, Theorem 10 can be reformulated in dynamical terms thus: the semidynamical system  $\pi$  admits a differential representation.

We conclude this section with an application of dynamical systems to control theory. The question treated may be indicated as follows: If  $x$  is any point on  $\partial\mathcal{R}(t)$  and  $t \geq s \geq 0$ , then obviously an optimal solution through  $x$  can be followed down till it meets  $\mathcal{R}(s)$ . The converse assertion, concerning backing out of  $\mathcal{R}(s)$  to reach  $\partial\mathcal{R}(t)$ , though plausible, does not have an obvious proof.

PROPOSITION 12. *Every optimal solution on  $[0, +\infty)$  can be extended to an optimal solution defined on  $R^1$ .*

*Proof.* First consider the special case that (S) is normal, so that semidynamical system theory may be applied. Consider any optimal solution  $x:I \rightarrow R^n$  which is inextensible. According to [1, p. 32], there are three possibilities. Either  $I = R^1$ , and there is nothing to prove; or

$$(9) \quad I = [\alpha, +\infty), \quad \alpha > -\infty, \quad x(\alpha) \text{ is a start point ;}$$

or finally,

$$(10) \quad \begin{aligned} I &= (\alpha, +\infty), \alpha > -\infty, \quad \text{and} \\ x(t) &\text{ has no cluster points in } \mathcal{R} \text{ as } t \rightarrow \alpha + . \end{aligned}$$

We proceed to eliminate the latter two alternatives. Since (S) is controllable,  $\mathcal{R}$  is open in  $R^n$ , and hence  $\pi$  is a semidynamical system on an  $n$ -manifold. According to [1, p. 124],  $\pi$  has no start points; thus (9) is impossible.

Finally, consider case (10). Since the right-hand side of (S) admits a linear estimate  $|A||x| + |B|$ , we cannot have  $x(t) \rightarrow \infty$  as  $t \rightarrow \alpha +$  (this is a minor modification of a standard result, see, e.g., [4, Theorem 3, Chap. 1]). Thus there is  $x(t_p) \rightarrow y$  for some  $t_p \rightarrow \alpha +$ ,  $y \in \mathcal{R} - \mathcal{R}$ . But then, from (6),

$$T(x(t_p)) = T(x(0)) - t_p \leq T(x_0) - \alpha < +\infty,$$

contradicting Corollary 2. This completes the proof in the case of normal (S).

Next, consider the general case, and let  $x:[0, +\infty) \rightarrow R^n$  be an optimal solution through  $x(0) = y_0 \in \mathcal{R}$ . Since normal systems are generic (see [3, p. 100], and recall that finite intersections of open dense sets are again such), there exist normal systems

$$(S_p) \quad \dot{x} = A_p x + B_p u, \quad A_p \rightarrow A, \quad B_p \rightarrow B.$$

From Lemma 4, there exist  $y_p \rightarrow y_0$  with  $y_p \in \mathcal{R}_p$ . We have just shown that  $(S_p)$  has an optimal solution  $x_p:R^1 \rightarrow R^n$  through  $y_p$ . Now, some subsequence of the  $x_p$ 's converges uniformly on compact subsets of  $(-\infty, 0]$ . Indeed, if the  $u_p$  are corresponding controls, then a subsequence converges weakly on compact subsets of  $(-\infty, 0]$ ; the limit is also an admissible control, and the  $B_p u_p$  converge weakly to  $Bu$ . Denote the limit by  $x$  again. Since  $x_p$  is optimal,

$$T_p(x_p(t)) = T_p(y_p) - t \quad \text{for } t \leq 0.$$

Take limits, and apply Theorem 5:

$$T(x(t)) = T(y_0) - t \quad \text{for } t \leq 0;$$

thus  $x$  is indeed an optimal solution of (S). This concludes the proof.

**5. Systems with uncoupled controls.** In this section only we make the following overall convention.

CONVENTION 13. *The control matrix  $B$  has rank  $B = n$ .*

A system is such if and only if it is linearly equivalent to a system with control matrix  $I$ , after possible omission of ineffective or duplicated controls; i.e., linearly

equivalent to a system with uncoupled controls

$$\dot{x}_i = \sum_j a_{ij}x_j + u_i, \quad 1 \leq i \leq n.$$

Note that automatically such systems are controllable, without any assumption on  $A$  (but not necessarily normal, e.g., if the minimal polynomial of  $A$  has degree  $< n$ ).

**PROPOSITION 14.** *The  $(n + 1)$ -dimensional reachable set  $\mathcal{R}^*$  contains a small cone with vertex at 0. More precisely, there exist  $\lambda > 0, \varepsilon > 0$ , such that*

$$x \in \mathcal{R}(\lambda|x|) \quad \text{whenever } |x| \leq \varepsilon.$$

Furthermore,  $\lambda$  may be taken arbitrarily close to

$$1/\min \left\{ \sum_i |b'_i c| : c \in R^n, |c| = 1 \right\},$$

where  $B = (b_1, \dots, b_m)$ .

*Proof.* First fix any  $\lambda > 0$ , and assume an  $x \in R^n$  with  $x \notin \mathcal{R}(\lambda|x|)$ . Since each  $\mathcal{R}(t)$  contains 0, some point  $y$  on the segment  $0x$  is on the boundary of  $\mathcal{R}(\lambda|x|)$ ; and, of course,  $|y| \leq |x| > 0$ .

Corresponding to  $y$ , there is a control, to reach  $y$  in time  $\lambda|x|$ , of a special type [2, p. 51]: for some  $c \in R^n$  with  $|c| = 1$ , if  $u(t) = \text{sgn}(B'e^{-A't}c)$ ,

$$y = \int_0^{\lambda|x|} e^{-As}Bu(s)ds.$$

It follows that

$$(c, y) = \int_0^{\lambda|x|} \sum_{i=1}^m |(B'e^{-A's}c)_i| ds.$$

The assumptions on magnitudes yield  $(c, y) \leq |x|$ ; thus

$$(11) \quad |x| \geq \lambda|x| \cdot \mu(\lambda|x|), \quad x \neq 0,$$

where  $\mu(\cdot)$  is defined by

$$\mu(t) = \inf \left\{ \sum_{i=1}^m |(B'e^{-A's}c)_i| : |c| = 1, 0 \leq s \leq t \right\}.$$

Evidently  $\mu$  is nonincreasing. To show that it is strictly positive, assume that

$$\sum_{i=1}^m |(B'e^{-A's_p}c_p)_i| \rightarrow 0 \quad \text{as } p \rightarrow \infty;$$

we may take  $s_p \rightarrow s, c_p \rightarrow c, |c| = 1$ . The limit relation shows that  $c \in R^n$  is perpendicular to all columns of  $e^{-As}B$ , contradicting

$$\text{rank } e^{-As}B = \text{rank } B = n, \quad c \neq 0.$$

Now choose any  $\varepsilon = \delta/\lambda, \delta > 0$ . Then, for all  $x$  with  $|x| \leq \varepsilon$ , we have  $\mu(\lambda|x|) \geq \mu(\delta) > 0$ , and (11) yields  $\lambda \leq 1/\mu(\delta)$ . In other words, for  $\lambda > 1/\mu(\delta)$ , necessarily no  $x \in R^n$  with  $|x| \leq \varepsilon$  can have  $x \notin \mathcal{R}(\lambda|x|)$ .

Evidently, we may take  $\lambda$  arbitrarily close to  $1/\mu(0)$ , at the cost of decreasing  $\varepsilon$ .

*Remark 9.* For Proposition 14 above, constancy of  $A, B$  may be weakened to:  $A$  is summable,  $B$  is continuous at 0,  $\text{rank } B(0) = n$ . However, subsequent results will require autonomy.

*Remark 10.* Obvious estimates yield that  $\mathcal{R}^*$  is also contained within a cone with vertex at 0 (locally at 0). This suggests that the numerical value of the actual opening, of  $\mathcal{R}^*$  at 0, might be a significant qualitative measure of performance of the system.

*Remark 11.* Consideration of the case  $A = 0$  suggests that our rank condition is also necessary; but this will not be needed here. The condition cannot be replaced by normality alone. Indeed, for the two-dimensional system

$$\dot{x} = y, \quad \dot{y} = u,$$

the three-dimensional reachable set intersects the  $(x, t)$ -plane in  $\{(x, t) : |x| \leq \frac{1}{4}t^2\}$ ; thus  $\mathcal{R}$  contains no small cone with vertex at 0, and the minimal-time function is not locally Lipschitzian at 0 (cf. Theorem 17). (Actually,  $T(x, y) = y + \sqrt{2y^2 + 4x}$ .)

*Remark 12.* It will be seen from subsequent results that good upper estimates of  $\lambda$  are of considerable interest (while  $\varepsilon$  is unimportant). For the one-dimensional system

$$(12) \quad \dot{x} = x + u,$$

$\mathcal{R}(t)$  is the segment  $[-(1 - e^{-t}), (1 - e^{-t})]$ ; thus the largest opening is 1. In this case also  $\min \{\dots\} = 1$ , so that our estimate of  $\lambda$  is exact. Somewhat more generally, it can be shown that, if  $B$  is square nonsingular, then  $\lambda$  may be taken arbitrarily close to  $|B^{-1}|$ .

Throughout this section,  $\lambda$  and  $\varepsilon$  will be taken as described in Proposition 14.

**COROLLARY 15.** For any  $x \in \mathcal{R}(t)$  and  $y$  with  $|y - x| \leq \varepsilon \exp -|A|t$  we have

$$y \in \mathcal{R}(t + |y - x|\lambda \exp |A|t).$$

Thus  $\mathcal{R}^*$  contains, with every  $x$ , a small cone with vertex at  $x$  and axis parallel to the  $t$ -axis.

*Proof.* This follows from Proposition 14 and the following elementary observation: if  $x \in \mathcal{R}(t)$ ,  $z \in \mathcal{R}(s)$ , then

$$x + e^{-At}z \in \mathcal{R}(t + s).$$

Next we turn to the study of the minimal-time function  $T$ .

**THEOREM 16.** For all  $x, y \in \mathcal{R}$ ,

$$(13) \quad |T(x) - T(y)| \leq |x - y|\lambda \exp (|A| \max (T(x), T(y))).$$

Consequently  $T: \mathcal{R} \rightarrow R^1$  is locally Lipschitzian, and the Lipschitz constant near  $x$  may be taken arbitrarily close to  $\lambda \exp |A|T(x)$ . Finally,  $\text{grad } T(x)$  exists almost everywhere in  $\mathcal{R}$ .

*Proof.* We first show that, for each  $x \in \mathcal{R}$ ,

$$T(y) - T(x) \leq |y - x|\lambda \exp |A|T(x) \quad \text{if } |y - x| < \varepsilon \exp (-|A|T(x)).$$

Indeed, choose any  $y$  as indicated, and then  $t > T(x)$  with  $|y - x| < \varepsilon \cdot \exp(-|A|t)$ . Then  $x \in \mathcal{R}(t)$ , and, from Corollary 15,

$$T(y) \leq t + |y - x|\lambda \exp |A|t;$$

finally, let  $t \rightarrow T(x)$ .

In proving (13) we may assume that  $T(y) < T(x)$ . Since  $\{x : T(x) \leq t\}$  is convex, the entire segment  $P$  from  $y$  to  $x$  has

$$T(z) \leq T(x) \quad \text{for } z \in P.$$

Now apply the previous estimate to a finite decomposition of  $P$  into segments of length  $< \varepsilon \exp(-|A|T(x))$ , noting that the Lipschitz constants are all  $\leq \lambda \cdot \exp |A|T(x)$ .

Finally, according to [6, p. 311, Theorem (14.2) (ii)] a Lipschitzian function has total differential almost everywhere.

PROPOSITION 17.  $\mathcal{R}$  contains the ball

$$U = \{x \in \mathbb{R}^n : |x| < 1/(\lambda|A|)\}$$

(with  $U = \mathbb{R}^n$  if  $A = 0$ ); in  $U$ ,

$$(14) \quad T(x) \leq \frac{1}{|A|} \log \frac{1}{1 - \lambda|A||x|},$$

so that, within any ball of radius  $\xi$  in  $U$ ,

$$|T(x) - T(y)| \leq |x - y| \frac{\lambda}{1 - \lambda|A|\xi}.$$

*Proof.* We obtain estimate (14) first. From Theorem 16, the function  $S$ ,

$$S(x) = \exp(-|A|T(x)) \quad \text{for } x \in \mathcal{R},$$

is Lipschitzian, with Lipschitz constant  $\lambda|A|$  (first locally only). Thus

$$|S(x) - S(0)| \leq \lambda|A||x|,$$

and the estimate follows.

To show that  $T$  is defined throughout  $U$  it is sufficient (since  $U$  is connected and  $\mathcal{R}$  open) to prove that  $\mathcal{R} \cap U$  is closed in  $U$ . Thus, let  $x_p \rightarrow x$ ,  $x_p \in \mathcal{R}$ ,  $x \in U$ . From the estimate, the  $T(x_p)$  are bounded, say by  $t$ . Thus  $x_p \in \mathcal{R}(t)$  and hence  $x \in \mathcal{R}(t) \subset \mathcal{R}$ .

Remark 13. For the example  $\dot{x} = x + u$  treated earlier, the estimates of both  $\mathcal{R}$  and  $T$  are sharp.



## REFERENCES

- [1] N. P. BHATIA AND O. HÁJEK, *Local Semi-Dynamical Systems*, Lecture Notes in Mathematics 90, Springer, Berlin, 1969.
- [2] H. HERMES AND J. P. LASALLE, *Functional Analysis and Time Optimal Control*, Academic Press, New York, 1969.
- [3] E. B. LEE AND L. MARKUS, *Foundations of Optimal Control Theory*, John Wiley, New York, 1967.
- [4] V. V. NIEMYCKIĀ AND V. V. STEPANOV, *Qualitative Theory of Differential Equations*, 2nd ed., Gostechizdat, Moscow, 1947.
- [5] E. ROXIN, *A geometric interpretation of Pontryagin's maximum principle*, Nonlinear Differential Equations and Nonlinear Mechanics, J. P. LaSalle and S. Lefschetz, eds., Academic Press, New York, 1963, pp. 303–324.
- [6] S. SAKS, *Theory of the Integral*, Monografie Matematyczne VII, Warszawa, 1937. (English transl.)

## ANOTHER PROOF OF THE LIAPUNOV CONVEXITY THEOREM\*

JAMES A. YORKE†

**Abstract.** A new proof of the Liapunov convexity theorem is presented.

**1. Introduction.** Let  $x$  be in  $R^n$  and let  $A(t)$  be a continuous  $n \times n$  matrix. Let  $b: [0, T] \rightarrow R^n$  be continuous. LaSalle's well-known bang-bang control theorem [5] for linear systems states that if  $X(t)$  is the fundamental matrix (with  $X(0) = I$ , the identity matrix) for  $x'(t) = A(t)x$ , then the "reachable set at time  $T$ " for

$$(1) \quad y'(t) = A(t)y(t) + b(t)u(t), \quad y(0) = 0$$

is compact and convex and the reachable set does not change if we restrict ourselves to bang-bang controls. That is:

Write  $E = [0, T]$ . Let  $\mathcal{F}$  be the set of integrable functions  $u: E \rightarrow [-1, +1]$ , and let  $\mathcal{F}^*$  be  $\{u \in \mathcal{F} : |u(t)| = 1 \text{ for almost all } t\}$ . If we let  $f(t) = X(T)X^{-1}(t)b(t)$ , then for  $u \in \mathcal{F}$ , the solution  $y(t)$  at time  $T$  satisfies  $y(T) = \int_E f(t)u(t) dt$ . Define

$$(2) \quad \mathcal{I}u = \int_E f(t)u(t) dt,$$

and for any set  $\mathcal{S}$  of functions  $u(\cdot)$  write  $\mathcal{I}\mathcal{S} = \{\mathcal{I}u : u \in \mathcal{S}\}$ . Then  $\mathcal{I}\mathcal{F}^*$  is the set of reachable points using bang-bang controls.

**THEOREM.**  $\mathcal{I}\mathcal{F}^*$  is compact and convex and equals  $\mathcal{I}\mathcal{F}$ .

Many proofs and extensions of this result have been given. The result (actually a minor variation) was first proved by A. M. Liapunov [3] in 1940. The shortest proof was given by Lindenstrauss [1] who used the Krein–Milman theorem. See Remark 2. Other proofs appear in [2], [4], [6]. The referee has pointed out that Olech used an induction argument in a way similar to that used here in Case 2. See [9, p. 89] and [10, p. 42].

**2. A proof of the theorem.** For a convex set  $U \subset R^n$ , the *dimension  $d$  of  $U$*  ( $\dim U$ ) may be defined to be the dimension of the smallest hyperplane  $H$  that contains  $U$ . For  $U \subset H \subset R^n$ , define  $\partial_H U = \overline{H - U} \cap \overline{U}$ .

**LEMMA.** Let  $U \subset H$  be convex and let  $H$  be a hyperplane with  $p \in \partial_H U$ . Assume the dimension of  $U$  equals the dimension of  $H$ . Then there is a linear functional  $h: R^n \rightarrow R$  which is not constant on  $H$  and satisfies  $\sup_{x \in U} h(x) = h(p)$ .

See, for example, [7], for the result that such an  $h: H \rightarrow R$  exists. The domain of  $h$  can then easily be extended to  $R^n$ .

---

\* Received by the editors September 8, 1970, and in revised form December 4, 1970.

† Institute for Fluid Dynamics and Applied Mathematics, University of Maryland, College Park, Maryland 20740. This research was supported in part by a Columbia University undergraduate scholarship, by the National Science Foundation with a Graduate Fellowship and under Grants GP-9347 and GP-17601.

*Proof of Theorem.* Let  $J$  be a measurable subset of  $E$  and let  $\phi : J \rightarrow \{-1, +1\}$  be measurable. We will say then that  $(\phi, J)$  is a *restriction*. Let

$$\begin{aligned} \mathcal{F}(\phi, J) &= \{u \in \mathcal{F} : u(t) = \phi(t) \text{ for almost all } t \in J\}, \\ \mathcal{F}^*(\phi, J) &= \{u \in \mathcal{F}^* : u(t) = \phi(t) \text{ for almost all } t \in J\}. \end{aligned}$$

We actually will prove the following stronger result.

**PROPOSITION.** *Let  $(\phi, J)$  be a restriction. Then  $\mathcal{I}\mathcal{F}(\phi, J) = \mathcal{I}\mathcal{F}^*(\phi, J)$ .*

The proposition reduces to the theorem when  $J$  is allowed to be trivial (i.e., have measure 0). Note that  $\mathcal{F}(\phi, J)$  is a convex closed bounded subset of  $L^2(E; \mathbb{R}^n)$  so  $\mathcal{F}(\phi, J)$  is weakly compact and its continuous image  $\mathcal{I}\mathcal{F}(\phi, J)$  is therefore convex and compact. ( $\mathcal{I}$  is continuous and linear so it is weakly continuous.)

Choose some restriction  $(\phi, J)$  and write  $\mathcal{F}(\phi, J) = \mathcal{F}_0$  and  $\mathcal{F}^*(\phi, J) = \mathcal{F}_0^*$ . Since  $\mathcal{F}_0^* \subset \mathcal{F}_0$ , it suffices to prove that  $\mathcal{I}\mathcal{F}_0^* \supset \mathcal{I}\mathcal{F}_0$ . Choose  $p \in \mathcal{I}\mathcal{F}_0$ . We now prove  $p \in \mathcal{I}\mathcal{F}_0^*$ . Let  $d$  be the dimension of  $\mathcal{I}\mathcal{F}_0$ . Then  $\mathcal{I}\mathcal{F}_0$  lies in a  $d$ -dimensional hyperplane  $H$ . If  $d = 0$ ,  $\mathcal{I}\mathcal{F}_0$  contains only one point, so  $\{p\} = \mathcal{I}\mathcal{F}_0^*$  and the result is true. We now prove the proposition by induction on  $d$ .

*Assumption.* The proposition is true if  $\dim \mathcal{I}\mathcal{F}(\phi, J) < d$ .

*Case 1.* Suppose  $p \in \partial_H \mathcal{I}\mathcal{F}_0$ . Letting  $U = \mathcal{I}\mathcal{F}_0$ , choose  $h$  to be the linear function in the lemma. Let  $J_0 = \{t \in E - J : h(f(t)) \neq 0\}$  and  $J_1 = J_0 \cup J$ . Let

$$(3) \quad \phi_1(t) = \begin{cases} \phi(t) & \text{for } t \in J, \\ \text{sgn } h(f(t)) & \text{for } t \in J_0. \end{cases}$$

Write  $\mathcal{F}_1 = \mathcal{F}(\phi_1, J_1)$  and  $\mathcal{F}_1^* = \mathcal{F}^*(\phi_1, J_1)$ . For  $u \in \mathcal{F}_0$ ,

$$\begin{aligned} (4) \quad h(\mathcal{I}u) &= \int_E u(t)h(f(t)) = \int_{J_1} u(t)h(f(t)) \\ &= \int_J \phi(t)h(f(t)) + \int_{J_0} u(t)h(f(t)) \\ &\leq \int_J \phi(t)h(f(t)) + \int_{J_0} |h(f(t))| = \int_E \phi_1(t)h(f(t)). \end{aligned}$$

It may be seen that we have equality in (4) for  $u \in \mathcal{F}_0$  if and only if  $u(t) = \text{sgn } h(f(t))$  for almost all  $t \in J_0$ . Since  $h(p) = \sup_{u \in \mathcal{F}_0} h(\mathcal{I}u)$ , we have  $h(p) \geq h(\mathcal{I}u)$  for all  $u \in \mathcal{F}_0$ . For  $u \in \mathcal{F}_0$ ,  $h(p) = h(\mathcal{I}u)$  if and only if  $u \in \mathcal{F}_1$ . Since  $h$  is nonconstant on  $H$ , the set  $H_p = \{x \in H : h(x) = h(p)\}$  has dimension  $d - 1$ . Since  $\mathcal{I}\mathcal{F}_1 \subset H_p$ ,  $\dim \mathcal{I}\mathcal{F}_1 < d$ . Since  $p \in \mathcal{I}\mathcal{F}_1$  we may use the assumption and conclude  $p \in \mathcal{I}\mathcal{F}_1^*$ , which proves the result because  $\mathcal{I}\mathcal{F}_1^* \subset \mathcal{I}\mathcal{F}_0^*$ .

*Case 2.* Suppose  $p \in \text{int}_H \mathcal{I}\mathcal{F}_0$ , where  $H$  is the smallest hyperplane containing  $\mathcal{I}\mathcal{F}_0$ . We now show that by extending the definition of  $\phi$  to  $\phi_\tau$  which is defined on the larger set  $J \cup [0, \tau]$ , the set  $\mathcal{I}\mathcal{F}(\phi_\tau, J^\tau)$  decreases continuously as  $\tau$  is increased until for some  $\sigma$ ,  $p \in \partial_H \mathcal{I}\mathcal{F}(\phi_\sigma, J^\sigma)$ , reducing the problem to Case 1. For  $\tau \in E$  let  $J^\tau$  denote  $J \cup [0, \tau]$  and define  $\phi_\tau$  on  $J^\tau$  by  $\phi_\tau(t) = \phi(t)$  for  $t \in J$  and  $\phi_\tau(t) = +1$  for  $t \in [0, \tau] - J$ . Consider the restriction  $(\phi_\tau, J^\tau)$  and write  $\mathcal{F}_\tau$  for  $\mathcal{F}(\phi_\tau, J^\tau)$ . If  $\tau_1 > \tau_2$ , then  $\mathcal{F}_{\tau_1} \subset \mathcal{F}_{\tau_2} \subset \mathcal{F}_0$ .  $\mathcal{I}\mathcal{F}_\tau$  is of course a convex set and we may let  $\sigma (\leq T)$  be the supremum of  $\tau \in E$  for which  $p \in \mathcal{I}\mathcal{F}_\tau$ . We claim  $p \in \partial_H \mathcal{I}\mathcal{F}_\sigma$ .

If  $\tau_1 > \tau_2$  and  $p \in \mathcal{I}\mathcal{F}_{\tau_2}$  and, say  $p = \mathcal{I}u_p$  for some  $u_p \in \mathcal{F}_{\tau_2}$ , then there is  $u_1 \in \mathcal{F}_{\tau_1}$  such that  $u_1 = u_p$  except perhaps on  $[\tau_2, \tau_1]$ ; hence  $\mathcal{I}u_1 \in \mathcal{I}\mathcal{F}_{\tau_1}$  and if  $\delta$  is the distance of  $p$  from  $\mathcal{I}\mathcal{F}_{\tau_1}$ ,

$$\delta \leq |\mathcal{I}u_1 - p| \leq \int_{\tau_2}^{\tau_1} |f(s)| |u_1(s) - u_p(s)| ds \leq 2|\tau_1 - \tau_2|M,$$

where  $M = \sup_E |f|$ . From this "continuity" of  $\mathcal{I}\mathcal{F}_{\tau}$ ,  $p \in \partial_H \mathcal{I}\mathcal{F}_{\sigma}$ . By considering the restriction  $(\phi_{\sigma}, J^{\sigma})$  instead of  $(\phi, J)$ , we have reduced the situation to Case 1. The proof of the theorem is complete since the induction argument for dimension  $d$  is complete.

*Remark 1.* It is known that the theorem is true if  $f(s)ds$  is replaced by a nonatomic vector measure but is false for more general measures. This proof uses the fact that  $f(s)ds$  is a nonatomic measure in Case 2 since otherwise  $\mathcal{I}\mathcal{F}(\phi_{\tau}, J^{\tau})$  does not have to be a continuous set function of  $\tau$  (in the Hausdorff metric).

*Remark 2.* This proof is similar to Lindenstrauss' proof in certain ways. It differs primarily in Case 2 where we have avoided the Krein–Milman theorem.

**Acknowledgment.** The author would like to thank H. Hermes and L. Markus for their helpful comments.

#### REFERENCES

- [1] JORAM LINDENSTRAUSS, *A short proof of Liapunoff's convexity theorem*, J. Math. Mech., 15 (1966), pp. 971–972.
- [2] H. HALKIN, *Some further generalizations of a theorem of Lyapunov*, Arch. Rational Mech. Anal., 17 (1964), pp. 272–277.
- [3] A. LIAPOUNOFF, *Sur les fonctions-vecteurs complètement additives*, Bull. Acad. Sci. URSS Sér. Math., 4 (1940), pp. 465–478.
- [4] P. R. HALMOS, *The range of a vector measure*, Bull. Amer. Math. Soc., 54 (1948), pp. 416–421.
- [5] J. P. LASALLE, *The time optimal control problem*, Contributions to the Theory of Nonlinear Oscillations, vol. 5, Princeton Univ. Press, Princeton, N.J., 1959, pp. 1–24.
- [6] C. OLECH, *On the range of an unbounded vector-valued measure*, Math. Systems Theory, 2 (1968), pp. 251–256.
- [7] E. B. LEE AND L. MARKUS, *Foundations of Optimal Control*, John Wiley, New York, 1967.
- [8] J. ÉTIENNE, *Sur une démonstration du bang-bang principe*, Bull. Soc. Roy. Sci. Liège, 37 (1968), no. 11–12, pp. 551–556.
- [9] C. OLECH, *Extremal solutions of a control system*, J. Differential Equations, 2 (1966), pp. 74–101.
- [10] A. BALAKRISHNAN AND LUCIEN W. NEUSTADT, eds., *Mathematical Theory of Control*, Academic Press, New York, 1967.

## ON THE SOLUTIONS OF A STOCHASTIC CONTROL SYSTEM\*

TYRONE DUNCAN† AND PRAVIN VARAIYA‡

**Abstract.** The control system considered in this paper is modeled by the stochastic differential equation

$$dx(t, \omega) = f(t, x(\cdot, \omega), u(t, \omega)) dt + dB(t, \omega),$$

where  $B$  is  $n$ -dimensional Brownian motion, and the control  $u$  is a nonanticipative functional of  $x(\cdot, \omega)$  taking its values in a fixed set  $U$ . Under various conditions on  $f$  it is shown that for every admissible control a solution is defined whose law is absolutely continuous with respect to the Wiener measure  $\mu$ , and the corresponding set of densities on the space  $C$  forms a strongly closed, convex subset of  $L^1(C, \mu)$ . Applications of this result to optimal control and two-person, zero-sum differential games are noted. Finally, an example is given which shows that in the case where only some of the components of  $x$  are observed, the set of attainable densities is not weakly closed in  $L^1(C, \mu)$ .

**1. Introduction and contents.** A stochastic control problem is defined by the specification of the stochastic differential equation which models the system dynamics, the information available to the controller and the corresponding set of admissible control laws, and the cost incurred by each control law. Of theoretical interest is the "existence" problem, which means determining in terms of the above three defining characteristics a class of control problems for which there exist control laws achieving minimum cost. Published results ([1], [2], [3], see especially the excellent survey article [4] of Fleming) differ from one another and are not usually comparable because either the models are different or the set of admissible control laws is different.

There are two basic steps involved in obtaining an existence result. The first step involves determining conditions which guarantee that a solution of the stochastic differential equation is defined for every admissible control law. The next step involves the search for a topology under which the set of solutions (or an equally good substitute) is compact, and the cost function is lower semicontinuous. Thus, for instance, Fleming and Nisio [1] consider stochastic differential equations of the form

$$dx(t) = f(t, x(\cdot)) u(t) dt + \sigma(t, x(\cdot), B(\cdot)) dB(t), \quad 0 \leq t < \infty,$$

where  $u(t)$  is any process taking values in the unit cube, and independent of future increments  $B(t_2) - B(t_1)$ ,  $t \leq t_1 \leq t_2$ , of the Brownian motion  $B$ . Various conditions on  $f$ ,  $\sigma$  are imposed to guarantee a solution for every admissible control. It is then shown that the set of laws of all the solutions of the differential equation

---

\* Received by the editors August 25, 1970, and in revised form January 11, 1971.

† Computer Information and Control Engineering Department, University of Michigan, Ann Arbor, Michigan 48104. The work of this author was done while he was visiting the Department of Electrical Engineering and Computer Sciences at the University of California, Berkeley, and was supported by the Office of Naval Research under Contract N0014-67-A-081-0021, NR041-376.

‡ Department of Electrical Engineering and Computer Sciences, Electronics Research Laboratory, University of California, Berkeley, California 94720. The work of this author was supported by the U.S. Army Research Office—Durham under Contract DAHC04-67-0046.

corresponding to the different control laws is compact in the Prohorov metric. Beneš [3] considers stochastic differential equations of the form

$$(1) \quad dx(t, \omega) = f(t, x(\cdot, \omega), u(t, \omega))dt + dB(t, \omega), \quad 0 \leq t \leq 1,$$

where  $f$  is measurable with respect to its arguments and continuous in  $u$ . The control law is any nonanticipative, measurable functional  $u(t, \omega) = \psi(t, x(\cdot, \omega))$  which takes values in a compact set  $U$ . He assumes that  $f$  satisfies a linear growth condition

$$|f(t, x(\cdot, \omega), u)|^2 \leq K(1 + \|x(\cdot, \omega)\|^2),$$

where  $\|x(\cdot, \omega)\| = \sup \{|x(t, \omega)| \mid 0 \leq t \leq 1\}$ .

The existence of a solution to (1) for every control law is guaranteed by a result of Girsanov [5] (see Corollary 3 below). The resulting law is absolutely continuous with respect to the Wiener measure  $\mu$  on the space  $C$  of all continuous functions from  $[0, 1]$  into  $R^n$ . Beneš shows that if  $f(t, x(\cdot, \omega), U)$  is convex for every  $t \in [0, 1]$  and  $x(\cdot, \omega) \in C$ , then the set of densities corresponding to all the admissible control laws is a convex and strongly closed (hence weakly compact) subset of  $L^1(C, \mu)$ .

In this paper, we show that the above result holds if the linear growth condition is replaced by the growth condition

$$(2) \quad |f(t, x(\cdot, \omega), u)| \leq f_0(\|x(\cdot, \omega)\|),$$

where  $f_0: R \rightarrow R$  is increasing, and the condition

$$(3) \quad \int_C \exp \left[ \int_0^1 \langle f(t, B, u), dB(t) \rangle - \frac{1}{2} \int_0^1 |f(t, B, u)|^2 dt \right] \mu(dB) = 1$$

for every admissible control law. An example is given to show that (2) does not imply (3). The linear growth condition implies (3) (see Corollary 3). Condition (3) also follows from (2), if the drift term  $f$  in (1) has a delay (see Corollary 4). Finally we show that in the important case where the control is allowed to depend only on some components of the state  $x$ , the set of densities is not always weakly closed in  $L^1(C, \mu)$ .

In § 2 we give some preliminary results and definitions, and in § 3 we present the main result on weak compactness of the attainable densities. In § 4 we give conditions which guarantee (3), in § 5 we present applications to optimal control and stochastic differential games, and in the final section we present the negative example for the problem with partial observations.

**2. Preliminaries.** In the main, we adopt the notations and definitions of Beneš [3]. Consider the stochastic differential equation

$$(1') \quad \begin{aligned} dx(t) &= f(t, x, u(t, x)) dt + dB(t), & 0 \leq t \leq 1, \\ x(0) &= 0, \end{aligned}$$

where  $B(t)$  is a standard  $n$ -dimensional Brownian motion process with continuous sample paths,  $x(t)$  is the state of the system and  $u(t, x)$  is the control law which takes values in a compact subset  $U$  of  $R^m$ . To state the precise conditions which  $f, u$  must satisfy we need the following definition.

DEFINITION 1. (a) Let  $C$  be the Banach space of all continuous functions  $z:[0, 1] \rightarrow R^n$  with norm  $\|z\| = \max \{|z(t)| | 0 \leqq t \leqq 1\}$ , where  $|y|$  is the Euclidean norm of  $y \in R^n$ .

(b) For each  $t \in [0, 1]$  let  $\mathcal{S}_t$  be the smallest  $\sigma$ -field of subsets of  $C$  which contains all sets of the form  $\{z | z(\tau) \in A\}$ , where  $\tau \in [0, t]$  and  $A$  is a Borel subset of  $R^n$ .

(c) Let  $\mathcal{S} = \mathcal{S}_1$ .

We shall define the solution of (1) in such a way that the sample paths of  $x$  are continuous (and have no explosions), so that  $f$  is a map from  $[0, 1] \times C \times U \rightarrow R^n$ . We impose throughout the following conditions on  $f$ .

C1.  $f$  is measurable with respect to the product  $\sigma$ -algebra  $\mathcal{B} \otimes \mathcal{S} \otimes \mathcal{B}_U$ , where  $\mathcal{B}$  ( $\mathcal{B}_U$ ) is the set of Borel measurable subsets of  $[0, 1]$  ( $U$ ).

C2. For fixed  $t \in [0, 1)$ ,  $f(t, \cdot, \cdot)$  is measurable with respect to the product  $\sigma$ -algebra  $\mathcal{S}_t \otimes \mathcal{B}_U$ .

C3. For fixed  $(t, z) \in [0, 1] \times C$ ,  $f(t, z, \cdot)$  is continuous on  $U$ .

C4. There exists an increasing function  $f_0: R \rightarrow R$  such that  $|f(t, z, u)| \leqq f_0(\|z\|)$  for all  $(t, z, u)$ .

C5.  $f(t, z, U)$  is closed and convex for every  $(t, z)$ .

DEFINITION 2. (a) An *admissible control (law)* is any map  $u:[0, 1] \times C \rightarrow U$  which is measurable with respect to  $\mathcal{B} \otimes \mathcal{S}$ , and for each fixed  $t \in [0, 1]$ ,  $u(t, \cdot)$  is measurable with respect to  $\mathcal{S}_t$ . Let  $\mathcal{U}$  be the set of all admissible control laws.

(b) For each  $u \in \mathcal{U}$ , the *drift corresponding to  $u$*  is the function  $g = g_u:[0, 1] \times C \rightarrow R^n$  defined by

$$g(t, z) = f(t, z, u(t, z)).$$

Let  $\mathcal{G} = \{g_u | u \in \mathcal{U}\}$ .

(c) For  $g \in \mathcal{G}$  and  $N \geqq 0$ , let  $g^N:[0, 1] \times C \rightarrow R^n$  be defined by

$$g^N(t, z) = \begin{cases} g(t, z) & \text{if } |z(\tau)| \leqq N \text{ for } \tau \leqq t, \\ 0 & \text{otherwise.} \end{cases}$$

Let  $\mathcal{G}^N = \{g^N | g \in \mathcal{G}\}$ .

DEFINITION 3. A function  $\psi:[0, 1] \times C \rightarrow R^n$  will be said to be *causal* if it is  $\mathcal{B} \otimes \mathcal{S}$  measurable, and if for each fixed  $t \in [0, 1]$ ,  $\psi(t, \cdot)$  is measurable with respect to  $\mathcal{S}_t$ .

From [6, Lemmas 1, 2] we can obtain the following useful characterization of  $\mathcal{G}$ . Condition C3 is needed only for Lemma 1. The reader should be warned that the proof of the "only if" part of Lemma 1 involves a nontrivial synthesis problem (lemma of Fillipov).

LEMMA 1. A causal function  $g:[0, 1] \times C \rightarrow R^n$  belongs to  $\mathcal{G}$  if and only if  $g(t, z) \in f(t, z, U)$  for all  $(t, z)$ .

It will prove convenient to work with sets larger than  $\mathcal{G}, \mathcal{G}^N$ .

DEFINITION 4. Let  $\Phi$  be the set of all causal maps  $\phi:[0, 1] \times C \rightarrow R^n$  such that  $|\phi(t, z)| \leqq f_0(\|z\|)$  for all  $(t, z)$ . Let  $\Phi^N = \{\phi | \phi \in \Phi, |\phi(t, z)| \leqq N \text{ for all } (t, z)\}$ .

Throughout the rest of this paper let  $\Omega$  be a fixed space and let  $\mathcal{A}_t, 0 \leqq t \leqq 1$ , be a fixed, increasing family of  $\sigma$ -fields of subsets of  $\Omega$ . Let  $\mathcal{A} = \mathcal{A}_1$ . We say that  $z(t)$  or  $z(t, \omega), 0 \leqq t \leqq 1$ , is a family of  $n$ -dimensional random variables on  $(\Omega, \mathcal{A}_t)$  if for each  $t, z(t, \cdot)$  is a map from  $\Omega$  into  $R^n$  which is measurable with respect to  $\mathcal{A}_t$ .

We shall need to consider various probability measures on  $\mathcal{A}$ . If  $z(t)$ ,  $0 \leq t \leq 1$ , is a family of  $n$ -dimensional random variables on  $(\Omega, \mathcal{A})$  and we wish to consider the stochastic process generated by  $z(t)$  corresponding to a particular probability measure  $P$  on  $\mathcal{A}$ , we will say that  $z(t)$ ,  $0 \leq t \leq 1$ , is an  $n$ -dimensional stochastic process on  $(\Omega, \mathcal{A}, P)$ . Finally let  $P_0$  be a distinguished probability measure on  $\mathcal{A}$ , and let  $x(t, \omega)$ ,  $0 \leq t \leq 1$ , be a fixed  $n$ -dimensional, Brownian motion process on  $(\Omega, \mathcal{A}, P_0)$  with almost all sample paths  $x(\cdot, \omega) \in C$ . We assume that the  $\sigma$ -fields  $\mathcal{A}_t$  are complete with respect to  $P_0$ .

DEFINITION 5. Let  $\psi : [0, 1] \times C \rightarrow R^n$  be a causal function such that

$$(4) \quad \int_0^1 |\psi(t, z)|^2 dt < \infty \quad \text{for all } z \in C.$$

Then  $\zeta^t(\psi)$ ,  $0 \leq t \leq 1$ , is the stochastic process on  $(\Omega, \mathcal{A}_t, P_0)$  with continuous sample paths, defined by

$$(5) \quad \zeta^t(\psi) = \int_0^t \langle \psi(\tau, x), dx(\tau) \rangle - \frac{1}{2} \int_0^t |\psi(\tau, x)|^2 d\tau.$$

For convenience, let  $\zeta(\psi) = \zeta^1(\psi)$ . (In (5), the first integral is to be interpreted as an Ito stochastic integral.)

The results of this section are immediate consequences of the work of Girsanov [5].

THEOREM 1 (Existence). Let  $\psi : [0, 1] \times C \rightarrow R^n$  be a causal function such that (4) holds.

(i) Then,

$$\int_{\Omega} \exp [\zeta(\psi)] P_0(d\omega) \leq 1.$$

(ii) Suppose

$$(6) \quad \int_{\Omega} \exp [\zeta(\psi)] P_0(d\omega) = 1,$$

and define the probability measure  $P_{\psi}$  on  $\mathcal{A}$  by

$$P_{\psi}(A) = \int_A \exp [\zeta(\psi)] P_0(d\omega), \quad A \in \mathcal{A}.$$

Then the stochastic process  $B(t)$  defined on  $(\Omega, \mathcal{A}_t, P_{\psi})$  by

$$B(t, \omega) = x(t, \omega) - \int_0^t \psi(\tau, x(\cdot, \omega)) d\tau, \quad 0 \leq t \leq 1,$$

is a Brownian motion.

(iii) If  $\psi$  is bounded, then (6) holds.

Proof. Parts (i), (ii) and (iii) are immediate consequences of Lemma 2, Theorem 1, and Lemma 1, respectively, of [5].



Theorem 1 immediately gives us a sufficient condition for the existence of a solution to (1). For, let  $u \in \mathcal{U}$ , and let  $g$  be the drift corresponding to  $u$ . If

$$\int_{\Omega} \exp [\zeta(g)] P_0(d\omega) = 1,$$

then the stochastic process  $x(t)$  on  $(\Omega, \mathcal{A}_t, P_g)$  satisfies the equation

$$x(t) = \int_0^t g(\tau, x) d\tau + \text{Brownian motion.}$$

LEMMA 2. Let  $\psi \in \Phi$ . Let  $y(t)$ ,  $0 \leq t \leq 1$ , be a stochastic process on  $(\Omega, \mathcal{A}_t, P)$  with continuous sample paths, such that the stochastic process  $B(t)$  on  $(\Omega, \mathcal{A}_t, P)$  defined by

$$(7) \quad B(t) = y(t) - \int_0^t \psi(\tau, y) d\tau, \quad 0 \leq t \leq 1,$$

is a Brownian motion. Then, the measure  $\nu$  induced by  $y$  on  $(C, \mathcal{S})$  is mutually absolutely continuous with respect to the Wiener measure  $\mu$ , and

$$(8) \quad \frac{d\mu}{d\nu}(y) = \exp \left[ - \int_0^1 \langle \psi(t, y), dB(t) \rangle - \frac{1}{2} \int_0^1 |\psi(t, y)|^2 dt \right].$$

Proof. Since  $|\psi(\cdot, z)| \leq f_0(\|z\|)$ , it follows from Lemma 7 of [5], that the measure  $\mu$  on  $(C, \mathcal{S})$  defined by

$$\mu(S) = \int_S \exp \left[ - \int_0^1 \langle \psi, dB \rangle - \frac{1}{2} \int_0^1 |\psi|^2 dt \right] d\nu$$

coincides with the Wiener measure. It is easy to see that

$$\exp \left[ - \int_0^1 \langle \psi, dB \rangle - \frac{1}{2} \int_0^1 |\psi|^2 dt \right] > 0,$$

$\nu$ -almost everywhere. The result follows.

COROLLARY 1. Let  $\psi \in \Phi$ , and let  $y(t)$ ,  $0 \leq t \leq 1$ , satisfy the hypothesis of Lemma 2. Then

$$(9) \quad \int_C \exp \left[ \int_0^1 \langle \psi(t, z), dz(t) \rangle - \frac{1}{2} \int_0^1 |\psi(t, z)|^2 dt \right] \mu(dz) = 1.$$

COROLLARY 2. Let  $\psi \in \Phi$ , and let  $y(t)$ ,  $0 \leq t \leq 1$ , satisfy the hypothesis of Lemma 2. Then the measure  $\nu$  on  $(C, \mathcal{S})$  induced by  $y$  is uniquely specified by  $\psi$  and is given by

$$(10) \quad \nu(S) = \int_S \exp \left[ \int_0^1 \langle \psi(t, z), dz(t) \rangle - \frac{1}{2} \int_0^1 |\psi(t, z)|^2 dt \right] \mu(dz).$$

Proof. The corollaries follow from (8) and the identity  $dB = dy - \psi(t, y) dt$ .

### 3. Main results.

DEFINITION 6. For any subset  $\Sigma \subset \Phi$ , let  $\mathcal{D}(\Sigma)$  be the subset of  $L^1(\Omega, \mathcal{A}, P_0)$  defined by

$$\mathcal{D}(\Sigma) = \{ \exp \zeta(\phi) | \phi \in \Sigma \}.$$

PROPOSITION 1.  $\mathcal{D}(\Phi^N)$  is a bounded subset of  $L^2(\Omega, \mathcal{A}, P_0)$ .

Proof. If  $\phi \in \Phi^N$ , then by definition,  $|\phi| \leq N$ . By Lemma 1 of [5], it follows that

$$\int_{\Omega} \exp 2\zeta^t(\phi) P_0(d\omega) \leq \exp tN^2.$$

For the rest of this paper let  $E_0$  denote expectation with respect to the probability measure  $P_0$ . Also if  $\gamma \in L^1(\Omega, \mathcal{A}, P_0)$ , then  $E_0(\gamma|\mathcal{A}_t)$  denotes the conditional expectation of  $\gamma$  with respect to  $\mathcal{A}_t$ .

The proofs of Lemmas 3 and 4 are simple modifications of the proofs of Theorems 4 and 3, respectively, of [3]. They are presented here for completeness and because we shall need to refer to parts of the proofs later.

LEMMA 3.  $\mathcal{D}(\Phi^N)$  is a closed subset of  $L^2(\Omega, \mathcal{A}, P_0)$ .

Proof. Let  $\phi_n, n = 1, 2, 3, \dots$ , be a sequence from  $\Phi^N$  and let  $\rho$  be such that

$$(11) \quad \lim_{n \rightarrow \infty} E_0|\rho - \exp \zeta(\phi_n)|^2 = 0$$

and

$$(12) \quad \lim_{n \rightarrow \infty} \exp \zeta(\phi_n) = \rho \quad \text{a.s. } P_0.$$

First of all  $\rho > 0$  a.s.  $P_0$ . Because, let  $A = \{\omega|\rho(\omega) = 0\}$ . Then from (12),

$$(13) \quad \lim_{n \rightarrow \infty} \zeta(\phi_n)(\omega) = -\infty \quad \text{for } \omega \in A.$$

Also,

$$\zeta(\phi_n) = \int_0^1 \langle \phi_n(s, x), dx(s) \rangle - \frac{1}{2} \int_0^1 |\phi_n(s, x)|^2 ds$$

and  $|\phi_n| \leq N$ , so that from (13),

$$\lim_{n \rightarrow \infty} \int_0^1 \langle \phi_n(s, x), dx(s) \rangle = -\infty \quad \text{on } A.$$

But

$$E_0 \left[ \int_0^1 \langle \phi_n, dx(s) \rangle \right]^2 = E_0 \int_0^1 |\phi_n(s)|^2 ds \leq N^2,$$

so that  $P_0(A) = 0$ . By Ito's representation [7], there is a causal map

$$\psi : [0, 1] \times C \rightarrow R^n$$

with

$$\int_0^1 |\psi(t, z)|^2 dt < \infty$$

for  $z$  in  $C$ , such that

$$\rho = 1 + \int_0^1 \langle \psi(t, x), dx(t) \rangle \quad \text{a.s. } P_0.$$

Let

$$\rho(t) = E_0(\rho|\mathcal{A}_t) = 1 + \int_0^t \langle \psi(s, x), dx(s) \rangle.$$

Then, by Jensen's inequality,

$$\int_0^1 E_0|\rho(t) - \exp \zeta^t(\phi_n)|^2 dt \leq \int_0^1 E_0|\rho - \exp \zeta(\phi_n)|^2 dt,$$

which converges to zero, so that taking subsequences if necessary we can assume that

$$(14) \quad \rho(t) = \lim_{n \rightarrow \infty} \exp \zeta^t(\phi_n) \quad \text{a.s. } l \otimes P_0,$$

where  $l$  denotes Lebesgue measure on  $[0, 1]$ . Next, by Ito's differentiation rule,

$$\exp \zeta(\phi_n) = 1 + \int_0^1 \exp \zeta^t(\phi_n) \langle \phi_n(t), dx(t) \rangle \quad \text{a.s. } P_0$$

so that

$$E_0 \int_0^1 |\exp \zeta^t(\phi_n) \phi_n(t) - \psi(t)|^2 dt = E_0 |\exp \zeta(\phi_n) - \rho|^2$$

converges to zero, and therefore, taking subsequences if necessary, we can assume that

$$\psi(t) = \lim_{n \rightarrow \infty} \exp \zeta^t(\phi_n) \phi_n(t) \quad \text{a.s. } l \otimes P_0.$$

Since  $\rho(t) > 0$  a.s.  $P_0$ , we see using (14) that

$$(15) \quad \psi(t)/\rho(t) = \lim_{n \rightarrow \infty} \phi_n(t) \quad \text{a.s. } l \otimes P_0.$$

It follows that there is a causal map  $\phi : [0, 1] \times C \rightarrow R^n$ ,

$$\phi(t, x(\cdot, \omega)) = \lim_{n \rightarrow \infty} \phi_n(t, x(\cdot, \omega)) \quad \text{a.s. } l \times P_0$$

and

$$\rho(t) = 1 + \int_0^t \rho(s) \langle \phi(s, x), dx(s) \rangle.$$

From Ito's differential rule we see that

$$d(\log \rho(t)) = \langle \phi(t), dx(t) \rangle - \frac{1}{2} |\phi(t)|^2 dt,$$

and hence,

$$\rho = \exp \zeta(\phi).$$

Because of (15) we can assume that  $|\phi| \leq N$ , so that the lemma is proved.

LEMMA 4.  $\mathcal{D}(\Phi^N)$  is a convex set.

Proof. Let  $\phi_i \in \Phi^N$ ,  $\lambda_i \geq 0$ ,  $i = 1, 2$ , with  $\lambda_1 + \lambda_2 = 1$ . By Ito's differentiation rule,

$$d[\exp \zeta^t(\phi_i)] = \exp \zeta^t(\phi_i) \langle \phi_i(t), dx(t) \rangle.$$

Define

$$\rho(t) = \lambda_1 \exp \zeta'( \phi_1 ) + \lambda_2 \exp \zeta'( \phi_2 ).$$

Then

$$d\rho(t) = \lambda_1 \exp \zeta'( \phi_1 ) \langle \phi_1(t), dx(t) \rangle + \lambda_2 \exp \zeta'( \phi_2 ) \langle \phi_2(t), dx(t) \rangle$$

which we can rewrite as

$$(16) \quad d\rho(t) = \rho(t) \langle \phi(t), dx(t) \rangle,$$

where

$$\phi(t) = \frac{\lambda_1 \exp \zeta'( \phi_1 )}{\sum_{i=1}^2 \lambda_i \exp \zeta'( \phi_i )} \phi_1(t) + \frac{\lambda_2 \exp \zeta'( \phi_2 )}{\sum_{i=1}^2 \lambda_i \exp \zeta'( \phi_i )} \phi_2(t).$$

Evidently  $\phi \in \Phi^N$  since  $\phi(t, z)$  is a convex combination of  $\phi_1(t, z)$  and  $\phi_2(t, z)$ . By Ito's differentiation rule from (16) we obtain

$$(17) \quad d(\log \rho(t)) = \langle \phi(t), dx(t) \rangle - \frac{1}{2} |\phi(t)|^2 dt$$

and  $\rho(0) = 1$  from (14) so that integrating (17) yields

$$\log \rho(t) = \int_0^t \langle \phi(s), dx(s) \rangle - \frac{1}{2} \int_0^t |\phi(s)|^2 ds.$$

Hence  $\rho(1) = \exp \zeta(\phi)$  and the lemma is proved.

We now state our main result and develop the proof through a sequence of lemmas.

**THEOREM 2.** (i)  $\mathcal{D}(\mathcal{G})$  is a convex set.

(ii) Let

$$\mathcal{G}^0 = \{g | g \in \mathcal{G}, E_0(\exp \zeta(g)) = 1\}.$$

Then,  $\mathcal{D}(\mathcal{G}^0)$  is a closed and convex subset of  $L^1(\Omega, \mathcal{A}, P_0)$ .

We shall develop the proof through a sequence of lemmas.

**LEMMA 5.**  $\mathcal{D}(\mathcal{G})$  is convex.

*Proof.* Let  $g_i(t, z) = f(t, z, u_i(t, z))$  with  $u_i \in \mathcal{U}_i$ ,  $i = 1, 2$ , and let  $\lambda_i \geq 0$ , with  $\lambda_1 + \lambda_2 = 1$ . By Ito's differentiation rule,

$$d(\exp \zeta'(g_i)) = \exp \zeta'(g_i) \langle g_i(t), dx(t) \rangle, \quad i = 1, 2.$$

Define

$$\rho(t) = \lambda_1 \exp \zeta'(g_1) + \lambda_2 \exp \zeta'(g_2).$$

Then if we repeat the proof of Lemma 4 we can conclude that (noting  $\rho(0) = 1$ )

$$\rho(1) = \lambda_1 \exp \zeta(g_1) + \lambda_2 \exp \zeta(g_2) = \exp \zeta(\phi),$$

where  $\phi(t, z)$  is a convex combination of  $g_1(t, z)$  and  $g_2(t, z)$ . Since  $g_i(t, z) \in f(t, z, U)$ , and since this set is convex by condition C5, we see that

$$\phi(t, z) \in f(t, z, U)$$

and hence  $\phi \in \mathcal{G}$  by Lemma 1. The lemma is proved.

LEMMA 6.  $\mathcal{D}(\mathcal{G}^0)$  is convex.

*Proof.* The set

$$\mathcal{R} = \{\rho | \rho \in L^1(\Omega, \mathcal{A}, P_0), \rho \geq 0, E_0\rho = 1\}$$

is convex, and

$$\mathcal{D}(\mathcal{G}^0) = \mathcal{D}(\mathcal{G}) \cap \mathcal{R},$$

so that the result follows from Lemma 5.

Next let  $g_n, n = 1, 2, \dots$ , be a sequence from  $\mathcal{G}^0$ , and let  $\rho$  be in  $L^1(\Omega, \mathcal{A}, P_0)$  such that

$$(18) \quad \lim_{n \rightarrow \infty} \exp \zeta(g_n) = \rho \quad \text{a.s. } P_0 \quad \text{and in } L^1(\Omega, \mathcal{A}, P_0).$$

For each positive integer  $N$ , let

$$g_n^N(t, z) = \begin{cases} g_n(t, z) & \text{if } |z(\tau)| \leq N \quad \text{for } \tau \leq t, \\ 0 & \text{otherwise,} \end{cases}$$

and for  $N = 1, 2, 3, \dots$ , inductively select subsequences  $g_k^N, k \in K_N$ , and  $\phi^N \in \Phi^N$  as follows:

For  $N = 1$ , let  $g_k^1, k \in K_1$ , be a subsequence of  $g_n^1, n = 1, 2, 3, \dots$ , and let  $\phi^1 \in \Phi^1$  be such that

$$\exp \zeta(\phi^1) = \text{w. lim}_{k \in K_1} \exp \zeta(g_k^1).$$

(Here and in the remainder w. lim means the weak limit in  $L^-(\Omega, \mathcal{A}, P_0)$ .) From Lemmas 4 and 5,  $\mathcal{D}(\Phi^N)$  is a weakly, sequentially compact subset of  $L^2(\Omega, \mathcal{A}, P)$  and  $g_n^N \in \Phi^N$  so that the above selection makes sense.

Suppose  $g_k^N, k \in K_N$ , and  $\phi^N \in \Phi^N$  are defined. Then let  $g_k^{N+1}, k \in K_{N+1}$ , be a sequence of  $g_n^N, k \in K_N$ , and let  $\phi^{N+1} \in \Phi^{N+1}$  be such that

$$\exp \zeta(\phi^{N+1}) = \text{w. lim}_{k \in K_{N+1}} \exp \zeta(g_k^{N+1}).$$

LEMMA 7. Let  $C_N^t = \{z | z \in C, |z(\tau)| \leq N \text{ for } \tau \leq t\}$ . Then for  $i \geq 0$ ,

$$\phi^{N+1}(t, z) = \phi^N(t, z) \quad \text{for } 0 \leq t \leq 1, \quad z \in C_N^1.$$

*Proof.* First of all from

$$\exp \zeta(\phi^N) = \text{w. lim}_{k \in K_N} \exp \zeta(g_k^N),$$

it is immediate that

$$E_0(\exp \zeta(\phi^N) | \mathcal{A}_t) = \text{w. lim}_{k \in K_N} E_0(\exp \zeta(g_k^N) | \mathcal{A}_t).$$

Secondly since

$$E_0(\zeta(\phi^N)) = 1, \quad E_0(\zeta(g_k^N)) = 1,$$

it follows that a.s.  $P_0$ ,

$$\exp \zeta^t(\phi^N) = E_0(\exp \zeta(\phi^N) | \mathcal{A}_t), \quad \exp \zeta^t(g_k^N) = E_0(\exp \zeta(g_k^N) | \mathcal{A}_t),$$

and hence,

$$(19) \quad \exp \zeta^i(\phi^N) = \text{w. lim}_{k \in K_N} \exp \zeta^i(g_k^N).$$

Next let

$$\Omega'_N = \{\omega | \omega \in \Omega, x(\cdot, \omega) \in C'_N\}.$$

By definition, for  $i \geq 0$ ,

$$g_k^N(\tau, x(\cdot, \omega)) = g_k^{N+i}(\tau, x(\cdot, \omega)) \quad \text{for } \tau \leq t, \quad \omega \in \Omega'_N,$$

so that from (19) for  $i \geq 0$ ,

$$\exp \zeta^i(\phi^{N+i})(\omega) = \exp \zeta^i(\phi^N)(\omega), \quad \tau \leq t, \quad \omega \in \Omega'_N.$$

The result now follows if we note that

$$\begin{aligned} 0 &= \int_{\Omega'_N} |\exp \zeta(\phi^{N+i}) - \exp \zeta(\phi^N)|^2 P_0(d\omega) = \int_{\Omega'_N} \left[ \int_0^1 \exp \zeta^i(\phi^{N+i}) \langle \phi^{N+i}(t), dx(t) \rangle \right. \\ &\quad \left. - \int_0^1 \exp \zeta^i(\phi^N) \langle \phi^N(t), dx(t) \rangle \right]^2 P_0(d\omega) \\ &= \int_{\Omega'_N} \left[ \int_0^1 \exp 2\zeta^i(\phi^N) |\phi^{N+i}(t, x(\cdot, \omega)) - \phi^N(t, x(\cdot, \omega))|^2 dt \right] P_0(d\omega) \end{aligned}$$

so that since  $\exp \zeta^i(\phi^N) > 0$  a.s.  $P_0$ , we must have

$$\int_{\Omega'_N} \left[ \int_0^1 |\phi^{N+i}(t, x(\cdot, \omega)) - \phi^N(t, x(\cdot, \omega))|^2 dt \right] P_0(d\omega) = 0,$$

and the lemma is proved.

Because of Lemma 7 we can define a causal function  $\phi : [0, 1] \times C \rightarrow R^n$  such that

$$(20) \quad \phi(t, z) = \phi^{N+i}(t, z) \quad \text{for } 0 \leq t \leq 1, \quad \|z\| \leq N, \quad i \geq 0.$$

From the proof of Lemma 8, and from (18) it follows that

$$\rho = \exp \zeta(\phi) \quad \text{a.s. } P_0.$$

Lemma 8 completes the proof of Theorem 2.

LEMMA 8.

$$\phi \in \mathcal{G}.$$

*Proof.* Because of (20) and Lemma 1 it is enough to show that

$$(21) \quad \phi^N(t, z) \in f(t, z, U) \quad \text{for } 0 \leq t \leq 1, \quad \|z\| \leq N.$$

Recall that

$$\exp \zeta(\phi^N) = \text{w. lim}_{k \in K_1} \exp \zeta(g_k^N).$$

From the properties of weak  $L^2$ -convergence it is known that there is a convex combination of the  $\exp \zeta(g_k^N)$  which converges to  $\exp \zeta(\phi^N)$  in the  $L^2$ -norm topology.

More precisely, for each  $n$ , there are nonnegative numbers  $\lambda_1^n, \dots, \lambda_n^n$  with  $\lambda_1^n + \dots + \lambda_n^n = 1$  such that

$$(22) \quad \lim_{n \rightarrow \infty} E_0 \left| \exp \zeta(\phi^N) - \sum_{i=1}^n \lambda_i^n \exp \zeta(g_i^N) \right|^2 = 0.$$

Let

$$h_n(t) = \sum_{i=1}^n \lambda_i^n \exp \zeta(g_i^N).$$

Repeating the proof of Lemma 4, we can conclude that

$$h_n(t) = \exp \zeta^t(\eta_n) \quad \text{a.s. } P_0,$$

where  $\eta_n(t, z)$  is a convex combination of  $g_1^N(t, z), \dots, g_n^N(t, z)$ . In particular, from the convexity of  $f(t, z, U)$  and the fact that  $g_i^N(t, z) = g_i(t, z) \in f(t, z, U)$  for  $\|z\| \leq N$ , it follows that for  $\|z\| \leq N$ ,

$$(23) \quad \eta_n(t, z) \in f(t, z, U).$$

Next

$$\begin{aligned} E_0 |\exp \zeta(\phi^N) - \exp \zeta(\eta_n)|^2 &= E_0 \left| \int_0^1 \exp \zeta^t(\phi^N) \langle \phi^N(t), dx(t) \rangle \right. \\ &\quad \left. - \int_0^1 \exp \zeta^t(\eta_n) \langle \eta_n(t), dx(t) \rangle \right|^2 \\ &= E_0 \int_0^1 |\exp \zeta^t(\phi^N) \phi^N(t) - \exp \zeta^t(\eta_n) \eta_n(t)|^2 dt \end{aligned}$$

converges to zero by (22). Taking subsequences if necessary we see that

$$(24) \quad \exp \zeta^t(\phi^N)(\omega) \phi^N(t, x(\cdot, \omega)) = \lim_{n \rightarrow \infty} \exp \zeta^t(\eta_n)(\omega) \eta_n(t, x(\cdot, \omega)) \quad \text{a.s. } l \otimes P_0,$$

where  $l$  denotes Lebesgue measure on  $[0, 1]$ .

Also a.s.  $P_0$ ,

$$\exp \zeta^t(\phi^N) = E_0(\zeta(\phi^N) | \mathcal{A}_t), \quad \exp \zeta^t(\eta_n) = E_0(\zeta(\eta_n) | \mathcal{A}_t),$$

so that from (22),

$$\lim_{n \rightarrow \infty} \int_0^1 E_0 |\exp \zeta^t(\phi^N) - \exp \zeta^t(\eta_n)|^2 dt = 0,$$

and hence taking subsequences if necessary, we have

$$\exp \zeta^t(\phi^N)(\omega) = \lim_{n \rightarrow \infty} \exp \zeta^t(\eta_n)(\omega) \quad \text{a.s. } l \otimes P_0.$$

Since  $\exp \zeta^t(\phi^N) > 0$  a.s.  $P_0$ , we conclude from (24) that

$$\phi^N(t, x(\cdot, \omega)) = \lim_{n \rightarrow \infty} \eta_n(t, x(\cdot, \omega)) \quad \text{a.s. } l \otimes P_0,$$

and hence from (23), and the fact that  $f(t, z, U)$  is closed, we see that

$$\phi^N(t, x(\cdot, \omega)) \in f(t, x(\cdot, \omega), U) \quad \text{for } \|x(\cdot, \omega)\| \leq N, \quad \text{a.s. } l \otimes P_0.$$

From Corollary 1 of Lemma 2 and from Theorem 2 we obtain Theorem 3.

**THEOREM 3.** *Suppose  $f$  satisfies C1–C5 of § 2.*

(i) *For an admissible control  $u \in \mathcal{U}$ , there exists a solution to (1) with continuous sample paths (without explosions) if and only if*

$$E_0 \exp \zeta(g_u) = 1.$$

(ii) *The set of densities  $\{\exp \zeta(g_u) | E_0 \exp \zeta(g_u) = 1\}$  is a convex set, which is closed in the norm topology of  $L^1(\Omega, \mathcal{A}, P_0)$ .*

**4. Sufficient conditions for  $E_0 \exp \zeta(\phi) = 1$ .**

**LEMMA 9.** *Let  $\phi : [0, 1] \times C \rightarrow R^n$  be a causal map such that  $\int_0^1 |\phi(t, z)|^2 dt < \infty$  for all  $z$  in  $C$ . Define  $T_\phi : C \rightarrow C$  by*

$$(25) \quad T_\phi(z)(t) = z(t) - \int_0^t \phi(\tau, z) d\tau.$$

*Suppose that for each  $N > 0$  there is  $M > 0$  such that  $\|T_\phi(z)\| \leq N$  implies  $\|z\| \leq M$ . Then,*

$$E_0 \exp \zeta(\phi) = 1.$$

*Proof.* The proof is immediate from Lemma 7 of [5].

As a consequence of Lemma 9, we can obtain the following sufficient conditions. The first result is due to Beneš [3].

**COROLLARY 3.** *Let  $\phi : [0, 1] \times C \rightarrow R^n$  be a causal map and suppose there is a constant  $K$  such that*

$$|\phi(t, z)| \leq K(1 + \max_{1 \leq \tau \leq t} |z(\tau)|).$$

*Then,*

$$E_0 \exp \zeta(\phi) = 1.$$

*Proof.* Let  $T_\phi(z)(t) = y(t)$ , and let  $\gamma(t) = \max_{0 \leq \tau \leq t} |z(\tau)|$ . Then, from (25),

$$\begin{aligned} \gamma(t) &\leq |y(t)| + \int_0^t K(1 + \gamma(\tau)) d\tau \\ &\leq (\|y\| + K) + \int_0^t K\gamma(\tau) d\tau. \end{aligned}$$

By the Bellman–Gronwall inequality,

$$\begin{aligned} \|z\| = \gamma(1) &\leq (\exp K)\gamma(0) + (\exp K)(\|y\| + K) \\ &\leq (\exp K)(2\|y\| + K), \end{aligned}$$

and the result follows from Lemma 9.

The next result is useful if we have a control system with delay.

**COROLLARY 4.** *Let  $\phi : [0, 1] \times C \rightarrow R^n$  be a causal map such that for some  $\delta > 0$ ,*

$$|\phi(t, z)| \leq f_0 \left( \max_{0 \leq \tau \leq t-\delta} |z(\tau)| \right),$$



where  $f_0 : R \rightarrow R$  is increasing. Then,

$$E_0 \exp \zeta(\phi) = 1.$$

*Proof.* Let  $y, \gamma$  be defined as in the previous proof. Then,

$$\begin{aligned} \gamma(\delta) &\leq \|y\| + f_0(\gamma(0)), \\ \gamma(2\delta) &\leq \|y\| + f_0(\gamma(\delta)) \leq \|y\| + f_0(\|y\| + f_0(\gamma(0))) \\ &= f_1(\|y\|, \gamma(0)) \quad \text{say.} \end{aligned}$$

By induction,

$$\gamma(i\delta) \leq f_i(\|y\|, \gamma(0)),$$

where  $f_i$  is increasing in each argument. Evidently if  $(m - 1)\delta < 1 \leq m\delta$ , we see that

$$\gamma(1) = \|z\| \leq f_m(\|y\|, |z(0)|),$$

and the result follows from Lemma 9.

*Remark.* McKean [8, p. 66] has shown that if  $\delta > 0$ , then all solutions of the one-dimensional diffusion equation

$$dx(t) = |x|^{1+\delta} dt + dB(t), \quad 0 \leq t < \infty,$$

explode with probability 1. It follows that condition (6) is a nontrivial restriction.

**5. Applications.** Consider a control system

$$dx(t, \omega) = f(t, x(\cdot, \omega), u(t, x(\cdot, \omega))) dt + dB(t, \omega),$$

where the control  $u$  takes values in a set  $U$  and  $f$  obeys the conditions C1–C5 of § 2. Let us impose an additional restriction.

C6. For every admissible  $u \in \mathcal{U}$ ,

$$E_0 \exp \zeta(g_u) = 1,$$

or equivalently (and directly in terms of  $u$ ) for

$$(26) \quad \rho_u(z) = \exp \left[ \int_0^1 \langle f(t, z, u(t, z)), dz(t) \rangle - \frac{1}{2} \int_0^1 |f(t, z, u(t, z))|^2 dt \right]:$$

C6'.

$$\int_C \rho_u(z) \mu(dz) = 1.$$

Instead we can limit ourself to the subset  $\mathcal{U}^0$  consisting of those  $u$  in  $\mathcal{U}$  which satisfy C6'.

Next let  $L : C \rightarrow R$  be a bounded function, measurable with respect to  $\mathcal{S}$ .  $L$  is the *cost function* and assigns to every  $u \in \mathcal{U}^0$  the cost

$$(27) \quad J(u) = \int_C L(z) \rho_u(z) \mu(dz).$$

**THEOREM 4.** *Suppose  $\mathcal{U}^0$  is nonempty. Then, there exists  $u^* \in \mathcal{U}^0$  such that*

$$J(u^*) \leq J(u) \text{ for all } u \in \mathcal{U}^0.$$

*Proof.* By Theorem 3, the set  $\{\rho_u | u \in \mathcal{U}^0\}$  is a strongly closed, convex subset of  $L^1(C, \mathcal{S}, \mu)$ . Hence it is weakly compact. Since  $\int_C L(z)\rho_u(z)\mu(dz)$  is linear and continuous in  $\rho_u$ , the result follows.

Let us note that a cost functional of the type (27) allows for a variable endpoint problem as follows. Let  $\mathcal{T}$  be a closed subset of  $[0, 1] \times R^n$  which includes the set  $\{1\} \times R^n$ . Let  $\lambda: [0, 1] \times C \rightarrow R^n$  be a bounded, causal function, and to each  $u \in \mathcal{U}^0$  assign the cost

$$J(u) = \int_C \left[ \int_0^{t(z)} \lambda(t, z) dt \right] \rho_u(z)\mu(dz),$$

where  $t(z) = \min \{ \tau | z(\tau) \in \mathcal{T} \}$ . The term in brackets is clearly of the form  $L(z)$  in (27).

If the cost also depends on the control  $u$ , then sometimes we can add an extra coordinate to the state vector and get an equivalent cost depending only on the state. See [3] for details.

As a second application consider a zero-sum stochastic differential game, with two players I and II, with controls  $u_1(t) \in U_1$  and  $u_2(t) \in U_2$  respectively, and dynamics given by

$$dx(t) = f(t, x, u_1(t), u_2(t)) dt + dB(t).$$

Suppose that  $f$  splits as

$$f(t, x, u_1, u_2) = \begin{pmatrix} f_1(t, x, u_1) \\ f_2(t, x, u_2) \end{pmatrix}.$$

Assume that  $f$  satisfies C1–C5 with C5 now restated as:  $f_1(t, z, U_1)$  and  $f_2(t, z, U_2)$  are closed and convex for each  $(t, z)$ . As before, we define the admissible controls for player  $i$ , as all causal maps  $u_i: [0, 1] \times C \rightarrow U_i, i = 1, 2$ . Let  $\mathcal{U}_i^0$  consist of those admissible controls  $u_i$  which satisfy

$$\int_C \rho_{u_i}^i(z)\mu(dz) = 1,$$

where

$$(28) \quad \rho_{u_i}^i = \exp \left[ \int_0^1 \langle f_i(t, z, u_i(t, z)), dz_i(t) \rangle - \frac{1}{2} \int_0^1 |f_i(t, z, u_i(t, z))|^2 dt \right].$$

Here we have split  $z = \begin{pmatrix} z_1 \\ z_2 \end{pmatrix}$  to be compatible with  $f = \begin{pmatrix} f_1 \\ f_2 \end{pmatrix}$ . Let  $L: C \rightarrow R$  be a bounded function, measurable with respect to  $\mathcal{S}$ , and to each pair  $(u_1, u_2) \in \mathcal{U}_1^0 \times \mathcal{U}_2^0$  assign to player I the payoff

$$(29) \quad J(u_1, u_2) = \int_C L(z)\rho_{(u_1, u_2)}(z)\mu(dz).$$

**THEOREM 5.** *Suppose  $\mathcal{U}_i^0$  is nonempty for  $i = 1, 2$ . Then, there exist  $u_i^* \in \mathcal{U}_i^0$ ,  $i = 1, 2$ , such that*

$$J(u_1, u_2^*) \leq J(u_1^*, u_2^*) \leq J(u_1^*, u_2) \quad \text{for all } u_i \in \mathcal{U}_i^0, \quad i = 1, 2.$$

*Proof.* From the definition (26) of  $\rho_{(u_1, u_2)}$  and the definition (28) we see that

$$(30) \quad J(u_1, u_2) = \int_C L(z) \rho_{u_1}^1(z) \rho_{u_2}^2(z) \mu(dz).$$

Next, from Theorem 3, the sets  $\{\rho_{u_i}^i | u_i \in \mathcal{U}_i^0\}$  are convex, closed subsets of  $L^1(C, \mathcal{S}, \mu)$ , hence weakly compact. Finally the integral in (30) is concave (in fact linear) and continuous in  $\rho_{u_1}^1$  for fixed  $\rho_{u_2}^2$ , and convex (in fact linear) and continuous in  $\rho_{u_2}^2$  for fixed  $\rho_{u_1}^1$ . Hence from the well-known results on two-person zero-sum games the existence of a saddle point  $(u_1^*, u_2^*)$  follows.

**6. Partial observations: A negative example.** Again consider the stochastic differential equation

$$(1) \quad dx(t) = f(t, x, u) dt + dB(t).$$

The conditions on  $f$  are as before, but now suppose that we consider the important case where the control  $u$  can only depend upon the past history of the last  $m$  ( $m < n$ ) components of  $x$ . More precisely, let  $\mathcal{Q}_t$  be the sub- $\sigma$ -algebra of  $\mathcal{S}_t$  generated by all sets of the form

$$\{z | z \in C, z_i(\tau) \in A\},$$

where  $\tau \leq t$ ,  $A$  is a Borel subset of  $R$  and  $n - m + 1 \leq i \leq n$ . Let  $\mathcal{U}_m$  be the set of all causal maps  $u: [0, 1] \times C \rightarrow U$  such that  $u(t, \cdot)$  is measurable with respect to  $\mathcal{Q}_t$ . Evidently,

$$\mathcal{U}_m \subset \mathcal{U} = \mathcal{U}_n.$$

First of all it should be clear from the proof of Lemma 4 that the set  $\{\exp \zeta(g_u) | u \in \mathcal{U}_m\}$  may fail to be convex. For, consider the two-dimensional system

$$dx_1 = u dt + dB_1, \quad dx_2 = dB_2,$$

where the control  $u$  is allowed to depend only on  $x_2$  and must take values in the set  $U = [-1, 1]$ . Now let  $u_1$  and  $u_2$  be control laws defined as follows:

$$u_1(t, x_2) \equiv 0, \\ u_2(t, x_2) = \begin{cases} 0, & 0 \leq t < \frac{1}{2}, \\ \text{sgn}(x_2(\frac{1}{2})), & \frac{1}{2} \leq t \leq 1. \end{cases}$$

It is easy to calculate that

$$\zeta^t(g_{u_1}) \equiv 0$$

so that  $\exp \zeta^t(g_{u_1}) \equiv 1$  and

$$\zeta^t(g_{u_2}) = \begin{cases} 0, & 0 \leq t < \frac{1}{2}, \\ \text{sgn}(x_2(\frac{1}{2}))(x_1(t) - x_1(\frac{1}{2})) - \frac{1}{2}(t - \frac{1}{2}), & t \geq \frac{1}{2}. \end{cases}$$

From the proof of Lemma 4 we know that the function  $u_\lambda$ , given by  $u_\lambda(t) \equiv 0$ ,  $0 \leq t \leq \frac{1}{2}$ , and

$$u_\lambda(t) = \frac{\lambda \exp \zeta'(g_{u_2})}{\lambda \exp \zeta'(g_{u_2}) + (1 - \lambda)} u_2(t), \quad \frac{1}{2} \leq t \leq 1,$$

satisfies

$$\exp \zeta'(g_{u_\lambda}) \equiv (1 - \lambda) \exp \zeta'(g_{u_1}) + \lambda \exp \zeta'(g_{u_2}).$$

Since  $d \exp \zeta'(g_{u_\lambda}) = \exp \zeta'(g_{u_\lambda}) \langle g_{u_\lambda}(t), dx(t) \rangle$ ,  $g_{u_\lambda}$  (and hence  $u_\lambda$ ) is uniquely defined by the previous equation. Since for  $0 < \lambda < 1$ ,  $u_\lambda$  depends on  $x_1$ , the assertion is proved.

Next, from Theorem 3 we see that the set  $\{\exp \zeta(g_u) | u \in \mathcal{U}_m\}$  is weakly compact in  $L^1(\Omega, \mathcal{A}, P_0)$  if and only if it is weakly closed. We give a simple example to show that in general we do not have weak closure.

Consider the two-dimensional system  $x = (x_1, x_2)$ , with  $u \in R$  depending only on  $x_2$ ,

$$dx_1(t) = f(t, x_1)u + dB_1(t), \quad dx_2(t) = dB_2(t),$$

where

$$f(t, x_1) = \begin{cases} 0, & t \leq \frac{1}{2}, \\ 2, & t > \frac{1}{2}, \quad x_1(\frac{1}{2}) > 0, \\ 1, & t > \frac{1}{2}, \quad x_1(\frac{1}{2}) \leq 0. \end{cases}$$

The control set is  $U = [-1, 1]$ . We shall define a sequence of control laws  $u_n(t, x_2)$  such that

$$u_n(t, x_2) = \begin{cases} 0, & t \leq \frac{1}{2}, \\ \gamma_n(x_2(\frac{1}{2})), & t > \frac{1}{2}, \end{cases}$$

where the functions  $\gamma_n$  are defined later. It follows that

$$\begin{aligned} \zeta_n = \zeta(g_{u_n}) &= \int_0^1 f(t, x_1)u_n(t) dx_1(t) - \frac{1}{2} \int_0^1 f^2(t, x_1)u_n^2(t) dt \\ &= \alpha\beta\gamma_n - \frac{1}{4}\beta^2\gamma_n^2, \end{aligned}$$

where

$$\alpha = x_1(1) - x_1(\frac{1}{2}), \quad \beta = \begin{cases} 2 & \text{if } x_1(\frac{1}{2}) > 0, \\ 1 & \text{if } x_1(\frac{1}{2}) \leq 0. \end{cases}$$

Therefore,

$$(31) \quad \exp \zeta_n = \exp(\alpha\beta\gamma_n) \exp(-\frac{1}{4}\beta^2\gamma_n^2).$$

We shall select  $\gamma_n$  such that  $|\gamma_n| \equiv 1$ , so that (31) simplifies to

$$(32) \quad \exp \zeta_n = \exp(\alpha\beta\gamma_n) \exp(-\frac{1}{4}\beta^2).$$

We define  $\gamma_n$  as follows:

Let  $\xi : R \rightarrow R$  be a measurable function such that under  $P_0$ ,  $\xi(x_2(\frac{1}{2}))$  is uniformly distributed over  $[0, 1]$ . For each integer  $n \geq 0$ , define  $\eta_n : [0, 1] \rightarrow \{-1, 1\}$  by

$$\eta_n(\xi) = \begin{cases} 1 & \text{if } \frac{2m}{2n} \leq \xi < \frac{2m+1}{2n}, & m = 0, 1, \dots, n-1, \\ -1 & \text{if } \frac{2m+1}{2n} \leq \xi < \frac{2m+2}{2n}, & m = 0, 1, \dots, n-1. \end{cases}$$

Finally, let

$$\gamma_n(x_2(\frac{1}{2})) = \eta_n(\xi(x_2(\frac{1}{2}))).$$

LEMMA 10.  $\exp \zeta_n$  converges to  $\frac{1}{2} [\exp(\alpha\beta) + \exp(-\alpha\beta)] \exp(-\frac{1}{4}\beta^2)$  in the weak topology of  $L^1(\Omega, \mathcal{A}, P_0)$ .

*Proof.* Let  $A_\alpha, A_\beta, A_\xi$  be Borel subsets of  $R$  and let  $I_A$  denote the indicator function of a set  $A$ . Let

$$\Pi_n = \int_{\Omega} I_{A_\alpha}(\alpha(\omega)) I_{A_\beta}(\beta(\omega)) I_{A_\xi}(\xi(\omega)) (\exp \zeta_n)(\omega) P_0(d\omega).$$

Now under  $P_0$  the random variables  $\alpha, \beta, \xi$  are independent, so that

$$(33) \quad \Pi_n = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} I_{A_\alpha}(\alpha) I_{A_\beta}(\beta) I_{A_\xi}(\xi) \exp(\alpha\beta\eta_n(\xi)) \exp\left(-\frac{1}{4}\beta^2\right) P_\xi(d\xi) P_\beta(d\beta) P_\alpha(d\alpha),$$

where  $P_\alpha, P_\beta, P_\xi$  are the marginal distributions of  $\alpha, \beta, \xi$  respectively. From the way  $\eta_n$  is defined and the fact that  $\xi$  is uniformly distributed on  $[0, 1]$  it follows that for fixed  $\alpha, \beta$ ,

$$\begin{aligned} & \lim_{n \rightarrow \infty} \int_{-\infty}^t \exp(\alpha\beta\eta_n(\xi)) \exp\left(-\frac{1}{4}\beta^2\right) P_\xi(d\xi) \\ &= \int_{-\infty}^t \frac{1}{2} [\exp(\alpha\beta) + \exp(-\alpha\beta)] \exp\left(-\frac{1}{4}\beta^2\right) P_\xi(d\xi) \end{aligned}$$

uniformly for  $t \in (-\infty, \infty)$ . It follows that  $\exp(\alpha\beta\eta_n(\xi)) \exp(-\frac{1}{4}\beta^2)$  converges to  $\frac{1}{2} [\exp(\alpha\beta) + \exp(-\alpha\beta)] \exp(-\frac{1}{4}\beta^2)$  weakly in  $L^1(R, P_\xi)$ . Since the integrands in (33) are uniformly integrable, it follows that

$$\begin{aligned} \lim_{n \rightarrow \infty} \Pi_n &= \int_{R^3} I_{A_\alpha}(\alpha) I_{A_\beta}(\beta) I_{A_\xi}(\xi) \frac{1}{2} [\exp(\alpha\beta) + \exp(-\alpha\beta)] \\ &\quad \cdot \exp\left(-\frac{1}{4}\beta^2\right) P_\xi(d\xi) P_\beta(d\beta) P_\alpha(d\alpha). \end{aligned}$$

From this it follows easily that  $\exp(\alpha\beta\eta_n(\xi)) \exp(-\frac{1}{4}\beta^2)$  converges to  $\frac{1}{2} [\exp(\alpha\beta) + \exp(-\alpha\beta)] \exp(-\frac{1}{4}\beta^2)$  weakly in  $L^1(R^3, P_\alpha \otimes P_\beta \otimes P_\xi)$  and the lemma is proved.

Next, by direct calculation we can show that the two-dimensional drift  $\hat{g}$ , where  $\hat{g}(t) = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$ ,  $0 \leq t \leq \frac{1}{2}$ , and

$$\hat{g}(t) = \begin{pmatrix} \beta \frac{\exp 2(x_1(t) - x_1(\frac{1}{2})) - 1}{\exp 2(x_1(t) - x_1(\frac{1}{2})) + 1} \\ 0 \end{pmatrix}, \quad \frac{1}{2} \leq t \leq 1,$$

satisfies

$$(34) \quad \exp \zeta(\hat{g}) = \frac{1}{2}[\exp(\alpha\beta) + \exp(-\alpha\beta)] \exp(-\frac{1}{4}\beta^2).$$

If we define  $\rho(t) = E_0(\exp \zeta(\hat{g}) | \mathcal{A}_t)$ , then

$$d\rho(t) = \rho(t)\langle \hat{g}(t), dx(t) \rangle, \quad 0 \leq t \leq 1,$$

so that (34) characterizes  $\hat{g}$  uniquely. Hence any control law  $\hat{u}$  such that  $\exp \zeta(g_u) = \exp \zeta(\hat{g})$  must satisfy  $g_u = \hat{g}$  so that  $\hat{u}$  must depend on  $x_1$ . Therefore the set of densities  $\exp \zeta(g_u)$  with  $u$  depending only on  $x_2$  is not weakly closed in  $L^1(\Omega, \mathcal{A}, P_0)$ .

Incidentally this example also shows that to guarantee weak closure, the convexity condition C5 is necessary, for even though  $u_n(t) \in \{-1, 1, 0\}$  for all  $t$ , it is not the case for  $\hat{u}(t)$ .

**Acknowledgments.** It should be obvious to anyone who has read [3] our great debt to V. E. Beneš. E. Wong clarified many subtle points about Ito's calculus and suggested the proofs in § 2. M. Davis caught many errors in earlier proofs. It is a pleasure to acknowledge our gratitude to them.

#### REFERENCES

- [1] W. H. FLEMING AND M. NISIO, *On the existence of optimal stochastic controls*, J. Math. Mech., 15 (1966), pp. 777-794.
- [2] H. J. KUSHNER, *On the existence of optimal stochastic controls*, this Journal, 3 (1965), pp. 463-474.
- [3] V. E. BENEŠ, *Existence of optimal stochastic control laws*, to appear.
- [4] W. H. FLEMING, *Optimal continuous-parameter stochastic control*, SIAM Rev., 11 (1969), pp. 470-509.
- [5] I. V. GIRSANOV, *On transforming a certain class of stochastic processes by absolutely continuous substitution of measures*, Theor. Probability Appl., 5 (1960), pp. 285-301.
- [6] V. E. BENEŠ, *Existence of optimal strategies based on specified information, for a class of stochastic decision problems*, this Journal, 8 (1970), pp. 179-188.
- [7] K. ITO, *Multiple Wiener integral*, J. Math. Soc. Japan, 3 (1951), pp. 158-161.
- [8] H. P. MCKEAN, JR., *Stochastic Integrals*, Academic Press, New York and London, 1969.

## ALMOST SURE BOUNDEDNESS OF RANDOMLY SAMPLED SYSTEMS\*

R. G. AGNIEL† AND E. I. JURY‡

**Abstract.** This paper discusses the almost sure boundedness of linear and nonlinear randomly sampled systems. It is shown that if an autonomous linear randomly sampled system exhibits almost sure asymptotic stability, then the system is almost surely bounded input–bounded output. Moreover, for a bounded input, the second moment of the output remains bounded and this bound is easily computable.

It is also found that linear or nonlinear systems which are almost surely asymptotically stable for a null input remain almost surely bounded when the input consists of an uncorrelated noise with finite variance.

**Introduction.** Recently, more and more attention has been devoted to sampled data systems with a randomly variable sampling period. The sampling intervals may be random for physical reasons, when they are contaminated by noise because of cheap equipment (“jitter”) or when the sampling is governed by a random process such as a sonar or radar echo. In certain cases, a randomly variable period is advantageous for economical reasons, for instance, when a time-shared computer controls several plants.

Moreover, randomly sampled systems were considered by many authors such as Kalman [3], Bergen [4], Bharucha [5] and others. In [3], Kalman repeatedly applies the Schwarz and Minkowski inequalities to prove that for an autonomous system which shows stability of the second moment of the output, this second moment remains bounded when a bounded input is applied to the system.

Following Bucy’s idea [6], [7] of using the submartingale convergence theorem, in this paper we will prove that not only the second moment remains bounded, but the output is almost surely bounded as well.

Moreover, it will be proved that a large class of linear or nonlinear randomly sampled systems, shown to be almost surely asymptotically stable with a null input, remain almost surely bounded when the input consists of an uncorrelated noise with finite variance.

**1. Stochastic stability concepts [8].** Let  $t_0 < t_1 < t_2 < t_3 < \dots < t_n < t_{n+1} < \dots$  be a discrete sequence of strictly ordered random times, and let  $\mathbf{X}(n, \mathbf{X}_0, n_0)$  be a discrete sequence of random vectors at time  $t_n$  such that  $\mathbf{X}(n_0, \mathbf{X}_0, n_0) = \mathbf{X}_0$  and satisfying the vector difference equation

$$\mathbf{X}(n + 1, \mathbf{X}_0, n_0) = f[\mathbf{X}(n, \mathbf{X}_0, n_0), \mathbf{r}_n, \mathbf{u}_n],$$

where the  $\mathbf{r}_n$ ’s are a sequence of random vectors and the  $\mathbf{u}_n$ ’s a sequence of input vectors.

---

\* Received by the editors December 12, 1969, and in final revised form December 3, 1970. This research was sponsored by the Air Force Office of Scientific Research, Office of Aerospace Research, under AFOSR Grant AF-AFOSR-68-1463.

† Direction Commerciale Groupe Tubes Electroniques, Thomson-CSF, Paris, France. Formerly with Electronics Research Laboratory, University of California, Berkeley, California.

‡ Department of Electrical Engineering and Computer Science, and Electronics Research Laboratory, University of California, Berkeley, California 94720.

Let  $\mathbf{0}$  be an equilibrium solution, i.e.,

$$f[\mathbf{0}, \mathbf{r}_\mu, \mathbf{0}] = \mathbf{0} \quad \text{for all } n.$$

The equilibrium solution  $\mathbf{X}_e \equiv \mathbf{0}$  is said :

(a) to be *almost surely stable* if given  $\epsilon, \epsilon' > 0$ , there exists  $\delta(\epsilon, \epsilon', n_0)$  such that  $\|\mathbf{X}_0\| < \delta$  implies

$$P[\sup_{n \geq n_0} \|\mathbf{X}(n, \mathbf{X}_0, n_0)\| > \epsilon'] < \epsilon,$$

where  $P[\cdot]$  denotes the probability and  $\|\cdot\|$  is a vector norm ;

(b) to be *almost surely asymptotically stable* if

- (i) it is almost surely stable,
- (ii) there exists  $\delta' > 0$  such that  $\|\mathbf{X}_0\| < \delta'$  implies

$$\|\mathbf{X}(n, \mathbf{X}_0, n_0)\| \xrightarrow{\text{a.s.}} 0 ;$$

(c) to show *stability of the p-th moment* if given  $\epsilon > 0$ , there exists  $\delta(\epsilon, n_0)$  such that  $\|\mathbf{X}_0\| < \delta$  implies

$$\sup_{n \geq n_0} E[\|\mathbf{X}(n, \mathbf{X}_0, n_0)\|^p] < \epsilon ;$$

(d) to show *asymptotic stability of the p-th moment* if

- (i) it shows stability of the pth moment,
- (ii) there exists  $\epsilon' > 0$  such that  $\|\mathbf{X}_0\| < \delta'$  implies

$$\lim_{n \rightarrow \infty} E[\|\mathbf{X}(n, \mathbf{X}_0, n_0)\|^p] = 0 ;$$

(e) to be *almost surely bounded input-bounded output* if for any  $\|\mathbf{X}_0\| < \infty$  and any almost surely bounded input,

$$\sup_n \|\mathbf{X}(n, \mathbf{X}_0, n_0)\| < \infty \quad \text{a.s.}$$

In what follows, any of these modes of stability will be said to hold absolutely if true for a null input.

**2. Stability of linear randomly sampled systems.**

**2.1. Description of the system.** Consider the closed loop sampled data system comprised of a random sampling and hold device followed by a time invariant plant and unity feedback (Fig. 1). This system has the following properties :

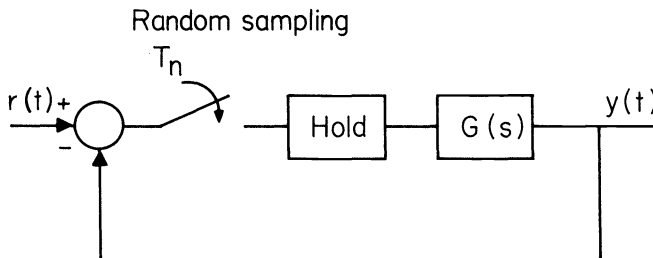


FIG. 1. Linear randomly sampled system



(a) *Sampling process*: The successive sampling intervals are stochastically independent, identically distributed, and are denoted by  $T_0, T_1, T_2, \dots, T_k, \dots$ . The sampling times are denoted by  $t_0, t_1, t_2, \dots, t_k, \dots$  (Fig. 2). The different distributions and their practical motivation are described by Agniel [8]. Since the successive periods are identically distributed, the expectation of any function of  $T_k$  is independent of  $k$ .

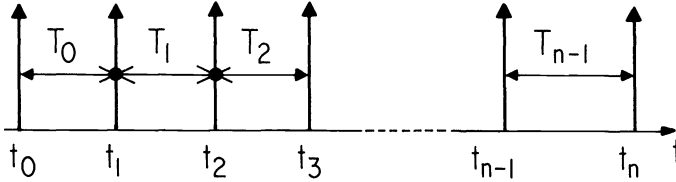


FIG. 2. *Sampling process*

(b) *Equation of the system*: Following Kalman and Bertram’s approach [9], we have the system equations

$$\begin{aligned}
 \mathbf{X}(k + 1) &= [G(k) + \mathbf{h}(k)\mathbf{a}']\mathbf{X}(k) + \mathbf{h}(k)r(t_k), \\
 (1) \quad y(k) &= \mathbf{a}'\mathbf{X}(k), \quad k = 0, 1, 2, 3, \dots, \\
 \mathbf{X}(0) &= \mathbf{X}_0,
 \end{aligned}$$

where  $\mathbf{X}(k)$  is the state vector at time  $t_k$ ,

$G(k) \triangleq G(T_k)$  is an  $n \times n$  matrix continuous in  $T_k$ ,

$\mathbf{h}(k) \triangleq \mathbf{h}(T_k)$  is an  $n$ -dimensional vector continuous in  $T_k$ ,

$\mathbf{a}$  is an  $n$ -dimensional constant vector,

$r(t_k)$  is the input at time  $t_k$ .

**2.2. Absolute stability of linear randomly sampled systems.** Suppose the input  $r(t) \equiv 0$  for all  $t < 0$ . Let

$$V_k = \mathbf{X}'(k)B\mathbf{X}(k),$$

where  $B$  is a symmetric positive definite matrix.

Taking  $V_k$  as a stochastic Lyapunov function and by a trivial application of Bucy’s results [6], [7], the linear system described in § 2.1(b) is almost surely asymptotically stable and shows asymptotic stability of the second moment if there exist positive definite matrices  $B$  and  $Q$  such that

$$(2) \quad E[\{G(k) + \mathbf{h}(k)\mathbf{a}'\}'B\{G(k) + \mathbf{h}(k)\mathbf{a}'\}] - B = -Q.$$

It was shown by Agniel and Jury [10] (or our Appendix B) that this can be done provided the eigenvalues of

$$E[\{G(k) + \mathbf{h}(k)\mathbf{a}'\} \otimes \{G(k) + \mathbf{h}(k)\mathbf{a}'\}]$$

lie inside the unit circle,  $\otimes$  denoting the first Krönercker product of matrices.

By using the same reasoning as in Kalman and Bertram [11], we see that condition (2) is equivalent to

$$(3) \quad E[\{G(k) + \mathbf{h}(k)\mathbf{a}'\}'B\{G(k) + \mathbf{h}(k)\mathbf{a}'\}] - \frac{1}{\eta}B \leq 0,$$

where  $\eta$  is a scalar defined by

$$(4) \quad \frac{1}{\eta} = \lambda_{\min}[E\{[(G(k) + \mathbf{h}(k)\mathbf{a}')B(G(k) + \mathbf{h}(k)\mathbf{a}')]B^{-1}\}],$$

where  $\eta > 1$ .

In conclusion, (2) implies for a null input that

$$(5) \quad E^{\mathcal{B}_k}V_{k+1} - V_k/\eta \leq 0, \quad \eta > 1, \quad k = 0, 1, 2, 3, \dots,$$

where  $E^{\mathcal{B}_k}(\cdot)$  denotes the conditional expectation given the  $\sigma$ -field  $\mathcal{B}_k$ ,<sup>1</sup> which is the  $\sigma$ -field generated by  $\mathbf{X}(1), \mathbf{X}(2), \dots, \mathbf{X}(k)$ . (See Bucy [6], [7] or Loève [12].)

**2.3. Stability for a bounded input.** Suppose the system is almost surely asymptotically stable and shows asymptotic stability of the second moment for a null input, i.e.,<sup>2</sup>

$$(6) \quad \sup_i |\lambda_i\{E[(G(k) + \mathbf{h}(k)\mathbf{a}') \otimes (G(k) + \mathbf{h}(k)\mathbf{a}')]]\}| < 1.$$

Suppose the system is now driven by a bounded deterministic input.

Consider the same Lyapunov function  $V_k = \mathbf{X}'(k)B\mathbf{X}(k)$ . We have

$$(7) \quad E^{\mathcal{B}_k}V_{k+1} - \frac{1}{\eta}V_k = E^{\mathcal{B}_k}\mathbf{X}'(k+1)B\mathbf{X}(k+1) - \frac{1}{\eta}\mathbf{X}'(k)B\mathbf{X}(k).$$

Since the sampling intervals are independent, and by properties of conditional expectations (see Loève [12, p. 350]), we get (using (1))

$$(8) \quad \begin{aligned} E^{\mathcal{B}_k}V_{k+1} - \frac{1}{\eta}V_k &= \mathbf{X}'(k)E[(G(k) + \mathbf{h}(k)\mathbf{a}')B(G(k) + \mathbf{h}(k)\mathbf{a}')] \mathbf{X}(k) \\ &\quad + 2E[\mathbf{h}'(k)B(G(k) + \mathbf{h}(k)\mathbf{a}')] \mathbf{X}(k)r(t_k) \\ &\quad + E[\mathbf{h}'(k)B\mathbf{h}(k)]r^2(t_k) - \frac{1}{\eta}\mathbf{X}'(k)B\mathbf{X}(k), \end{aligned}$$

which implies

$$(9) \quad E^{\mathcal{B}_k}V_{k+1} - \frac{1}{\eta}V_k - E[\mathbf{h}'(k)B\mathbf{h}(k)]r^2(t_k) = \mathbf{X}'(k)\theta\mathbf{X}(k) + 2\gamma'\mathbf{X}(k)r(t_k)$$

with

$$(10) \quad \theta = E[(G(k) + \mathbf{h}(k)\mathbf{a}')B(G(k) + \mathbf{h}(k)\mathbf{a}')] - \frac{1}{\eta}B,$$

$$(11) \quad \gamma' = E[\mathbf{h}'(k)B(G(k) + \mathbf{h}(k)\mathbf{a}')].$$

Equation (9) can be rewritten as

$$(12) \quad E^{\mathcal{B}_k}V_{k+1} - \frac{1}{\eta}V_k - [E\mathbf{h}'(l)B\mathbf{h}(k) + \rho]r^2(t_k) = [\mathbf{X}'(k), r(t_k)] \begin{bmatrix} \theta\gamma \\ \gamma' - \rho \end{bmatrix} \begin{bmatrix} \mathbf{X}(k) \\ r(t_k) \end{bmatrix}.$$

<sup>1</sup> In this case,  $\mathcal{B}_k$  is the  $\sigma$ -field generated by  $\mathbf{X}(k)$  only. The derivation remains valid for sampling periods which are not independent.

<sup>2</sup> Condition (6) is necessary and sufficient for absolute stability of the second moment [5]. However, it is only sufficient for a.s. absolute stability [5], [10].

Using the results of Gantmacher [13, pp. 46, 306], we see that the right-hand side of (12) is negative definite if

(i)  $\theta$  is negative definite, which is always verified by our assumption (the system is absolutely asymptotically stable, i.e., inequality (3) is satisfied);

(ii)  $\rho > -\gamma'\theta^{-1}\gamma > 0$ .

Therefore,

$$(13) \quad E^{\theta_k} V_{k+1} - \frac{1}{\eta} V_k - C_k < 0 \quad \text{a.s. for all } k$$

with  $C_k = [E\{h'(k)Bh(k)\} + \rho]r^2(t_k) \geq 0$ , and since  $|r(t_k)| \leq M < \infty$ , we have

$$(14) \quad E^{\theta_k} V_{k+1} - \frac{1}{\eta} V_k - C < 0 \quad \text{a.s. for all } k$$

with  $C = \sup_k C_k \leq [E\{h'(k)Bh(k)\} + \rho]M^2 < \infty$ .

It is not necessary to compute  $\theta^{-1}$  since a close bound for  $\rho$  can be found by writing

$$0 < \rho \leq -[\lambda_{\min}(\theta)]^{-1}\gamma'\gamma.$$

The properties of systems satisfying inequality (14) will be studied in § 4 and the following theorem will be proved.

**THEOREM.** *For the system described in § 2.1 :*

1.  $V_k$  is almost surely bounded for all  $k$ .
2.  $E(V_k)$  is bounded for all  $k$  and

$$EV_k \leq \sup \left[ C \frac{\eta}{\eta - 1}, V_0 \right].$$

3.  $\|\mathbf{X}(k)\|$  is almost surely bounded for all  $k$ .
4.  $E[\|X(k)\|^2]$  is bounded for all  $k$  and

$$E[\|\mathbf{X}(k)\|^2] \leq [\min_i \lambda_i(B)]^{-1} \sup \left[ C \frac{\eta}{\eta - 1}, \mathbf{X}'_0 B \mathbf{X}_0 \right].$$

It will now be shown that linear or nonlinear absolutely stable systems satisfy inequality (14) when the input is an uncorrelated noise with finite variance.

**3. Randomly sampled systems with a noise as an input.**

**3.1. Description of the input.** Randomly sampled systems where the input is uncorrelated noise with zero mean and finite variance will now be investigated. Such noise is different from white noise since the variance is finite, and then does not require infinite energy.

This input noise will be denoted by  $r(t)$  and it satisfies the following conditions :

- (i)  $r(t)$  has zero mean, i.e.,  $E[r(t)] = 0$  for all  $t$ ;
- (ii)  $r(t)$  is uncorrelated, i.e.,  $E[r(t)r(t + \tau)] = 0$  for all  $t$  and all  $\tau > 0$ ;
- (iii)  $r(t)$  has finite variance, i.e.,

$$0 \leq E[r^2(t)] \leq C_1 \leq \infty \quad \text{for all } t,$$

$C_1$  being a positive constant;

- (iv) the input and the sampling process are independent.

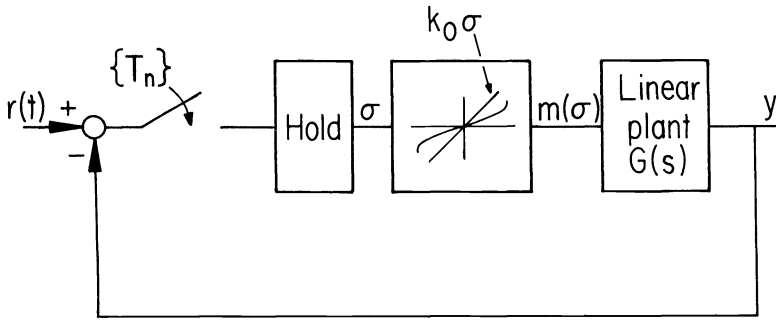


FIG. 3. Nonlinear randomly sampled system

CLAIM. For linear or nonlinear randomly sampled systems shown to be almost surely asymptotically stable for a null input, (14) is verified when the input consists of the noise described in § 3.1.

This claim will be proved for a particular configuration. Consider a closed loop sampled data system with a random sample and hold element followed by a time invariant nonlinearity and a plant with unity feedback (Fig. 3). The system satisfies the following assumptions:

(i) the sampling intervals are independent and identically distributed;

(ii) the nonlinearity is differentiable, of sector type and has bounded slope, i.e., if  $m(\sigma)$  is the output of the nonlinearity for an input  $\sigma$ ,

$$(15) \quad |dm/d\sigma| < K,$$

$$(16) \quad 0 \leq \sigma m(\sigma) \leq k_0 \sigma^2;$$

(iii) the input noise is described in § 3.1;

(iv) the linear plant is time invariant, i.e. (see for instance Kalman and Bertram [9]),

$$(17) \quad \mathbf{X}(k+1) = G(k)\mathbf{X}(k) + \mathbf{h}(k)m(\sigma_k)$$

with the same notation as in § 2.1, and

$$(18) \quad m(\sigma_k) = m[\mathbf{a}'\mathbf{X}(k) + r(t_k)].$$

Consider the stochastic Lyapunov function  $V_k = \mathbf{X}'(k)B\mathbf{X}(k)$ , where  $B$  is a positive definite matrix.

If the system has been proven to be almost surely asymptotically stable for a null input, then there exists<sup>3</sup> a positive definite  $B$  such that (see Agniel and Jury [10])

$$(19) \quad E^{\mathcal{B}_k} V_{k+1} - \frac{1}{\eta} V_k \leq 0 \quad \text{a.s. for all } k, \eta > 1,$$

<sup>3</sup>The necessity of the existence of a Lyapunov function can be proved by an analogy with the Massera theorem (see [11, p. 379]).

i.e.,

$$(20) \quad \mathbf{X}'(k) \left\{ E[G'(k)BG(k)] - \frac{1}{\eta} \mathbf{B} \right\} \mathbf{X}(k) + 2E[h'(k)BG(k)]\mathbf{X}(k)m(\mathbf{a}'\mathbf{X}(k)) + E[\mathbf{h}'(k)\mathbf{Bh}(k)]m^2(\mathbf{a}'\mathbf{X}(k)) \leq 0 \quad \text{for all } k \quad \text{a.s.}$$

If the input  $r(t)$  is as described in § 3.1, then

$$(21) \quad E^{\mathcal{B}_k} V_{k+1} = E^{\mathcal{B}_k} \{ [\mathbf{X}'(k)G'(k)BG(k)\mathbf{X}(k)] + 2\mathbf{h}'(k)BG(k)\mathbf{X}(k)m[\mathbf{a}'\mathbf{X}(k) + r(t_k)] + \mathbf{h}'(k)\mathbf{Bh}(k)m^2[\mathbf{a}'\mathbf{X}(k) + r(t_k)] \}.$$

The matrix  $G'(k)BG(k)$ , the vector  $\mathbf{h}'(k)BG(k)$  and  $\mathbf{h}'(k)\mathbf{Bh}(k)$  are Borel functions of  $T_k$ ; moreover, the  $\sigma$ -field induced by  $T_k$  and  $\mathcal{B}_k$  are independent by assumption. By properties of conditional expectations (see Loève [12, pp. 349, 350]),

$$(22) \quad E^{\mathcal{B}_k} V_{k+1} = \mathbf{X}'(k)E[G'(k)BG(k)]\mathbf{X}(k) + 2E[\mathbf{h}'(k)BG(k)]\mathbf{X}(k)E^{\mathcal{B}_k}[m(\mathbf{a}'\mathbf{X}(k) + r(t_k))] + E[\mathbf{h}'(k)\mathbf{Bh}(k)]E^{\mathcal{B}_k}[m^2(\mathbf{a}'\mathbf{X}(k) + r(t_k))] \quad \text{a.s.}$$

Since the nonlinearity is differentiable, we can apply MacLaurin's formula to get

$$(23) \quad m[\mathbf{a}'\mathbf{X}(k) + r(t_k)] = m[\mathbf{a}'\mathbf{X}(k)] + \dot{m}(\theta)r(t_k)$$

with  $|\dot{m}| = |dm/d\sigma| \leq K < \infty$  by assumption (from (15)) and  $\theta \in [\mathbf{a}'\mathbf{X}(k), \mathbf{a}'\mathbf{X}(k) + r(t_k)]$ .

Since the sequences  $\{r(t_k)\}$  and  $\{T_k\}$  are independent by assumption, and by using the properties of  $r(t_k)$  and the fact that  $m(\mathbf{a}'\mathbf{X}(k))$  is  $\mathcal{B}_k$ -measurable, we obtain (see Loève [12, pp. 348, 349, 350])

$$(24) \quad E^{\mathcal{B}_k} \{ m^2[\mathbf{a}'\mathbf{X}(k) + r(t_k)] \} \leq m^2[\mathbf{a}'\mathbf{X}(k)] + K^2 C_1 \quad \text{a.s.,}$$

$$(25) \quad E^{\mathcal{B}_k} [m(\mathbf{a}'\mathbf{X}(k) + r(t_k))] = m(\mathbf{a}'\mathbf{X}(k)) \quad \text{a.s.;}$$

then

$$(26) \quad E^{\mathcal{B}_k} V_{k+1} - \frac{1}{\eta} V_k \leq \mathbf{X}'(k) \left\{ E[G'(k)BG(k)] - \frac{1}{\eta} \mathbf{B} \right\} \mathbf{X}(k) + 2E[\mathbf{h}'(k)BG(k)]\mathbf{X}(k)m[\mathbf{a}'\mathbf{X}(k)] + E[\mathbf{h}'(k)\mathbf{Bh}(k)][m^2(\mathbf{a}'\mathbf{X}(k)) + K^2 C_1] \quad \text{for all } k, \quad \text{a.s.,}$$

and from (19),

$$(27) \quad E^{\mathcal{B}_k} V_{k+1} - \frac{1}{\eta} V_k - E[\mathbf{h}'(k)\mathbf{Bh}(k)]K^2 C_1 < 0 \quad \text{for all } k, \quad \text{a.s.}$$

We obtain an equation of type (14) with

$$C = E[\mathbf{h}'(k)\mathbf{Bh}(k)]K^2 C_1 > 0.$$

The same procedure can be applied successfully to a large class of system configurations and Lyapunov functions; in particular,

- linear randomly sampled systems,
- nonlinear randomly sampled systems with nonlinear element in the feedback loop,
- nonlinear randomly sampled systems with a stochastic Lyapunov function in the Luré form as considered by Agniel and Jury [10], Kushner and Tobias [15].

It can be shown that the constant  $C$  in inequality (14) is proportional to  $C_1$  in most cases, and the physical interpretation of this remark will be illustrated in § 5.

For linear systems, this property remains true if the system has an output consisting of a bounded signal with a noise as described in § 3.1. Of course, this property is not true for nonlinear systems, but for some simple inputs, one can sometimes circumvent this difficulty by considering the deterministic signal as an initial condition.

It will be shown in the next section that for the class of systems considered which have been shown to satisfy equation (14):

1. The sequence  $V_k$  ( $k = 0, 1, \dots$ ) is almost surely bounded;
2.  $EV_k$  is bounded for all  $k$  and  $EV_k \leq \sup [C(\eta/\eta - 1), V_0]$ ;
3.  $\|X(k)\|$  remains almost surely bounded.

**4. Properties of systems described by  $E^{\mathcal{B}_k}V_{k+1} - (1/\eta)V_k - C < 0$  a.s.  $k = 0, 1, 2, \dots, n$ .**

**4.1. Properties of  $V_k$ .**

**THEOREM 1.**  $V_k$  is almost surely bounded.

*Proof.* As it is,  $\{V_k\}$  is not supermartingale. Consider  $\{W_k\}$  defined by

$$(28) \quad W_k = V_k - C \frac{\eta}{\eta - 1}.$$

Obviously,  $W_k$  and  $V_k$  are measurable with respect to the same  $\sigma$ -field  $\mathcal{B}_k$ . Since conditional expectations are a.s. linear operators and since the conditional expectation of a constant is a.s. equal to this constant (see Loève [12, pp. 347, 348]), inequality (14) becomes

$$(29) \quad E^{\mathcal{B}_k}W_{k+1} - \frac{1}{\eta}W_k < 0 \quad \text{a.s.,}$$

i.e.,  $W_k$  is a supermartingale satisfying an equation similar to (5). It is deduced that

$$(30) \quad E(-W_0) \leq E(-W_1) < \dots \leq E(-W_k) \leq \dots,$$

and recalling that  $V_k$  is positive for all  $k$ , then

$$(31) \quad E(-W_k) = E(W_k^- - W_k^+) = C \frac{\eta}{\eta - 1} - EV_k,$$

$$(32) \quad E(|W_k|) = E(W_k^+ + W_k^-) \leq C \frac{\eta}{\eta - 1} + EV_k.$$

Summing up, one has

$$(33) \quad E(-W_k)^+ \leq C \frac{\eta}{\eta - 1} < \infty \quad \text{since } \eta > 1.$$

The submartingale convergence theorem (see Appendix A) can now be applied to the submartingale  $\{-W_k\}$  to obtain

$$(34) \quad -W_k = C \frac{\eta}{\eta - 1} - V_k \xrightarrow{\text{a.s.}} W < \infty,$$

i.e.,

$$(35) \quad V_k \xrightarrow{\text{a.s.}} V < \infty \quad \text{a.s.},$$

and from (28) and (30),  $V_k < \infty$  a.s. for all  $k$ .

**THEOREM 2.**  $EV_k$  remains bounded for all  $k$  and

$$EV_k \leq \sup \left[ C \frac{\eta}{\eta - 1}, V_0 \right].$$

*Proof.* Inequality (14) implies (see Loève [12, p. 341])

$$(36) \quad \eta EV_k < EV_{k-1} + \eta C,$$

and by induction,

$$(37) \quad EV_k < \frac{V_0}{\eta^k} + C \frac{\eta}{\eta - 1} \left( 1 - \frac{1}{\eta^{k+1}} \right).$$

By a straightforward optimization procedure, (37) implies

$$(38) \quad EV_k \leq \sup \left[ C \frac{\eta}{\eta - 1}, V_0 \right] \quad \text{for all } k.$$

#### 4.2. Properties of $\mathbf{X}(k)$ .

**THEOREM 3.**  $\|\mathbf{X}(k)\|$  is almost surely bounded for all  $k$ .

The proof is trivial since it is known that  $V_k$  is a stochastic Lyapunov function if there exist nondecreasing functions  $\alpha(\cdot)$  and  $\beta(\cdot)$  such that  $\alpha(0) = \beta(0) = 0$ ,

$$(39) \quad \alpha(\|\mathbf{X}(k)\|) \leq V_k \leq \beta(\|\mathbf{X}(k)\|).$$

If  $V_k = \mathbf{X}'(k)B\mathbf{X}(k)$ , then  $(\mathbf{X}'(k)B\mathbf{X}(k))^{1/2}$  is a norm for  $\mathbf{X}(k)$ , and then  $V_k = \|\mathbf{X}(k)\|^2 < \infty$  for all  $k$ , a.s.

**THEOREM 4.**  $E\|\mathbf{X}(k)\|$  is bounded for all  $k$ .

In the general case, since  $V_k$  is positive definite, an  $\alpha(\cdot)$  which is convex can be found such that

$$(40) \quad \alpha(\|\mathbf{X}(k)\|) \leq V_k \quad \text{for all } k.$$

By applying Jensen's inequality,

$$(41) \quad \alpha(E\|\mathbf{X}(k)\|) \leq EV_k \leq \sup \left[ V_0, C \frac{\eta}{\eta - 1} \right] < \infty,$$

and the procedure is as in Theorem 3.

However, in most cases of interest,  $V_k$  is a quadratic form or a Luré form (see Agniel and Jury [10]); for a quadratic form, if  $V = \mathbf{X}'\mathbf{B}\mathbf{X}$ ,

$$(42) \quad E[\|\mathbf{X}(k)\|^2] \leq [\min_i \lambda_i(\mathbf{B})]^{-1} \sup \left[ C \frac{\eta}{\eta - 1}, \mathbf{X}'_0 \mathbf{B} \mathbf{X}_0 \right],$$

i.e., a bound has been found for the *second moment* of  $\|\mathbf{X}(k)\|$ . The same kind of result can be obtained for a Luré form.

**5. Conclusions.**

**5.1. Linear randomly sampled systems.** For periodically sampled systems, it is well known that if (see (1))

$$(43) \quad |\lambda_i\{G(T)\}| < 1,$$

then :

- (i) the unforced system is asymptotically stable in the large,
- (ii) there exists a unique solution to the matrix equation

$$(44) \quad G'(T)BG(T) - B = -Q,$$

i.e., given  $Q$  positive definite, there exists a unique positive definite  $B$  such that

- (iii) the system is bounded input–bounded output stable.

For randomly sampled systems with independent identically distributed sampling intervals, it is shown that if

$$(45) \quad |\lambda_i\{E[G(k) \otimes G(k)]\}| < 1,$$

then :

- (i) the unforced system is almost surely asymptotically stable and also shows stability of the second moment (see Kalman [3], Bharucha [5]);
- (ii) the matrix equation

$$(46) \quad E[G'(k)BG(k)] - B = -Q$$

has a unique solution and  $Q$  positive definite implies  $B$  positive definite (see Agniel and Jury [10], or our Appendix B);

- (iii) the system is almost surely bounded input–bounded output stable and also shows boundedness of the second moment for a bounded input (see § 2 and § 4).

In other words, with respect to almost sure stability and stability of the first and second moments, condition (45) is similar to condition (43) for periodically sampled systems.

It was shown by Agniel and Jury [10] that condition (45) is a necessary assumption to show the stability of nonlinear randomly sampled systems. It will be proved in a forthcoming paper that this condition is of interest to show almost sure stability by a Popov-type method.

**5.2. Nonlinear randomly sampled systems.** It was shown that provided the autonomous system was proven to be almost surely stable by a Lyapunov function approach, a noise with zero mean and finite energy introduced bounded disturbances.

The bound for  $EV_k$  is  $\sup [C(\eta/(\eta - 1)), V_0]$ . This introduces an interesting physical interpretation, i.e., the bound is given by the major cause of disturbances,



$V_0$  representing the initial conditions and  $C(\eta/(\eta - 1))$  the size of the input. Moreover, it is not difficult to show (see Agniel [8]) that the bigger  $\eta$  is, the faster the system converges to its equilibrium for a null input and the smaller is the parameter domain of stability. Here, the bound  $C(\eta/(\eta - 1))$  is continuously decreasing when we increase  $\eta > 1$ , i.e., when we make the system more and more stable.

**5.3. Related work.** In his book, Kushner [14] gives a comprehensive study of equations of type (14) for finite time stability and excursion times.

**Appendix A.**

SUBMARTINGALE CONVERGENCE THEOREM (see Loève [12, p. 393]). *Let the random variables  $X_n$  form a submartingale sequence.*

(i) *If  $\sup EX_n^+ < \infty$ , then  $X_n \xrightarrow{a.s.} X < \infty$  with*

$$EX \leq \sup EX_n^+,$$

$$E|X| \leq \sup E|X_n|.$$

(ii)  $X_n \xrightarrow{r} X$  where  $r \geq 1$  if and only if the  $|X_n|^r$  are uniformly integrable, and then  $X_n \xrightarrow{a.s.} X$ .

**Appendix B.** We wish to solve

$$(B.1) \quad E[G'(k)BG(k)] - B = -Q.$$

*Problem.* Given a positive definite matrix  $Q$ , under what conditions can we find a positive definite matrix  $B$  satisfying (B.1)?

*Existence and uniqueness.* Although  $G(k)$  is not necessarily symmetric,  $G'(k)BG(k)$  is symmetric and it does not affect the generality of the proof to assume  $B$  and  $Q$  symmetric (see Gantmacher [13, p. 294]). If  $b_{ij}$  and  $q_{ij}$  are the general elements of  $B$  and  $Q$  respectively, (B.1) can be rewritten as

$$(B.2) \quad [(E[G(k) \otimes G(k)]) - I] \begin{bmatrix} b_{11} \\ b_{12} \\ \vdots \\ b_{1n} \\ b_{22} \\ b_{23} \\ \vdots \\ b_{2n} \\ b_{33} \\ \vdots \\ b_{3n} \\ \vdots \\ b_{nn} \end{bmatrix} = - \begin{bmatrix} q_{11} \\ q_{12} \\ \vdots \\ q_{1n} \\ q_{22} \\ q_{23} \\ \vdots \\ q_{2n} \\ q_{33} \\ \vdots \\ q_{3n} \\ \vdots \\ q_{nn} \end{bmatrix},$$

where  $G(k) \otimes G(k)$  is the first Krönecker product of  $G(k)$  by itself. From (B.2) we see that given the  $q_{ij}$ 's we can determine the  $b_{ij}$ 's provided

$$(B.3) \quad \lambda_i [E\{G(k) \otimes G(k)\}] \neq 1.$$

Since  $B$  and  $Q$  are taken symmetric, we can determine  $B$  completely from  $Q$ . Condition for a positive definite  $B$ . Consider the linear system described by

$$(B.4) \quad \begin{aligned} \mathbf{X}(k+1) &= G(k)\mathbf{X}(k), \\ \mathbf{X}(0) &= \mathbf{X}_0, \end{aligned} \quad k = 0, 1, 2, 3, \dots,$$

where  $\mathbf{X}(k)$  is an  $n$ -dimensional vector with components  $x_1(k), x_2(k), \dots, x_n(k)$  and  $G(k)$  is an  $n \times n$  matrix continuous in  $T_k$ . If  $B$  is positive definite, condition (B.1) implies that the system (B.4) shows stability of the second moment, and it is necessary that (see Kalman [3], Bharucha [5])

$$(B.5) \quad \sup_i |\lambda_i \{E[G(k) \otimes G(k)]\}| < 1;$$

then,

$$(B.6) \quad E[x_i(k)x_j(k)] \rightarrow 0 \quad \text{as } k \rightarrow \infty \quad \text{for all } i, j = 1, 2, \dots, n.$$

Conversely, assume that (B.5) is satisfied and  $B$  is not positive definite. Then there exists a state  $\mathbf{X}_0$  such that  $\mathbf{X}'_0 B \mathbf{X}_0 < 0$ . This implies that

$$E[\mathbf{X}'(1)B\mathbf{X}(1)] = \mathbf{X}'_0 B \mathbf{X}_0 - \mathbf{X}'_0 Q \mathbf{X}_0 < 0$$

since  $Q$  is positive definite.

By induction,

$$E[\mathbf{X}'(k)B\mathbf{X}(k)] < \mathbf{X}'_0 B \mathbf{X}_0 < 0 \quad \text{for all } k = 0, 1, 2, \dots,$$

which contradicts (B.6). Therefore  $B$  must be positive definite.

In conclusion, given a positive definite matrix  $Q$ , there exists a unique positive definite matrix  $B$  satisfying (B.1) provided (B.5) is satisfied.

#### REFERENCES

- [1] G. A. BEKEY, J. M. BIDDLE AND A. J. JACOBSON, *The effect of a random sampling interval on a sampled-data model of the human operator*, Proc. 3rd Conference on Manual Control, Rep. SP-144, NASA, 1967.
- [2] J. D. FORESTER, *A stochastic revised sampled data model of the human operator*, S.M. Thesis Rep. MVT-68-2, Manned Vehicle Laboratory, Center for Space Research, MIT, Cambridge, Mass., 1968.
- [3] R. E. KALMAN, *Analysis and synthesis of linear dynamical systems operating on randomly sampled data*, Doctoral thesis, Columbia University, New York, 1957.
- [4] A. R. BERGEN, *Stability of systems with randomly time varying parameters*, IRE Trans. Automatic Control, AC-5 (1960), pp. 265-269.
- [5] B. H. BHARUCHA, *On the stability of randomly varying systems*, Doctoral thesis, University of California, Berkeley, 1961.
- [6] R. S. BUCY, *Stability and positive supermartingales*, J. Differential Equations, 1 (1965), pp. 151-155.
- [7] R. S. BUCY AND P. D. JOSEPH, *Filtering for Stochastic Processes with Application to Guidance*, Interscience, New York, 1968.
- [8] R. G. AGNIEL, *Stability of randomly sampled systems*, Doctoral thesis, University of California, Berkeley, 1969.

- [9] R. E. KALMAN AND J. E. BERTRAM, *A unified approach to the theory of sampling systems*, J. Franklin Institute, 267 (1959), pp. 405–436.
- [10] R. G. AGNIEL AND E. I. JURY, *Stability of nonlinear randomly sampled systems*, Allerton Conference, Urbana, Illinois, 1969, pp. 710–719.
- [11] R. E. KALMAN AND J. E. BERTRAM, *Control system analysis and design via the second method of Lyapunov. I: Continuous time systems; II: Discrete time systems*, Trans. ASME Ser. D. J. Basic Engrg., 82 (1960), pp. 371–400.
- [12] M. LOÈVE, *Probability Theory*, 3rd ed., Van Nostrand, Princeton, N.J., 1963.
- [13] F. R. GANTMACHER, *The Theory of Matrices*, Chelsea, New York, 1959.
- [14] H. J. KUSHNER, *Stochastic Stability and Control*, Academic Press, New York, 1967.
- [15] H. J. KUSHNER AND L. TOBIAS, *On the stability of randomly sampled systems*, IEEE Trans. Automatic Control, AC-14 (1969), pp. 219–324.

## THE SOLUTION OF A QUADRATIC PROGRAMMING PROBLEM USING SYSTEMATIC OVERRELAXATION\*

COLIN W. CRYER†

**Abstract.** Let  $\mathbf{A}$  be a real symmetric positive definite  $n \times n$  matrix and  $\mathbf{b}$  a real column  $n$ -vector. We consider the following problem: Find real column  $n$ -vectors  $\mathbf{x}$  and  $\mathbf{y}$  such that

$$\begin{aligned}\mathbf{Ax} &= \mathbf{b} + \mathbf{y}, \\ \mathbf{x}^T \mathbf{y} &= 0, \quad \mathbf{x} \geq 0, \quad \mathbf{y} \geq 0.\end{aligned}$$

Problems of this type occur when the method of Christopherson is used to solve free boundary problems for journal bearings. In such cases,  $\mathbf{A}$  is a "finite difference" matrix.

We present a method for solving the above problem which is a modification of systematic overrelaxation. This method is particularly suitable when  $\mathbf{A}$  is a finite difference matrix.

**1. Introduction.** Let  $\mathbf{A} = (a_{ij})$  be a real symmetric positive definite  $n \times n$  matrix and  $\mathbf{b} = (b_i)$  a real column  $n$ -vector. We shall be concerned with the following problem.

*Problem 1.* Find real column  $n$ -vectors  $\mathbf{x} = (x_i)$  and  $\mathbf{y} = (y_i)$  such that

$$(1.1) \quad \mathbf{Ax} - \mathbf{y} = \mathbf{b},$$

$$(1.2) \quad \mathbf{x}^T \mathbf{y} = 0,$$

$$(1.3) \quad \mathbf{x} \geq 0, \quad \mathbf{y} \geq 0.$$

It is known that Problem 1 is equivalent to a quadratic programming problem, Problem 2 (see § 2). Both Problems 1 and 2 have been extensively studied from the viewpoint of linear and quadratic programming (Cottle and Dantzig [1], Hadley [5, p. 212], and Lemke [6]) and there are many methods available for solving these problems.

Our interest in Problem 1 arose because problems of this type occur when the method of Christopherson is used to solve free boundary problems for journal bearings (Cryer [3]). In such cases, Problem 1 has certain features which are unusual in nonlinear programming problems:

- (i)  $\mathbf{A}$  is a large matrix, perhaps a  $10,000 \times 10,000$  matrix.
- (ii)  $\mathbf{A}$  is a "finite difference" matrix. Typically, each row of  $\mathbf{A}$  will have no more than five nonzero elements. However,  $\mathbf{A}^{-1}$  is a full matrix.
- (iii) Because of the physical significance of the solution vector  $\mathbf{x}$ , most of the components  $x_i$  may be expected to be positive.

When these features are present, the conventional methods for solving Problems 1 and 2 have substantial disadvantages.

In § 3, we introduce a method for solving Problem 1 which is particularly suitable when  $\mathbf{A}$  is a "finite difference" matrix, since the method is a modified version of S.O.R. (systematic overrelaxation). In § 3, we prove that this method

\* Received by the editors February 9, 1970, and in revised form November 30, 1970.

† Computer Sciences Department, University of Wisconsin, Madison, Wisconsin 53706. This work was supported by the Office of Naval Research under Contract N00014-67-A-0004.

converges, and in § 4, we study how the rate of convergence depends upon the relaxation parameter.

**2. Existence and uniqueness of solution.** It is convenient to introduce the following quadratic programming problem.

*Problem 2.* Find a column  $n$ -vector  $\mathbf{x}$  which maximizes

$$(2.1) \quad f(\mathbf{x}) = \mathbf{b}^T \mathbf{x} - (\mathbf{x}^T \mathbf{A} \mathbf{x})/2,$$

subject to the constraints

$$(2.2) \quad \mathbf{x} \geq 0.$$

**THEOREM 2.1.** *Problems 1 and 2 are equivalent: if  $\{\mathbf{x}, \mathbf{y}\}$  is a solution of Problem 1, then  $\mathbf{x}$  is a solution of Problem 2; if  $\mathbf{x}$  is a solution of Problem 2, then  $\{\mathbf{x}, \mathbf{A}\mathbf{x} - \mathbf{b}\}$  is a solution of Problem 1.*

*There exists a unique solution to Problem 2 (and hence to Problem 1).*

*Proof.* Since  $\mathbf{A}$  is positive definite, the equivalence of Problems 1 and 2 follows from the Kuhn–Tucker theory (Hadley [5, pp. 212–214]).

Since  $\mathbf{A}$  is positive definite,  $f(\mathbf{x})$  is strictly concave (Hadley [5, p. 213]). Hence, since  $\mathbf{x} = \mathbf{0}$  is a “feasible” solution of Problem 2, there exists a unique solution to Problem 2. The proof of the theorem is therefore complete.

**3. Application of S.O.R.** We study the following algorithm for solving Problem 1.

**ALGORITHM 1.** Choose a column  $n$ -vector  $\mathbf{x}^{(0)} = (x_i^{(0)})$ , where  $\mathbf{x}^{(0)} \geq 0$ . Choose a *relaxation parameter*  $\omega$ , where  $0 < \omega < 2$ .

Generate a sequence of column  $n$ -vectors  $\mathbf{x}^{(k)} = (x_i^{(k)})$ ,  $\mathbf{r}^{(k)} = (r_i^{(k)})$ ,  $\mathbf{y}^{(k)} = (y_i^{(k)})$ ,  $k = 1, 2, \dots$ , using the equations,

$$(3.1) \quad r_i^{(k+1)} = b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(k+1)} - \sum_{j=i}^n a_{ij} x_j^{(k)},$$

$$(3.2) \quad x_i^{(k+1)} = \max \{0, x_i^{(k)} + \omega r_i^{(k+1)} / a_{ii}\},$$

$$(3.3) \quad y_i^{(k+1)} = -r_i^{(k+1)} + a_{ii} x_i^{(k+1)} - x_i^{(k)}.$$

(We remind the reader that we have assumed that  $\mathbf{A}$  is positive definite so that  $a_{ii} > 0$  for  $1 \leq i \leq n$ .)

Algorithm 1 is a generalization of methods used by Christopherson [2] and Gnanadoss and Osborne [4]; a brief account of the history of the algorithm is given in Cryer [3].

Algorithm 1 can be interpreted in two ways. On the one hand, Algorithm 1 consists of applying S.O.R. to the equations  $\mathbf{A}\mathbf{x} = \mathbf{b}$  with the proviso that the vectors  $\mathbf{x}^{(k)}$  should be nonnegative. On the other hand, as will be seen in the proof of Theorem 3.1,  $f(\mathbf{x}^{(k+1)}) \geq f(\mathbf{x}^{(k)})$ , so that Algorithm 1 is a method for maximizing  $f(\mathbf{x})$  subject to the restraint that  $\mathbf{x} \geq 0$ . Of course, it is not surprising that two interpretations of Algorithm 1 exist, since it has been known for a long time (Temple [10]) that there is a connection between relaxation methods and the minimization of quadratic forms.

**THEOREM 3.1.** *Let  $\mathbf{x}^{(k)}$  and  $\mathbf{y}^{(k)}$  be generated by Algorithm 1. Then  $\mathbf{x}^{(k)} \rightarrow \mathbf{x}$  and  $\mathbf{y}^{(k)} \rightarrow \mathbf{y}$ , where  $\{\mathbf{x}, \mathbf{y}\}$  is the solution of Problem 1.*

*Proof.* The method of proof is similar to that used by Schechter [8], [9].

For any column  $n$ -vector  $\mathbf{u}$  let

$$(3.4) \quad G(\mathbf{u}) = -2f(\mathbf{u}) = \mathbf{u}^T \mathbf{A} \mathbf{u} - 2\mathbf{u}^T \mathbf{b}.$$

Then, direct computation shows that if  $\mathbf{u}$  and  $\mathbf{v}$  are column  $n$ -vectors,

$$(3.5) \quad G(\mathbf{u}) - G(\mathbf{v}) = (\mathbf{u} - \mathbf{v})^T \mathbf{A} (\mathbf{u} - \mathbf{v}) + 2(\mathbf{u} - \mathbf{v})^T (\mathbf{A} \mathbf{v} - \mathbf{b}).$$

It is convenient to introduce the vectors  $\mathbf{x}^{(k,l)} = (x_i^{(k,l)})$ , where

$$(3.6) \quad x_i^{(k+1,l)} = \begin{cases} x_i^{(k+1)} & \text{if } 1 \leq i \leq l, \\ x_i^{(k)} & \text{if } l < i \leq n \end{cases}$$

for  $k \geq 0$  and  $0 \leq l \leq n$ . Then,

$$(3.7) \quad \mathbf{x}^{(k+1,0)} = \mathbf{x}^{(k)}, \quad \mathbf{x}^{(k+1,n)} = \mathbf{x}^{(k+1)},$$

and, from (3.1),

$$(3.8) \quad r_i^{(k+1)} = [\mathbf{b} - \mathbf{A} \mathbf{x}^{(k+1,i-1)}]_i.$$

Let

$$(3.9) \quad \omega_{k+1,i} = \begin{cases} [x_i^{(k+1)} - x_i^{(k)}] a_{ii} / r_i^{(k+1)} & \text{if } r_i^{(k+1)} \neq 0, \\ \omega & \text{if } r_i^{(k+1)} = 0. \end{cases}$$

Then, noting (3.2),

$$(3.10) \quad x_i^{(k+1)} = x_i^{(k)} + \omega_{k+1,i} r_i^{(k+1)} / a_{ii},$$

and

$$(3.11) \quad 0 \leq \omega_{k+1,i} \leq \omega.$$

Using (3.5), (3.6), (3.8) and (3.10), we find that

$$\begin{aligned} G(\mathbf{x}^{(k+1,i)}) - G(\mathbf{x}^{(k+1,i-1)}) &= [\mathbf{x}^{(k+1,i)} - \mathbf{x}^{(k+1,i-1)}]^T \mathbf{A} [\mathbf{x}^{(k+1,i)} - \mathbf{x}^{(k+1,i-1)}] \\ &\quad + 2[\mathbf{x}^{(k+1,i)} - \mathbf{x}^{(k+1,i-1)}]^T [\mathbf{A} \mathbf{x}^{(k+1,i)} - \mathbf{b}], \end{aligned}$$

and

$$\begin{aligned} (3.12) \quad G(\mathbf{x}^{(k+1,i)}) - G(\mathbf{x}^{(k+1,i-1)}) &= a_{ii} [x_i^{(k+1)} - x_i^{(k)}]^2 - 2[x_i^{(k+1)} - x_i^{(k)}] r_i^{(k+1)} \\ &= -\omega_{k+1,i} (2 - \omega_{k+1,i}) [r_i^{(k+1)}]^2 / a_{ii}. \end{aligned}$$

Remembering that  $0 < \omega < 2$ , it follows from (3.11) and (3.12) that  $G(\mathbf{x}^{(k+1,i)}) \leq G(\mathbf{x}^{(k+1,i-1)})$ . Therefore, the sequence  $\{G(\mathbf{x}^{(k,i)})\}$  is monotone decreasing. But  $\mathbf{A}$  is positive definite so that  $G(\mathbf{u})$  is strictly convex and hence bounded below.

Consequently, there is a constant,  $G_\infty$  say, such that

$$(3.13) \quad G(\mathbf{x}^{(k,l)}) \downarrow G_\infty.$$

Next, we prove that

$$(3.14) \quad |x_i^{(k+1)} - x_i^{(k)}| \leq [a(-1 + 2/\omega)]^{1/2} [G(\mathbf{x}^{(k+1, i-1)}) - G(\mathbf{x}^{(k+1, i)})]^{1/2},$$

where

$$(3.15) \quad a = \min_i a_{ii}.$$

If  $x_i^{(k+1)} = x_i^{(k)}$ , then (3.14) is trivially true, so that we need only consider the case when  $x_i^{(k+1)} \neq x_i^{(k)}$ . But then, from (3.10),  $\omega_{k+1, i} \neq 0$ , so that, from (3.10), (3.11), (3.12) and (3.15),

$$(3.16) \quad \begin{aligned} G(\mathbf{x}^{(k+1, i-1)}) - G(\mathbf{x}^{(k+1, i)}) &= (-1 + 2/\omega_{k+1, i})a_{ii}[x_i^{(k+1)} - x_i^{(k)}]^2 \\ &\geq (-1 + 2/\omega)a[x_i^{(k+1)} - x_i^{(k)}]^2. \end{aligned}$$

Inequality (3.14) follows immediately from (3.16).

Noting (3.13), it follows from (3.14) that

$$(3.17) \quad x_i^{(k+1)} - x_i^{(k)} \rightarrow 0 \quad \text{as } k \rightarrow \infty, \quad 1 \leq i \leq n.$$

Now, let  $\mathbf{x}$  be any limit point of the sequence  $\{\mathbf{x}^{(k)}\}$ . Then there is an increasing sequence of integers  $\{k_p\}$ ,  $p = 1, 2, \dots$ , such that

$$(3.18) \quad \mathbf{x}^{(k_p)} \rightarrow \mathbf{x} \quad \text{as } p \rightarrow \infty.$$

From (3.1), (3.3) and (3.18), we have that, as  $p \rightarrow \infty$ ,

$$(3.19) \quad \begin{aligned} \mathbf{r}^{(k_p)} &\rightarrow \mathbf{r} = \mathbf{b} - \mathbf{A}\mathbf{x}, \\ \mathbf{y}^{(k_p)} &\rightarrow \mathbf{y} = -\mathbf{r}. \end{aligned}$$

We assert that

$$(3.20) \quad \mathbf{x} \geq 0, \quad \mathbf{r} \leq 0.$$

That  $\mathbf{x} \geq 0$  follows immediately from the fact that  $\mathbf{x}^{(k)} \geq 0$  for all  $k$ . To prove that  $\mathbf{r} \leq 0$ , suppose that this is not the case. Then there is an  $\varepsilon > 0$  and integers  $i_0$  and  $k'_0$  such that  $r_{i_0}^{(k_p)} \geq \varepsilon$  for  $k_p \geq k'_0$ . Hence, from (3.2),

$$x_{i_0}^{(k_p)} - x_{i_0}^{(k_p-1)} \geq \varepsilon\omega/a_{i_0i_0}$$

for  $k_p \geq k'_0$ . But this contradicts (3.17).

Next, we show that

$$(3.21) \quad \mathbf{r}^T \mathbf{x} = 0.$$

Suppose that this is not the case. Then, noting (3.20), we see that there is an  $\varepsilon > 0$  and integers  $i_0$  and  $k'$  such that  $r_{i_0}^{(k_p)} \leq -\varepsilon$  and  $x_{i_0}^{(k_p)} \geq \varepsilon$  for  $k_p \geq k'$ . It follows from (3.2) that if  $k_p \geq k'$ , then  $x_{i_0}^{(k_p-1)} \geq x_{i_0}^{(k_p)}$  and

$$|x_{i_0}^{(k_p)} - x_{i_0}^{(k_p-1)}| \geq \omega\varepsilon/a_{i_0i_0}.$$

But this contradicts (3.17).

From (3.19), (3.20) and (3.21), it follows that  $\{\mathbf{x}, \mathbf{y}\}$  satisfies (1.1) through (1.3) and is the (unique) solution of Problem 1.

To complete the proof of the theorem we must show that the sequence  $\{\mathbf{x}^{(k)}\}$  has at least one limit point. But this is a consequence of the fact that (see (3.13))  $\mathbf{x}^{(k)} \in R$  for all  $k$ , where  $R$  is the compact set

$$R = \{\mathbf{x}; G(\mathbf{x}) \leq G(\mathbf{x}^{(0)})\}.$$

**4. Determination of the optimum relaxation parameter.** It is natural to ask how the convergence of Algorithm 1 depends upon  $\omega$ , and whether there is a value of  $\omega$  for which the rate of convergence is maximized. In this section, we partially answer these questions.

Since we make use of the theory of S.O.R., we first summarize this theory.

Let

$$(4.1) \quad \tilde{\mathbf{A}}\tilde{\mathbf{x}} = \tilde{\mathbf{b}},$$

and

$$(4.2) \quad \tilde{\mathbf{A}} = \tilde{\mathbf{D}} - \tilde{\mathbf{E}} - \tilde{\mathbf{F}},$$

where  $\tilde{\mathbf{x}}$  is a column  $m$ -vector,  $\tilde{\mathbf{D}}$  is a diagonal  $m \times m$  matrix, and  $\tilde{\mathbf{E}}$  and  $\tilde{\mathbf{F}}$  are respectively strictly upper and lower triangular  $m \times m$  matrices. Let

$$(4.3) \quad \mathfrak{Q}_\omega(\tilde{\mathbf{A}}) = (\tilde{\mathbf{D}} - \omega\tilde{\mathbf{E}})^{-1}\{(1 - \omega)\tilde{\mathbf{D}} + \omega\tilde{\mathbf{F}}\}.$$

For a given relaxation parameter  $\omega$  and initial guess  $\tilde{\mathbf{x}}^{(0)}$ , let  $\tilde{\mathbf{x}}^{(k)}$ ,  $k = 1, 2, \dots$ , denote the iterates generated by S.O.R. applied to (4.1). Let

$$(4.4) \quad \tilde{\mathbf{e}}^{(k)} = \tilde{\mathbf{x}}^{(k)} - \tilde{\mathbf{x}}.$$

Then (Varga [11, p. 59]),

$$(4.5) \quad \tilde{\mathbf{e}}^{(k+1)} = \mathfrak{Q}_\omega(\tilde{\mathbf{A}})\tilde{\mathbf{e}}^{(k)}.$$

From (4.5) it can be seen that  $\tilde{\mathbf{e}}^{(k)}$  depends upon  $\tilde{\mathbf{A}}$ ,  $\tilde{\mathbf{e}}^{(0)}$  and  $\omega$ . The asymptotic rate of convergence corresponding to  $\tilde{\mathbf{A}}$  and  $\omega$  is (Varga [11, p. 67])

$$(4.6) \quad R_\infty[\mathfrak{Q}_\omega(\tilde{\mathbf{A}})] = -\log [\rho[\mathfrak{Q}_\omega(\tilde{\mathbf{A}})]],$$

where  $\rho[\mathfrak{Q}_\omega(\tilde{\mathbf{A}})]$  denotes the spectral radius of  $\mathfrak{Q}_\omega(\tilde{\mathbf{A}})$ . Equivalently,

$$(4.7) \quad R_\infty[\mathfrak{Q}_\omega(\tilde{\mathbf{A}})] = -\log [\sup_{\tilde{\mathbf{e}}^{(0)}} \limsup_{k \rightarrow \infty} \|\tilde{\mathbf{e}}^{(k)}\|^{1/k}].$$

The optimum relaxation parameter,  $\omega_b = \omega_b(\tilde{\mathbf{A}})$ , is defined by means of the relation (Varga [11, p. 109])

$$(4.8) \quad R_\infty[\mathfrak{Q}_{\omega_b}(\tilde{\mathbf{A}})] = \max_{0 < \omega < 2} R_\infty[\mathfrak{Q}_\omega(\tilde{\mathbf{A}})].$$

For certain classes of matrices  $\tilde{\mathbf{A}}$ , notably 2-cyclic consistently ordered matrices,  $\omega_b(\tilde{\mathbf{A}})$  is known in terms of the eigenvalues of the Jacobi matrix corresponding to  $\tilde{\mathbf{A}}$  (Varga [11, p. 110]).

Next, we introduce some notation. We set

$$(4.9) \quad Z = \{1, 2, \dots, n\}.$$

Let  $T \subset Z$ ,  $\mathbf{B} = (B_{ij})$  be an  $n \times n$  matrix and  $\mathbf{z} = (z_i)$  be an  $n$ -vector. Then  $|T|$  denotes the number of elements of  $T$ ;  $\mathbf{B}(T)$  is the  $|T| \times |T|$  submatrix of  $\mathbf{B}$  obtained



by deleting those elements  $B_{ij}$  for which  $i \notin T$  or  $j \notin T$ ; and  $\mathbf{z}(T)$  is the  $|T| \times 1$  subvector of  $\mathbf{z}$  obtained by deleting those elements  $z_i$  for which  $i \notin T$ . Finally,

$$(4.10) \quad Z(\mathbf{z}) = \{i \in Z; z_i \neq 0\}.$$

We are now ready to consider Algorithm 1. Let  $\{\mathbf{x}, \mathbf{y}\}$  be the solution of Problem 1, and let  $\{\mathbf{x}^{(k)}, \mathbf{y}^{(k)}\}$  be generated using Algorithm 1. We set

$$(4.11) \quad \mathbf{e}^{(k)} = \mathbf{x}^{(k)} - \mathbf{x}.$$

From (3.1) through (3.3) and (4.11) we see that  $\mathbf{e}^{(k)}$  depends upon  $\mathbf{A}, \mathbf{b}, \mathbf{e}^{(0)}$ , and  $\omega$ . By Theorem 3.1,  $\mathbf{e}^{(k)} \rightarrow 0$  as  $k \rightarrow \infty$  for any  $\mathbf{e}^{(0)}$ .

Corresponding to (4.7) and following Ortega and Rockoff [7] we define the asymptotic rate of convergence corresponding to  $\mathbf{A}, \mathbf{b}$ , and  $\omega$ , to be

$$(4.12) \quad R(\mathbf{A}, \mathbf{b}, \omega) = -\log \left\{ \sup_{\mathbf{e}^{(0)}} \limsup_{k \rightarrow \infty} \|\mathbf{e}^{(k)}\|^{1/k} \right\}.$$

LEMMA 4.1. *Given  $\mathbf{x}^{(0)}$ , there is an integer  $k_0$  such that for  $k \geq k_0$ ,*

$$(4.13) \quad \begin{aligned} x_i^{(k)} &> 0 && \text{if } i \in X, \\ x_i^{(k)} &= 0 \quad \text{and} \quad y_i^{(k)} > 0 && \text{if } i \in Y, \end{aligned}$$

where  $X = Z(\mathbf{x})$  and  $Y = Z(\mathbf{y})$ .

*Proof.* Let  $\mathbf{x}^{(0)}$  be given. Let  $\mathbf{r}^{(k)}$  be as in Algorithm 1. Since  $\mathbf{x}^{(k)} \rightarrow \mathbf{x}$  and  $\mathbf{r}^{(k)} \rightarrow -\mathbf{y}$ , it follows that there is an  $\varepsilon > 0$  and an integer  $k_1$  such that if  $k \geq k_1$ , then

$$\begin{aligned} x_i^{(k)} &> 0 && \text{if } i \in X, \\ y_i^{(k)} &> 0 \quad \text{and} \quad r_i^{(k)} < -\varepsilon && \text{if } i \in Y. \end{aligned}$$

Noting (3.2), we see that (4.13) holds if  $k_0 \geq k_1 + 1 + [\max_i a_{ii} x_i^{(k_1)}] / (\varepsilon \omega)$ .

THEOREM 4.2. *Let  $\mathbf{A}$  and  $\mathbf{b}$  be such that*

$$(4.14) \quad |x_i| + |y_i| > 0, \quad 1 \leq i \leq n,$$

where  $\{\mathbf{x}, \mathbf{y}\}$  is the solution of Problem 1. Let  $X = Z(\mathbf{x})$ .

Then,

$$(4.15) \quad R(\mathbf{A}, \mathbf{b}, \omega) = \begin{cases} \infty & \text{if } X \text{ is empty,} \\ R_\infty(\mathfrak{Q}_\omega[\mathbf{A}(X)]) & \text{otherwise.} \end{cases}$$

*Proof.* Let  $\mathbf{x}^{(0)}$  be given. Then it follows from Lemma 4.1, (3.1), (3.2) and (4.14), that, for  $k \geq k_0$ ,

$$\begin{aligned} \mathbf{e}^{(k+1)}(X) &= \mathfrak{Q}_\omega(\mathbf{A}(X))\mathbf{e}^{(k)}(X), \\ \mathbf{e}^{(k+1)}(Y) &= 0. \end{aligned}$$

The theorem follows from (4.7) and (4.12).

Condition (4.14) is satisfied by "almost all"  $\mathbf{A}$  and  $\mathbf{b}$ , and the following theorem covers an important subclass of the remaining problems.

THEOREM 4.3. *Let*

$$(4.16) \quad 0 < \omega \leq 1$$

and

$$(4.17) \quad a_{ij} \leq 0 \quad \text{for } i \neq j.$$

Then,  $R[\mathbf{A}, \mathbf{b}, \omega] \leq R_\infty[\mathfrak{Q}_\omega(\mathbf{A}[T])]$ , where  $T = Z - Z(\mathbf{y})$ .

*Proof.* Let  $\mathbf{x}^{(0)}$  be given. Using an idea due to Gnanadoss and Osborne [4], we see from Lemma 4.1, (3.1) and (3.2), that for  $k \geq k_0$ ,

$$(4.18) \quad \begin{aligned} \mathbf{e}^{(k+1)}(T) &= \mathbf{C}^{(k+1)}\mathbf{e}^{(k)}(T), \\ \mathbf{e}^{(k+1)}(Y) &= 0. \end{aligned}$$

Here,  $\mathbf{C}^{(k+1)}$  is a  $|T| \times |T|$  matrix such that

$$(4.19) \quad \mathbf{C}^{(k+1)} = \prod_{l=1}^{|T|} (\mathbf{H}^{(k+1,l)}\mathbf{L}^{(l)});$$

$\mathbf{L}^{(l)}$  are  $|T| \times |T|$  matrices such that

$$(4.20) \quad \mathfrak{Q}_\omega[\mathbf{A}(T)] = \prod_{l=1}^{|T|} \mathbf{L}^{(l)};$$

and  $\mathbf{H}^{(k+1,l)}$  is a  $|T| \times |T|$  diagonal matrix with diagonal elements equal to either 0 or 1. In particular, when  $T = Z$ , then  $\mathbf{L}^{(l)} = (L_{ij}^{(l)})$  and  $\mathbf{H}^{(k+1,l)} = \text{diag}(H_{ii}^{(k+1,l)})$ , where

$$L_{ij}^{(l)} = \begin{cases} 1 & \text{if } i = j \text{ and } i \neq l, \\ 1 - \omega & \text{if } i = j = l, \\ -\omega a_{lj}/a_{ll} & \text{if } i = l \text{ and } j \neq l, \\ 0 & \text{otherwise,} \end{cases}$$

$$H_{ii}^{(k+1,l)} = \begin{cases} 1 & \text{if } i \neq l, \\ 1 & \text{if } i = l \text{ and } x_i^{(k+1)} > 0, \\ 0 & \text{otherwise.} \end{cases}$$

From (4.16) and (4.17), it follows that  $\mathbf{L}^{(l)} \geq 0$ ; that is, the elements of  $\mathbf{L}^{(l)}$  are nonnegative. Hence, we see from (4.19) and (4.20) that  $0 \leq \mathbf{C}^{(k+1)} \leq \mathfrak{Q}_\omega[\mathbf{A}(T)]$ . The theorem follows immediately.

On the basis of Theorems 4.2 and 4.3 we make the following conjecture.

CONJECTURE 4.4.  $R[\mathbf{A}, \mathbf{b}, \omega] \leq R_\infty[\mathfrak{Q}_\omega(\mathbf{A}[T])]$ , where  $T = Z - Z(\mathbf{y})$ , and  $\{\mathbf{x}, \mathbf{y}\}$  is the solution of Problem 1.

Theorems 4.2 and 4.3 provide some help in choosing  $\omega$  so as to maximize the rate of convergence of Algorithm 1.

If (4.14) holds, then we see from Theorem 4.2 that we should set  $\omega = \omega_{\text{opt}}$ , where

$$(4.21) \quad \omega_{\text{opt}} = \omega_b[\mathbf{A}(X)].$$

Of course, (4.21) does not give  $\omega_{\text{opt}}$  explicitly, since, in general, neither  $X$  nor  $\omega_b[\mathbf{A}(X)]$  is known explicitly. However:

(i) There are several methods of estimating an optimum overrelaxation parameter  $\omega_b(\tilde{\mathbf{A}})$ . One approach is to obtain rigorous a priori bounds on the

spectral radius of the Jacobi matrix corresponding to  $\tilde{\mathbf{A}}$ . Another approach is to choose a value for the overrelaxation parameter  $\omega$ , iterate a number of times, and use the information gained to estimate  $\omega_b(\tilde{\mathbf{A}})$ . Details and references are given by Varga [11, p. 283].

(ii) If  $\mathbf{A}$  is 2-cyclic, consistently ordered, and satisfies (4.17), then it follows from the Perron–Frobenius theory for nonnegative matrices (Varga [11, p. 26]) that, for any  $X \subset Z$ ,  $\omega_b(\mathbf{A}[X]) \leq \omega_b(\mathbf{A})$ . Remembering that it is in general better to overestimate  $\omega_{\text{opt}}$  rather than underestimate  $\omega_{\text{opt}}$  (Varga [11, p. 114]), we can estimate  $\omega_{\text{opt}}$  by  $\omega_b(\mathbf{A})$ .

(iii) After performing  $k$  iterations of Algorithm 1, the set  $X$  can be estimated, according to Lemma 4.1, by the set  $Z(\mathbf{x}^{(k)})$ .

In [3], some of these ideas were used to estimate  $\omega_{\text{opt}}$  for Christopherson's problem for  $n = 64$ . The estimated value of  $\omega_{\text{opt}}$  was  $\omega^* = 1.906$ , and with  $\omega = \omega^*$  it was found that 146 iterations were needed to reduce the residuals  $r_i^{(k+1)}$  (see (3.1)) to less than  $10^{-7}$ . The iterations were also performed with  $\omega = 1.0, 1.1, 1.2, \dots, 1.9$ ; the maximum number of iterations needed was 811 (for  $\omega = 1.0$ ), and the minimum number needed was 70 (for  $\omega = 1.8$ ).

If (4.14) does not hold, then we can say much less about the choice of  $\omega$ . However, if  $\mathbf{A}$  is 2-cyclic and consistently ordered and (4.17) is satisfied, then, from the Perron–Frobenius theory for nonnegative matrices, it follows that for any  $T \subset Z$ ,  $R_\infty[\mathcal{L}_\omega(\mathbf{A}(T))]$  is a monotone decreasing function of  $\omega$  for  $0 < \omega \leq 1$ . Hence, the results of Theorem 4.3 suggest that we should choose  $\omega \geq 1$ .

**Acknowledgment.** It is a pleasure to thank Professor J. B. Rosen for his help.

#### REFERENCES

- [1] R. W. COTTLE AND G. B. DANTZIG, *Complementary pivot theory of mathematical programming*, Mathematics of the Decision Sciences, Part I, G. B. Dantzig and A. F. Veinott, Jr., eds., American Mathematical Society, Providence, R.I., 1968.
- [2] D. G. CHRISTOPHERSON, *A new mathematical method for the solution of film lubrication problems*, Proc. Inst. Mech. Engrs., 146 (1941), pp. 126–135.
- [3] C. W. CRYER, *The method of Christopherson for solving free boundary problems for infinite journal bearings by means of finite differences*, Tech. Rep. 72, Computer Sciences Dept., University of Wisconsin, Madison, 1969.
- [4] A. A. GNANADOSS AND M. R. OSBORNE, *The numerical solution of Reynolds' equation for a journal bearing*, Quart. J. Mech. Appl. Math., 17 (1964), pp. 241–246.
- [5] G. HADLEY, *Nonlinear and Dynamic Programming*, Addison-Wesley, Reading, Mass., 1964.
- [6] C. E. LEMKE, *On complementary pivot theory*, Mathematics of the Decision Sciences, Part I, G. B. Dantzig and A. F. Veinott, Jr., eds., American Mathematical Society, Providence, R.I., 1968.
- [7] J. M. ORTEGA AND M. L. ROCKOFF, *Nonlinear difference equations and Gauss-Seidel type iterative methods*, SIAM J. Numer. Anal., 3 (1966), pp. 497–513.
- [8] S. SCHECHTER, *Iteration methods for nonlinear problems*, Trans. Amer. Math. Soc., 104 (1962), pp. 179–189.
- [9] ———, *Relaxation methods for convex problems*, SIAM J. Numer. Anal., 5 (1968), pp. 601–612.
- [10] G. TEMPLE, *The general theory of relaxation methods applied to linear systems*, Proc. Roy. Soc. Ser. A., 169 (1939), pp. 476–500.
- [11] R. S. VARGA, *Matrix Iterative Analysis*. Prentice-Hall, Englewood Cliffs, N.J., 1962.

## PERTURBATIONS OF LINEAR CONTROL SYSTEMS\*

JERALD P. DAUER†

**Abstract.** The purpose of this paper is to consider the controllability of linear control systems which are perturbations (with respect to  $L^p$ -norms) of controllable linear systems. We show that the set of all completely controllable linear systems is open and dense in the set of all linear control systems.

**1. Introduction.** The problem of controllability of a linear control system  
 (1)  $\dot{x} = A(t)x + B(t)u \quad (\dot{x} = dx/dt)$

defined on a bounded interval has been studied by several authors [1]–[4]. Kalman, Ho and Narendra [1], Weiss [2], and Silverman and Meadows [3] gave three basic necessary and sufficient conditions when  $A$  and  $B$  are continuous. Conti [4] studied the problem for  $A \in L^1$  and  $B \in L^p$ ,  $1 \leq p < \infty$ . Related results were obtained by Youla [5], who studied realizations of weighting patterns with  $A, B \in L^2$  (see also [6], [7]). In addition, when  $A$  and  $B$  are autonomous, Lee and Markus [8, p. 100] showed that the set of completely controllable linear autonomous systems is open and dense in the set of all linear autonomous control systems.

The object of this paper is to generalize this result of Lee and Markus to non-autonomous linear systems. In §3 we show openness with respect to  $L^p$ -norms for linear systems defined on bounded intervals. We give a density theorem in §4 which also proves an assertion of Kalman, Ho and Narendra [1]. They commented that intuitively linear systems are completely controllable. In fact, as they point out, there are many linear systems that are not completely controllable, such as  $\dot{x}(t) = 2$ . However, if the system represents a physical process that involves approximated parameters, the result of §4 shows that we can assume that the system is completely controllable.

**2. Preliminaries.** Suppose system (1) is defined on a bounded interval  $I = [t_0, t_1]$ . The state vector function  $x$  is  $n$ -dimensional, the control vector function  $u$  is  $m$ -dimensional, and  $A$  and  $B$  are  $n \times n$  and  $n \times m$  matrix functions which are (Lebesgue) integrable on  $I$ . Let  $X$  be the fundamental matrix solution of  $\dot{z} = A(t)z$  such that  $X(t_0)$  is the identity matrix.

If  $f$  is a real-valued measurable function defined on  $I$  and  $r$  is such that  $1 \leq r < \infty$ , then define

$$\|f\|_r = \left( \int_I |f(s)|^r ds \right)^{1/r}.$$

For  $r = \infty$ , define  $\|f\|_r = \text{ess sup } \{|f(s)| : s \in I\}$ . Let  $L^r$  be the set of all such functions  $f$  satisfying  $\|f\|_r < \infty$ .

System (1) is said to be *completely controllable in  $L^r$*  if for every  $x_0, x_1 \in E^n$ , Euclidean  $n$ -space, there exists a control function  $u \in L^r$  such that the solution  $x$  of

$$\dot{x} = A(t)x + B(t)u(t),$$

$$x(t_0) = x_0$$

satisfies  $x(t_1) = x_1$ .

\* Received by the editors June 9, 1970, and in revised form December 3, 1970.

† Department of Mathematics, University of Kansas, Lawrence, Kansas. Now at Department of Mathematics, University of Nebraska, Lincoln, Nebraska 68508.

Let  $p$  be given,  $1 \leq p \leq \infty$ , and let  $q$  satisfy  $1/p + 1/q = 1$ . For  $y \in E^n$ , let  $y^*$  denote the transpose of  $y$ .

The following result is due to Conti [4].

**THEOREM 1.** *Let  $1 < q \leq \infty$  and assume  $B \in L^p$ . System (1) is completely controllable in  $L^q$  if and only if*

$$\inf \left\{ \left( \int_I |y^* X(t_1) X^{-1}(s) B(s)|^p ds \right)^{1/p} : |y^*| = 1 \right\} > 0.$$

A set of vector functions  $\{x_1, \dots, x_k\}$  is said to be linearly independent over  $I$  if for every nonzero vector  $(a_1, \dots, a_k) \in E^k$  there exists a subset  $J$  of  $I$  with positive (Lebesgue) measure such that

$$\sum_{i=1}^k a_i x_i(t) \neq 0 \quad \text{for all } t \in J.$$

Weiss [2] and Kalman, Ho and Narendra [1] proved results similar to the following two theorems. Their proofs can be extended to these results.

**THEOREM 2.** *Assume  $B \in L^\infty$ . System (1) is completely controllable in  $L^1$  if and only if the rows of the matrix function  $\{X^{-1}(t)B(t)\}$  are linearly independent over  $I$ .*

**THEOREM 3.** *Let  $2 \leq p \leq \infty$  and assume  $B \in L^p$ . System (1) is completely controllable in  $L^p$  if and only if*

$$W = \int_I X^{-1}(s)B(s)B(s)^* X^{-1}(s)^* ds$$

*is positive definite.*

In the next section, we will use the following two lemmas in connection with the preceding two theorems.

**LEMMA 1.** *Let  $\{x_1, \dots, x_k\}$  be a set of measurable vector functions in  $E^n$  which is linearly independent over  $I$ . Then there exists  $\varepsilon > 0$  such that any set of vector functions  $\{y_1, \dots, y_k\}$ , satisfying  $|y_i(t) - x_i(t)| < \varepsilon$  for  $1 \leq i \leq k$  and almost every  $t \in I$ , is linearly independent over  $I$ .*

*Proof.* Let the function  $f$  be defined by

$$f(a, G(t)) = \sum_{i=1}^k a_i x_i(t),$$

where  $a = (a_1, \dots, a_k) \in E^k$  and  $G(t) = (x_1(t), \dots, x_k(t))$  is an  $n \times k$  matrix function.

Fix an arbitrarily chosen  $a \in \partial S_1(0)$ , the boundary of the open ball with center 0 and radius 1. Then there exists a set  $J_a \subseteq I$  with positive measure and a number  $b_a > 0$  such that  $|f(a, G(t))| \geq b_a$  for all  $t \in J_a$ . Further, by Lusin's theorem, there exists a compact set  $K_a \subseteq I$  such that  $G$  is continuous on  $K_a$  and such that the set  $K_a \cap J_a$  has positive measure.

Since  $f$  is uniformly continuous on the compact set

$$\mathfrak{A} = \partial S_1(0) \times \{n \times k \text{ matrices } H : |H - G(t)| \leq 1 \text{ for some } t \in K_a\},$$

there exists  $\varepsilon_a > 0$  such that if  $t \in K_a$ ,  $(c, H_c) \in \mathfrak{A}$  and  $|a - c| + |H_c - G(t)| < \varepsilon_a$ ,

then

$$|f(a, G(t)) - f(c, H_c)| < b_a/2.$$

Hence if  $t \in K_a \cap J_a$ , we have  $|f(c, H_c)| \geq b_a/2$ .

Covering the compact set  $\partial S_1(0)$  with the open balls  $S_{\varepsilon_a/2}(a)$ , for  $a \in \partial S_1(0)$ , gives us a finite subcover with centers  $a_1, \dots, a_r$ . Let  $b_{a_i} = b_i$ ,  $\varepsilon_{a_i} = \varepsilon_i$  and  $\varepsilon = (1/2k) \min \{\varepsilon_i : 1 \leq i \leq r\}$ .

Without loss of generality, suppose  $|y_i(t) - x_i(t)| < \varepsilon$  for all  $t \in I$ . Let  $H(t) = (y_1(t), \dots, y_k(t))$ .

Then, given any  $a \in \partial S_1(0)$ , there exists  $i \in \{1, \dots, r\}$  such that  $|a - a_i| < \varepsilon_i/2$ . Hence,  $|a - a_i| + |H(t) - G(t)| < \varepsilon_i$  for all  $t \in K_{a_i} \cap J_{a_i}$ . Therefore,  $|f(a, H(t))| \geq b_i/2 > 0$  for all  $t \in K_{a_i} \cap J_{a_i}$ .

Taking  $a \in E^k$ ,  $a \neq 0$ , we have  $a/|a| \in \partial S_1(0)$ . Using the above argument, there exists  $i \in \{1, \dots, r\}$  such that

$$\frac{1}{|a|} \sum_{j=1}^k a_j y_j(t) \geq \frac{b_i}{2} > 0 \quad \text{for all } t \in K_{a_i} \cap J_{a_i}.$$

Hence  $\{y_1, \dots, y_k\}$  is linearly independent.

The following lemma can be easily shown.

LEMMA 2. *Suppose  $W$  is a positive definite  $n \times n$  matrix. There exists  $\varepsilon > 0$  such that if the  $n \times n$  matrix  $V$  satisfies  $|W - V| < \varepsilon$ , then  $V$  is positive definite.*

**3. Openness.** In this section, we will assume that system (1) is completely controllable in various  $L^r$  spaces. We then examine the controllability of linear systems which are perturbations of system (1). Throughout this section,  $C$  and  $D$  will denote measurable  $n \times n$  and  $n \times m$  matrix functions.

The first result considers the natural setting of controllability in  $L^q$ , where  $B \in L^p$ .

THEOREM 4. *Suppose  $A \in L^1$  and  $B \in L^p$ ,  $1 \leq p \leq \infty$ , and assume that system (1) is completely controllable in  $L^q$ . There exists  $\varepsilon > 0$  such that if*

$$\|A - C\|_p + \|B - D\|_p < \varepsilon,$$

then the system

$$(2) \quad \dot{y} = C(t)y + D(t)u$$

is completely controllable in  $L^q$ .

Since the interval  $I$  is bounded, we have the following result.

COROLLARY. *Assume the conditions of Theorem 4. For every  $r$ ,  $p \leq r \leq \infty$ , there exists  $\varepsilon > 0$  such that if*

$$\|A - C\|_r + \|B - D\|_r < \varepsilon,$$

then system (2) is completely controllable in  $L^q$ .

*Proof of Theorem 4.* We first show that if  $\|A - C\|_1$  is sufficiently small, then the system

$$(3) \quad \dot{y} = C(t)y + B(t)u$$

is completely controllable in  $L^q$ .

Let  $Y$  be the fundamental matrix solution of  $\dot{z} = C(t)z$  such that  $Y(t_0)$  is the identity matrix. Then

$$\dot{Y}(t) = A(t)Y(t) + [C(t) - A(t)]Y(t).$$

Applying the variation of parameters formula, we have

$$Y(t) = X(t)X^{-1}(t_0)Y(t_0) + X(t) \int_{t_0}^t X^{-1}(s)[C(s) - A(s)]Y(s) ds.$$

Letting  $M = \max_{t \in I} \|X(t)\|$  and  $N = \max_{t \in I} \|X^{-1}(t)\|$ , we have

$$\|Y(t) - X(t)\| \leq M^2N \|A - C\|_1 + MN \int_{t_0}^t \|A(s) - C(s)\| \|Y(s) - X(s)\| ds.$$

Thus, by Grönwall's inequality, we have

$$(4) \quad \|Y(t) - X(t)\| \leq M^2N \|A - C\|_1 e^{(MN \|A - C\|_1)t}.$$

Similarly, using adjoint equations,

$$(5) \quad \|Y^{-1}(t) - X^{-1}(t)\| \leq MN^2 \|A - C\|_1 e^{(MN^2 \|A - C\|_1)t}.$$

Assume  $q > 1$ , and let

$$\delta = \inf \left\{ \left( \int_{t_0}^{t_1} |y^* X(t_1) X^{-1}(s) B(s)|^p ds \right)^{1/p} : |y^*| = 1 \right\}.$$

Theorem 1 states that  $\delta > 0$  and that if  $|y^*| = 1$ , we have

$$\|y^* X(t_1) X^{-1} B\|_p \geq \delta > 0.$$

Hence  $\|B\|_p \neq 0$ . Therefore, using (4) and (5), we can choose  $\varepsilon > 0$  such that if  $\|A - C\|_1 < \varepsilon$ , then

$$\begin{aligned} \frac{\delta}{2\|B\|_p} &> M \|Y^{-1} - X^{-1}\|_\infty + \|Y - X\|_\infty N \\ &+ \|X - Y\|_\infty \|Y^{-1} - X^{-1}\|_\infty. \end{aligned}$$

Hence, for every  $y^*$  such that  $|y^*| = 1$ , we have

$$\|y^* Y(t_1) Y^{-1} B\|_p \geq \|y^* X(t_1) X^{-1} B\|_p - \frac{\delta}{2} \geq \frac{\delta}{2} > 0.$$

Theorem 1 shows that system (3) is completely controllable in  $L^q$  for  $1 < q \leq \infty$ .

If  $q = 1$ , then  $B \in L^\infty$ . By Theorem 2, the rows of the matrix function  $\{X^{-1}(t)B(t)\}$  are linearly independent over  $I$ . Let  $\tilde{\varepsilon} > 0$  be given by Lemma 1 for this function. Then, for sufficiently small  $\|A - C\|_1$ , equation (5) gives

$$\|Y^{-1}(t) - X^{-1}(t)\| < \frac{\tilde{\varepsilon}}{\|B\|_\infty} \quad \text{for all } t \in I.$$

Using Lemma 1 and Theorem 2, we have that system (3) is completely controllable in  $L^q$  for  $q = 1$ .

Next, we show that if  $\|B - D\|_p$  is sufficiently small, then the system

$$(6) \quad \dot{y} = A(t)y + D(t)u$$

is completely controllable in  $L^q$ .

Assume that  $q > 1$  and let  $M, N$  and  $\delta$  be as above. Let  $\varepsilon = \delta/(2MN)$ . Then, for  $\|B - D\|_p < \varepsilon$  and every  $y^*$  with  $|y^*| = 1$ , we have

$$\|y^*X(t_1)X^{-1}D\|_p \geq \|y^*X(t_1)X^{-1}B\|_p - \frac{\delta}{2} \geq \frac{\delta}{2} > 0.$$

Theorem 1 shows that system (6) is completely controllable in  $L^q$  for  $1 < q \leq \infty$ .

The proof of the result for  $q = 1$  is similar to that of the corresponding result in the first part of the proof. Here, we need

$$\|B - D\|_\infty < \frac{\tilde{\varepsilon}}{\|X^{-1}\|_\infty}.$$

Again, we use Lemma 1 and Theorem 2 to show system (6) completely controllable. Combining the results on systems (3) and (6) gives Theorem 4.

The following theorem considers controllability in a smaller set of controls than the preceding result.

**THEOREM 5.** *Suppose  $A \in L^1$  and  $B \in L^p$ ,  $2 \leq p \leq \infty$ , and assume that system (1) is completely controllable in  $L^p$ . There exists  $\varepsilon > 0$  such that if*

$$\|A - C\|_q + \|B - D\|_q < \varepsilon,$$

then the system

$$(2) \quad \dot{y} = C(t)y + D(t)u$$

is completely controllable in  $L^p$ .

**COROLLARY.** *Assume the conditions of Theorem 5. For every  $r, q \leq r \leq \infty$  (in particular for  $r = p$ ), there exists  $\varepsilon > 0$  such that if*

$$\|A - C\|_r + \|B - D\|_r < \varepsilon,$$

then system (2) is completely controllable in  $L^p$ .

*Proof of Theorem 5.* As in the proof of Theorem 4, we will first consider system (3).

Let  $Y, M$  and  $N$  be as in the proof of Theorem 4. Define  $V$  by

$$V = \int_{t_0}^{t_1} Y^{-1}(s)B(s)B(s)^*Y^{-1}(s)^* ds.$$

If  $W$  is defined as in Theorem 3, then we have

$$|W - V| \leq (2N + \|Y^{-1} - X^{-1}\|_\infty)\|Y^{-1} - X^{-1}\|_\infty \int_{t_0}^{t_1} |B(s)|^2 ds.$$

Since  $I$  is bounded, we have  $\|B\|_p < \infty$  implies  $\|B\|_2 < \infty$ . Thus, using (5), we can make  $|W - V|$  as small as we wish by making  $\|A - C\|_1$  sufficiently small. Since  $W$  is positive definite, Lemma 2 and Theorem 3 show that system (3) is completely controllable for  $\|A - C\|_1$  sufficiently small.



Now, consider system (6). Define  $P(t) = B(t) - D(t)$ , let  $W$  be as in Theorem 3, and let

$$U = W - \int_{t_0}^{t_1} X^{-1}(s)B(s)P(s)^*X^{-1}(s)^* ds - \int_{t_0}^{t_1} X^{-1}(s)P(s)B(s)^*X^{-1}(s)^* ds.$$

Then

$$|W - U| = 2\|X^{-1}\|_\infty^2\|B\|_p\|P\|_q.$$

Since  $W$  is positive definite, it follows from Lemma 2 that if  $\|B - D\|_q = \|P\|_q$  is sufficiently small, then  $U$  is positive definite. Noting that

$$\int_{t_0}^{t_1} X^{-1}(s)P(s)P(s)^*X^{-1}(s)^* ds$$

is positive semidefinite, we have that if  $\|B - D\|_q$  is sufficiently small, then

$$\begin{aligned} &\int_{t_0}^{t_1} X^{-1}(s)D(s)D(s)^*X^{-1}(s)^* ds \\ &= U + \int_{t_0}^{t_1} X^{-1}(s)P(s)P(s)^*X^{-1}(s)^* ds \end{aligned}$$

is positive definite. Theorem 3 shows that system (6) is completely controllable for  $\|B - D\|_q$  sufficiently small. Combining these results gives us Theorem 5.

**4. Density.** The following is a result on the denseness of completely controllable systems. The generality of this result is exhibited in the corollary.

**THEOREM 6.** Consider system (1), where  $A, B \in L^1$ . For each  $\varepsilon > 0$  there exists a measurable  $n \times m$  matrix function  $D$  such that

$$|B(t) - D(t)| < \varepsilon \quad \text{for all } t \in I$$

and such that the system

$$\dot{y} = A(t)y + D(t)u$$

is completely controllable in  $L^\infty$ .

Further, given any  $\delta > 0$ , we can choose  $D$  such that the (Lebesgue) measure of  $\{t \in I : D(t) \neq B(t)\}$  is less than  $\delta$ . Also, if  $B$  is (piecewise) continuous, we can choose  $D$  to be (piecewise) continuous. If  $B$  is constant we can choose  $D$  constant.

*Proof.* Let  $\varepsilon, \delta > 0$  be given. Since  $B$  is measurable, there exists a closed set  $T_1 \subseteq I$  of positive measure such that  $B$  is continuous on  $T_1$ . Hence there exists  $T \subseteq T_1$  and  $\bar{t} \in T$  such that the measure of  $T$  is positive, but less than  $\delta$ , and such that

$$|B(t) - B(\bar{t})| < \varepsilon/2 \quad \text{for all } t \in T.$$

Let  $D_1$  be a matrix whose entries are algebraically independent over the rational numbers and which satisfies

$$|B(\bar{t}) - D_1| < \varepsilon/2.$$

Thus (the determinant)  $\det(D_1) \neq 0$  and

$$|B(t) - D_1| < \varepsilon \quad \text{for all } t \in T.$$

Define the matrix function  $D$  by

$$D(t) = \begin{cases} B(t) & \text{for } t \in I \setminus T, \\ D_1 & \text{for } t \in T. \end{cases}$$

Then  $\det(X(t_1)X^{-1}(s)D(s)) = \det(X(t_1)X^{-1}(s)) \det(D(s)) \neq 0$  for all  $s \in T$ . Therefore, for each  $s \in T$  the rows of the matrix  $\{X(t_1)X^{-1}(s)D(s)\}$  are linearly independent. Hence, for each  $y^*, |y^*| = 1$  and each  $s \in T$ , we have

$$y^* X(t_1)X^{-1}(s)D(s) \neq 0.$$

Since the function  $f(s, y^*) = y^* X(t_1)X^{-1}(s)D(s)$  is uniformly continuous on the compact set  $\mathfrak{A} = T \times \{y^* : |y^*| = 1\}$ , there exists  $(\bar{s}, \bar{y}^*) \in \mathfrak{A}$  such that

$$0 < \alpha = |f(\bar{s}, \bar{y}^*)| \leq |f(s, y^*)| \quad \text{for all } (s, y^*) \in \mathfrak{A}.$$

Hence, for each  $y^*, |y^*| = 1$ , we have

$$\int_T |y^* X(t_1)X^{-1}(s)D(s)| ds \geq \alpha \cdot (\text{measure of } T) > 0.$$

Since  $D \in L^1$ , Theorem 1 gives the desired result.

If  $B$  is (piecewise) continuous, we can assume that  $T = [t_2, t_3]$  for some  $t_2, t_3 \in I$ . By defining  $D_1$  on a subinterval of  $T$  we can easily make  $D$  a (piecewise) continuous function. If  $B$  is constant, define  $D(t) = D_1$ .

**COROLLARY.** Consider system (1) with  $A \in L^1$  and  $B \in L^p, 1 \leq p \leq \infty$ . For each pair  $r, v$ , with  $1 \leq r, v \leq \infty$ , and each  $\varepsilon > 0$  there exists a measurable matrix function  $D$  such that

$$\|B - D\|_r < \varepsilon$$

and such that the system

$$\dot{y} = A(t)y + D(t)u$$

is completely controllable in  $L^v$ .

**5. Remark.** Suppose  $A, B \in L^1$ , and let completely controllable mean completely controllable in  $L^\infty$ . Theorems 4 and 6 show that the set of completely controllable linear systems is open and dense in the set of all linear control systems with respect to the  $L^p$ -norm for  $1 \leq p \leq \infty$ . This statement is also true if we replace the words linear systems by continuous linear systems, piecewise continuous linear systems or constant linear systems.

**Acknowledgment.** This paper is a chapter of the author's doctoral thesis written under the guidance of Professor Fred S. Van Vleck at the University of Kansas.

REFERENCES

[1] R. E. KALMAN, Y. C. HO AND K. S. NARENDRA, *Controllability of linear dynamical systems*, Contributions to Differential Equations, 1 (1963), pp. 189–213.

- [2] L. WEISS, *The concepts of differential controllability and differential observability*, J. Math. Anal. Appl., 10 (1965), pp. 442–449.
- [3] L. M. SILVERMAN AND A. E. MEADOWS, *Controllability and observability in time-variable linear systems*, this Journal, 5 (1967), pp. 64–73.
- [4] R. CONTI, *Contributions to linear control theory*, J. Differential Equations, 1 (1965), pp. 427–445.
- [5] D. C. YOULA, *The synthesis of linear dynamical systems from prescribed weighting patterns*, SIAM J. Appl. Math., 14 (1966), pp. 527–549.
- [6] R. E. KALMAN, *Mathematical description of linear dynamical systems*, this Journal, 1 (1963), pp. 152–192.
- [7] L. WEISS AND R. E. KALMAN, *Contributions to linear system theory*, Internat. J. Engrg. Sci., 3 (1965), pp. 141–171.
- [8] E. B. LEE AND L. MARKUS, *Foundations of Optimal Control Theory*, John Wiley, New York, 1967.

## BOUNDARY VALUE CONTROL THEORY OF THE HIGHER-DIMENSIONAL WAVE EQUATION, PART II\*

DAVID L. RUSSELL†

**Abstract.** The present work extends controllability results obtained for the wave equation in two or more space variables in an earlier article [13]. In the earlier paper approximate controllability was established for time  $T > 2T_0$ , where  $T_0$  is a constant related to the wave propagation speeds in the medium, but only for three or fewer space dimensions. In the present article we establish this result for an arbitrary space dimension. We also examine the controllability problem for  $T = 2T_0$ , the "critical time," and show that here controllability depends upon certain relationships between the coefficients of the partial differential equation and the shape of the spatial domain under consideration.

**1. Introduction.** This paper is a sequel to [13], where we began our study of the approximate controllability of the higher-dimensional wave equation with boundary value controls. There, and here, we let  $\Omega$  be a bounded, open, connected domain in  $R^n$  whose boundary,  $\Gamma$ , is an analytic (or  $C^\infty$  and piecewise analytic)  $(n - 1)$ -dimensional surface in  $R^n$ . We parametrize  $\Gamma$  with an  $(n - 1)$ -dimensional vector variable  $s$  and indicate points on  $\Gamma$  by  $x(s)$ . Integrals over  $\Omega$  are written as  $\int_\Omega (\cdot) dx$ , while integrals over  $\Gamma$  are written  $\int_\Gamma (\cdot) ds$ . Taking  $\tilde{\Gamma}$  to be a relatively open subset of  $\Gamma$  and  $T$  a positive number, we define an admissible control to be a function  $f: \Gamma \otimes [0, T] \rightarrow R^1$  such that  $f \in C^\infty(\Gamma \otimes [0, T])$  and  $f$  vanishes identically outside a compact subset of  $\tilde{\Gamma} \otimes (0, T)$ .

For all such admissible controls  $f$  we let  $w^f(x, t)$  solve the linear hyperbolic mixed initial boundary value problem

$$(1.1) \quad \rho(x)w_{tt}^f - \sum_{i,j=1}^n (\alpha_{ij}(x)w_i^f)_j = 0 \quad \text{in } \Omega \otimes [0, T],$$

$$(1.2) \quad w_x^f(x(s), t)A(x(s))\eta(x(s)) = f(s, t) \quad \text{on } \Gamma \otimes [0, T],$$

$$(1.3) \quad w^f(x, 0) \equiv w_x^f(x, 0) \equiv 0, \quad x \in \Omega.$$

The subscripts  $t$  and  $i$  denote partial differentiation with respect to  $t$  and  $x^i$  (the  $i$ th component of  $x \in R^n$ ) respectively. The subscript  $x$  indicates the gradient vector of the function to which it is applied. The vector  $\eta(x(s))$  is the outward unit normal to  $\Gamma$  at  $x(s) \in \Gamma$ . The real analytic functions  $\rho(x), \alpha_{ij}(x), i, j = 1, 2, \dots, n$ , are such that

$$\begin{aligned} \alpha_{ij}(x) &= \alpha_{ji}(x), \\ \rho(x) &\geq \rho_0 > 0, \\ v'A(x)v &\geq \delta_0 \|v\|^2, \quad \delta_0 > 0, \end{aligned}$$

\* Received by the editors September 16, 1970.

† Departments of Mathematics and Computer Sciences, University of Wisconsin, Madison, Wisconsin 53706 and Honeywell, Inc., St. Paul, Minnesota. This research was supported in part by the National Science Foundation under Grant GP-20858.

in some open set which includes  $\Omega \cup \Gamma$ . Here  $A(x)$  is the  $n \times n$  symmetric matrix whose entries are  $\alpha_{ij}(x)$ .

From [3] and [8] we learn that (1.1), (1.2), (1.3) has a unique  $C^\infty$ -solution in  $\Omega \otimes [0, T]$ . Thus we may let  $R_T$  denote the set of all terminal states  $(w^f(\cdot, T), w_t^f(\cdot, T))$ . The set  $R_T$  is a subspace of the Hilbert space  $H_E(\Omega)$  of finite energy states with inner product

$$\langle (u, u_t); (v, v_t) \rangle_E = \int_{\Omega} [\rho(x)u_t(x)v_t(x) + u_x(x)A(x)v_x(x)'] dx$$

(here  $'$  denotes the transpose of a vector) and norm

$$\|(v, v_t)\|_E = (\langle (v, v_t); (v, v_t) \rangle_E)^{1/2}.$$

The gradients  $u_x, v_x$  are defined in the sense of the theory of distributions. To avoid an indefinite inner product, two states which differ by  $(c, 0)$ , where  $c$  is a constant function on  $\Omega$ , are identified. However, we will continue to speak of elements of  $H_E(\Omega)$  as "states" rather than as "equivalence classes of states."

The control system (1.1), (1.2) is said to be approximately controllable in time  $T$  if  $R_T$  is dense in  $H_E(\Omega)$ , i.e., if the validity of the equation

$$\langle (w^f(\cdot, T), w_t^f(\cdot, T)); (\hat{v}, \hat{v}_t) \rangle_E = 0$$

for all  $f \in R_T$  implies that  $(\hat{v}, \hat{v}_t) = (c, 0)$ , a zero energy state in  $H_E(\Omega)$ .

In [13] we showed that  $\Omega, \Gamma, \rho$  and  $A$  determine a positive number  $T_0$  such that:

- (i) if  $T < 2T_0$ , the system (1.1), (1.2), (1.3) is not approximately controllable in time  $T$ ;
- (ii) if  $T > 2T_0$  and  $n \leq 3$ , then the system is approximately controllable in time  $T$ .

We will refer to  $2T_0$  as the *critical time*. When  $n = 1$  it is known (see [5], [14], [15], e.g.) that approximate controllability continues to hold for  $T = 2T_0$ .

The purpose of the present paper is two-fold. First, we show in §2 that if  $T > 2T_0$ , approximate controllability holds without any restriction on the dimension  $n$ . Second, we show in the remaining sections that if  $n \geq 2$ , approximate controllability may or may not hold for  $T = 2T_0$ , the critical time, depending on certain relationships between  $\Gamma, \rho$  and  $A$ . Because the proofs for  $T = 2T_0$  are very detailed, they are given only for special examples. In the concluding remarks we describe the form which a general theory of critical time approximate controllability would take.

**2. A new proof of approximate controllability for  $T > 2T_0$ .** The theorem which we will prove in this section replaces Theorem 4 in [13]. The new result has the advantage of being valid for all positive integers  $n$ . Many of the details of the proof are the same as in the earlier result. Therefore we will concentrate on the essential differences, referring the reader to [13] for complete treatment of parts common to both proofs.

Let  $(\hat{v}, \hat{v}_t)$  be a finite energy state, i.e.,  $\|(\hat{v}, \hat{v}_t)\|_E < \infty$ , and assume that  $(\hat{v}, \hat{v}_t)$  is orthogonal to all states  $(w^f(\cdot, T), w_t^f(\cdot, T))$  in  $R_T$  relative to the energy inner

product. Thus

$$(2.1) \quad \langle (w^f(\cdot, T), w_t^f(\cdot, T)); (\hat{v}, \hat{v}_t) \rangle_E = \int_{\Omega} [\rho(x)w_t^f(x, T)\hat{v}_t(x) + w_x^f(x, T)A(x)\hat{v}'_x(x)] dx = 0$$

for all admissible controls  $f$ . We let  $v(x, t)$  be the generalized solution of the mixed problem

$$(2.2) \quad \rho(x)v_{tt} - \sum_{i,j=1}^n (\alpha_{ij}(x)v_{ij}) = 0 \quad \text{in } \Omega \otimes [0, T],$$

$$(2.3) \quad v_x(x(s), t)A(x(s))\eta(x(s)) = 0, \quad (x(s), t) \in \Gamma \otimes [0, T],$$

$$(2.4) \quad v(x, T) = \hat{v}(x), \quad v_t(x, T) = \hat{v}_t(x).$$

The existence of such a solution is proved, e.g., in [8] and [10], where it is likewise shown that  $v(\cdot, t)$  and  $v_t(\cdot, t)$  define continuous functions from  $[0, T]$  into  $H^1(\Omega)$  and  $H^0(\Omega) = L^2(\Omega)$ , respectively. (Recall that if  $m$  is a nonnegative integer, then  $H^m(\Omega)$  consists of real functions  $u(x)$  whose derivatives of order  $\leq m$ , taken in the sense of the theory of distributions, lie in  $L^2(\Omega)$ .  $H^m(\Omega)$  is a Hilbert space with inner product

$$(u, \hat{u})_{H^m(\Omega)} = \sum_{0 \leq \|p\| \leq m} \int_{\Omega} [D^p u(x) D^p \hat{u}(x)] dx.$$

Here  $p$  is an  $n$ -vector with nonnegative integer components  $p_1, p_2, \dots, p_n$ ,  $\|p\| = p_1 + p_2 + \dots + p_n$ , and  $D^p$  denotes  $\partial^{\|p\|} / (\partial x^1)^{p_1} (\partial x^2)^{p_2} \dots (\partial x^n)^{p_n}$ .

As in [13] we smooth the solution  $v(x, t)$  by a process of antidifferentiation and formation of finite differences. The innovation lies in the way in which the antiderivatives are defined. We consider the elliptic operator

$$Bu = \frac{1}{\rho(x)} \sum_{i,j=1}^n (\alpha_{ij}(x)u_{ij})$$

which is defined on functions  $u \in C^2(\bar{\Omega})$  ( $\bar{\Omega} = \Omega \cup \Gamma$ ) satisfying the boundary conditions

$$u(x(s))A(x(s))\eta(x(s)) = 0, \quad x(s) \in \Gamma.$$

This unbounded operator is symmetric with respect to the inner product  $(u, \hat{u})_{\rho} = \int_{\Omega} u(x)\hat{u}(x)\rho(x) dx$  and has an unbounded self-adjoint (with respect to that inner product) extension, which we shall still call  $B$ , defined on a domain  $D$  dense in  $L^2(\Omega)$ . (See, e.g., [4], [6].) Moreover, if  $(u, 1)_{\rho} = 0$ , then there is a positive number  $\lambda_0$ , the smallest eigenvalue of  $B$  except 0, such that

$$\|Bu\|_{\rho} \geq \lambda_0 \|u\|_{\rho}.$$

From this it follows that if we let  $\hat{B}$  denote the restriction of  $B$  to  $D \cap \{u \in L^2(\Omega) | (u, 1)_{\rho} = 0\}$ , then  $\hat{B}^{-1}$  is defined, bounded and self-adjoint on  $\{u \in L^2(\Omega) | (u, 1)_{\rho} = 0\}$ , which we shall call  $\hat{D}$ .

From the work of Lions–Magenes [9, p. 165 ff.] it is known that if  $g \in \hat{D} \cap H^m(\Omega)$ ,  $m \geq 0$ , then  $\hat{B}^{-1}g \in \hat{D} \cap H^{m+2}(\Omega)$  and the mapping  $g \in \hat{D} \cap H^m(\Omega) \rightarrow \hat{B}^{-1}g \in \hat{D} \cap H^{m+2}(\Omega)$  is continuous with respect to the norms  $\|\cdot\|_{H^m(\Omega)}$ ,  $\|\cdot\|_{H^{m+2}(\Omega)}$ .

We return to  $(\hat{v}, \hat{v}_t)$  and let  $c_1$  and  $c_2$  be real constants such that

$$(2.5) \quad \int_{\Omega} (\hat{v}(x) - c_1)\rho(x) dx = \int_{\Omega} (\hat{v}_t(x) - c_2)\rho(x) dx = 0.$$

Then

$$\tilde{v}(x, t) = v(x, t) - c_1 - c_2(t - T)$$

satisfies (2.2) and (2.3) and  $\tilde{v}(\cdot, t) \in \hat{D} \cap H^1(\Omega)$ ,  $t \in [0, T]$ . Likewise,

$$\tilde{v}_t(x, t) = v_t(x, t) - c_2$$

is such that  $\tilde{v}_t(\cdot, t) \in \hat{D} \cap H^0(\Omega)$ ,  $t \in [0, T]$ . We define, for each nonnegative integer  $k$ ,

$$D^{-2k}\tilde{v} = \hat{B}^{-k}\tilde{v}, \quad D^{-2k+1}\tilde{v} = \hat{B}^{-k}\tilde{v}_t$$

and conclude from the above cited work in [9] that for a nonnegative integer  $m$ ,

$$D^{-m}\tilde{v}(\cdot, t) \in \hat{D} \cap H^{m+1}(\Omega), \quad t \in [0, T],$$

and that  $D^{-m}\tilde{v}(\cdot, t)$  is a continuous function of  $t$  relative to the norm  $\|\cdot\|_{H^{m+1}(\Omega)}$ . Since  $\tilde{v}(\cdot, t)$  is a generalized solution of  $\tilde{v}_{tt} = B\tilde{v}$  (i.e., (2.2)) one can verify without difficulty that  $D^{-m}\tilde{v}$  satisfies the same equation (in the strict sense if  $m > 0$ ) and that

$$\frac{d^m}{dt^m}(D^{-m}\tilde{v}(\cdot, t)) = \tilde{v}(\cdot, t).$$

Next we define

$$(2.6) \quad D^{-m}v(\cdot, t) = D^{-m}\tilde{v}(\cdot, t) + c_1 \frac{(t - T)^m}{m!} + c_2 \frac{(t - T)^{m+1}}{(m + 1)!}$$

and verify that

$$\frac{d^m}{dt^m}(D^{-m}v(\cdot, t)) = v(\cdot, t).$$

It is not in general true that  $D^{-m}v(\cdot, t)$  is a solution of  $v_{tt} = Bv$ . But since  $c_1$  and  $c_2$  are constants, it is clear that we still have

$$(2.7) \quad D^{-m}v(\cdot, t) \in H^{m+1}(\Omega), \quad t \in [0, T], \quad m \geq 0.$$

We now refer to the theorem of Sobolev (see, e.g., [1, p. 32]) which states that if  $v \in H^m(\Omega)$  and if  $l$  is a positive integer strictly less than  $m - n/2$ , then  $v \in C^l(\Omega)$ . Moreover, there is a constant  $K$ , independent of  $v$ , such that

$$(2.8) \quad \|v\|_{C^l(\bar{\Omega})} \leq K\|v\|_{H^m(\Omega)}.$$

We choose  $m = 2k$  to be a positive integer such that  $m - n/2 > 1$ . Then from (2.7) and the Sobolev theorem we have

$$D^{-m}v(\cdot, t) \in C^2(\bar{\Omega}), \quad D^{-m+1}v(\cdot, t) = D^{-m}v_t(\cdot, t) \in C^1(\bar{\Omega}).$$

The continuity of  $D^{-m}v(\cdot, t)$ ,  $D^{-m+1}v(\cdot, t)$ , as functions of  $t$ , with respect to  $\|\cdot\|_{H^{m+1}(\Omega)}$ ,  $\|\cdot\|_{H^m(\Omega)}$ , respectively, combined with (2.8), then shows that  $D^{-m}v(\cdot, t) \in C^2(\bar{\Omega} \otimes [0, T])$ .

Now, for  $\delta > 0$ , we define

$$\begin{aligned} \Delta(D^{-m}v(\cdot, t)) &= D^{-m}v(\cdot, t + \delta) - D^{-m}v(\cdot, t), & t \in [0, T - \delta], \\ \Delta^k(D^{-m}v(\cdot, t)) &= \Delta(\Delta^{k-1}(D^{-m}v(\cdot, t))), & t \in [0, T - k\delta]. \end{aligned}$$

Noting (2.6), the fact that  $D^{-m}\tilde{v}$  solves  $\tilde{v}_{tt} = B\tilde{v}$ , and the fact that  $v \in C^2(\bar{\Omega} \otimes [0, T])$ , we see that  $\Delta^m(D^{-m}v(\cdot, t))$ , which we will call  $v^\delta(\cdot, t)$ , is such that  $v^\delta(x, t)$  is a  $C^2$ -solution of (2.2), (2.3) in  $\bar{\Omega} \otimes [0, T]$ .

The rest of the proof proceeds much as in [13] and we will give an outline only. The interested reader should consult the earlier paper for details, noting that there  $\tilde{\Gamma}$  of this paper was called  $\Gamma_1$ .

Using the divergence theorem one shows that (2.1) implies (with  $D$  denoting  $\partial/\partial t$ )

$$(2.9) \quad \int_{\tilde{\Gamma} \otimes [0, T]} [D^{-m+1}v(x(s), t)D^m f(s, t)] dx dt = 0$$

for all admissible controls  $f$ . This implies that  $D^{-m+1}v(x(s), t) = (D^{-m}v(x(s), t))_t$  is a polynomial in  $t$  of degree at most  $m - 1$  whose coefficients are  $C^1$ -functions of  $x(s)$ , for  $(x(s), t) \in \tilde{\Gamma} \otimes [0, T]$ . Then

$$(\Delta^m(D^{-m}v(x(s), t)))_t = \Delta^m((D^{-m}v(x(s), t))_t) \equiv 0, \quad (x(s), t) \in \tilde{\Gamma} \otimes [0, T - m\delta].$$

This, combined with the fact that  $\Delta^m(D^{-m}v)$  satisfies the boundary condition (2.3), enables us to use the Holmgren–Fritz John uniqueness theorem [7] to show that  $(\Delta^m(D^{-m}v))_t$  must vanish identically for  $(x, t) \in K(\tilde{\Gamma}, 0, T - m\delta)$ , the intersection of the forward cone of influence of  $\tilde{\Gamma}$  at time 0 with the backward cone of influence of  $\tilde{\Gamma}$  at time  $T - m\delta$ . If  $T > 2T_0$ , the set  $K(\tilde{\Gamma}, 0, T - m\delta)$  includes a set  $\bar{\Omega} \otimes [(T/2) - \varepsilon, (T/2) + \varepsilon]$  for some  $\varepsilon > 0$ , provided  $\delta > 0$  is sufficiently small. (See figures in [13].) Thus,

$$(\Delta^m(D^{-m}v(x, t)))_t \equiv 0, \quad (x, t) \in \bar{\Omega} \otimes [(T/2) - \varepsilon, (T/2) + \varepsilon],$$

which clearly implies

$$(\Delta^m(D^{-m}v(x, t)))_{tt} \equiv 0, \quad (x, t) \in \bar{\Omega} \otimes [(T/2) - \varepsilon, (T/2) + \varepsilon].$$

Since  $\Delta^m(D^{-m}v)$  is a  $C^2$ -solution of (2.2), (2.3) we conclude that

$$v(x, t) \equiv v(x), \quad (x, t) \in \bar{\Omega} \otimes [(T/2) - \varepsilon, (T/2) + \varepsilon],$$

where  $v(x)$  is a  $C^2$ -solution of the elliptic boundary value problem

$$(2.10) \quad \sum_{i,j=1}^n (\alpha_{ij}(x)u_{ij}) = 0, \quad x \in \Omega,$$

$$(2.11) \quad u_x(x(s))A(x(s))\eta(x(s)) = 0, \quad x(s) \in \Gamma.$$

But the only solutions of (2.10), (2.11) have the form

$$u(x) = c, \quad \text{a constant, } x \in \Omega.$$

Thus

$$\Delta^m(D^{-m}v(x, t)) = c, \quad (x, t) \in \bar{\Omega} \otimes [(T/2) - \varepsilon, (T/2) + \varepsilon],$$



so that  $D^{-m}v(x, t)$  is a polynomial in  $t$  of degree at most  $m$  whose coefficients are  $C^2$ -functions of  $x$  for  $x \in \Omega$ . Then  $v(x, t) = D^m(D^{-m}v(x, t))$  is a constant in  $\Omega \otimes [(T/2) - \varepsilon, (T/2) + \varepsilon]$ . In particular,

$$(v(\cdot, T/2), v_t(\cdot, T/2)) = (c, 0),$$

a zero energy state. Applying the conservation of energy principle, which is valid for generalized solutions of (2.2), (2.3), we infer that

$$(v(\cdot, T), v_t(\cdot, T)) = (v, v_t) = (c, 0).$$

We see therefore that if (2.1) holds for all admissible controls  $f$ , so that  $(v, v_t)$  is orthogonal, relative to the energy inner product  $\langle \cdot ; \cdot \rangle_E$ , to every state in  $R_T$ , then  $\|(v, v_t)\|_E = 0$  and  $(v, v_t)$  is the null element in  $H_E(\Omega)$ . We have proved this without making any special assumptions on  $n$ , the dimension of the space in which  $\Omega$  lies. Thus Theorem 4 of [13] can be replaced by the following stronger theorem.

**THEOREM 4a.** *The system (1.1), (1.2) is approximately controllable in time  $T$  if  $T > 2T_0$ .*

Combined with Theorem 2 of [13], which states that the system (1.1), (1.2) is not approximately controllable in time  $T$  if  $T < 2T_0$ , we see that we are justified in referring to  $2T_0$  as the critical time. We will see in the sequel that, if  $n \geq 2$ , critical time approximate controllability is a rather delicate question.

**3. The critical time control problem.** We are going to study the problem for a particular partial differential equation in certain special domains. In § 6 we will indicate a more general theory.

In  $R^n, n \geq 2$ , we consider "rectangles"  $\Sigma_r, r = 1, 2, \dots, n$ , of dimension  $n - r$ , defined by

$$\Sigma_r = \{x = (x^1, x^2, \dots, x^n) \in R^n | x^i = 0, i = 1, \dots, r, 0 \leq x^j \leq 1, j = r + 1, \dots, n\}.$$

Of course,  $\Sigma_n$  is just the origin in  $R^n$ . For all real  $\xi$  we define

$$\tilde{\rho}(\xi) = \exp(1 - 1/\xi^2) \quad (\equiv 0 \text{ if } \xi = 0),$$

and for all  $x = (x^1, x^2, \dots, x^n) \in R^n$  we put

$$\rho(x) = \tilde{\rho}(x^1) + \tilde{\rho}(x^2) + \dots + \tilde{\rho}(x^n).$$

We define domains  $\Omega_r \subseteq R^n$  as follows:

$$\Omega_r = \{x \in R^n | \inf_{y \in \Sigma_r} \rho(x - y) < 1\}.$$

Then  $\Omega_r$  is an open, bounded, simply connected region in  $R^n$  whose boundary

$$\Gamma_r = \{x \in R^n | \inf_{y \in \Sigma_r} \rho(x - y) = 1\}$$

is an  $n$ -dimensional surface of class  $C^\infty$  which is piecewise analytic.

In  $\Omega_r$  we consider a boundary value control problem for the ordinary wave equation :

$$(3.1) \quad w_{tt}^f - \sum_{i=1}^n w_{ii}^f = 0 \quad \text{in } \Omega_r \otimes [0, T],$$

$$(3.2) \quad w_x^f(x(s), t)\eta(x(s)) = f(s, t), \quad (x(s), t) \in \Gamma_r \otimes [0, T],$$

$$(3.3) \quad w^f(x, 0) \equiv w_t^f(x, 0) \equiv 0, \quad x \in \Omega_r.$$

We take  $\tilde{\Gamma} = \Gamma_r$ , i.e., admissible controls are  $C^\infty$ -functions whose supports are compact subsets of the interior of  $\Gamma_r \otimes [0, T]$ . Thus control forces operate over the whole boundary of  $\Omega_r$ .

For (3.1) there is a universal wave propagation speed, 1. Thus, given an instant  $t_0$ , the forward cone of influence of  $\Gamma_r$  at time  $t_0$  is given by

$$(3.4) \quad K^+(\Gamma_r, t_0) = \{(x, t) \in \Omega_r \otimes [t_0, +\infty] \mid \inf_{y \in \Gamma_r} \|x - y\| \leq t - t_0\}$$

and the backward cone of influence of  $\Gamma_r$  at time  $t_0$  is

$$K^-(\Gamma_r, t_0) = \{(x, t) \in \Omega_r \otimes (-\infty, t_0] \mid (x, 2t_0 - t) \in K^+(\Gamma_r, t_0)\}.$$

(In (3.4)  $\|\cdot\|$  denotes the Euclidean norm in  $R^n$ .) We define, for  $T > 0$ ,

$$K(\Gamma_r, 0, T) = K^+(\Gamma_r, 0) \cap K^-(\Gamma_r, T).$$

As shown in [13, § 3], the critical time  $T_0$  has the property that

$$\Omega_r \otimes \{T_0\} \subseteq K(\Gamma_r, 0, 2T_0)$$

but  $\Omega_r \otimes \{T/2\}$  is not a subset of  $K(\Gamma_r, 0, T)$  if  $T < 2T_0$ . In the present case it follows that  $T_0 = 1$ , and hence the critical time is  $T = 2$ , because

$$\sup_{x \in \Omega_r} \{ \inf_{y \in \Gamma_r} (\|x - y\|) \} = 1.$$

We shall prove two theorems regarding approximate controllability of (3.1), (3.2) in the critical time  $T = 2$ . We give these theorems the numbers 5 and 6 since they complement the four theorems proved in [13] and § 2 of the present paper.

**THEOREM 5.** *If  $r = 1$ , the system (3.1), (3.2) is not approximately controllable in the critical time  $T = 2$ .*

**THEOREM 6.** *If  $2 \leq r \leq n$ , the system (3.1), (3.2) is approximately controllable in the critical time  $T = 2$ .*

The reader should be aware that these theorems apply for  $n \geq 2$  only. When  $n = 1$  the analogue of Theorem 5 is not true, for it has already been shown in [5], [14], [15] that we do have critical time approximate controllability in this case.

In order to prove Theorems 5 and 6 we need certain results from the theory of distributions.

**4. Distributions in  $H^{-1}(\Omega_r)$  with support in  $\Sigma_r$ .** As in § 2, we denote by  $H^1(\Omega_r)$  real-valued functions  $v(x)$  defined on  $\Omega_r$  which lie in  $H^0(\Omega_r) = L^2(\Omega_r)$  and have first order partial derivatives, defined in the sense of the theory of

distributions, which also lie in  $H^0(\Omega_r)$ . With the inner product

$$(u, v)_{H^1(\Omega_r)} = \int_{\Omega_r} \left[ u(x)v(x) + \sum_{i=1}^n u_i(x)v_i(x) \right] dx$$

(again the subscript  $i$  refers to partial differentiation with respect to  $x^i$ )  $H^1(\Omega_r)$  is a Hilbert space. We have

$$H^1(\Omega_r) \subseteq H^0(\Omega_r)$$

and for each  $v \in H^1(\Omega_r)$ ,

$$\|v\|_{H^1(\Omega_r)} \geq \|v\|_{H^0(\Omega_r)},$$

which shows that the injection mapping of  $H^1(\Omega_r)$  into  $H^0(\Omega_r)$  is continuous.

We will now indicate the construction of a third Hilbert space  $H^{-1}(\Omega_r)$  with

$$H^1(\Omega_r) \subseteq H^0(\Omega_r) \subseteq H^{-1}(\Omega_r),$$

and the injection of  $H^0(\Omega_r)$  into  $H^{-1}(\Omega_r)$  is likewise continuous. To begin, let  $u \in H^0(\Omega_r)$ . We define a continuous linear functional on  $H^0(\Omega_r)$ :

$$(4.1) \quad l_u(v) = (u, v)_{H^0(\Omega_r)}, \quad v \in H^0(\Omega_r).$$

Now if  $v \in H^1(\Omega_r)$ ,

$$|l_u(v)| \leq \|u\|_{H^0(\Omega_r)} \|v\|_{H^0(\Omega_r)} \leq \|u\|_{H^0(\Omega_r)} \|v\|_{H^1(\Omega_r)}$$

and we conclude that (4.1) also defines  $l_u$  as a continuous linear functional on  $H^1(\Omega_r)$ . It follows that there is a unique element  $\hat{u} \in H^1(\Omega_r)$ , such that

$$(4.2) \quad l_u(v) = (\hat{u}, v)_{H^1(\Omega_r)}.$$

We define

$$(4.3) \quad \|u\|_{H^{-1}(\Omega_r)} = \|\hat{u}\|_{H^1(\Omega_r)}.$$

Now for all  $u \in H^0(\Omega_r)$ ,

$$\begin{aligned} \|u\|_{H^{-1}(\Omega_r)} &= \sup_{\substack{v \in H^1(\Omega_r) \\ v \neq 0}} \frac{|(\hat{u}, v)_{H^1(\Omega_r)}|}{\|v\|_{H^1(\Omega_r)}} \\ &= \sup_{\substack{v \in H^1(\Omega_r) \\ v \neq 0}} \frac{|(u, v)_{H^0(\Omega_r)}|}{\|v\|_{H^1(\Omega_r)}} \leq \sup_{\substack{v \in H^1(\Omega_r) \\ v \neq 0}} \frac{|(u, v)_{H^0(\Omega_r)}|}{\|v\|_{H^0(\Omega_r)}} \\ &= \sup_{\substack{v \in H^0(\Omega_r) \\ v \neq 0}} \frac{|(u, v)_{H^0(\Omega_r)}|}{\|v\|_{H^0(\Omega_r)}} = \|u\|_{H^0(\Omega_r)}, \end{aligned}$$

the second last equality being true because  $H^1(\Omega_r)$  is dense in  $H^0(\Omega_r)$  relative to the topology induced by the norm  $\|\cdot\|_{H^0(\Omega_r)}$ .

We define  $H^{-1}(\Omega_r)$  to be the completion of  $H^0(\Omega_r)$  relative to the norm  $\|\cdot\|_{H^{-1}(\Omega_r)}$ . Now  $\|u\|_{H^{-1}(\Omega_r)} = \|\hat{u}\|_{H^1(\Omega_r)}$  holds for  $u \in H^0(\Omega_r)$ , which is clearly dense in  $H^{-1}(\Omega_r)$ , and this relationship extends (see [12]) to an isometry  $u \leftrightarrow \hat{u}$

between  $H^{-1}(\Omega_r)$  and  $H^1(\Omega_r)$ . The space  $H^{-1}(\Omega_r)$  is a Hilbert space with

$$(4.4) \quad (u, v)_{H^{-1}(\Omega_r)} = (\hat{u}, \hat{v})_{H^1(\Omega_r)}.$$

The elements  $\phi$  of  $H^{-1}(\Omega_r)$  correspond to distributions  $l_\phi$  of order at most 1 (see [16]) on  $\Omega_r$ .

We are now ready to prove two lemmas which will be of great importance in the proofs of Theorems 5 and 6.

LEMMA 1. *If  $n \geq 2$ , there exists a nontrivial element  $\phi \in H^{-1}(\Omega_1)$  such that:*

- (i) *the support of  $l_\phi$  is a subset of  $\Sigma_1$ ;*
- (ii) *if  $c$  is a constant function on  $\Omega_1$ , then  $l_\phi(c) \equiv (\hat{\phi}, c)_{H^1(\Omega_1)} = 0$ .*

LEMMA 2. *If  $2 \leq r \leq n$ , there is no nontrivial distribution in  $H^{-1}(\Omega_r)$  with support in  $\Sigma_r$ .*

*Proof of Lemma 1.* Let  $\psi$  denote a real-valued function of  $n - 1$  variables  $x^2, x^3, \dots, x^n$  such that, with  $\tilde{\Sigma}_1$  defined by

$$\tilde{\Sigma}_1 = \{ \tilde{x} = (x^2, \dots, x^n) \in R^{n-1} \mid (0, \tilde{x}) \in \Sigma_1 \},$$

$\psi \in C^2(\tilde{\Sigma}_1)$  vanishes outside a compact subset of the interior of  $\tilde{\Sigma}_1$ , and

$$(4.5) \quad \int_{\tilde{\Sigma}_1} \psi(\tilde{x}) d\tilde{x} = 0, \quad \int_{\tilde{\Sigma}_1} (\psi(\tilde{x}))^2 d\tilde{x} \neq 0.$$

For positive integers  $k = 4, 5, 6, \dots$  define

$$(4.6) \quad \tilde{\theta}_k(\xi) = \begin{cases} 0, & -1 \leq \xi \leq -3/4, \\ -(1/2)(\xi + 3/4)^2, & -3/4 \leq \xi \leq -1/4, \\ -(1/2)\xi - 1/4, & -1/4 \leq \xi \leq 1/k, \\ (k/4)\xi^2 + 1/(4k) - 1/4, & -1/k \leq \xi \leq 1/k, \\ \tilde{\theta}_k(-\xi), & 1/k \leq \xi \leq 1. \end{cases}$$

Then, for  $x \in \Omega_1$ , put

$$(4.7) \quad \theta_k(x) = \theta_k(x^1, \tilde{x}) = \begin{cases} 0 & \text{if } \tilde{x} \notin \tilde{\Sigma}_1, \\ \tilde{\theta}_k(x^1)\psi(\tilde{x}) & \text{if } x \in \tilde{\Sigma}_1. \end{cases}$$

Then  $\theta_k$  is defined as a function of class  $C^2$  in  $\Omega_1$  for  $k = 4, 5, 6, \dots$ . Now compute, for any  $v \in H^1(\Omega_1)$ ,

$$\begin{aligned} \int_{\Omega_1} \frac{\partial \theta_k(x)}{\partial x^1} \frac{\partial v(x)}{\partial x^1} dx &= - \int_{\Omega_1} \frac{\partial^2 \theta_k(x)}{(\partial x^1)^2} v(x) dx \\ &= - \left[ \int_{[-3/4, -1/4] \otimes \tilde{\Sigma}} -\psi(\tilde{x})v(x) dx + \int_{[-1/k, 1/k] \otimes \tilde{\Sigma}} \frac{k}{2} \psi(\tilde{x})v(x) dx \right. \\ &\quad \left. + \int_{[1/4, 3/4] \otimes \tilde{\Sigma}} -\psi(\tilde{x})v(x) dx \right]. \end{aligned}$$

(In each case the set described beneath the integral is the domain of integration

for that integral.) Thus

$$\int_{[-1/k, 1/k] \otimes \tilde{\Sigma}} \frac{k}{2} \psi(\tilde{x})v(x) dx = \int_{[-3/4, -1/4] \otimes \tilde{\Sigma}} \psi(\tilde{x})v(x) dx + \int_{[1/4, 3/4] \otimes \tilde{\Sigma}} \psi(\tilde{x})v(x) dx - \int_{\Omega_1} \frac{\partial \theta_k(x)}{\partial x^1} \frac{\partial v(x)}{\partial x^1} dx$$

and, for  $k = 4, 5, 6, \dots, j = 4, 5, 6, \dots,$

$$\int_{[-1/k, 1/k] \otimes \tilde{\Sigma}} \frac{k}{2} \psi(\tilde{x})v(x) dx - \int_{[-1/j, 1/j] \otimes \tilde{\Sigma}} \frac{j}{2} \psi(\tilde{x})v(x) dx = \int_{\Omega_1} \left( \frac{\partial \theta_j(x)}{\partial x^1} - \frac{\partial \theta_k(x)}{\partial x^1} \right) \frac{\partial v(x)}{\partial x^1} dx.$$

Applying the Schwarz inequality, we have

$$\begin{aligned} & \left| \int_{[-1/k, 1/k] \otimes \tilde{\Sigma}} \frac{k}{2} \psi(\tilde{x})v(x) dx - \int_{[-1/j, 1/j] \otimes \tilde{\Sigma}} \frac{j}{2} \psi(\tilde{x})v(x) dx \right| \\ (4.8) \quad & \leq \left\| \frac{\partial \theta_j}{\partial x^1} - \frac{\partial \theta_k}{\partial x^1} \right\|_{H^0(\Omega_1)} \left\| \frac{\partial v}{\partial x^1} \right\|_{H^0(\Omega_1)} \\ & \leq \|\theta_j - \theta_k\|_{H^1(\Omega_1)} \|v\|_{H^1(\Omega_1)}. \end{aligned}$$

An inspection of (4.6), (4.7) readily shows that

$$\|\theta_j - \theta_k\|_{H^1(\Omega)} = \varepsilon_{jk},$$

where

$$\lim_{\substack{j \rightarrow \infty \\ k \rightarrow \infty}} \varepsilon_{jk} = 0.$$

Let us put

$$(4.9) \quad \phi_k(x) = \phi_k(x^1, \tilde{x}) = \begin{cases} \frac{k}{2} \psi(\tilde{x}), & (x^1, \tilde{x}) \in \left[-\frac{1}{k}, \frac{1}{k}\right] \otimes \tilde{\Sigma}, \\ 0 & \text{otherwise.} \end{cases}$$

Then (4.8) shows that the continuous linear functionals  $l_{\phi_k}$ , defined on  $H^1(\Omega_1)$  as in (4.1), (4.2), have the property that

$$|l_{\phi_k}(v) - l_{\phi_j}(v)| \leq \varepsilon_{jk} \|v\|_{H^1(\Omega_1)}$$

which implies (cf. (4.2)) that

$$\|\hat{\phi}_k - \hat{\phi}_j\|_{H^1(\Omega_1)} \leq \varepsilon_{jk},$$

and therefore, from (4.3),

$$\|\phi_k - \phi_j\|_{H^{-1}(\Omega_1)} \leq \varepsilon_{jk}.$$

Thus, in  $H^{-1}(\Omega_1)$ ,  $\{\phi_k\}$  is a Cauchy sequence and has a limit  $\phi \in H^{-1}(\Omega_1)$ . It remains only to show that  $\phi$  has properties (i) and (ii) stated in Lemma 1.

Let  $v \in C^\infty(\Omega_1)$  have support  $K$  which is a compact subset of  $\Omega_1 - \Sigma_1$ . Then, for all sufficiently large  $k$ ,  $K \cap ([-1/k, 1/k] \otimes \tilde{\Sigma})$  is empty and

$$l_{\phi_k}(v) = (\phi_k, v)_{H^0(\Omega_1)} = (\hat{\phi}_k, v)_{H^1(\Omega_1)} = 0.$$

Since  $\phi_k$  converges to  $\phi$  in  $H^{-1}(\Omega_1)$ ,  $\hat{\phi}_k$  converges to  $\hat{\phi}$  in  $H^1(\Omega_1)$ , by virtue of the isometry discussed just prior to (4.4). Therefore

$$l_\phi(v) = (\hat{\phi}, v)_{H^1(\Omega_1)} = \lim_{k \rightarrow \infty} (\hat{\phi}_k, v) = 0.$$

Thus  $l_\phi$  vanishes when applied to  $v \in C^\infty(\Omega_1)$  with support  $K$  not meeting  $\Sigma_1$ , and we have shown that the support of  $l_\phi$  must be a subset of  $\Sigma_1$ .

Similarly, for  $k = 4, 5, 6, \dots$ ,  $c$  constant,

$$\begin{aligned} l_{\phi_k}(c) &= (\phi_k, c)_{H^0(\Omega_1)} = \int_{[-1/k, 1/k] \otimes \tilde{\Sigma}_1} \frac{k}{2} \psi(\tilde{x}) c \, dx^1 \, d\tilde{x} \\ &= c \int_{\tilde{\Sigma}_1} \psi(\tilde{x}) \, d\tilde{x} = 0, \end{aligned}$$

from (4.5). Thus part (ii) of Lemma 1 is proved. The second part of (4.5) readily shows that  $\phi$  is nontrivial, and the proof of Lemma 1 is complete.

*Proof of Lemma 2.* For  $p > 0$  we define

$$(4.10) \quad h_p(x^1, x^2, \dots, x^r) = 1 - [(x^1)^2 + (x^2)^2 + \dots + (x^r)^2]^{1/p}.$$

We compute

$$(4.11) \quad \sum_{i=1}^r \left( \frac{\partial h_p}{\partial x^i} \right)^2 = \frac{4}{p^2} [(x^1)^2 + (x^2)^2 + \dots + (x^r)^2]^{2/p-1}.$$

Integrating (4.11) over the unit ball in  $R^r$  we obtain the integral

$$4\omega_{r-1}/(4p + (r - 2)p^2),$$

where  $\omega_{r-1}$  is the integral of 1 over the  $(r - 1)$ -dimensional sphere of radius 1. Thus we see that if  $B$  is any bounded open set in  $R^r$ , then  $h_p \in H^1(B)$  for  $p > 0$  and

$$(4.12) \quad \lim_{p \rightarrow +\infty} \|h_p\|_{H^1(B)} = 0.$$

(Note that  $r \geq 2$  is necessary for these conclusions.)

Given  $x = (x^1, \dots, x^r, x^{r+1}, \dots, x^n) \in R^n$ , let us set  $y = (x^1, \dots, x^r)$ ,  $z = (x^{r+1}, \dots, x^n)$ . Each distribution  $l$  defined on  $R^{n-r}$  has a natural extension to a distribution  $\hat{l}$  defined on  $R^n$ . If  $\hat{v} (= \hat{v}(y, z)) \in C^\infty(R^n)$ , we let  $\hat{v}$  be defined on  $R^{n-r}$  by  $v(z) = \hat{v}(0, z)$ . Then  $\hat{l}(\hat{v}) = l(v)$  defines the extension  $\hat{l}$  of  $l$ .

A result in [9, p. 78] shows that if  $\phi \in H^{-1}(\Omega_r)$  then the distribution  $l_\phi$  associated with  $\phi$  can be expressed in the form

$$l_\phi(u) = (g_0, u)_{H^0(\Omega_r)} + \sum_{i=1}^n \left( g_i, \frac{\partial u}{\partial x^i} \right)_{H^0(\Omega_r)}, \quad u \in C^\infty(\Omega_r),$$

where  $g^i \in H^0(\Omega_r)$ ,  $i = 0, 1, \dots, n$ . This shows that  $l_\phi$  is a distribution of order at most 1 (i.e., if the  $v_k$  are  $C^\infty$ -functions converging to zero in  $C^1(\Omega_r)$  as  $k \rightarrow \infty$ , then  $\lim_{k \rightarrow \infty} l_\phi(v_k) = 0$ ).

Thus if we take  $\phi \in H^{-1}(\Omega)$ ,  $l_\phi$  has order at most 1. A theorem in [16, p. 99 ff.] then shows that if the support of  $l_\phi$  is a subset of  $\Sigma_r$ , there are distributions  $l_0, l_1, l_2, \dots, l_r$  defined on  $R^{n-r}$  with support in  $\tilde{\Sigma}_r = \{z|(0, z) \in \Sigma_r\}$  such that

$$l_\phi = \hat{l}_0 + \sum_{i=1}^r \frac{\partial \hat{l}_i}{\partial x^i}.$$

Let  $\psi(z) \in C^\infty(R^{n-r})$  have support  $K$  contained in some small neighborhood of  $\tilde{\Sigma}_r$  in  $R^{n-r}$ . Define  $v(= v(y, z))$  in  $\Omega_r$  by

$$(4.13) \quad v(x) = v(y, z) = h_4(y)\psi(z),$$

where  $h_4$  is given by (4.10). Then  $v \in H^1(\Omega_r)$ . Since  $\phi \in H^{-1}(\Omega_r)$ , the linear functional  $l_\phi$  can be defined on all of  $H^1(\Omega_r)$  and we have

$$\begin{aligned} l_\phi(v) &= \hat{l}_0(v) + \sum_{i=1}^r \left( \frac{\partial \hat{l}_i}{\partial x^i} \right) (v) \\ &= \hat{l}_0(v) - \sum_{i=1}^r \hat{l}_i \left( \frac{\partial v}{\partial x^i} \right). \end{aligned}$$

Let  $e_j = (0, \dots, 0, 1, 0, \dots, 0)$  in  $R^r$ , the 1 being in the  $j$ th position. Define

$$v_\varepsilon(x) = v_\varepsilon(y, z) = v(y + \varepsilon e_j, z).$$

One verifies without difficulty that

$$(4.14) \quad \lim_{\varepsilon \rightarrow 0} \|v - v_\varepsilon\|_{H^1(\Omega_r)} = 0.$$

On the other hand,

$$\begin{aligned} l_\phi(v_\varepsilon) &= \hat{l}_0(v_\varepsilon) - \sum_{i=1}^r \hat{l}_i \left( \frac{\partial v_\varepsilon}{\partial x^i} \right) \\ &= h_4(0 + \varepsilon e_j)l_0(\psi) - \sum_{i=1}^r \frac{\partial h_4}{\partial x^i}(0 + \varepsilon e_j)l_i(\psi) \\ &= h_4(\varepsilon e_j)l_0(\psi) - \frac{\partial h_4}{\partial x^i}(\varepsilon e_j)l_j(\psi) \\ &= (1 - \varepsilon^{1/2})l_0(\psi) + \frac{1}{2}\varepsilon^{-1/2}l_j(\psi). \end{aligned}$$

Thus if  $l_j(\psi)$  is different from zero,

$$\lim_{\varepsilon \rightarrow 0} l_\phi(v_\varepsilon) = +\infty$$

and then (4.14) shows that  $l_\phi$  cannot be a continuous linear functional on  $H^1(\Omega_r)$ , contrary to our assumption that  $\phi \in H^{-1}(\Omega_r)$ . We conclude that  $l_j(\psi) = 0$ . Since this is true for all such  $\psi$  and the support of  $l_j$  is a subset of  $\tilde{\Sigma}_r$ , we conclude that  $l_j = 0$ . We can do this for  $j = 1, 2, \dots, r$  and conclude that

$$l_\phi = \hat{l}_0.$$

Now, for  $p > 0$ , define  $v_p(x)$  as in (4.13), replacing 4 by  $p$ . Compute

$$l_\phi(v_p) = h_p(0)l_0(\psi) = l_0(\psi)$$

for all  $p > 0$ , since  $h_p(0) = 1$ . But (4.12) is easily seen to imply that

$$\lim_{p \rightarrow \infty} \|v_p\|_{H^1(\Omega_r)} = 0$$

and thus  $l_\phi$  cannot be a continuous linear functional on  $H^1(\Omega_r)$  unless  $l_0(\psi) = 0$ . Therefore, since  $l_\phi$  is a continuous linear functional on  $H^1(\Omega_r)$ ,  $l_0(\psi) = 0$  for all  $\psi$  of the form prescribed above. Then the fact that  $l_0$  has support in  $\tilde{\Sigma}_r$  shows that  $l_0 = 0$ . We have now shown that

$$l_\phi = 0,$$

and Lemma 2 has been proved.

*Remarks.* Some readers may find the dual role of  $l_\phi$ , as a distribution of order  $\leq 1$  and as a linear functional on  $H^1(\Omega_r)$ , slightly confusing. Given  $\phi \in H^{-1}(\Omega_r)$  there is associated with it a unique element  $\hat{\phi} \in H^1(\Omega_r)$  and for all  $v \in H^1(\Omega_r)$ ,

$$l_\phi(v) = (\hat{\phi}, v)_{H^1(\Omega_r)}.$$

This also defines  $l_\phi$  as a continuous linear functional on  $C^\infty(\Omega_r)$ , since convergence in  $C^\infty(\Omega_r)$  implies convergence in  $H^1(\Omega_r)$ . Thus  $l_\phi$  is also a distribution in the sense of Schwartz [16].

One can easily see that Lemma 1, part (i) continues to hold for  $n = 1$ . (Just put  $\phi = \delta$ , the Dirac distribution.) But (ii) cannot hold for  $n = 1$ . The function  $\psi$  cannot be constructed as in (4.5). This explains why Theorem 5 is true for  $n \geq 2$  but not for  $n = 1$ .

**5. Proof of Theorem 5.** A result in Lions–Magenes [9, p. 202] states that if  $\tilde{\phi} \in H^{-1}(\Omega)$  satisfies (ii) of Lemma 1, then there is a unique function  $\tilde{v} \in H^1(\Omega)$  with  $\int_{\Omega_1} \tilde{v}(x) dx = 0$  such that, in the sense of the theory of distributions,

$$(5.1) \quad \sum_{i=1}^n \tilde{v}_{ii} = \tilde{\phi} \quad \text{in } \Omega_1,$$

$$(5.2) \quad \tilde{v}_{,x}(x(s))\eta(x(s)) = 0, \quad x(s) \in \Gamma_1.$$

The sense in which (5.2) holds is also explained in [9]. In our applications  $\tilde{v}$  is harmonic outside a compact subset of  $\Omega_1$  and (5.2) holds in the classical sense. Moreover, there is a constant  $M > 0$  such that

$$(5.3) \quad \|\tilde{v}\|_{H^1(\Omega_1)} \leq M \|\phi\|_{H^{-1}(\Omega_1)}.$$

Let the functions  $\phi_k$  be defined on  $\Omega_1$  as in (4.9) and let  $\tilde{v}^k$  be the corresponding solutions of (5.1), (5.2) with  $\tilde{\phi}$  replaced by  $\phi_k$ . Also, let  $\tilde{v}$  satisfy (5.1), (5.2) with  $\tilde{\phi}$  replaced by the element  $\phi \in H^{-1}(\Omega_1)$  constructed in Lemma 1. Since  $\lim_{k \rightarrow \infty} \|\phi - \phi_k\|_{H^{-1}(\Omega)} = 0$ , (5.3) implies that

$$\lim_{k \rightarrow \infty} \|\tilde{v} - \tilde{v}^k\|_{H^1(\Omega_1)} = 0.$$

It is clear that  $\tilde{v}$  cannot be a constant on  $\Omega_1$ ; therefore  $(\tilde{v}, 0)$  is a nonzero energy state. We let  $v(x, t), v^k(x, t), k = 4, 5, 6, \dots$ , be generalized solutions in  $\Omega_1 \times [0, 2]$  of

$$(5.4) \quad v_{tt} - \sum_{i=1}^n v_{ii} = 0,$$



$$(5.5) \quad v_x(x(s), t)\eta(x(s)) = 0, \quad (x(s), t) \in \Gamma_1 \otimes [0, 2],$$

satisfying

$$v(x, 1) \equiv \tilde{v}(x), \quad v_t(x, 1) \equiv 0, \quad v^k(x, 1) \equiv \tilde{v}^k(x), \quad v_t^k(x, 1) \equiv 0.$$

By the principle of conservation of energy,  $(v(\cdot, 2), v_t(\cdot, 2))$  is also a nonzero energy state.

Let  $f$  be an admissible control. Then the support of  $f$  lies in a set  $\Gamma_1 \otimes [\delta, 2 - \delta]$  for some  $\delta > 0$ . Since the support of  $\phi_k$  is  $[-1/k, 1/k] \otimes \tilde{\Sigma}_1$ ,  $\tilde{v}^k$  is harmonic in  $\Omega_1 - ([-1/k, 1/k] \otimes \tilde{\Sigma}_1)$ . Then, by a familiar uniqueness result in the theory of hyperbolic partial differential equations (see, e.g., [2]),

$$v^k(x, t) \equiv \tilde{v}^k(x), \quad v_t^k(x, t) \equiv 0 \quad \text{for } |t - 1| \leq \inf_{y \in [-1/k, 1/k] \otimes \tilde{\Sigma}_1} \{\|x - y\|\}.$$

Thus, for sufficiently large  $k$ ,  $v_t^k(x(s), t) \equiv 0$ ,  $(x(s), t) \in \Gamma_1 \otimes [\delta, 2 - \delta]$  and an application of the divergence theorem (cf. Theorem 1 in [13]) shows that

$$(5.6) \quad \int_{\Omega_1} \left[ w_t^f(x, 2)v_t^k(x, 2) + \sum_{i=1}^n w_i^f(x, 2)v_i^k(x, 2) \right] dx = \int_{\Gamma_1 \otimes [0, 2]} v_t^k(x(s), t)f(s, t) ds = 0.$$

(The solution  $w^f \in C^\infty(\Omega_1 \otimes [0, 2])$  and it is proved in [10] that  $v^k \in H^2(\Omega^1 \otimes [0, 2])$ . This enables one to use the divergence theorem without difficulty.)

Noting that

$$\begin{aligned} \lim_{k \rightarrow \infty} \|v_t(\cdot, 2) - v_t^k(\cdot, 2)\|_{H^0(\Omega_1)} &= 0, \\ \lim_{k \rightarrow \infty} \|v_i(\cdot, 2) - v_i^k(\cdot, 2)\|_{H^0(\Omega_1)} &= 0, \quad i = 1, 2, \dots, n, \end{aligned}$$

we conclude from (5.6) that

$$\int_{\Omega_1} \left[ w_t^f(x, 2)v_t(x, 2) + \sum_{i=1}^n w_i^f(x, 2)v_i(x, 2) \right] dx = 0.$$

Since  $f$  is an arbitrary admissible control we have shown that the nonzero energy state  $(v(\cdot, 2), v_t(\cdot, 2))$  lies in  $R_2^\perp$  and thus that  $R_2$  is not dense in  $H_E(\Omega_1)$  relative to the norm  $\|\cdot\|_E$ . Thus Theorem 5 is proved.

**6. Proof of Theorem 6.** Much of the work necessary to prove Theorem 6 has already been done in § 2 in the proof of Theorem 4a. We again assume that  $(v, v_t)$  is a finite energy state which satisfies (2.1) (with  $\rho \equiv 1$ ,  $A \equiv I$  and  $\Omega = \Omega_r$ ) and we let  $v(x, t)$  be the generalized solution of (5.4), (5.5) satisfying the terminal conditions (2.4). The solution  $v$  is smoothed by the same process of forming antiderivatives and finite time differences as described in (2.5)–(2.8) ff. The divergence theorem can again be used to obtain (2.9) (with  $\tilde{\Gamma}$  replaced by  $\Gamma_r$ ), and thus, via the Holmgren–Fritz John uniqueness theorem [7] to prove that  $(\Delta^m(D^{-m}v))_t$  must vanish identically for  $(x, t) \in K(\Gamma_r, 0, T - m\delta)$ , the intersection of the forward cone of influence of  $\Gamma_r$  at time 0 with the backward cone of influence of  $\Gamma_r$  at time  $T - m\delta$ .

The essential difference between the proof of Theorem 6 and that of Theorem 4a lies in the fact that when  $T = 2$ , the critical time,  $K(\Gamma_r, 0, 2 - m\delta)$  does not include any set  $\bar{\Omega}_r \otimes [1 - \varepsilon, 1 + \varepsilon]$  for any  $\varepsilon > 0$ , no matter how small we take  $\delta > 0$  to be.

If  $\delta > 0$  is small, the functions  $\Delta^m(D^{-m}v(x, 1))$  are defined and twice continuously differentiable for  $x \in \bar{\Omega}_r$ . Now the operator  $\Delta$  depends on  $\delta$ , and we define

$$v^\delta(x) = \delta^{-m} \Delta^m(D^{-m}v(x, 1)), \quad x \in \bar{\Omega}_r.$$

The continuity of  $v(\cdot, t)$  as a mapping from  $R^1$  into  $H^1(\Omega_r)$  enables one to show by elementary means that

$$(6.1) \quad \lim_{\delta \rightarrow 0} \|v^\delta(x) - v(x, 1)\|_{H^1(\Omega_r)} = 0.$$

Now  $\delta^{-m} \Delta^m(D^{-m}v(x, t)) \equiv v^\delta(x, t)$  is twice continuously differentiable in  $\bar{\Omega}_r \otimes [0, 2 - m\delta]$  and there satisfies

$$\sum_{i=1}^n v_{ii}^\delta = v_{ii}^\delta$$

and boundary conditions of the form (5.5). Thus the functions

$$g^\delta(x) = v_{ii}^\delta(x, T_0)$$

are, for  $\delta > 0$ , continuous in  $\bar{\Omega}_r$  and we have

$$\sum_{i=1}^n v_{ii}^\delta(x) \equiv g^\delta(x), \quad x \in \bar{\Omega}_r.$$

Now  $v_i^\delta(x, t) (= \delta^{-m} \Delta^m(D^{-m}v)_i(x, t))$  has been shown to vanish in  $K(\Gamma_r, 0, 2 - \delta)$ , which implies that  $v_{ii}^\delta(x, t)$  vanishes there also. We conclude therefore that

$$(6.2) \quad g^\delta(x) \equiv v_{ii}^\delta(x, T_0) \equiv 0, \quad x \in \Omega_r^\delta,$$

where

$$(6.3) \quad \Omega_r^\delta = \{x \in \Omega_r | (x, 1) \in K(\Gamma_r, 0, 2 - m\delta) \cap (\Omega_r \otimes \{1\})\}.$$

The sets  $\Omega_r^\delta$  are monotone increasing as  $\delta \rightarrow 0$  with the property

$$(6.4) \quad \bigcap_{\delta > 0} (\Omega_r - \Omega_r^\delta) = \Sigma_r.$$

Let  $u \in H^1(\Omega_r) \subseteq H^0(\Omega_r)$ . Since  $g^\delta \in C^0(\bar{\Omega}_r) \subseteq H^0(\Omega_r)$  we can form the inner product  $(g^\delta, u)_{H^0(\Omega_r)}$ . Integrating by parts we find that

$$\begin{aligned} |(g^\delta, u)_{H^0(\Omega_r)}| &\equiv \left| \left( \sum_{i=1}^n v_{ii}^\delta, u \right)_{H^0(\Omega_r)} \right| \\ &= \left| - \sum_{i=1}^n (v_i^\delta, u_i)_{H^0(\Omega_r)} \right| \leq \|v^\delta\|_{H^1(\Omega_r)} \|u\|_{H^1(\Omega_r)}. \end{aligned}$$

Thus  $g^\delta$  is an element of  $H^0(\Omega_r)$  which defines, via  $(g^\delta, u)_{H^0(\Omega_r)}$ , a continuous linear functional  $l_{g^\delta}$  on  $H^1(\Omega_r)$  for which

$$\|l_{g^\delta}\| \leq \|v^\delta\|_{H^1(\Omega_r)}.$$

There is an element  $\hat{g}^\delta \in H^1(\Omega_r)$  such that

$$l_{g^\delta}(u) = (\hat{g}^\delta, u)_{H^1(\Omega_r)}, \quad u \in H^1(\Omega_r).$$

Then, reasoning as in § 4, we have

$$\|g^\delta\|_{H^{-1}(\Omega_r)} = \|\hat{g}^\delta\|_{H^1(\Omega_r)} \leq \|v^\delta\|_{H^1(\Omega_r)}.$$

Similarly for  $\delta_1 > 0, \delta_2 > 0,$

$$\|g^{\delta_1} - g^{\delta_2}\|_{H^{-1}(\Omega_r)} \leq \|v^{\delta_1} - v^{\delta_2}\|_{H^1(\Omega_r)}.$$

Now if we take a sequence  $\{\delta_k\}$  of positive numbers with  $\lim_{k \rightarrow \infty} \delta_k = 0,$  we have

$$\lim_{k \rightarrow \infty} \|v^{\delta_k} - v(\cdot, 1)\|_{H^1(\Omega_r)} = 0$$

from (6.1). Thus

$$\lim_{\substack{k \rightarrow \infty \\ j \rightarrow \infty}} \|g^{\delta_k} - g^{\delta_j}\|_{H^{-1}(\Omega_r)} = \lim_{\substack{k \rightarrow \infty \\ j \rightarrow \infty}} \|v^{\delta_k} - v^{\delta_j}\|_{H^1(\Omega_r)} = 0$$

and we see that  $\{g^{\delta_k}\}$  is Cauchy in  $H^{-1}(\Omega_r),$  converging to an element  $g \in H^{-1}(\Omega_r).$

Let  $l_g$  be the distribution (also linear functional on  $H^1(\Omega_r)$ ) associated with  $g.$  We claim that the support of  $g$  is contained in  $\Sigma_r.$  For, if  $u \in C^\infty(\Omega_r)$  has support  $K$  which does not meet  $\Sigma_r,$  then (6.4) shows that  $K$  is a subset of  $\Omega_r^\delta$  for sufficiently small  $\delta > 0.$  Then

$$l_g(u) = \lim_{k \rightarrow \infty} l_{g^{\delta_k}}(u) = (g^{\delta_k}, u)_{H^0(\Omega_r)} = 0,$$

as we see from (6.2). Thus  $l_g(u)$  vanishes whenever the support of  $u \in C^\infty(\Omega_r)$  does not meet  $\Sigma_r$  and we conclude that the support of  $l_g$  lies in  $\Sigma_r.$

In § 4 we showed that if  $g \in H^{-1}(\Omega_r)$  and  $l_g$  has support in  $\Sigma_r, 2 \leq r \leq n,$  then  $g = 0.$  Thus

$$0 = \|g\|_{H^{-1}(\Omega_r)} = \lim_{k \rightarrow \infty} \|g^{\delta_k}\|_{H^{-1}(\Omega_r)},$$

and for every  $u \in H^1(\Omega_r),$

$$(6.5) \quad \lim_{k \rightarrow \infty} l_{g^{\delta_k}}(u) = (g^{\delta_k}, u)_{H^0(\Omega_r)} = 0.$$

Setting  $u = -v(\cdot, 1)$  in (6.5), we have

$$0 = \lim_{k \rightarrow \infty} (g^{\delta_k}, -v(\cdot, 1))_{H^0(\Omega_r)} = \lim_{k \rightarrow \infty} \left( \sum_{i=1}^n (v_i^{\delta_k}, v_i(\cdot, 1))_{H^0(\Omega_r)} \right).$$

Since  $v^{\delta_k}$  converges to  $v(\cdot, 1)$  in  $H^1(\Omega_r)$  we have

$$0 = \lim_{k \rightarrow \infty} \left( \sum_{i=1}^n (v_i^{\delta_k}, v_i(\cdot, 1))_{H^0(\Omega_r)} \right) = \sum_{i=1}^n (v_i(\cdot, 1), v_i(\cdot, 1))_{H^0(\Omega_r)}$$

and we conclude that there is a constant  $c$  such that

$$v(x, T_0) \equiv c, \quad x \in \Omega_r.$$

Since  $v_t(\cdot, t)$  is a continuous mapping from  $[0, 2]$  into  $H^0(\Omega_r)$  one can show by elementary means that

$$(6.6) \quad \lim_{\delta \rightarrow 0} \|v_t^\delta(\cdot, 1) - v_t(\cdot, 1)\|_{H^0(\Omega_r)} = 0.$$

But  $v_t^\delta(x, 1) \equiv 0$  for  $x \in \Omega_r^\delta$  as we see from (6.3) and the fact that  $v_t^\delta(x, t) \equiv 0$  in  $K(\Gamma_r, 0, 2 - m\delta)$ . Combined with (6.4) this shows that

$$\lim_{\delta \rightarrow 0} v_t^\delta(x, 1) = 0 \quad \text{a.e. in } \Omega_r,$$

and then (6.6) shows that

$$v_t(\cdot, T_0) = 0 \quad \text{in } H_0(\Omega_r).$$

Thus  $(v(\cdot, T_0), v_t(\cdot, T_0)) = (c, 0)$  is a zero energy state. The conservation of energy principle then shows that  $(v(\cdot, 2), v_t(\cdot, 2))$  is likewise a zero energy state and, reasoning as in the proof of Theorem 4a, we see that (3.1), (3.2) is approximately controllable in time  $T = 2$ . Thus Theorem 6 is proved.

*Remark.* The Holmgren–Fritz John uniqueness theorem [7] cited here and in § 2 was originally proved under the assumption that the boundary  $\Gamma$  of  $\Omega$  is analytic. The boundaries  $\Gamma_r$  of the sets  $\Omega_r$  constructed in § 3 do not have this property—they are  $C^\infty$  and piecewise analytic. However, the results of [7] can be extended to such boundaries with very little difficulty. If  $\Gamma_r = \Gamma_r^1 \cup \dots \cup \Gamma_r^s$ , where the  $\Gamma_r^k$  are relatively closed in  $\Gamma_r$  with disjoint relative interiors  $\overset{\circ}{\Gamma}_r^k$ , and if each  $\overset{\circ}{\Gamma}_r^k$  is an analytic surface, then  $(\Delta^m(D^{-m}v))_t \equiv 0$  on  $\overset{\circ}{\Gamma}_r^k$  implies, via [7], that this identity continues to hold in  $K(\overset{\circ}{\Gamma}_r^k, 0, 2 - m\delta)$ . But the interior of  $K(\overset{\circ}{\Gamma}_r^k, 0, 2 - m\delta)$  is included in the set  $\bigcup_k K(\overset{\circ}{\Gamma}_r^k, 0, 2 - m\delta)$  and thus the continuous function  $(\Delta^m(d^{-m}v))_t$  also vanishes in  $K(\Gamma_r, 0, 2 - m\delta)$ , as we need for our proof.

**7. Concluding remarks.** While Theorems 5 and 6 are stated for special domains  $\Omega$ , and a special hyperbolic partial differential equation, it is not difficult to extrapolate these results to systems of the form (1.1), (1.2) in more general domains  $\Omega$  with boundary  $\Gamma$  which includes a relatively open subset  $\tilde{\Gamma}$  whereon control is exercised.

Given the critical time  $T = 2T_0$  one forms sets  $K(\Gamma_1, 0, 2T_0 - m\delta)$  as in the proof of Theorem 6. (See [13] for complete description.) Then we form the sets

$$\Omega^\delta = \{x|(x, T_0) \in K(\Gamma_1, 0, 2T_0 - m\delta) \cap [\Omega \otimes \{T_0\}]\}.$$

As  $\delta$  tends to zero the sets  $\Omega^\delta$  increase. The complementary sets  $\Omega - \Omega^\delta$  decrease and we put

$$\Sigma = \bigcap_{\delta > 0} (\Omega - \Omega^\delta).$$

The dimension of  $\Sigma$  is what is critical. If  $\Sigma$  contains a smooth manifold of dimension  $n - 1$ , the system will not be controllable in time  $T = 2T_0$ , for one can construct a distribution  $\phi \in H^{-1}(\Omega)$  with properties (i) and (ii) of Lemma 1, solve

$$\sum_{i,j=1}^n (\alpha_{ij}(x)\tilde{v}_i)_j = \phi,$$

set

$$(7.1) \quad v(x, T_0) \equiv \tilde{v}(x), \quad v_i(x, T_0) \equiv 0$$

and then let  $v(x, t)$  be the generalized solution of (2.2), (2.3) satisfying (7.1). The state  $(v(\cdot, 2T_0), v_i(\cdot, 2T_0))$  will then lie in  $R_{2T_0}^\perp$  relative to the energy inner product in  $H_E(\Omega)$ . If  $\Sigma$  has dimension  $n - 2$  or less one can show, as in Lemma 2, that  $\Sigma$  cannot be the support of a nontrivial distribution in  $H^{-1}(\Omega)$  and prove critical time controllability as in Theorem 6.

It is clear that in the "typical" case  $\Sigma$  will have dimension less than  $n - 1$ . In fact,  $\Sigma$  will be a single point in many instances. It seems reasonable to conjecture that  $\Sigma$  cannot have dimension greater than  $n - 2$  if  $\Gamma$  is an analytic surface. Thus critical time approximate controllability is the rule, not the exception.

The results of [13] and the present paper leave the theory of *approximate* boundary value controllability of systems (1.1), (1.2) in a fairly satisfactory state. However, much remains to be done. Perhaps the most important task is that of characterizing all finite energy states which can be reached (from a zero initial state) in a time  $T \geq 2T_0$  using controls  $f \in L^2(\Gamma_1 \otimes [0, T])$ . A first step is to consider  $f \in C^\infty(\Gamma_1 \otimes [0, T])$  as in the present paper and try to bound  $\|f\|_{L^2(\Gamma_1 \otimes [0, T])}$  in terms of  $w^f(\cdot, T)$  and its derivatives. Similar work has already been done in [5], [15] for the wave equation in one space dimension. Results in this direction would enable one to undertake a systematic study of the applicability of the quadratic criterion to hyperbolic boundary value control problems, as has been done, e.g., in [11] for the case of spatially distributed controls.

#### REFERENCES

- [1] S. AGMON, *Lectures on Elliptic Boundary Value Problems*, Van Nostrand, Princeton, N.J., 1965.
- [2] R. COURANT AND D. HILBERT, *Methods of Mathematical Physics, vol. II, Partial Differential Equations*, Interscience, New York, 1962.
- [3] G. F. D. DUFF, *Mixed problems for hyperbolic equations of general order*, *Canad. J. Math.*, 11 (1959), pp. 195–221.
- [4] N. DUNFORD AND J. T. SCHWARTZ, *Linear Operators, Part II: Spectral Theory*, Interscience, New York, 1963.
- [5] J. J. GRAINGER, *Boundary-value control of distributed systems characterized by hyperbolic differential equations*, Doctoral thesis, Dept. of Electrical Engineering, University of Wisconsin, Madison, 1967.
- [6] G. HELLWIG, *Differential Operators of Mathematical Physics*, Addison-Wesley, Reading, Mass., 1964.
- [7] FRITZ JOHN, *On linear partial differential equations with analytic coefficients—Unique continuation of data*, *Comm. Pure Appl. Math.*, 2 (1949), pp. 209–253.
- [8] P. D. LAX, *On Cauchy's problem for hyperbolic equations and the differentiability of solutions of elliptic equations*, *Ibid.*, 8 (1955), pp. 615–633.
- [9] J. L. LIONS AND E. MAGENES, *Problèmes aux limites non homogènes et applications*, vol. 1, Dunod, Paris, 1968.
- [10] ———, *Problèmes aux limites non homogènes et applications*, vol. 2, Dunod, Paris, 1968.
- [11] D. L. LUKES AND D. L. RUSSELL, *The quadratic criterion for distributed systems*, *this Journal*, 7 (1969), pp. 101–121.
- [12] P. H. RABINOWITZ, *Periodic solutions of nonlinear hyperbolic partial differential equations*, *Comm. Pure Appl. Math.*, 20 (1967), pp. 145–204.
- [13] D. L. RUSSELL, *Boundary value control of the higher dimensional wave equation*, *this Journal*, 9 (1971), pp. 29–42.

- [14] ———, *On boundary- value controllability of linear symmetric hyperbolic systems*, Mathematical Theory of Control, A. V. Balakrishnan and L. W. Neustadt, eds., Academic Press, New York, 1967, pp. 312–321.
- [15] ———, *Nonharmonic Fourier series in the control theory of distributed parameter systems*, J. Math. Anal. Appl., 18 (1967), pp. 542–560.
- [16] L. SCHWARTZ, *Théorie des distributions*, Hermann, Paris, 1950.

## RECURSIVE IDENTIFICATION OF LINEAR SYSTEMS\*

J. RISSANEN†

**Abstract.** Let the three matrices  $\Sigma(N) = (G(N), F(N), H(N))$  define a linear constant system of least degree which realizes the set of numbers  $f_1, \dots, f_N$  regarded as a partial impulse response of a system. An algorithm has been developed for recursively calculating the minimal partial realizations for each  $N = 1, 2, \dots$  such that

$$\dots \Sigma(N - k) \subseteq \Sigma(N) \subseteq \dots$$

This algorithm differs from the previous ones, such as that of B. L. Ho's, in that there is a recursion on  $N$  as well. Because of this, no a priori guess of the order of the system is required. Moreover, an addition of terms to the initial sequence causes the computation of only a few new elements. When combined with another algorithm for factoring covariance matrices the given algorithm permits a recursive identification of linear random systems. No earlier recursive identification methods seem to appear in the literature. Finally, a categorical description of the abstract realizations is given.

**1. Introduction.** Let  $A_1, A_2, \dots$  be the impulse response of a linear constant system of the type

$$(1.1) \quad \begin{aligned} x(k + 1) &= Fx(k) + Gu(k), \\ y(k) &= Hx(k), \end{aligned}$$

where  $x(k) \in X = R^n$ ,  $u(k) \in U = R^p$  and  $y(k) \in Y = R^m$  for some integers  $n, p$ , and  $m$ . In other words,

$$(1.2) \quad A_i = HF^{i-1}G, \quad i = 1, 2, \dots$$

The so-called realization or identification problem is one of recovering the maps or matrices  $(G, F, H)$  from the sequence  $A_1, A_2, \dots$ . The problem is classical, but the first good solution was not found until 1965 when B. L. Ho, in his doctoral dissertation, presented a neat algorithmic or nearly algorithmic solution to the problem (see [1]). At about the same time related algorithms were described also by L. Silverman and D. Youla [2], [10].

Subsequently, the problem was extended to the interesting case where only a finite sequence  $A_1, \dots, A_N$  was given. The same algorithm was shown to give a solution [3], [4]. In that form the problem was a generalization of the Padé approximation problem, and the algorithm provided a solution which in many respects is superior to the classical Cauchy–Jacoby formulas.

In this paper we consider the following more general variant of the same problem: Given the partial sequence  $A_1, \dots, A_N$ , for each  $N = 1, 2, 3, \dots$  find a sequence of minimal partial realizations

$$(1.3) \quad \Sigma(N) = (G(N), F(N), H(N))$$

such that

$$\dots \Sigma(N') \subseteq \Sigma(N) \subseteq \dots \quad \text{if } N' < N.$$

---

\* Received by the editors November 3, 1970.

† IBM, Inc., San Jose, California 95114.

Here the inclusion symbol means that the matrices in  $\Sigma(N')$  appear as submatrices of the corresponding ones in  $\Sigma(N)$ . Thus, not only do we look for an algorithmic solution to the partial realization problem for some  $N$ , but we want a recursion on  $N$  as well. Moreover, we look for a solution where each extension of the sequence of the  $A_i$ 's can be met by calculating just a few new elements in the corresponding realization.

We point out that if Ho's algorithm is used in this case, then each partial realization must be calculated anew. Even when applied to a single partial sequence our algorithm turns out to be simpler than that of Ho's.<sup>1</sup>

In §§ 2 and 3 we give a self-contained abstract discussion of the realization problem. We give a categorical characterization of the fundamental notions of the minimal state space and the associated canonical factorization and realization.

In the final section we present an application where the problem naturally leads to finding the partial realizations for an indefinite number of values  $N$ .

The discussion throughout is restricted to the case where the input and the output spaces are one-dimensional. This is done to avoid introducing messy irrelevant notations and indices. The reader should have little trouble extending the results to the general case.

**2. Linear input-output systems.** We begin by giving a condensed but self-contained exposition of how linear systems are characterized by their input-output properties. Our approach differs from that of Kalman's [3] above all in that the important notions of canonical factorization and the associated minimal state space are defined in categorical terms. This better emphasizes the universality of these constructs; see also [6].

Let  $U$  and  $Y$  be one-dimensional vector spaces over the field  $R$  of the real numbers; i.e., both of them may be identified with  $R$ . Let  $T_-$  denote the set of nonpositive integers, and let  $\Omega = \text{hom}_c(T_-, U)$  denote the linear space of all sequences  $\omega: T_- \rightarrow U$  of finite support; i.e.,  $\omega$  has only finitely many nonzero components. Finally, let  $\Gamma = \text{hom}(T_+, Y)$  denote the linear space of all functions or sequences  $\gamma: T_+ \rightarrow Y$ , where  $T_+$  is the set of positive integers.

The linear space  $L = \Omega \oplus \Gamma$  admits a ring structure with convolution as the product:

$$(2.1) \quad (a * b)_n = \sum_{i \in T} a_{n-i} b_i, \quad n \in T = T_- \cup T_+, \quad a, b \in L.$$

This product is well-defined since the sum contains only finitely many nonzero terms for each  $n$ . Moreover, the product is commutative, and the ring has the element  $(\dots, 0, 1, 0, \dots) = e, 1$  in the zeroth position, as the identity.

Any element  $f$  of  $\Gamma$  defines a linear mapping

$$(2.2) \quad f: \Omega \rightarrow \Gamma$$

---

<sup>1</sup> When this paper was written Professor Kalman called our attention to a recent paper by Zeiger, *Some computational aspects of Ho's algorithm* (we have not found it published), in which Ho's algorithm was improved by a special factorization of the so-called Hankel matrix. That factorization is related to the one discussed in § 4. However, we exploit the factorization in a different way and the result is an altogether new algorithm with features not obtained by that of Zeiger's which still basically remains similar to Ho's algorithm.



by the ring product

$$(2.3) \quad (f(\omega))_n = (f * \omega)_n, \quad n \in T_+.$$

We take any such mapping  $f$  as an input-output description of a linear constant system. This is justified, for if  $f_i = HF^{i-1}G$ , then any system (1.1) defines an element of  $\Gamma$ , and as will be shown shortly, any  $f$  in  $\Gamma$  admits a representation (1.1).

The element  $e_1 = (\dots, 0, 1, 0)$  of  $\Omega$  defines a left shift, say  $\sigma_\Omega$ :

$$(2.4) \quad \sigma_\Omega: \Omega \rightarrow \Omega, \quad \sigma_\Omega(\omega) = e_1 * \omega.$$

Similarly, it defines a left shift  $\sigma_\Gamma$  on  $\Gamma$ :

$$(2.5) \quad \sigma_\Gamma: \Gamma \rightarrow \Gamma, \quad (\sigma_\Gamma(\gamma))_n = (e_1 * \gamma)_n, \quad n \in T_+.$$

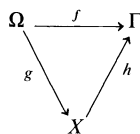
The commutativity of the ring implies that a linear mapping  $f$  of  $\Gamma$  satisfies

$$(2.6) \quad f\sigma_\Omega = \sigma_\Gamma f,$$

which in [3] was taken as the definition of linear constant systems.

We now turn to the question of how to associate a state space, above all a minimal one, to the mapping  $f$ . This question is central to the whole theory of systems described by their input-output properties, and we shall give two equivalent characterizations of the minimal state space and the associated canonical factorization.

Consider any factorization of  $f$ :

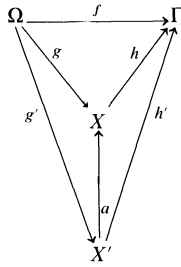


DIAG. 2.7

If  $g$  is surjective and  $h$  injective, the factorization is said to be *canonical* with a *minimal* state space  $X$ . Such factorizations obviously exist;  $X = \Omega/\ker f$ , with  $g$  the natural projection, gives one. Moreover, all canonical factorizations are equivalent in the sense that  $hg = h'g'$  implies the existence of a unique isomorphism  $q: X \rightarrow X'$  such that  $g' = qg$  and  $h = h'q$  (see [3], [5]).

These facts suggest an alternative definition of the canonical factorization and the associated state space which incorporates the universal properties of these constructs virtually without any further proving. We form the category  $\mathcal{L}_s$  of linear systems as follows. The objects are Diag. 2.7 with  $g$  surjective;  $\Omega, f$ , and  $\Gamma$  are the same in all the objects. An object may then be denoted by  $(g, X, h)$ . As the set of morphisms,  $\text{hom}[(g', X', h'), (g, X, h)]$ , take all linear maps such that Diag. 2.8 commutes. Compositions and the identities are the obvious ones, and the category axioms [9] for  $\mathcal{L}_s$  are satisfied. Define a *terminal object* in this category to be a canonical factorization through the associated minimal state space.

We could give a standard construct as a limit for the terminal objects, which would prove their existence. But since we must anyway show that this definition



DIAG. 2.8

gives the canonical factorization in the same sense as the previous one, we avoid such constructs by proving the following theorem.

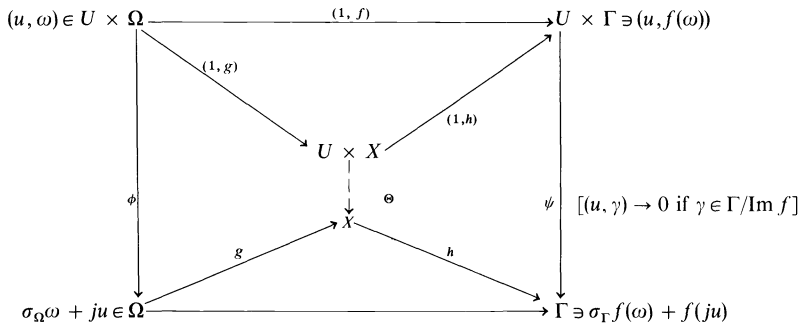
**THEOREM 2.9.** *The factorization  $f = hg$  through  $X = \Omega/\ker f$  gives a terminal object in  $\mathcal{L}_s$ .*

*Proof.* We must show that there exists a unique  $a$  in Diag. 2.8, where  $X$  now is taken as  $\Omega/\ker f$  and  $g$  the associated natural projection. Let  $z$  belong to  $X'$ . Define  $a: z \rightarrow g(\omega)$ , where  $\omega$  is any point in  $\Omega$  such that  $g'(\omega) = z$ ; one such exists since  $g'$  is surjective. This definition gives a function, since if also  $g'(\omega') = z$ , then  $f(\omega) = h'g'(\omega) = h'g'(\omega') = f(\omega')$ , and  $\omega - \omega' \in \ker f$ . Hence,  $g(\omega) = g(\omega')$ . That  $a$  is linear and satisfies  $ag' = g$  is clear.  $a$  is the only such mapping since  $g'$  is surjective, which means that  $g = ag' = a'g'$  implies  $a = a'$ .

The case where the dimension of the minimal state space is finite is of particular interest. This is the case which we shall be primarily concerned with.

**3. Minimal realizations.** By a minimal or canonical realization of the system described by an input-output mapping (2.2) is meant the intrinsic state description (1.1) of the same system characterized by the three mappings  $\Sigma = (G, F, H)$  and the associated spaces,  $U, Y, X$  with  $X$  minimal. Before describing how such a realization is constructed recursively, we quickly define them in abstract terms.

The problem is to indicate in Diag. 2.7 how input sequences are successively built into longer ones and to show how this induces the state transitions. Let  $j: U \rightarrow \Omega$  by  $u \rightarrow (\dots, 0, u)$  denote the natural injection. Consider Diag. 3.1, where  $X$  is minimal.



DIAG. 3.1

The mapping  $\phi$  induces the mappings  $\Theta$  and  $\psi$ ,  $\Theta$  unique, such that the diagram commutes. The latter is defined in the diagram, and the former results from the fact that  $h$  is injective:

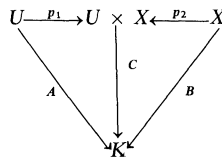
$$(3.2) \quad \Theta = k\psi(1, h),$$

where  $k:\Gamma \rightarrow X$  satisfies  $kh = 1_X$ .

$\Theta$  can be written as

$$(3.3) \quad \Theta(u, x) = Fx + Gu.$$

To see this, recall first that by  $U \times X$  is meant the direct sum  $U \xrightarrow{p_1} U \times X \xleftarrow{p_2} X$ , where  $p_1$  and  $p_2$  are the canonical injections. Hence, there is a bijection:  $\text{hom}(U \times X, K) \leftrightarrow \text{hom}(U, K) \times \text{hom}(X, K)$  in Diag. 3.4



DIAG. 3.4

for any linear space  $K$ . Taking  $K = X$  and  $C = \Theta$  gives (3.3).

More explicitly, the mapping  $G$  is given by Diag. 3.1 as

$$(3.5) \quad G(u) = \Theta(u, 0) = gi(u)$$

and  $F$  as

$$(3.6) \quad F(x) = \Theta(0, x) = k\sigma_\Gamma h(x).$$

If  $P_1:\Gamma \rightarrow Y$ ,  $P_1(y_1, y_2, \dots) = y_1$ , then the last mapping in (1.1) is given by

$$(3.7) \quad H = P_1 h: X \rightarrow Y.$$

By Theorem 2.9 any other canonical realization is isomorphic to the one given above. The preceding formulas form the basis for the algorithm to be described in the next section.

**4. Realization algorithm.** We need a few preliminary results. Let  $b_1 = f(e) = (f_1, f_2, \dots) \in \Gamma$  and  $b_i = \sigma_\Gamma^{i-1} b_1$ . We arrange the entries in the so-called Hankel matrix:

$$(4.1) \quad A = \begin{bmatrix} f_1 & f_2 & \cdots & f_n & | & \cdots \\ f_2 & f_3 & \cdots & f_{n+1} & | & \cdots \\ \text{---} & \text{---} & \text{---} & \text{---} & | & \cdots \\ f_m & f_{m+1} & \cdots & f_{m+n-1} & | & \cdots \\ \vdots & \vdots & \vdots & \vdots & | & \vdots \end{bmatrix}.$$

Let  $A(m, n)$  denote the dashed submatrix of  $A$ .

Observe first from Diag. 2.7 that since  $h$  is injective in a canonical factorization, the dimension of  $X$ , or the order of the system, is the same as the dimensionality of the image of  $f$ . Since the set  $\{\sigma_{\Omega}^i e\}$  spans  $\Omega$ , the image of  $f$  is spanned by the set  $\{b_i\}$ . The order of the system is then the cardinality of the linearly independent vectors in the set  $\{b_i\}$ . We need the following lemma.

LEMMA 4.2. *If  $\dim \{b_1, b_2, \dots\} = n$ , then  $b_1, b_2, \dots, b_n$  are the linearly independent vectors in  $\{b_i\}$ . In addition,  $A(n, n)$  is nonsingular and  $A(n, m)$ ,  $m \geq n$ , has rank  $n$ .*

*Proof.* Suppose that  $b_{k+1}$  is linearly dependent on  $b_1, \dots, b_k$ ; i.e.,  $(\sigma_{\Gamma}^k + \sum_1^k c_i \sigma_{\Gamma}^{i-1})b_1 = 0$ . Then,  $\sigma_{\Gamma}$  of that same expression also vanishes, or  $b_{k+2} + c_k b_{k+1} + \dots + c_1 b_2 = 0$ , and also  $b_{k+2}$  is linearly dependent on  $b_1, \dots, b_k$ . Hence, the first  $n$  must be linearly independent.

Consider the last row in  $A(n+1, n)$ . Since  $b_{n+1}$  is linearly dependent on  $b_1, \dots, b_n$ , the rank of  $A(n+1, n)$  is still  $n$ . By symmetry the same is true about  $A(n, n+1)$ . Hence the ranks of  $A(n, m)$  for  $m \geq n$  are no greater than  $n$ . But since they clearly can neither be less than  $n$  they all must be equal to  $n$ .

The algorithm is based on a factorization of  $A(n, m)$  of the following type:

$$(4.3) \quad A(n, m) = P(n, n)Q(n, m), \quad m \geq n, \quad \text{rank } A(n, m) \geq n - 1,$$

where  $P(n, n)$  is lower triangular with 1's on the diagonal; i.e.,

$$(4.4) \quad \begin{bmatrix} f_1 & \cdots & f_i & \cdots & f_{i+1} & \cdots & f_m \\ f_2 & \cdots & f_{i+1} & \cdots & f_{i+2} & \cdots & \cdots \\ \vdots & & \vdots & & \vdots & & \vdots \\ f_n & \cdots & f_{i+n-1} & \cdots & f_{i+n} & \cdots & f_{n+m-1} \end{bmatrix} = \begin{bmatrix} 1 & & & & & & \\ p_{21} & 1 & & & & & \\ \vdots & & \vdots & & \vdots & & \vdots \\ p_{n1} & \cdots & p_{n,n-1} & & 1 & & \end{bmatrix} \begin{bmatrix} q_{11} & \cdots & q_{1i} & q_{1,i+1} & \cdots & q_{1m} \\ q_{21} & \cdots & q_{2i} & & & \cdots \\ \vdots & & \vdots & & & \vdots \\ q_{n1} & \cdots & q_{ni} & q_{n,i+1} & \cdots & q_{n,m} \end{bmatrix}.$$

The factors are not unique, a fact which we shall take advantage of to obtain factors with further desired features. By setting certain elements  $q_{ij} = 0$  we shall be able to calculate the  $p_{ij}$ 's recursively one by one. Moreover, an addition of rows and columns to  $A(n, m)$  will not change the numbers already calculated. A further property of  $Q(n, m)$  is that the first  $n - 1$  rows are linearly independent and the last row is zero if the rank of  $A(n, m) = n - 1$ . With a minor modification the algorithm would work regardless of what the rank of  $A(n, m)$  is. Since we took  $U$  and  $Y$  to be one-dimensional we shall have no need for that case, however.

The factorization algorithm runs as follows.

Step 1. Set  $q_{1i} = f_i$  for all  $i$ . If  $n = 1$  we are done:  $P(1, 1) = (1)$ .

Step 2. In the other event, proceeding recursively, we have at the  $i$ th step or row all the  $p_{jk}$ 's and  $q_{jk}$ 's,  $j = 0, 1, \dots, i - 1$ , determined. Let  $s(j)$  be the least

integer such that  $q_{j,s(j)} \neq 0, j < n. s(j)$  exists because of the rank condition on  $A(n, m)$ . Set  $q_{k,s(j)} = 0$  for  $k > j$ . Equation (4.4), then, leads to a set of  $i - 1$  equations, one for each column  $s(j), j = 1, \dots, i - 1$ . Because of the previous conditions, the unknowns  $p_{i1}, \dots, p_{i,i-1}$  can be solved recursively one by one from these equations. The submatrix  $P(i, i)$  with (4.4) determine the remaining elements of the  $i$ th row of  $Q(n, m)$ , which completes the cycle.

As an example, consider

$$(4.5) \quad A(4, 5) = \begin{bmatrix} 1 & 1 & 1 & 2 & 1 \\ 1 & 1 & 2 & 1 & 3 \\ 1 & 2 & 1 & 3 & 2 \\ 2 & 1 & 3 & 2 & 3 \end{bmatrix}.$$

We have  $s(1) = 1$ . Set  $q_{i1} = 0, i > 1$ . Then

$$p_{21} \cdot 1 + 1 \cdot q_{21} = 1, \quad p_{21} = 1.$$

Further,  $q_{22} = 0, q_{23} = 1, q_{24} = -1, q_{25} = 2$ . Then, since  $s(2) = 3, q_{i3} = 0, i > 2$ , and the first and the third column give

$$p_{31}q_{11} = 1 \quad \text{or} \quad p_{31} = 1, \\ 1 \cdot q_{13} + p_{32}q_{23} = f_5 = 1 \quad \text{or} \quad p_{32} = 0.$$

Continuing we get the result

$$(4.6) \quad A(4, 5) = \begin{bmatrix} 1 & & & & \\ 1 & 1 & & & \\ 1 & 0 & 1 & & \\ 2 & 1 & -1 & 1 & \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 & 2 & 1 \\ 0 & 0 & 1 & -1 & 2 \\ 0 & 1 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}.$$

Define

$$(4.7) \quad P_*(n - 1, n - 1) = \begin{bmatrix} p_{21} & 1 & & & \\ p_{31} & p_{32} & 1 & & \\ \vdots & & & \ddots & \\ p_{n1} & \dots & & & p_{n,n-1} \end{bmatrix}, \quad G(n - 1) = \begin{bmatrix} q_{11} \\ q_{21} \\ \vdots \\ q_{n-1,1} \end{bmatrix},$$

$H(n - 1) = (1, 0, \dots, 0), n - 1$  elements.

**THEOREM 4.8.** *Given  $f_1, f_2, \dots$ , let  $A(n - 1, n - 1)$  be nonsingular and  $m$  any integer such that  $A(n, m)$  has rank  $n - 1$ . Then  $\Sigma(n - 1) = (G(n - 1), F(n - 1), H(n - 1)), X(n - 1) = R^{n-1}$ , where*

$$(4.9) \quad F(n - 1) = P^{-1}(n - 1, n - 1)P_*(n - 1, n - 1)$$

is a minimal realization of the sequence  $f_1, \dots, f_{n+m-1}$ .

*Proof.* By the factoring algorithm the last row of  $Q(n, m)$  is zero. Hence by writing the equality between the dashed columns of  $A(n, m)$  in (4.4) we get with (4.9)

$$(4.10) \quad F(n - 1) \begin{bmatrix} q_{1i} \\ \vdots \\ q_{n-1,i} \end{bmatrix} = \begin{bmatrix} q_{1,i+1} \\ \vdots \\ q_{n-1,i+1} \end{bmatrix} \quad \text{for all } i.$$

Denote the elements of  $R^{n-1}$  by  $x = \text{col}(x_1, \dots, x_{n-1})$  and consider the equations

$$(4.11) \quad \begin{aligned} x(i+1) &= F(n-1)x(i) + G(n-1)u(i), \\ y(i) &= H(n-1)x(i). \end{aligned}$$

With  $x(0) = 0$  and  $u(0) = 1, u(i) = 0, i > 0$ , the first equation in (4.11) describes the consecutive states which are just the columns of  $Q(n-1, m)$  in (4.4), or (4.10), and the second gives  $y(i) = q_{1i} = f_i$  for  $i = 1, \dots, m$ . We still have to show that (4.11) realizes the rest of the  $f_i$ 's. This will follow from the special form of  $F(n-1)$  which it inherits from  $P_*(n-1)$ , namely,

$$(4.12) \quad F(n-1) = \begin{bmatrix} f_{11} & 1 & 0 & \cdots & 0 \\ f_{21} & f_{22} & 1 & 0 & \cdots & 0 \\ \vdots & & & & \vdots & \\ & & & & & 1 \\ f_{n-1,1} & \cdots & & & & f_{n-1,n-1} \end{bmatrix}.$$

Indeed, applying formula (4.11) for  $i = m+1, \dots, m+n-1$  extends  $Q(n, m)$  to  $\bar{Q}(n, m+n-1)$  (the last row is extended as a zero row). Multiplying the result by  $P(n, n)$  we get the corresponding extension of  $A(n, m)$  to  $\bar{A}(n, m+n-1)$ , say. But due to the special shape of  $F(n-1)$ ,  $\bar{A}(n, m+n-1)$  will have the following elements:

$$\begin{bmatrix} f_1 & \cdots & f_m & \cdots & f_{m+n-1} \\ \vdots & & & & \cdot \\ f_n & \cdots & f_{m+n-1} & & \end{bmatrix}.$$

This means that  $y(i) = q_{1i} = f_i$  for all  $i \leq m+n-1$ .

Any realization of the sequence  $f_1, \dots, f_{m+n-1}$  extends it indefinitely. If such a realization has order  $k < n-1$ , then  $A(n-1, n-1)$  has rank  $k$  by Lemma 4.2 which contradicts the assumptions. Hence, (4.11) is minimal.

We are now ready to describe the realization algorithm. It works under the hypothesis that from some source we can pick the numbers  $f_1, \dots, f_N$  for any  $N = 1, 2, \dots$ .

*Step 1.* Let  $k$  be the least integer for which  $f_k \neq 0$ . Take  $N = 2k+1$  and form  $A(k+1, k+1)$ . It has rank  $\geq k$ .

*Step 2.* Apply the factoring algorithm and find  $P(k+1, k+1)$  and  $Q(k+1, k+1)$  (equation (4.4)). If the last row of  $Q$  is nonzero, the rank of  $A(k+1, k+1)$  is  $k+1$ . Increase  $N$  by 2, form  $A(k+2, k+2)$  and continue the factorization. If the last row of  $Q(k+2, k+2)$  is still nonzero, increase  $N$  by 2 and repeat until, say, for  $N = 2n-1$ , the last row of  $Q(n, n)$  is zero. Such an  $n$  exists by Lemma 4.2 if the numbers  $f_i$  admit a finite order realization.

*Step 3.* From formulas (4.7) and (4.9) calculate the partial realization  $\Sigma(n-1)$ . Observe that the inverse of  $P(n-1, n-1)$  can also be calculated recursively since this matrix together with its inverse is lower triangular.

*Step 4.* Increase  $N$  by 1. Continue the factorization for  $A(n, n + 1)$ . If the last row of  $Q(n, n + 1)$  remains zero, increase  $N$  by 1, and repeat. If the last row remains zero for all  $m$ , we have found the realization. (This cannot, of course, be decided, so that the algorithm would never stop. A stopping rule is introduced by setting an upper limit for  $m$ .)

*Step 5.* If for some  $m(>n)$  the last element in the last row of  $Q(n, m)$  becomes nonzero, the last picked element  $f_{m+n-1}$  is not realized by the partial realization  $\Sigma(n - 1)$ . In this case, pick one new point,  $f_{m+n}$ , and form  $A(n + 1, m)$ . Continue the factorization, pick one new point and repeat until either  $Q(n', m)$  for some least  $n' \leq m$  has last row zero or  $n' = m$  and the last row is nonzero. In the previous case go to Step 3. In the latter case go to Step 2.

*Comments.* Due to the special form of the matrices (4.7) and (4.12) this algorithm has the crucial “nesting” property (for  $k$ , see (4.13)):

$$\dots \Sigma(n - k) \subseteq \Sigma(n) \subseteq \dots, \quad \text{if } k < 1;$$

where the inclusion sign means that the matrices in  $\Sigma(n - k)$  are submatrices of the corresponding ones in  $\Sigma(n)$ . Hence, at each step only a few new elements as a function of the old ones and the new picked elements  $f_i$  need be calculated. In fact, we have established a recursion of the type,

$$\begin{aligned} \Sigma(n + k) &= \Pi_1(n, \Sigma(n), f_{2n+1}, \dots, f_m), \\ (4.13) \quad m &= \Pi_2(n, \Sigma(n), f_{2n+1}, \dots, f_{m-1}), \\ k &= \Pi_3(n, \Sigma(n), f_{2n+1}, \dots, f_m). \end{aligned}$$

Another advantage of this algorithm is that the formulas (4.7) and (4.9) are actually simpler than those in the Ho algorithm [3].

To illustrate the algorithm we give the sequence of partial realizations corresponding to the example in (4.5) and (4.6). In Step 2,  $n = 2$ , and in Step 3 we get  $\Sigma(1) = ((1), (1), (1))$ . In Step 4 an addition of  $f_4 = 2$  makes  $q_{23} = 1 \neq 0$ . In Step 5 pick  $f_5 = 1$ , and form  $A(3, 3)$ . Returning to Step 2, we pick two new points,  $f_6 = 3$  and  $f_7 = 2$ . This time the last row of  $Q(4, 4)$  is zero, and we calculate  $\Sigma(3)$  following Step 3:

$$G(3) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \quad F(3) = \begin{bmatrix} 1 & 1 & 0 \\ 0 & -1 & 1 \\ 1 & 0 & -1 \end{bmatrix}, \quad H(3) = (1, 0, 0),$$

which realizes all the numbers given in  $A(4, 5)$ .

**5. Identification of random systems.** In this section we briefly describe an application of the realization algorithm where its recursive nature proves to be particularly useful.

Consider a system whose output forms a stationary random process  $\{y_t\}$ ,  $t = 1, 2, \dots$ . Take  $E(y_t) = 0$  for all  $t$ , and denote the covariances by  $r_t = E(y_t y_{t-i})$ . Form the infinite symmetric covariance matrix,  $R = (r_{ij})$ ,  $r_{ij} = r_{i-j}$ .

Assuming the process to be of full rank, i.e.,  $R > 0$ , we can factor  $R$  recursively [7], [8] into the product of a lower triangular matrix  $B$  and its transpose:

$$(5.1) \quad R = BB',$$

where

$$(5.2) \quad B = \begin{bmatrix} b_{00} & & & & \\ b_{10} & b_{11} & & & \\ & & & & \\ b_{n0} & b_{n1} & \cdots & b_{nn} & \\ \dots & \dots & \dots & \dots & \dots \end{bmatrix}.$$

The recursive equations for the  $b_{ij}$ 's are obtained by comparing elements in (5.1) starting on the top row and progressing down row by row. Observe that the result is just a special case of the factoring algorithm in § 4.

**THEOREM 5.3.** *If  $R$  satisfies Szego's criterion,*

$$\int_{|z|=1} \log \left( \sum_{-\infty}^{\infty} r_i z^i \right) dz > -\infty,$$

then  $b_{ij} \rightarrow f_{i-j}$  as  $i \rightarrow \infty$ .

We omit the proof, which is a simple modification of the proof of the related theorem, Theorem 6 in [8].

Applying the realization algorithm to the sets of numbers  $f_1, \dots, f_N$ ,  $N = 1, 2, \dots$ , we get a series of realizations. If the process is generated by a finite order system, this series leads to the minimal realization.

Since we actually must use the approximations  $b_{ij}$  for  $f_{i-j}$ , the realizations will be approximations, too. Because of this, the precise rank condition in the algorithm should be replaced by some rule such as: the last row of  $Q(n, m)$  is considered as zero if  $|q_{ni}|/\|q_i\| \leq \epsilon$ , where  $q_i$  is the  $i$ th column of  $Q(n-1, m)$ . Since the last row represents the  $n$ th components of the state vectors, the error so introduced admits a clear interpretation.

The above described scheme seems to compare favorably with the commonly used procedures in that all of the computations progress in an orderly, recursive fashion without any trial and error involved.

The outlined procedure has the defect that the computation of (5.1) involves a growing amount of memory even when  $R(z) = \sum_{-\infty}^{\infty} r_i z^i$  is a rational function. We sketch an approach which in some respect is an improved version of a related one by P. Faure (see, e.g., [11]); see also [13].

One first applies the realization algorithm to the sequences of numbers  $r_0, r_1, \dots, r_N$ . This gives the realization  $(K, F, H)$  of order  $n$ , say, and the rational function

$$R(z) = H(zI - F)^{-1}FK + K'F'(z^{-1}I - F')^{-1}H' + HK = \frac{P(z)}{Q(z)Q(z^{-1})}.$$

The polynomial  $P(z) = \sum_{-m}^m p_i z^i$ ,  $m \leq n$ , satisfies (5.3). Hence, the factorization obtained by replacing  $R$  in (5.1) by the Toeplitz matrix  $P$  associated with  $P(z)$  which has only  $2m + 1$  nonzero diagonals converges and gives in the limit

$$P(z) = P_1(z)P_1(z^{-1}).$$

The transfer function of the system is  $P_1(z)/Q(z)$ . We leave the details of this



procedure, above all the proof that the result is a “minimum phase” system, to another context. Such a result, in effect, was announced in [12], which reference was communicated to us by H. Aasnaes.

**Acknowledgment.** We are indebted to Professor R. Kalman, Dr. F. Palermo and Professor T. Kailath for stimulating discussions during the preparation of this paper.

#### REFERENCES

- [1] B. L. HO AND R. E. KALMAN, *Effective construction of linear state-variable models from input/output functions*, Proc. Third Allerton Conference, Urbana, Illinois, 1966, pp. 449–459.
- [2] D. C. YOULA, *The synthesis of linear dynamic systems from prescribed weighting patterns*, SIAM J. Appl. Math., 14 (1966), pp. 527–549.
- [3] R. E. KALMAN, P. L. FALB AND M. A. ARBIB, *Topics in Mathematical System Theory*, McGraw-Hill, New York, 1968.
- [4] ———, *Lectures on Controllability and Observability*, Centro Internazionale Matematico Estivo, Bologna, Italy, 1968.
- [5] H. P. ZEIGER, *Ho's algorithm, commutative diagrams, and the uniqueness of minimal linear systems*, Information and Control, 11 (1967), pp. 71–79.
- [6] M. A. ARBIB AND H. P. ZEIGER, *On the relevance of abstract algebra to control theory*, Automatika, 5 (1969), pp. 589–606.
- [7] J. RISSANEN, *An algebraic approach to the problems of linear prediction and identification*, Res. Rep. RJ468, IBM, Yorktown Heights, N.Y., 1967.
- [8] J. RISSANEN AND L. BARBOSA, *Properties of infinite covariance matrices and stability of optimum predictors*, Information Sci., 1 (1969), pp. 221–236.
- [9] S. MACLANE AND G. BIRKHOFF, *Algebra*, Macmillan, New York, 1967.
- [10] L. SILVERMAN, *Representation and realization of time-variable linear systems*, Tech. Rep. 94, Dept. of Electrical Engineering, Columbia Univ., New York, 1966.
- [11] P. FAURRE AND J. P. MARMORAT, *Un algorithme de realization stochastique*, C. R. Acad. Sci. Paris, Sér. A, 268 (1969), pp. 978–981.
- [12] F. BAUER, *Ein directes Iterationsverfahren zur Hurwitz-Zerlegung eines Polynoms*, Mitt. Math. Inst. Tech. Hochschule München, A.E.Ü., 9 (1955), pp. 285–290.
- [13] T. KAILATH AND R. GEESEY, *An innovations approach to least squares estimation, Part IV: Recursive estimation given the covariance functions*, IEEE Trans. Automatic Control, to appear.

## GLOBAL VARIATION CRITERIA FOR STABILITY OF LINEAR TIME-VARYING SYSTEMS\*

Y. V. VENKATESH†

**Abstract.** The system under study is a linear feedback system consisting of a time-invariant block  $G$  in the forward path and a time-varying gain  $k(t)$  in the feedback path. Improved sufficient conditions for its asymptotic stability are derived by a combination of Brockett's factorization technique of generating Lyapunov functions and the Krasovskii-Corduneanu theorem. The resulting bound on  $((dk/dt)/k)$  resembles but is distinct from the average criterion of Freedman and Zames [1].

**1. Introduction.** Consider the feedback system of Fig. 1 which is governed by the linear differential equation

$$(1) \quad p(D)y + k(t)q(D)y = 0 \quad \text{on the interval } [t_0, \infty),$$

where

$$p(D) = D^n + p_{n-1}D^{n-1} + \dots + p_0,$$

$$q(D) = q_m D^m + q_{m-1}D^{m-1} + \dots + q_0$$

are constant coefficient differential operators with the order  $n$  of  $p(D)$  at least one higher than the order  $m$  of  $q(D)$ .

Let  $y = x_1$ ,  $x_2 = dx_1/dt$ ,  $\dots$ ,  $x_n = dx_{n-1}/dt$ ; and  $\mathbf{x} = \text{col } [x_1, x_2, \dots, x_n]$ . Then (1) can be written as the vector differential equation

$$(2) \quad \frac{d\mathbf{x}}{dt} = A_0\mathbf{x} - k(t)\mathbf{b}\mathbf{c}'\mathbf{x} \triangleq A(t)\mathbf{x},$$

where  $A_0$  is a stable matrix having the form:

$$A_0 = \begin{bmatrix} 0 & 1 & 0 & 0 & \dots & 0 \\ 0 & 0 & 1 & 0 & \dots & 0 \\ \cdot & \cdot & \cdot & \cdot & \dots & \cdot \\ 0 & 0 & 0 & 0 & \dots & 1 \\ -p_0 & -p_1 & -p_2 & -p_3 & \dots & -p_{n-1} \end{bmatrix},$$

and  $\mathbf{b}$ ,  $\mathbf{c}$  are  $n$ -vectors given by

$$\mathbf{b} = \text{col } [0, 0, \dots, 1],$$

$$\mathbf{c} = \text{col } [-q_0, -q_1, \dots, -q_m, \dots, 0].$$

The gain  $k(t)$  is assumed to be absolutely continuous on the interval  $[t_0, \infty)$ . Let  $G(s)$  be the transfer function of the forward block, i.e.,  $G(s) = q(s)/p(s)$ .

\* Received by the editors January 2, 1969, and in final revised form December 7, 1970.

† Department of Electrical Engineering, Indian Institute of Science, Bangalore 12, India.

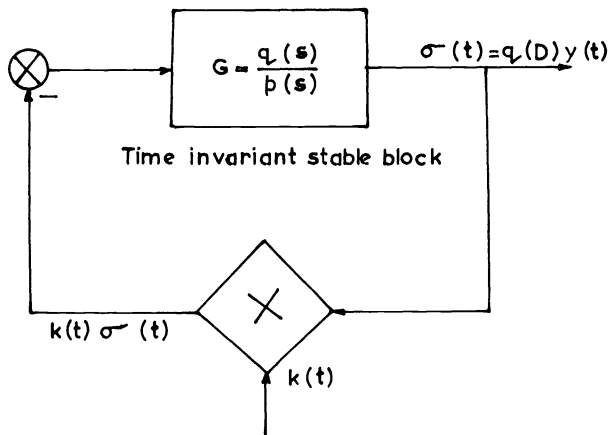


FIG. 1. A linear time-varying feedback system

*Remark.* The stability properties of the null solution (n.s.) of (1) are identical with the stability properties of the n.s. of (2) and vice versa.

**ASSUMPTION.** The n.s. of (1) is asymptotically<sup>1</sup> stable (a.s.) for every constant function  $k(t) = K$  in  $[0, \infty)$ .

**PROBLEM.** Find conditions for the n.s. of (1) to be a.s. for every (absolutely) continuous function  $k(t)$  with values in  $[0, \infty)$ .

Literature on the preceding problem is too vast to receive due acknowledgment here. For all unbounded  $dk/dt$ , the circle criteria of Zames [2], Narendra and Goldwyn [3] and Sandberg [4] are well known. The circle criterion is applicable even when the feedback is nonlinear and time-varying. Using a positive multiplier approach similar to that of Zames [2], Brockett and Forys [5] and others introduced a bound on the rate of variation  $((dk/dt)/k)$  permitting the hypothesis on  $G(s)$  to be weakened. Freedman and Zames [1] obtained an explicit stability condition free of multipliers in terms of a bound on an average of  $((dk/dt)/k)$ , namely,

$$(3) \quad \sup_{t \geq 0} \frac{1}{T} \int_t^{t+T} \left| \frac{dk(\tau)}{d\tau} / k(\tau) \right| d\tau < 4\beta_{sh}$$

for some  $T > 0$ , where  $\beta_{sh}$  is a constant determined from the Nyquist diagram of  $G(s)$ . What is noteworthy here is that the use of an implicit multiplier has been avoided. However, condition (3) is not necessary and weaker conditions may be possible.

**2. Solution of the main problem.** The present objective is to explore a route towards obtaining sufficient conditions for the a.s. of (1) based on Brockett's technique [9] of generating a Lyapunov function and the Krasovskii–Corduneanu theorem [6, pp. 56–57], [7]. In fact, a new condition on  $((dk/dt)/k)$  is derived below which is distinct from the modulus condition (3) of Freedman and Zames [1]. In order to state the main results of the paper, the following definitions and lemmas will be needed.

<sup>1</sup> See Definitions 3 and 4 below.

DEFINITION 1. A complex-valued function  $Z(s)$  of a complex variable  $s$  is called a *positive real (strictly positive real) function of the argument  $s$* , if  $Z(s)$  is real for real values of  $s$ , and for  $\text{Re } s > 0$  (where  $\text{Re}$  denotes the real part) is analytic and satisfies the inequality

$$\text{Re } Z(s) \geq 0 \quad (>0).$$

The following notation will be used:  $\|\mathbf{x}\|$  denotes the norm of  $\mathbf{x}$ , where  $\|\mathbf{x}\|^2 = \mathbf{x}'\mathbf{x}$ ;  $\mathbf{x}_0$  denotes  $\mathbf{x}(t_0)$ ;  $\mathbf{x}(t; t_0, \mathbf{x}_0)$  denotes the solution of (2) which takes the value  $\mathbf{x}_0$  for  $t = t_0$ .

DEFINITION 2. The null solution of (2) is said to be *stable* if for every  $\varepsilon > 0$  there exists  $\delta(\varepsilon; t_0)$  such that if  $\|\mathbf{x}_0\| < \delta(\varepsilon; t_0)$ , then  $\|\mathbf{x}(t; t_0, \mathbf{x}_0)\| < \varepsilon$  for  $t \geq t_0$ .

DEFINITION 3. The null solution of (2) is said to be *asymptotically stable* if it is stable and, in addition, there exists a  $\delta_0(t_0) > 0$  with the property that if  $\|\mathbf{x}_0\| < \delta_0$ , then  $\lim_{t \rightarrow \infty} \mathbf{x}(t; t_0, \mathbf{x}_0) = 0$ .

DEFINITION 4. The null solution of (2) is said to be *exponentially stable* if there exist positive constants  $\varepsilon_1, \varepsilon_2$ , such that, for  $t \geq t_0$ ,

$$\|\mathbf{x}(t; t_0, \mathbf{x}_0)\| < \varepsilon_2 \|\mathbf{x}_0\| \exp[-\varepsilon_1(t - t_0)].$$

LEMMA 1. *The asymptotic stability of (1) is equivalent to the asymptotic stability of*

$$(4) \quad p(D)n(D)y + k(t)q(D)n(D)y = 0,$$

*provided  $n(D)z = 0$  represents an asymptotically stable system.*

*Proof.* Since the system described by  $n(D)z = 0$  is asymptotically stable,  $n(s) = 0$  has all its zeros in the half-plane  $\text{Re } s < 0$ .

If  $\phi(t)$  is a solution of (4), then  $n(D)\phi(t)$  is a solution of (1). Further, if  $n(D)\phi(t)$  tends to zero as  $t \rightarrow \infty$ , so does  $\phi(t)$  showing that the stability of (1) implies the stability of (4).

To prove the converse, note that if  $\psi(t)$  is a solution of (1), then there is a function  $y(t)$ , satisfying  $n(D)y = \psi(t)$  and such that  $y(t)$  is a solution of (4). Now, because of the a.s. of the system described by (4),  $\psi(t)$  tends to zero as  $t \rightarrow \infty$ . But  $n(s) = 0$  has all its zeros in the half-plane  $\text{Re } s < 0$ . Consequently,  $y(t)$ , treated as the output of a system with the transfer function  $1/n(s)$ , also tends to zero as  $t \rightarrow \infty$ .

LEMMA 2. *A real function of a complex variable  $Z(s) = m(s)/n(s)$ , where  $m(s)$  and  $n(s)$  are finite polynomials in  $s$ , is positive real if and only if:*

- (i)  $n(s) + m(s)$  has no zeros in the closed right half-plane ( $\text{Re } s \geq 0$ ), and
- (ii)  $\text{Re } Z(j\omega) \geq 0$  for all real  $\omega$ .

*Proof.* See Weinberg and Slepian [8].

A. If  $u(s)$  is a polynomial with real coefficients, let  $\text{Ev } u(s) = \{u(s) + u(-s)\}/2$  denote its even part. Note that for real  $\omega$ ,  $\text{Ev } u(j\omega) = \text{Re } u(j\omega)$ . If  $u(s)$  is even and  $\text{Re } u(j\omega) \geq 0$  for all real  $\omega$ , then according to a theorem of Wiener, there exists a unique polynomial  $v(s)$  with real positive coefficients such that  $v(s)v(-s) = u(s)$ , and  $v(s)$  has zeros only in the closed left half-plane  $\text{Re } s \leq 0$ . Let  $[u(s)]^{(+)}$  stand for  $v(s)$  and  $[u(s)]^{(-)}$  for  $v(-s)$ . The latter is also known as the right half-plane spectral factor, or, at times, as the negative spectral factor of the even polynomial  $u(s)$ . Observe that  $[u(s)]^{(-)}$  has zeros only in the closed right half-plane  $\text{Re } s \geq 0$ .

B. If  $u_1(s)$  and  $u_2(s)$  are two polynomials with real coefficients and  $\text{Re} [u_1(j\omega)/u_2(j\omega)] \geq 0$  for all real  $\omega$ , then clearly  $\text{Re} [u_1(j\omega)u_2(-j\omega)] \geq 0$  for all real  $\omega$ , and there is a unique factor  $[\text{Ev } u_1(s)u_2(-s)]^{(-)}$ .

Similarly, if  $\beta$  is a real constant and  $\text{Re } Z(j\omega - \beta)G(j\omega - \beta) \geq 0$  for all real  $\omega$ , then

$$(5) \quad r_1(s) = \{\text{Ev } m(s - \beta)q(s - \beta)n(-s - \beta)p(-s - \beta)\}^{(-)}$$

is uniquely defined. Moreover, if  $\alpha$  is a real constant and  $\text{Re } Z(j\omega - \alpha) \geq 0$  for all real  $\omega$ , then

$$(6) \quad r_2(s) = \{\text{Ev } m(s - \alpha)n(-s - \alpha)\}^{(-)}$$

is also uniquely defined. The polynomials  $r_1(s)$  and  $r_2(s)$  have their zeros in the closed right half-plane only.

C. Let  $f(\mathbf{x}(t))$  be a given (scalar) function of trajectories  $\mathbf{x}(t)$  on  $[0, \infty)$ , for example, of solutions of (2). Suppose the integral

$$(7) \quad I = \int_{t_1}^{t_2} f(\mathbf{x}(\tau)) d\tau$$

exists and is identical for all trajectories  $\mathbf{x}(t)$  satisfying  $\mathbf{x}(t_1) = \mathbf{x}_1$  and  $\mathbf{x}(t_2) = \mathbf{x}_2$ , where  $\mathbf{x}_1$  and  $\mathbf{x}_2$  are constant vectors. Then the integral (7) will be represented by the notation

$$I = \int_{t(\mathbf{x}_1)}^{t(\mathbf{x}_2)} f(\mathbf{x}(\tau)) d\tau.$$

Now suppose that at time  $t$ , the state is  $\mathbf{x}$  and at  $(t + \Delta t)$ , the state is  $\mathbf{x}_0$ . Then

$$I(\mathbf{x}) - I(\mathbf{x}_0) = \int_t^{t+\Delta t} f(\mathbf{x}(\tau)) d\tau.$$

The time derivative of  $I$  along the solutions of (2) is obtained by allowing  $\Delta t \rightarrow 0$  in  $\{I(\mathbf{x}) - I(\mathbf{x}_0)\}/\Delta t$  and substituting for  $\mathbf{x}(t)$  the solution of (2):

$$\left. \frac{dV}{dt} \right|_{(2)} = f(\mathbf{x}(t)),$$

where now  $\mathbf{x}(t)$  is the solution of (2).

D. Let  $Z(s - \beta)G(s - \beta)$  be assumed to be positive real for some  $\beta \geq 0$ . Then by Lemma 2, the solution  $z'(t)$  of the equation

$$(8) \quad m(D - \beta)q(D - \beta)z' + n(D - \beta)p(D - \beta)z' = 0$$

is asymptotically stable. The order  $\eta$  of this equation is greater than or equal to  $n$ . Let  $z'(t) = z(t) \exp(\beta t)$ . Then  $z(t)$  satisfies the  $\eta$ th order equation

$$(9) \quad m(D)q(D)z + n(D)p(D)z = 0,$$

and is asymptotically stable. Further, let

$$z_1(t) = z(t), \quad z_2(t) = \frac{dz_1}{dt}, \quad z_3(t) = \frac{dz_2}{dt}, \dots, z_\eta(t) = \frac{dz_{\eta-1}}{dt},$$

and

$$\mathbf{z}(t) = \text{col} [z_1, \dots, z_\eta].$$

Note that the order of (4) is also  $\eta$ . The state vector  $\mathbf{x}(t)$  of (1) is a subspace of  $\mathbf{x}^*(t)$ , the state vector of (4).

LEMMA 3. If  $Z(s - \beta)G(s - \beta)$  is positive real for some  $\beta \geq 0$ , and  $m(s - \beta)q(s - \beta)$  and  $n(s - \beta)p(s - \beta)$  have no zeros on the imaginary axis, then there exists a positive definite function  $V_1(\mathbf{z}, t)$  quadratic in  $\mathbf{z}$  defined by

$$(10) \quad V_1(\mathbf{z}, t) = e^{-2\beta t} \int_{t(\mathbf{0})}^{t(\mathbf{z})} \{ [m(D - \beta)q(D - \beta)(ze^{\beta t})][n(D - \beta)p(D - \beta)(ze^{\beta t})] - [r_1(D)ze^{\beta t}]^2 \} d\tau,$$

where  $z(t)$  is any solution of (9) on  $[0, \infty)$ .

Proof. As shown in [9], the integral in (10) is path independent. It remains to show that  $V_1(\mathbf{z}, t)$  is positive definite. To this end, rewrite (10) letting  $z' \triangleq ze^{\beta t}$  to obtain

$$(11) \quad V_1(\mathbf{z}, t) = e^{-2\beta t} \int_{t(\mathbf{0})}^{t(\mathbf{z})} \{ [m(D - \beta)q(D - \beta)z'] [n(D - \beta)p(D - \beta)z'] - [r_1(D)z']^2 \} d\tau.$$

But  $z'$  satisfies (8), from which (11) takes the form

$$(12) \quad V_1(\mathbf{z}, t) = -e^{-2\beta t} \int_{t(\mathbf{0})}^{t(\mathbf{z})} \{ [m(D - \beta)q(D - \beta)z']^2 + [r_1(D)z']^2 \} d\tau$$

or

$$(13) \quad V_1(\mathbf{z}, t) = -e^{-2\beta t} \int_{t(\mathbf{0})}^{t(\mathbf{z})} \{ [n(D - \beta)p(D - \beta)z']^2 + [r_1(D)z']^2 \} d\tau.$$

Further, by the asymptotic stability of (8), (12) and (13) are valid for  $t(\mathbf{0}) = \infty$ ,  $t(\mathbf{z}) = 0$ . Hence,

$$V_1(\mathbf{z}, t) = e^{-2\beta t} \int_0^\infty \{ [m(D - \beta)q(D - \beta)z']^2 + [r_1(D)z']^2 \} d\tau$$

or

$$V_1(\mathbf{z}, t) = e^{-2\beta t} \int_0^\infty \{ [n(D - \beta)p(D - \beta)z']^2 + [r_1(D)z']^2 \} d\tau,$$

from which it is evident that  $V_1$  is positive semidefinite. Suppose  $V_1$  is not positive definite. Then

$$m(D - \beta)q(D - \beta)z' = -n(D - \beta)p(D - \beta)z' = r_1(D)z' = 0$$

for some nonzero  $z'$ . That is,  $m(s - \beta)q(s - \beta)$ ,  $n(s - \beta)p(s - \beta)$  and  $r_1(s)$  have common factors. This can only happen if  $m(\cdot)q(\cdot)$  and  $n(\cdot)p(\cdot)$  have imaginary zeros, a situation ruled out by hypothesis. This proves the positive definiteness of  $V_1(\mathbf{z}, t)$ .

E. Observing that  $m(D - \beta)q(D - \beta)(ze^{\beta t}) = e^{\beta t}m(D)q(D)z(t)$ , we can write (10) as

$$(14) \quad V_1(\mathbf{z}, t) = e^{-2\beta t} \int_{t(0)}^{t(\mathbf{z})} e^{2\beta\tau} \{ [m(D)q(D)z][n(D)p(D)z] - [r_1(D + \beta)z]^2 \} d\tau.$$

The time derivative of  $V_1(\mathbf{z}, t)$  is then given by

$$(15) \quad \frac{dV_1(\mathbf{z}, t)}{dt} = -2\beta V_1(\mathbf{z}, t) + [m(D)q(D)z][n(D)p(D)z] - [r_1(D + \beta)z]^2.$$

In order to find the value of  $dV_1(\mathbf{z}, t)/dt$  when  $z(t)$  is the solution of (4), replace  $\mathbf{z}(t)$  by  $\mathbf{x}^*(t)$  and  $z(t)$  by  $y(t)$  in (15) and use (4) (after multiplication by  $m(D)q(D)y$ ) to obtain

$$(16) \quad \left. \frac{dV_1(\mathbf{x}^*, t)}{dt} \right|_{(4)} = -2\beta V_1(\mathbf{x}^*, t) - k(t)[m(D)q(D)y][n(D)q(D)y] - [r_1(D + \beta)y]^2.$$

F. Let it be assumed that  $Z(s - \alpha)$  is positive real (p.r.) for some  $\alpha \geq 0$ . Then by Lemma 2, the solution  $W(t)$  of the equation

$$(17) \quad m(D - \alpha)(e^{\alpha t}q(D)W) + n(D - \alpha)(e^{\alpha t}q(D)W) = 0$$

is asymptotically stable. Let  $\eta_1$  be the order of (17). It is easy to verify that  $\eta_1 \geq n$  and the order of (8) is equal to  $\eta_1$  or, at the most, greater than  $\eta_1$  by 1. Further, let  $\mathbf{W}(t)$  be the state vector for (17).

LEMMA 4. *If  $Z(s - \alpha)$  is p.r. for some  $\alpha \geq 0$ , then there exists a positive semi-definite quadratic form in  $\mathbf{W}$  defined by*

$$(18) \quad V_2(\mathbf{W}, t) = e^{-2\alpha t} \int_{t(0)}^{t(\mathbf{W})} \{ [m(D - \alpha)(e^{\alpha\tau}q(D)W)][n(D - \alpha)(e^{\alpha\tau}q(D)W)] - [r_2(D)(e^{\alpha\tau}q(D)W)]^2 \} d\tau.$$

*Proof.* The proof is similar to the proof of Lemma 3.

G. The time derivative of (18) is given by

$$(19) \quad \frac{dV_2(\mathbf{W}, t)}{dt} = -2\alpha V_2(\mathbf{W}, t) + [m(D)q(D)W][n(D)q(D)W] - [r_2(D + \alpha)q(D)W]^2.$$

Its value along the trajectories of (4) is obtained by substituting  $y(t)$ , the solution of (4), for  $W(t)$  in (19):

$$(20) \quad \left. \frac{dV_2(\mathbf{W}, t)}{dt} \right|_{(4)} = -2\alpha V_2(\mathbf{x}^*, t) + [m(D)q(D)y][n(D)q(D)y] - [r_2(D + \alpha)q(D)y]^2.$$

It can be verified that

$$(21) \quad \begin{aligned} \frac{d}{dt}[k(t)V_2(\mathbf{W}, t)] \Big|_{(4)} &= \left\{ \left( \frac{dk}{dt} / k \right) - 2\alpha \right\} k(t)V_2(\mathbf{x}^*, t) \\ &+ k(t)[m(D)q(D)y][n(D)q(D)y] - k(t)[r_2(D + \alpha)q(D)y]^2. \end{aligned}$$

The proof of the following lemma is obvious.

LEMMA 5. Let  $\delta_1(t), \delta_2(t)$  be bounded real functions on  $[0, \infty)$ . If  $V_1(\mathbf{x}^*, t) \geq 0, V_2(\mathbf{x}^*, t) \geq 0$  and  $V_{12}(\mathbf{x}^*, t) = \delta_1(t)V_1(\mathbf{x}^*, t) + \delta_2(t)V_2(\mathbf{x}^*, t)$ , then

$$V_{12}(\mathbf{x}^*, t) \leq \sup_{t \geq 0} [\delta_1(t), \delta_2(t)](V_1 + V_2)(\mathbf{x}^*, t).$$

H. Consider

$$[r_1(D + \beta)y]^2 = \{[Ev m(D)q(D)n(-D)p(-D)]^{(-)}y\}^2,$$

which is nonnegative and quadratic in  $\mathbf{x}^*$ . From Lemma 3,  $V_1(\mathbf{x}, t)$  is positive definite and quadratic in  $\mathbf{x}^*$ . By a well-known property of quadratic forms, there exists a nonnegative constant  $\gamma_1$  such that

$$[r_1(D + \alpha)y]^2 \geq \gamma_1 V_1(\mathbf{x}^*, t).$$

Similarly, there exists a nonnegative constant  $\gamma_2$  such that

$$[r_2(D + \alpha)q(D)y]^2 \geq \gamma_2 V_2(\mathbf{x}^*, t).$$

Let  $\gamma$  be equal to the minimum of the two numbers  $\gamma_1, \gamma_2$ ; and let  $V(\mathbf{x}^*, t) = V_1(\mathbf{x}^*, t) + k(t)V_2(\mathbf{x}^*, t)$ . Therefore, in view of the fact that  $k(t)$  is nonnegative,

$$(22) \quad [r_1(D + \beta)y]^2 + k(t)[r_2(D + \alpha)q(D)y]^2 \geq \gamma V(\mathbf{x}^*, t).$$

I. The time derivative of  $V(\mathbf{x}^*, t)$  along the solution of (4) is obtained by adding (16) and (21) to give

$$(23) \quad \begin{aligned} \frac{dV(\mathbf{x}^*, t)}{dt} \Big|_{(4)} &= -2\beta V_1(\mathbf{x}^*, t) + \left[ \left( \frac{dk}{dt} / k \right) - 2\alpha \right] k(t)V_2(\mathbf{x}^*, t) \\ &- [r_1(D + \beta)y]^2 - k(t)[r_2(D + \alpha)q(D)y]^2. \end{aligned}$$

Let  $\theta = ((dk/dt)/k)$  and

$$\sup_{t \geq 0} (-2\beta - \gamma, \theta(t) - 2\alpha - \gamma) = -\xi(t).$$

Then, using Lemma 5, from (23) one obtains

$$(24) \quad \frac{dV(\mathbf{x}^*, t)}{dt} \Big|_{(4)} \leq -\xi(t)V(\mathbf{x}^*, t).$$

J. DEFINITION. The solutions of (2) are said to have Property K-C if there exist a positive constant  $\eta_0$  and a real function  $\mu(t)$  on  $[t_0, \infty)$  such that

$$(25) \quad \|\mathbf{x}(t)\| \leq \|\mathbf{x}(t_0)\| \eta_0 e^{-[\mu(t) - \mu(t_0)]/2}, \quad t \geq t_0.$$

LEMMA 6 (The Krasovskii-Corduneanu theorem). The solutions of (2) have Property K-C if there exist a positive definite and decrescent quadratic form



$v(\mathbf{x}, t) = \mathbf{x}'P(t)\mathbf{x}$ , and a real-valued function  $\lambda(t)$  on  $[t_0, \infty)$  such that the derivative of  $v(\mathbf{x}, t)$  along the solutions of (2) satisfies the inequality

$$(26) \quad \left. \frac{dv}{dt} \right|_{(2)} \leq -\lambda(t)v.$$

*Proof.* Because  $v(\mathbf{x}, t)$  is positive definite and decrescent, there exist positive constants  $\alpha_1$  and  $\alpha_2$  such that

$$\alpha_1 \mathbf{x}'\mathbf{x} \leq \mathbf{x}'P(t)\mathbf{x} \leq \alpha_2 \mathbf{x}'\mathbf{x}.$$

Integration of (26) gives

$$v(\mathbf{x}, t) \leq v(\mathbf{x}_0, t_0) \exp \left( - \int_{t_0}^t \lambda(\tau) d\tau \right).$$

Consequently,

$$\begin{aligned} \alpha_1 \mathbf{x}'\mathbf{x} &\leq v(\mathbf{x}, t) \leq v(\mathbf{x}_0, t_0) \exp \left( - \int_{t_0}^t \lambda(\tau) d\tau \right) \\ &\leq \alpha_2 \|\mathbf{x}_0\|^2 \exp \left( - \int_{t_0}^t \lambda(\tau) d\tau \right), \end{aligned}$$

from which

$$\|\mathbf{x}\|^2 \leq \frac{\alpha_2}{\alpha_1} \|\mathbf{x}_0\|^2 \exp \left( - \int_{t_0}^t \lambda(\tau) d\tau \right).$$

Therefore, the solutions of (2) have Property K–C with  $\eta_0 = \sqrt{\alpha_2/\alpha_1}$  and  $\lambda(t)$  equal to the derivative of  $\mu(t)$ .

**COROLLARY 1.** *If  $\int_{t_0}^t \lambda(\tau) d\tau$  increases without bound as  $t \rightarrow \infty$ , then the system (2) is asymptotically stable.*

**COROLLARY 2.** *If  $(1/T)\int_{t_0}^{t_0+T} -\lambda(\tau) d\tau \leq -\gamma$  for some constants  $\gamma > 0$  and  $T > 0$ , then*

$$\|\mathbf{x}\| \leq \|\mathbf{x}_0\| \eta_0 \exp [-\gamma(t - t_0)/2], \quad t \geq t_0,$$

and the system is exponentially stable.

**K.** Let  $\zeta(t)$  be a nonnegative (integrable and bounded) function on  $[t_0, \infty)$ , and  $h(t) = \exp -\int_{t_0}^t \zeta(\tau) d\tau$ . Assume that the integral  $\int_{t_0}^t \zeta(\tau) d\tau \leq M < \infty$  for all  $t$  in  $[t_0, \infty)$ , and

$$0 < \varepsilon \leq \lim_{t \rightarrow \infty} \int_{t_0}^t \zeta(\tau) d\tau \leq M < \infty.$$

Then  $h(t)$  is a bounded positive function. Note that

$$\left( \frac{dh(t)}{dt} / h(t) \right) = -\zeta(t),$$

which is nonpositive.

Let

$$(27) \quad V_0(\mathbf{x}^*, t) = h(t)\{V_1(\mathbf{x}^*, t) + k(t)V_2(\mathbf{x}^*, t)\},$$

where  $V_1(\mathbf{x}^*, t)$ ,  $V_2(\mathbf{x}^*, t)$  are as defined in Lemmas 3 and 4, respectively. The time derivative of  $V_0(\mathbf{x}^*, t)$  along the solutions of (4) is given by

$$\begin{aligned}
 \left. \frac{dV_0(\mathbf{x}^*, t)}{dt} \right|_{(4)} &= -\zeta(t)V_0(\mathbf{x}^*, t) + h(t)\{-2\beta V_1(\mathbf{x}^*, t) \\
 &\quad + [\theta(t) - 2\alpha]k(t)V_2(\mathbf{x}^*, t) - [r_1(D + \beta)y]^2 \\
 &\quad - k(t)[r_2(D + \alpha)q(D)y]^2\} \\
 (28) \qquad \qquad \qquad &\leq \sup_{t \geq 0} \{[-2\beta - \gamma, \theta(t) - 2 - \gamma] - \zeta(t)\} V_0(\mathbf{x}^*, t).
 \end{aligned}$$

**3. Main results.**

**THEOREM 1.** *If*

- (a)  $Z(s)$  (the “multiplier”) is a function of the complex variable  $s$  such that  $Z(s - \alpha)$  is positive real for some  $\alpha \geq 0$ ;
- (b)  $Z(s - \beta)G(s - \beta)$  is positive real for some  $\beta \geq 0$ , and  $m(s - \beta)q(s - \beta)$ ,  $n(s - \beta)p(s - \beta)$  have no imaginary zeros; and
- (c) for some positive constant  $\gamma$ ,

$$(29) \qquad - \int_{t_0}^t \sup_{\tau \geq 0} \left[ -2\beta - \gamma, \left( \frac{dk}{d\tau} / k \right) - 2\alpha - \gamma \right] d\tau$$

increases without bound as  $t \rightarrow \infty$ , then the null solution of (1) is asymptotically stable.

**THEOREM 2.** *If the hypotheses (a) and (b) of Theorem 1 are valid,  $M$  is a positive constant,  $\theta(t)$  and  $[\theta(t) + 2(\beta - \alpha)]^+$  denote respectively  $((dk/dt)/k)$  and the positive values of  $[\theta(t) + 2(\beta - \alpha)]$ , and*

(d)

$$\begin{aligned}
 (i) \qquad \qquad \qquad & \frac{1}{T} \int_{t_0}^{t_0+T} [\theta(\tau) + 2(\beta - \alpha)]^+ d\tau \leq M < \infty \\
 (30) \qquad \qquad \qquad & \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \text{for all finite } T > 0, \\
 (ii) \qquad \qquad \qquad & \lim_{T \rightarrow \infty} \frac{1}{T} \int_{t_0}^{t_0+T} [\theta(\tau) + 2(\beta - \alpha)]^+ d\tau \leq 2\beta + \gamma - \nu \\
 & \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \text{for some } \nu > 0,
 \end{aligned}$$

then the null solution of (1) is exponentially stable.

*Proof of Theorem 1.* As a Lyapunov function candidate for (4) and hence for (1), choose

$$V(\mathbf{x}^*, t) = V_1(\mathbf{x}^*, t) + k(t)V_2(\mathbf{x}^*, t),$$

where  $V_1(\mathbf{x}^*, t)$  and  $V_2(\mathbf{x}^*, t)$  are defined by (10) and (18) respectively. Clearly,  $V(\mathbf{x}^*, t)$  is positive definite, radially unbounded, has continuous first partial derivatives, and satisfies decrescent conditions by virtue of the boundedness of  $k(t)$ . Its time derivative along the solutions of (4) was shown to satisfy inequality (24). Invoking the Krasovskii–Corduneanu theorem (Corollary 1), we conclude that hypothesis (c) implies asymptotic stability.

*Proof of Theorem 2.* As a Lyapunov function candidate for (4), choose

$$V_0(\mathbf{x}^*, t) = h(t)\{V_1(\mathbf{x}^*, t) + k(t)V_2(\mathbf{x}^*, t)\},$$

where  $V_1(\mathbf{x}^*, t)$ ,  $V_2(\mathbf{x}^*, t)$  are defined by (10), (18) respectively, and  $h(t)$  is defined in subsection K above.  $V_0(\mathbf{x}^*, t)$  is positive definite, radially unbounded, has continuous partial derivatives, and satisfies decrescent conditions by virtue of the boundedness of  $k(t)$  and  $h(t)$ . Its time derivative along the solutions of (4) was shown to satisfy inequality (28). Invoking Corollary 2 of the Krasovskii–Corduneanu theorem, we conclude that hypothesis (d) implies exponential stability.

*Remarks.* (a) The present global condition (29) for asymptotic stability of (1) differs from the Freedman and Zames average condition (3) in that the integrand of (29) may assume negative values, whereas the integrand of (3) may not.

(b) (Note that  $\alpha \geq \beta$ .) The average condition (30) for exponential stability of (1) allows larger positive variations of  $\theta(t)$  over a finite interval than (3) and the negative lobes of  $\theta(t)$  do not enter into the integrand of (30).

(c) The condition on  $G(s)$  in Theorems 1 and 2 can be replaced by an explicit geometric condition, free of multipliers (for details see [1] and [9]).

#### REFERENCES

- [1] M. FREEDMAN AND G. ZAMES, *Logarithmic variation criteria for the stability of systems with time-varying gains*, this Journal, 6 (1968), pp. 487–507.
- [2] G. ZAMES, *On the stability of nonlinear time-varying feedback systems*, Proc. NEC, 20 (1964), pp. 725–730.
- [3] K. S. NARENDRA AND R. M. GOLDWYN, *A geometrical criterion for the stability of certain nonlinear nonautonomous systems*, IEEE Trans. Circuit Theory, CT-11 (1964), pp. 406–407.
- [4] I. W. SANDBERG, *A frequency domain condition for the stability of systems containing a single time varying nonlinear element*, Bell System Tech. J., 43 (1964), pp. 1601–1638.
- [5] R. W. BROCKETT AND L. J. FORYS, *On the stability of systems containing a time-varying gain*, Proc. Second Allerton Conference, University of Illinois, Urbana, 1964, pp. 413–430.
- [6] N. N. KRASOVSKII, *Stability of Motion*, Stanford University Press, Stanford, Calif., 1963.
- [7] C. CORDUNEANU, *The application of differential inequalities to the theory of stability*, An. Sti. Univ. “Al. I. Cuza” Iasi Sect. I (N.S.), 6 (1960), pp. 47–58.
- [8] L. WEINBERG AND P. SLEPIAN, *Positive real matrices*, J. Math. Mech., 9 (1960), pp. 71–83.
- [9] R. W. BROCKETT AND J. L. WILLEMS, *Frequency-domain stability criteria, I*, IEEE Trans. Automatic Control, AC-10 (1965), pp. 255–261.

## A NOTE ON THE NECESSARY CONDITIONS FOR OPTIMAL STRATEGIES IN A CLASS OF NONCOOPERATIVE $N$ -PERSON DIFFERENTIAL GAMES\*

I. G. SARMA AND U. R. PRASAD†

**Abstract.** In this note, we show that the main results of our previous paper [1], the properties of the value function vector and the equilibrium point principle stated in Theorems 1 and 2 of [1], are extendable to a slightly larger but considerably more realistic class of games.

In general, an  $N$ -person game can have multiple equilibrium points. The equilibrium points of a two-person zero-sum game (called the saddle points), if nonunique, are automatically both interchangeable and equivalent. This is *not* the case for a general  $N$ -person game. However, if all the equilibrium points of an  $N$ -person game are interchangeable (without necessarily being equivalent), then all of them constitute the Nash noncooperative solution of the game. These notions are discussed by Nash [2] and Luce and Raiffa [3].

In our original paper [1], the solution of an  $N$ -person game is characterized by (2.4) and the resulting payoff function is assumed single-valued on  $\mathcal{R}$ . This is equivalent to assuming the twin conditions of interchangeability and equivalence which are valid in general *only* for two-person zero-sum games. To emphasize that two-person zero-sum ideas are not applicable when dealing with  $N$ -person games, we introduce the more realistic concept known as the Nash noncooperative solution for these games. This solution also satisfies (2.4) but the resulting payoff function can be discontinuous across certain well-defined manifolds in  $\mathcal{R}$ . The following simple example bears out these points.

*Example.* The state of the game satisfies the following differential equations:

$$(1) \quad \begin{aligned} \dot{x}_1 &= -x_2, \\ \dot{x}_2 &= u^1 + cu^2 \end{aligned}$$

with  $u^1$  and  $u^2$  constrained as

$$(2) \quad |u^1| \leq 1; \quad |u^2| \leq 1.$$

The terminal surface  $\mathcal{T}_1 \cup \mathcal{T}_2$  is given in terms of the parameter  $\sigma = (x_{2f}, t_f)$  as

$$(3) \quad \begin{aligned} \mathcal{T}_1 &= \left\{ (x_{1f}, x_{2f}) : x_{1f} = \frac{x_{2f}^2}{2(1+c)}; x_{2f} \geq 0 \right\}, \\ \mathcal{T}_2 &= \left\{ (x_{1f}, x_{2f}) : x_{1f} = \frac{x_{2f}^2}{2}; x_{2f} \geq 0 \right\}. \end{aligned}$$

The playing space is in between  $\mathcal{T}_1$  and  $\mathcal{T}_2$  as shown in Fig. 1. The payoff functionals

\* Received by the editors April 7, 1970, and in revised form December 10, 1970.

† Department of Electrical Engineering, Indian Institute of Technology, Kanpur, India.

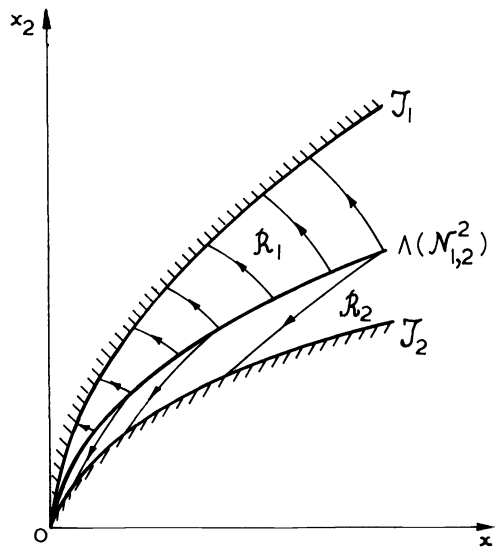


FIG. 1. Nash noncooperative solution of example

of the players are given by

$$\begin{aligned}
 (4) \quad J^1[x_0, \mathbf{u}] &= \phi^1(\sigma) + \int_0^{t_f} dt, \\
 J^2[x_0, \mathbf{u}] &= \phi^2(\sigma) + \int_0^{t_f} \{|u^1| + b|u^2|\} dt,
 \end{aligned}$$

where

$$\begin{aligned}
 (5) \quad \phi^1(\sigma) &= \begin{cases} \frac{x_{2f}}{(1+c)}, & x_f \in \mathcal{T}_1, \\ x_{2f}, & x_f \in \mathcal{T}_2, \end{cases} \\
 \phi^2(\sigma) &= \begin{cases} \frac{x_{2f}(1+b)}{(1+c)}, & x_f \in \mathcal{T}_1, \\ x_{2f}, & x_f \in \mathcal{T}_2. \end{cases}
 \end{aligned}$$

The constants  $b$  and  $c$  in (1) and (4) are related by

$$(6) \quad 2 < c < b.$$

The optimal control actions should necessarily minimize the Hamiltonians for the players. Hence, we have

$$\begin{aligned}
 (7) \quad H^1 &= 1 - \lambda_1^1 x_2 + \lambda_2^1 (u^1 + cu^2), \\
 H^2 &= |u^1| + b|u^2| - \lambda_1^2 x_2 + \lambda_2^2 (u^1 + cu^2),
 \end{aligned}$$

and

$$\begin{aligned}
 (8) \quad u^{1*} &= -\text{sgn } \lambda_2^1, \\
 u^{2*} &= -\text{dez}(\lambda_2^2 c/b),
 \end{aligned}$$

where we define

$$(9) \quad \begin{aligned} \operatorname{sgn} z &= \begin{cases} +1, & z > 0, \\ -1, & z < 0, \end{cases} \\ \operatorname{dez} z &= \begin{cases} +1, & z > 1, \\ 0, & |z| < 1, \\ -1, & z < -1. \end{cases} \end{aligned}$$

The adjoint equations are given for  $l = 1, 2$  by

$$(10) \quad \begin{aligned} \dot{\lambda}_1^l &= 0, \\ \dot{\lambda}_2^l &= \lambda_1^l. \end{aligned}$$

As  $u^{1*}$  and  $u^{2*}$  assume the only values of  $\pm 1$  and 0, in view of (9), the term  $H_{u^m}^l U_x^m$  corresponding to the other player's optimal strategy is absent in (10).

Applying the transversality condition (4.4) at  $\mathcal{T}_1$ , we have

$$(11) \quad \begin{aligned} \frac{1}{1+c} - \lambda_1^1 \left( \frac{x_{2f}}{1+c} \right) - \lambda_2^1 &= 0, \\ 1 - \lambda_1^1 x_{2f} + \lambda_2^1 (u^1 + cu^2) &= 0, \\ \frac{1+b}{1+c} - \lambda_1^2 \left( \frac{x_{2f}}{1+c} \right) - \lambda_2^2 &= 0, \\ |u^1| + b|u^2| - \lambda_1^2 x_{2f} + \lambda_2^2 (u^1 + cu^2) &= 0 \end{aligned}$$

with the quantities referred to time  $t_f$ . Solving (11) to be consistent with (6), (8) and (10) yields

$$(12) \quad \begin{aligned} \lambda_1^1(t_f) &= \frac{1}{x_{2f}}; & \lambda_2^1(t_f) &= 0, \\ \lambda_1^2(t_f) &= \frac{2+c+b}{(2+c)x_{2f}}; & \lambda_2^2(t_f) &= \frac{b}{2+c} \end{aligned}$$

and

$$(13) \quad u^1(t) = 1; \quad u^2(t) = 0 \quad \text{for } t < t_f.$$

By a similar application at  $\mathcal{T}_2$ , we have

$$(14) \quad u^1(t) = 1; \quad u^2(t) = -1 \quad \text{for } t < t_f.$$

The resulting paths are shown in Fig. 1. It is clear that the strategy of player 1 is continuous on  $\mathcal{R}$ . The second player's strategy is discontinuous giving rise to his dispersal surface  $\Lambda$ . This surface divides the playing space into two regions,  $\mathcal{R}_1$  and  $\mathcal{R}_2$ . (Because of this, this switching surface could be denoted by  $\mathcal{N}_{1,2}^2$ .) For starting points on  $\Lambda$ , the value to the second player  $W^2$  must be the same

whether the paths reach  $\mathcal{T}_1$  or  $\mathcal{T}_2$ . By actually equating the values, we get

$$(15) \quad \Lambda = \left\{ (x_1, x_2) : x_1 = \frac{\eta}{2} x_2^2, x_2 \geq 0 \right\},$$

where  $\eta$  is given by

$$(16) \quad \begin{aligned} & \frac{b+1}{c-1} - \left( \frac{2+b-c}{c-1} \right) \left( \frac{-1+(c-1)\eta}{c-2} \right)^{1/2} \\ & = -1 + (2+c+b) \left( \frac{1+\eta}{(1+c)(2+c)} \right)^{1/2}. \end{aligned}$$

The optimality of this solution can be conclusively shown by verifying the Bellman equations for  $W^1$  and  $W^2$ . It can be further verified easily that for starting points on  $\Lambda$ , the value function  $W^1$  of player 1 is greater for the paths reaching  $\mathcal{T}_2$  in comparison with the paths ending on  $\mathcal{T}_1$ . This result has no counterpart in two-person zero-sum differential games, for obvious reasons.

Now, we make the necessary changes in the results of our original paper to include these situations. Let  $\mathbf{U}^*$  be a playable  $N$ -tuple which introduces a regular decomposition on  $\mathcal{R}$ . The hypersurface  $\mathcal{N}_{i_1, i_2}^l$  which (if it exists) separates  $\mathcal{R}_{i_1}$  and  $\mathcal{R}_{i_2}$  and across which  $U^{l*}$  is discontinuous, is called the dispersal surface of player  $l$  (R. Isaacs, [4]). The payoff  $\mathbf{P}(t, x, \mathbf{U}^*)$  is single-valued in  $\mathcal{R}$  except on the dispersal surfaces. On the dispersal surface  $\mathcal{N}_{i_1, i_2}^l$  of player  $l$ , the  $l$ th component of  $\mathbf{P}(t, x, \mathbf{U}^*)$  is assumed to be independent of the two optimal paths branching into  $\mathcal{R}_{i_1}$  and  $\mathcal{R}_{i_2}$ . Then  $\mathbf{U}^*$  is said to be the Nash noncooperative solution relative to the normal form defined by the strategy sets  $U_i^l, l = 1, 2, \dots, N$ , if the addition  $\mathbf{P}(t, x, \mathbf{U}^*)$  satisfies (2.4) of [1] with the interpretation that suitable one-sided limits are taken whenever necessary. The rest of the assumptions on the optimal paths remain the same as in the paper. The solution  $\mathbf{U}^*$  in this sense consists of interchangeable (but not necessarily equivalent) equilibrium points, whenever multiple paths arise.

For the enlarged class of games having the Nash noncooperative solution, Theorem 1 holds with the difference that the value function  $\mathbf{W}$  need not be continuous across the  $N$  manifolds. As explained earlier,  $W^l$  will however be continuous across  $\mathcal{N}_{i_1, i_2}^l$ . Analogous to [5], it is straightforward to show that for any  $(t, x)$  on  $\mathcal{N}_{i_1, i_2}^l$ , the following condition holds:

$$(17) \quad \begin{aligned} & H^l(t, x, (\mathbf{U}^*; U_{i_1}^l), \lambda_{i_1}^l) dt - \lambda_{i_1}^l dx \\ & = H^l(t, x, (\mathbf{U}^*; U_{i_2}^l), \lambda_{i_2}^l) dt - \lambda_{i_2}^l dx, \end{aligned}$$

where  $dt$  and  $dx$  are any differentials on the manifold and the subscripts  $i_1, i_2$  indicate the appropriate one-sided limits.

It may be noted that in (4.3), it is tacitly assumed that the terminal surface has the parametric representation

$$(18) \quad t = T_{ij_i}(\sigma); \quad x = X_{ij_i}(\sigma).$$

Similarly, in (4.6), the following parametric representation is assumed for the

manifold  $\mathcal{M}_{ik}$ :

$$(19) \quad t = T_{ik}(\sigma); \quad x = X_{ik}(\sigma).$$

These are clear from the context.

#### REFERENCES

- [1] I. G. SARMA, R. K. RAGADE AND U. R. PRASAD, *Necessary conditions for optimal strategies in a class of noncooperative  $N$ -person differential games*, this Journal, 7 (1969), pp. 637–644.
- [2] J. NASH, *Non-cooperative games*, Ann. of Math., 54 (1951), pp. 286–295.
- [3] R. O. LUCE AND H. RAIFFA, *Games and Decisions*, John Wiley, New York, 1957.
- [4] RUFUS ISAACS, *Differential Games*, John Wiley, New York, 1965.
- [5] L. D. BERKOVITZ, *Necessary conditions for optimal strategies in a class of differential games and control problems*, this Journal, 5 (1967), pp. 1–24.



## EXISTENCE OF OPTIMAL STOCHASTIC CONTROL LAWS\*

V. E. BENEŠ†

**Abstract.** We give an approach to optimal control of systems described by stochastic functional-differential equations of Itô's type:  $dx(t) = f(t, x, u(t, x)) dt + dw(t)$ ,  $0 \leq t \leq 1$ , with a cost functional  $k(u) = E \int_0^1 c(t, x, u(t, x)) dt$  to be minimized. Here  $w(\cdot)$  is Brownian motion,  $f$  and  $c$  are nonanticipative functionals describing system dynamics and cost rate respectively, and  $u(t, \cdot)$  is a causal control law, to be chosen, taking values in a space  $\Gamma$  of control points, and depending at  $t$  at most on given information about the past  $\{x(s), s \leq t\}$ .

A key role is played by the information available for control. This is described by giving, for each  $t$ , a sub- $\sigma$ -algebra  $G_t$  of the  $\sigma$ -algebra  $S_t$  over the continuous functions  $C = C[0, 1]$  generated by sets  $\{y: y(s) \in A\}$  with  $0 \leq s \leq t$  and  $A$  Borel in  $R^d$ .  $S_t$  is the  $\sigma$ -algebra corresponding to knowing the whole past of the trajectory prior to  $t$ . The set  $\mathcal{U}$  of admissible control laws consists of functions  $u: [0, 1] \times C \rightarrow \Gamma$  which are Lebesgue in  $t$ , and  $G_t$ -measurable in  $y$  for each  $t$ . There is a  $\sigma$ -algebra  $G$  over  $[0, 1] \times C$  such that admissibility is equivalent to  $G$ -measurability.

Our formulation is based on a result of Girsanov (Teoriya Veroyatnostei, 5 (1960), p. 285): For  $\varphi$  a nonanticipative functional of Brownian motion  $w$ , the transformed measure  $d\tilde{P} = \exp \zeta(\varphi) dP$  with

$$\zeta(\varphi) = \int_0^1 \varphi dw - \frac{1}{2} \int_0^1 |\varphi|^2 dt$$

makes the functions  $w(\cdot) - \int_0^\cdot \varphi dt$  a Wiener process, provided  $E \exp \zeta(\varphi) = 1$ . This result suggests and justifies taking, for the solution  $x(\cdot)$  of the system equations, the process determined by Girsanov's device with  $\varphi = f(t, w, u(t, w))$ ; in this case we say that  $u$  attains the density  $\exp \zeta(\varphi)$ .

The control problem is reformulated as a search for admissible  $u$  that achieve  $\inf_{u \in \mathcal{U}} E \exp \zeta(\varphi) \int_0^1 c dt$ . That is, with each  $u \in \mathcal{U}$  we associate, as the solution of the system equations to be considered for the purpose of our minimization of  $k(\cdot)$ , the functions  $w(\cdot)$  under the measure  $\exp \zeta(\varphi) dP$ , with the justification that under this measure

$$w(t) - \int_0^t f(s, w, u(s, w)) ds$$

is a Wiener process.

Novelty of the approach lies in these features: (i) control is closed loop; (ii) admissible controls need not be smooth; (iii) the Radon–Nikodym derivative used by Girsanov directly gives a measure corresponding to a solution of the system equations.

As a principal result, we prove that if  $\Gamma$  is compact metric, if  $f(t, y, u)$  grows at most linearly with  $y(t)$ , and if  $G_t = S_t$  (i.e., if the whole past is available for control), then the set of densities  $\{\exp \zeta(\varphi): \varphi = f(t, w, u(t, w)), u \in \mathcal{U}\}$  attainable by the admissible control laws is convex, and there exists an optimal control law  $u^* \in \mathcal{U}$  achieving  $\inf_{u \in \mathcal{U}} E \exp \zeta(\varphi) \int_0^1 c dt$ .

**1. Introduction.** We give a new approach to two aspects of the optimal control of systems described by stochastic functional-differential equations. These aspects are the *formulation* of the control problem, and the *existence of* optimal control laws. Particular emphasis is placed on the role of the information about the past of the trajectory that is available to the controller as a basis for control decisions.

Novelty of the approach lies in four features: (i) A generalized notion of what is a solution of the stochastic system equations is used: a Radon–Nikodym

---

\* Received by the editors May 19, 1970, and in revised form October 19, 1970.

† Bell Telephone Laboratories, Inc., Murray Hill, New Jersey 07974.

derivative functional used by Girsanov [1] directly gives a measure corresponding to a solution. (ii) Control is closed loop. (iii) No restriction of the admissible control laws to be (for example) Lipschitz is imposed; a control law is admissible if it is unambiguous, nonanticipative, suitably measurable, and restricted in its dependence to the available information; it does not have to be smooth. (iv) A general description of the information available to the controller is assumed; to our knowledge, only the case of complete information about the past carries with it sufficient convexity to yield closure and existence results by the usual methods.

The systems of principal interest here are those governed by a stochastic functional-differential equation

$$(1) \quad dx(t) = f(t, x, u(t, x)) dt + dw(t), \quad t \in [0, 1],$$

where  $w(\cdot)$  is a  $d$ -dimensional Brownian motion,  $f$  is a  $d$ -vector-valued nonanticipative functional, and the “control”  $u$  is a function (with values in a prescribed space  $\Gamma$ ) whose value at  $t$  may depend at most on specified information about the past of  $x(\cdot)$  prior to  $t$ . The tasks to be tackled are first, (because  $u$  need not be smooth) to formulate a suitable meaning for the “existence” of solutions to (1), second, to pose a minimization problem for cost functionals

$$k(u) = E \int_0^1 c(t, x, u(t, x)) dt, \quad c \geq 0,$$

under the constraint (1) and the information restriction on  $u$ , and third, to prove within this formulation that optimal admissible controls exist. Cost criteria involving final values of smooth functions and random stopping times can also be handled by our approach, as the reader can verify.

We first informally describe our approach to the existence problem by considering the system (1) with  $u(\cdot, \cdot)$  chosen and fixed, and we set

$$h(t, y) = f(t, y, u(t, y))$$

for  $y$  in the space  $C$  of continuous  $R^d$ -valued functions. We also assume that a Brownian motion  $w(\cdot)$  is given on a probability space of points  $\omega$ . In many cases of interest, e.g., if  $h(t, \cdot)$  is Lip in  $y$  uniformly in  $t$ , (1) has a unique solution [2], [3], a stochastic process  $x(t, \omega)$  with continuous sample paths; the measure  $\nu$  which  $x(\cdot, \omega)$  induces on  $C$  is absolutely continuous with respect to Wiener measure  $\mu$ , with Radon–Nikodym derivative [3]:

$$(2) \quad \frac{d\nu}{d\mu}(y) = \exp \left\{ \int_0^1 h(t, y) dy(t) - \frac{1}{2} \int_0^1 |h(t, y)|^2 dt \right\}, \quad y \in C.$$

In such a case we do not need to solve (1) to find the cost of using control law  $u(\cdot, \cdot)$ ; the cost is simply

$$(3) \quad k(u) = \int_C \int_0^1 c(t, y, u(t, y)) dt \frac{d\nu}{d\mu}(y) d\mu(y).$$

These circumstances suggest directly using a Radon–Nikodym derivative functional like (2) to *define* a solution of (1), to single it out as the relevant one,

and thereby to give the definite form (3) to the cost of using the control  $u(\cdot, \cdot)$ . That this is possible was shown by the late I. V. Girsanov [1] in a 1960 paper that is only now attracting the attention it deserves.

Girsanov considered a nonanticipating Brownian functional  $\varphi(t, \omega)$ , square-integrable almost surely, and defined the functional

$$(4) \quad e^{\zeta_s(\varphi)} = \exp \left\{ \int_s^t \varphi(t, \omega) dw(t) - \frac{1}{2} \int_s^t |\varphi(t, \omega)|^2 dt \right\}.$$

He then showed that if  $E \exp \zeta_s^t(\varphi) = 1$ , then the transformed measure  $d\tilde{P} = \exp \zeta_s^t(\varphi) dP$  makes the function

$$w(\cdot, \omega) - w(s, \omega) - \int_s^\cdot \varphi(u, \omega) du$$

a Wiener process on  $[s, t]$ , i.e., it makes  $w(\cdot, \omega) - w(s, \omega)$  a solution of the equation

$$w(t) - w(s) - \int_s^t \varphi(s, \omega) ds = \text{Wiener process.}$$

In other words, the functional “solves” the stochastic differential equation  $dx = \varphi(t, \omega) dt + dw$  by giving a transformation of the basic measure into one corresponding to a solution. Girsanov’s result is a generalization of Cameron and Martin’s translation theorem [4] in which  $\varphi$  was not random.

We propose to use Girsanov’s theorem to accomplish our first task, the clarification of the existence of solutions, as follows: to obtain the stochastic process corresponding to use of  $u(\cdot, \cdot)$  as control law, we put

$$(5) \quad \varphi(t, \omega) = f(t, w(\omega), u(t, w(\omega))), \quad \zeta(\varphi) = \zeta_0^1(\varphi),$$

in Girsanov’s functional (4), and we take the functions  $w(\cdot, \omega)$  under the measure  $\exp \zeta(\varphi) dP$  as the required process. This procedure gives the conveniently explicit form

$$(6) \quad E e^{\zeta(\varphi)} \int_0^1 c(t, w(\omega), u(t, w(\omega))) dt$$

for the cost of using  $u(\cdot, \cdot)$ , and obviates solving the system equation (1), salutary effects indeed; the cost can be computed by numerical integration over Wiener space!

Our second task, of posing a minimization problem, is now straightforward. We describe the information available for control in terms of a  $\sigma$ -algebra  $G$  on  $[0, 1] \times C$ , and we can characterize the admissible controls as the  $G$ -measurable functions taking values in the control space  $\Gamma$ . The control problem becomes this: To minimize (6) over admissible controls  $u(\cdot, \cdot)$ , with the understanding (5).

With a definite minimization problem at hand, we can pass to the third task, proving that optimal admissible control laws exist. Our basic result is that if the entire past of the trajectory is available for control, then optimal admissible control laws exist, provided that  $f(t, y, u)$  is continuous in the control variable  $u$ , that it does not grow faster than linearly in  $y(t)$ , that  $f(t, y, \Gamma)$  is convex, and that  $\Gamma$  is compact metric. In this case it is possible to construct, out of an arbitrary

minimizing sequence of control laws, a new minimizing sequence for which the corresponding stochastic processes converge to a process obtained by using an optimal control law. Closure and existence theorems, analogous to those used in deterministic optimal control, turn out to be difficult to prove, harder than in the nonstochastic case. This additional difficulty is due in large part to the role of the information available to the controller. As in the nonstochastic problems, convexity and continuity are basic in existence proofs, as is a version of Filippov's lemma [5].

**2. Historical remarks.** The problem of the existence of optimal stochastic control laws for processes of diffusion type has been studied by Kushner [6] for Lipschitz closed loop controls and measurable open loop control entering linearly. Fleming and Nisio [2] considered open loop controls entering multiplicatively, and proved existence theorems using Prokhorov's topology. The case of incomplete information was broached by Fleming [7] in a paper on control of partially observable diffusions: only some of the state-vector components are observed.

The functionals  $e^{\zeta}$  originated with the work [4] of Cameron and Martin on Wiener measure. After Itô extended stochastic integrals to random but non-anticipative integrands, they reappeared in the work of Ventcel [8] on additive functionals, in that of Skorokhod [3] on differentiability of measures corresponding to diffusion processes, and in Girsanov's work. They were first used in control theory by Mortensen [9]; Kailath [10], Duncan [11], and Kallianpur and Striebel [12] have noted their relevance to estimation and filtering. They are implicit in the work [13] of Stroock and Varadhan on diffusion processes with continuous coefficients, and explicit in McKean's exposition [14] of stochastic differentials and integrals.

**3. Formulation.** Let  $\Gamma$  be a compact metric space of control points, and let  $C = C[0, 1]$  denote the space of continuous functions  $y(\cdot)$  with  $y: [0, 1] \rightarrow R^d$ . The sets  $\{y(\cdot) \in C: y(s) \in A\}$  for  $0 \leq s \leq t \leq 1$  and  $A$  Borel in  $R^d$ , generate a  $\sigma$ -algebra  $S_t$  of  $C$ -subsets. This is the  $\sigma$ -algebra representing knowledge of the past from 0 to  $t$ . We shall suppose that the system dynamics are given by a function  $f: [0, 1] \times C \times \Gamma \rightarrow R^d$  with these properties:

- (i)  $f(t, y, \cdot)$  is continuous on  $\Gamma$  for each  $t, y$ ;
- (ii)  $f$  is nonanticipative in the strong sense that  $f(t, \cdot, u)$  is measurable with respect to  $S_t$  for each  $t, u \in [0, 1] \times \Gamma$ ;
- (iii)  $|f(t, y, u)|^2 \leq \kappa(1 + |y(t)|^2)$ ,  $\kappa$  a constant, for every  $t, y, u \in [0, 1] \times C \times \Gamma$ ,
- (iv)  $f(\cdot, y, u)$  is Lebesgue measurable for  $y, u \in C \times \Gamma$ .

We take the view that an admissible control law is a function  $u: [0, 1] \times C \rightarrow \Gamma$  with the interpretation that  $u(t, \cdot)$  indicates what point of the control space  $\Gamma$  is to be exercised as control at time  $t$ , and with the proviso that  $u(t, \cdot)$  depend only on whatever information (about the past of the trajectory) the controller is allowed to know, remember, and use at time  $t$ . The restrictions on  $u(\cdot, \cdot)$  representing the pattern of available information will be described mathematically by the concept of measurability with respect to a  $\sigma$ -algebra. Roughly speaking, if  $\mathcal{A}_1$  and  $\mathcal{A}_2$  are two  $\sigma$ -algebras over the same space, and if  $\mathcal{A}_1 \subseteq \mathcal{A}_2$ , then the functions measurable on  $\mathcal{A}_1$  are less complicated, ("depend on less") than those which are measurable with respect to  $\mathcal{A}_2$  but not  $\mathcal{A}_1$ .

The information pattern will be described by specifying, for each  $t \in [0, 1]$ , a  $\sigma$ -algebra  $G_t \subseteq S_t$ , and by imposing the condition that for each  $t$   $u(t, \cdot)$  be  $G_t$ -measurable.

Thus the *admissible controls* will be all those functions  $u: [0, 1] \times C \rightarrow \Gamma$  such that

- (i)  $u(\cdot, y)$  is Lebesgue for  $y \in C$ ,
- (ii)  $u(t, \cdot)$  is  $G_t$ -measurable for  $t \in [0, 1]$ .

It is possible and convenient to express the property of being an admissible control in terms of a single  $\sigma$ -algebra. This is done as follows: consider the measurable subsets  $E$  of  $[0, 1] \times C$  such that

- (i) every  $t$ -section of  $E$  is a  $G_t$ -set, for  $t \in [0, 1]$ ,
- (ii) every  $y$ -section of  $E$  is a Lebesgue set, for  $y \in C$ .

It is easy to verify that  $G$  is an algebra; since  $G$  is closed under monotone limits, it is a  $\sigma$ -algebra; it can then be proved that a function  $h$  on  $[0, 1] \times C$  is  $G$ -measurable if  $h(t, \cdot)$  is  $G_t$ -measurable for fixed  $t \in [0, 1]$  and  $h(\cdot, y)$  is Lebesgue measurable for  $y \in C$ . Thus measurability with respect to  $G$  concisely expresses the requirement of admissibility.

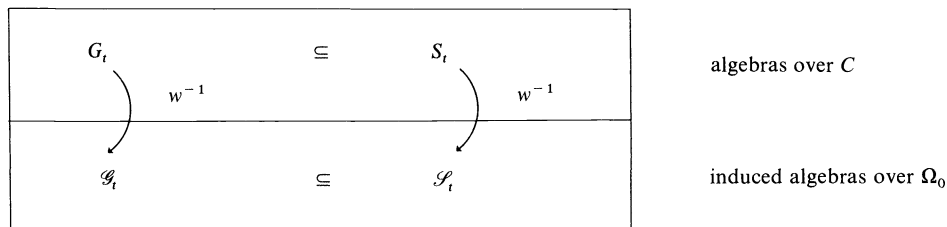
The preceding assumptions have concerned the information available for control and the system dynamics. The reader should note, though, that although  $\sigma$ -algebras over  $C$  were used, no probabilistic machinery has been introduced yet; this is now done.

We assume as given a probability space  $(\Omega, P, \mathcal{B})$  of points  $\omega \in \Omega$ , with  $P(A)$  the probability of a  $\mathcal{B}$ -set  $A$ . On this space is defined a measurable separable Brownian motion  $\{w(t, \omega), 0 \leq t \leq 1, \omega \in \Omega\}$  taking values in  $R^d$ , with continuous sample paths.

There is a set  $\Omega_0 \in \mathcal{B}$  of full measure such that  $w(\cdot, \omega) \in C$  for  $\omega \in \Omega_0$ . The process induces a measurable map  $w: \Omega_0 \rightarrow C$  according to the formula  $w(\omega) = w(\cdot, \omega)$ . This is seen as follows: every set  $\{y \in C: y(t) \in A\}$  for  $A$  Borel is in  $S_1$ ; every set  $\{\omega: w(t, \omega) \in A\}$  for  $A$  Borel is in  $\mathcal{B}$ ; but

$$\{\omega: w(t, \omega) \in A\} \cap \Omega_0 = w^{-1}\{y \in C: y(t) \in A\},$$

and so  $w^{-1}S_1 \subseteq \mathcal{B}$ . The classes  $w^{-1}G_t, t \in [0, 1]$ , and  $w^{-1}S_t, t \in [0, 1]$ , are all  $\sigma$ -algebras; they will provide us with a way of doing all our work in the probability space  $(\Omega, P, \mathcal{B})$  and then returning, for our control laws, to the space  $C$ . Setting  $\mathcal{G}_t = w^{-1}G_t$  and  $\mathcal{S}_t = w^{-1}S_t$ , we have the following diagram:



express what is available for control

express the total past

Let  $\varphi:[0, 1] \times \Omega_0 \rightarrow [0, 1] \times C$  according to the formula  $\varphi(t, \omega) = t, w(\omega)$ . Define  $\mathcal{F} = \varphi^{-1}G$ , where  $G$  is the  $\sigma$ -algebra expressing admissibility.  $\mathcal{F}$  is a  $\sigma$ -algebra of  $[0, 1] \times \Omega_0$  sets, and will be used to express the requirement of admissibility in terms of functions of  $t, \omega$  rather than  $t, y$  for  $y \in C$ . Indeed, we shall prove the existence of optimal control laws by first expressing them as  $\mathcal{F}$ -measurable  $t, \omega$  functions, and then (properly, now) as  $G$ -measurable  $t, y$  functions, using the following result.

LEMMA. If  $h:[0, 1] \times \Omega_0 \rightarrow \Gamma$  compact metric is  $\mathcal{F}$ -measurable, then there exists a  $G$ -measurable  $u:[0, 1] \times C \rightarrow \Gamma$  such that

$$u(t, w(\omega)) = h(t, \omega), \quad \omega \in \Omega_0.$$

This result is elementary (see [15, p. 185]).

The following convenient terminology is used:

The set  $\mathcal{A}$  of admissible drifts consists of all functions  $g:[0, 1] \times \Omega \rightarrow R^d$  of the form

$$g(t, \omega) = f(t, w(\omega), u(t, w(\omega))), \quad u \in \mathcal{U},$$

i.e., where  $u:[0, 1] \times C \rightarrow \Gamma$  is an admissible control law.

The set  $\mathcal{A}$  of admissible drifts consists of all functions  $g:[0, 1] \times \Omega \rightarrow R^d$  form

$$\zeta(\omega) = e^{\zeta(g)}, \quad g \in \mathcal{A},$$

i.e., where  $g$  is an admissible drift.

It is to be noted that admissible drifts are random processes and attainable densities are random variables.

$(\Omega, \mathcal{B}, P)$  is the probability space on which is defined a Wiener process (Brownian motion)  $\{w(t, \omega), \omega \in \Omega\}$  in  $d$ -dimensions. We propose to cut through all the questions of existence and uniqueness of solutions to (1), and at the same time obtain and use a representation of the criterion value achieved by a given admissible control law  $u$ , by taking, as "the solution" of (1), the process obtained by taking the functions  $w(\cdot, \omega)$  under the measure

$$d\tilde{P} = e^{\zeta(g)} dP, \quad g = f(t, w, u(t, w)),$$

where it is assumed (and for suitable  $f$ , proved) that  $\tilde{P}(\Omega) = 1$ , and where

$$\begin{aligned} \zeta(g)_\omega &= \int_0^1 f(t, w(\omega), u(t, w(\omega))) dw(t) \\ &\quad - \int_0^1 |f(t, w(\omega), u(t, w(\omega)))|^2 dt. \end{aligned}$$

This procedure provides a solution of (1) in the sense that under the transformed measure  $\tilde{P}$ ,

$$w(t, \omega) - \int_0^t f(s, w(\omega), u(s, w(\omega))) ds$$

is a Wiener process. If we define

$$W(t, \omega) = w(t, \omega) - \int_0^t f(s, w(\omega), u(s, w(\omega))) ds,$$

then  $W(\cdot, \omega)$  is a Wiener process under  $\tilde{P}$ . If in this context we give the new name  $x(t, \omega)$  to  $w(t, \omega)$ , then the definition of  $W(\cdot, \omega)$  tells us that

$$x(t, \omega) = \int_0^t f(s, x(s, \omega), u(s, x(s, \omega))) ds + W(t, \omega).$$

With measure  $\tilde{P}$  this asserts that  $x$  is the solution of an equation with drift  $f(s, x, u(s, x))$  and Brownian noise  $W$ , which is what we wanted to find. The correspondence between  $\Omega$  and the solution curves is the same as before, but now  $\tilde{P}$  gives their distribution. This procedure does not assign, to each sample function  $w(\cdot, \omega)$  of the original Wiener process, a “solution function”  $x(\cdot, \omega)$  constructed out of  $w(\omega)$ .

It is possible and useful to eliminate the dependence of the criterion on the control law through the cost function  $c(\cdot, \cdot, \cdot)$ . In the nonstochastic case this is done by adding one more differential equation

$$dx_0(t) = c(t, x, u(t, x)) dt$$

and minimizing  $x_0(1)$ . In our stochastic case we can make an analogous simplification, provided that  $c$  is nonanticipative, a reasonable condition. Heuristically: we replace the  $d$ -vector function  $f$  by the  $(1 + d)$ -vector function  $c, f$ ; we add another dimension  $w_0$  of Brownian motion, independent of  $w$ , to  $w$  to get a  $(1 + d)$ -dimensional motion

$$z = (w_0, w) = w_0, w_1, \dots, w_d;$$

and we swap the integral

$$\int_0^1 c dt$$

for the end-value  $w_0(1)$ . If now  $\xi = \xi_{c,f}$  is defined by

$$\xi = \int_0^1 h(t, z) dz(t) - \frac{1}{2} \int_0^1 |h(t, z)|^2 dt$$

with  $h(t, z) = c(t, w(\omega), u(t, w(\omega))), f(t, w(\omega), u(t, w(\omega)))$ , then the functions  $z(\cdot)$  under the measure  $\exp \xi dP$  will, if  $E e^\xi = 1$ , have the property that

$$z(t) - \int_0^t h(s, z) ds$$

forms a Wiener process of  $d + 1$  dimensions. Replacement of the time-integral by  $w_0(1)$  is justified by the fact that

$$E \int_0^1 c(t, w(\omega), u(t, w(\omega))) dt e^{\xi(f)} = E w_0(1) e^\xi.$$

To see this, note that

$$\xi = \zeta(f) + \int_0^1 c(t, w, u(t, w)) dw_0(t) - \frac{1}{2} \int_0^1 |c(t, w, u(t, w))|^2 dt$$

and that under the measure  $\exp \xi dP$  the function

$$w_0(t) - \int_0^t c(s, w, u(s, w)) ds$$

is a Wiener process in one dimension; hence,

$$\begin{aligned} 0 &= E \left( w_0(1) - \int_0^1 c(s, w, u(s, w)) ds \right) e^\xi \\ &= E w_0(1) e^\xi - E \int_0^1 c ds e^{\xi(f)} E \{ e^{\xi - \xi(f)} | \mathcal{S}_1 \}. \end{aligned}$$

If  $c(\cdot, \cdot, \cdot)$  satisfies growth conditions similar to those on  $f(\cdot, \cdot, \cdot)$ , and  $\mathcal{S}_1$  denotes the  $\sigma$ -algebra corresponding to the  $w$ -process over  $[0, 1]$ , the conditional expectation on the right will be 1.

Thus, we may and shall assume henceforth that the cost rate  $c(t, y, u)$  is simply the zero-index component  $f_0(t, y, u)$  of a new  $(1 + d)$ -dimensional "right-hand side"  $f(t, y, u)$ , with  $f: [0, 1] \times C \times \Gamma \rightarrow R^{1+d}$ .

Finally we discuss the question of a prescribed initial value  $x(0) = a$  for our solution  $x(\cdot)$  of

$$dx = f(t, x, u(t, x)) dt + dw.$$

As presently formulated this initial condition is  $x(0) = 0$ , since the Wiener process starts at 0, and we took for  $x(\cdot)$  the functions  $w(\cdot)$  under the measure  $e^{\xi(g)} dP$ . The initial condition  $x(0) = a \neq 0$  can be realized by taking instead the functions  $w(\cdot) + a$  under the same measure; if desired,  $g$  can be redefined so as to be a functional of  $w(\cdot) + a$  instead of  $w(\cdot)$ ; this is a simple shift. As will be shown, the important thing is that  $\exp \{ (w(t) + a) \cdot \theta - \theta \cdot \int_0^t g ds \}$  be a martingale with respect to  $e^{\xi(g)} dP$  for all  $\theta$ ; since this process differs from the previous one in having a factor  $\exp a \cdot \theta$ , the result is clear. In view of these facts we shall assume  $x(0) = 0$  henceforth.

We thus arrive at this formulation of an optimal control problem for stochastic functional-differential equations: To minimize the integral

$$\int w_0(1) e^{\xi(g)} dP = E w_0(1) e^{\xi(g)}$$

subject to the condition that  $g$  be an admissible drift, i.e., have the form

$$g(t, \omega) = f(t, w(\omega), u(t, w(\omega))),$$

where  $w(\cdot)$  is a Wiener process,  $u(\cdot, \cdot)$  is an admissible control, and

$$f_0(t, y, v) = c(t, y, v), \quad t, y, v \in [0, 1] \times C \times \Gamma,$$

i.e., the zero component of  $f$  is the cost rate  $c$ . In this form of the problem we minimize the average of the value of  $x_0(\cdot)$  at the endpoint 1; the functional  $\exp \xi(g)$  determines what this averaging is.

The functions  $z(t) = E \{ e^{\xi(\varphi)} | \mathcal{S}_t \}$  have many interesting properties, described by McKean [14]; among them is the fact that they are the unique solution of the



stochastic differential equation

$$(7) \quad dz(t) = z(t)\varphi(t) dw(t), \quad 0 \leqq t \leqq 1.$$

This allows us to rephrase the control problem as that of minimizing  $E x_0(1)z(1)$  subject to (7) and the same conditions on  $\varphi$  as before.

**4. Martingale proof of Girsanov's theorem.** We have proposed basing an approach to optimal stochastic control on an absolute continuity theorem of Girsanov [1, Theorem 1]. In this section we prove this theorem by a direct martingale method, and note its applicability under our assumptions. Girsanov proved that if  $\varphi$  is a nonanticipative Brownian functional with  $\int_0^1 |\varphi|^2 dt < \infty$  a.s., and if  $E \exp \zeta(\varphi) = 1$ , then the process

$$w(t) - \int_0^t \varphi ds$$

is a Wiener process under the measure  $\exp \zeta(\varphi) dP$ . We shall put his result into a more symmetric form.

**THEOREM 1.** *Let  $\varphi$  be a nonanticipative Brownian functional with  $\int_0^1 |\varphi|^2 dt < \infty$  a.s. The following are equivalent:*

- (i)  $w(t) - \int_0^t \varphi ds$  is a Wiener process under  $e^{\zeta(\varphi)} dP$ ;
- (ii)  $E e^{\zeta(\varphi + \theta)} = 1$  for every constant  $\theta \in R^d$ ;
- (iii)  $E e^{\zeta(\varphi)} = 1$ .

Proof of (i) implies (ii):  $d\tilde{P} = e^{\zeta(\varphi)} dP$  makes  $W(t) = w(t) - \int_0^t \varphi ds$  a Wiener process, so for  $\theta \in R^d$ ,

$$\tilde{E} \exp \{ \theta \cdot W(1) - \frac{1}{2} |\theta|^2 \} = 1.$$

But

$$\begin{aligned} \zeta(\varphi + \theta) &= \zeta(\varphi) + \zeta(\theta) - \int_0^1 \varphi \cdot \theta ds \\ &= \zeta(\varphi) + \theta \cdot w(1) - \frac{1}{2} |\theta|^2 - \theta \cdot \int_0^1 \varphi ds \\ &= \zeta(\varphi) + \theta \cdot W(1) - \frac{1}{2} |\theta|^2. \end{aligned}$$

So

$$E \exp \{ \zeta(\varphi + \theta) \} = \tilde{E} \{ \theta \cdot W(1) - \frac{1}{2} |\theta|^2 \} = 1.$$

Proof of (ii) implies (i): Let  $E e^{\zeta(\varphi)} = 1$ . Then since  $e^{\zeta_b(\varphi)}$  is a supermartingale [14, p. 25] we have a.s., for  $0 \leqq s \leqq t \leqq 1$ ,

$$\begin{aligned} E \{ e^{\zeta_b(\varphi)} | \mathcal{L}_s \} &\leqq e^{\zeta_b(\varphi)}, \\ E e^{\zeta_b(\varphi)} &= EE \{ e^{\zeta_b(\varphi)} | \mathcal{L}_s \} \leqq E e^{\zeta_b(\varphi)}. \end{aligned}$$

Thus  $E e^{\zeta_b(\varphi)} = 1$ . Further for  $s < t$ ,

$$0 = 1 - E e^{\zeta_b(\varphi)} = E e^{\zeta_b(\varphi)} E \{ e^{\zeta_s(\varphi)} - 1 | \mathcal{L}_s \}.$$

Since  $e^{\zeta_b(\varphi)} > 0$  a.s. we have

$$E \{ e^{\zeta_s(\varphi)} | \mathcal{L}_s \} = 1 \quad \text{a.s.}$$

In particular (ii) implies that for any  $\theta \in R^d$ ,

$$E\{e^{\zeta_s^t(\varphi+\theta)}|\mathcal{L}_s\} = 1 \quad \text{a.s.}$$

It follows from Doob [22, Theorem 11.9, p. 384] that if a continuous process  $y(\cdot)$  has the property that

$$e^{\theta \cdot y(t) - |\theta|^2 t/2}$$

is a martingale for every  $\theta \in R^d$ , then  $y(\cdot)$  is a Wiener process. Thus it is enough to show that

$$\exp\left\{\theta \cdot w(t) - \theta \cdot \int_0^t \varphi \, ds - \frac{1}{2}|\theta|^2 t\right\}$$

is a martingale with respect to  $e^{\zeta(\varphi)} dP$ , i.e., that for  $0 \leq s \leq t \leq 1$ ,

$$\tilde{E}\left\{\exp\left\{\theta \cdot [w(t) - w(s)] - \theta \cdot \int_s^t \varphi \, du - |\theta|^2(t - s)/2\right\}\middle|\mathcal{L}_s\right\} = 1 \quad \text{a.s.}$$

The conditional expectation above is an  $\mathcal{L}_s$ -measurable function  $\psi(\omega)$ ,  $\omega \in \Omega$ , such that for  $B \in \mathcal{L}_s$ ,

$$\begin{aligned} \int_B \psi(\omega) e^{\zeta(\varphi)} dP &= \int_B \exp\left\{\theta \cdot [w(t) - w(s)] - \theta \cdot \int_s^t \varphi \, du - \frac{1}{2}|\theta|^2(t - s)\right\} e^{\zeta(\varphi)} dP \\ &= \int_B e^{\zeta_t^t(\varphi) + \zeta_s^s(\varphi+\theta) + \zeta_\delta^\delta(\varphi)} dP \\ &= \int_B E\{e^{\zeta_t^t(\varphi)}|\mathcal{L}_t\} e^{\zeta_s^s(\varphi+\theta) + \zeta_\delta^\delta(\varphi)} dP \\ &= \int_B E\{e^{\zeta_s^s(\varphi+\theta)}|\mathcal{L}_s\} e^{\zeta_\delta^\delta(\varphi)} dP \\ &= \int_B e^{\zeta_\delta^\delta(\varphi)} dP \\ &= \int_B e^{\zeta_\delta^\delta(\varphi)} dP = \tilde{P}(B), \end{aligned}$$

because the conditional expectations above all equal 1 a.e. It follows that  $\psi = 1$  a.e.

To encompass Girsanov's theorem and complete the proof of Theorem 1 we must show that (iii) implies (ii). Accordingly let  $E \exp \zeta(\varphi) = 1$ , and let  $\varphi_N$  bounded be such that

$$\int_0^1 |\varphi_N - \varphi|^2 dt \rightarrow 0 \quad \text{a.e.}, \quad N \rightarrow \infty.$$

Then also

$$\int_0^1 |\varphi_N + \theta - \varphi - \theta|^2 dt \quad \text{a.e.},$$

and thus  $\zeta(\varphi_N + \theta) \rightarrow \zeta(\varphi + \theta)$  in probability. Together with the hypothesis  $E \exp \zeta(\varphi) = 1$ , this implies by Girsanov's elementary Lemma 5 [1] that

$$e^{\zeta(\varphi_N + \theta)} \rightarrow e^{\zeta(\varphi + \theta)} \quad \text{in } L_1.$$

Each  $\varphi_N + \theta$  is bounded, so

$$1 = E e^{\zeta(\varphi_N + \theta)} \rightarrow E e^{\zeta(\varphi + \theta)}, \quad N \rightarrow \infty.$$

This completes the proof of Theorem 1.

Girsanov's theorem can be viewed as stating that the Wiener measure  $\mu$  on  $C$  is a fixed point of the composition of two transformations of measures, the first corresponding to weighting by  $e^{\zeta(\varphi)}$ , the second to translation by  $-\int_0 \varphi ds$ . Since  $e^{\zeta(\varphi)}$  is  $\mathcal{S}_1$ -measurable and integrable, it corresponds to an  $S_1$ -measurable  $\mu$ -integrable functional  $\eta$  on  $C$  according to the formula  $\int_A \eta d\mu = \int_{w^{-1}(A)} e^{\zeta(\varphi)} dP$ ,  $A \in S_1$ . For  $y \in C$  set  $Ty(t) = y(t) - \int_0^t \varphi(s, y) ds$ . For measures  $\nu$  on  $C$  define  $W\nu$  and  $\nu T^{-1}$  for  $A \in S_1$  by

$$\begin{aligned} W\nu(A) &= \int_A \eta d\nu, \\ \nu T^{-1}(A) &= \nu(T^{-1}A). \end{aligned}$$

Then Girsanov's theorem states that  $\mu = (W\mu)T^{-1}$ .

**5. Some preliminary results.** In view of the role of condition  $E \exp \zeta(\varphi) = 1$  in application of Girsanov's theorem, it is particularly important to clarify the conditions under which it obtains for attainable densities  $\exp \zeta(g)$ ,  $g \in \mathcal{A}$ . We prove the following lemma.

LEMMA 0.  $E \exp \zeta(\alpha g) = 1$  for all  $g \in \mathcal{A}$  and  $\alpha \geq 0$ .

*Proof.* The proof to be given is the same for all  $\alpha \geq 0$ , so only the case  $\alpha = 1$  is described. Let  $g \in \mathcal{A}$ , and

$$y(t, \omega) = w(t, \omega) - \int_0^t g(s, \omega) ds$$

so that in Girsanov's terminology [1]  $y(\cdot, \omega)$  is an Itô process with respect to  $w(\cdot, \omega)$  corresponding to the matrix  $I$  and the vector  $g(\cdot, \omega)$ . The result will follow from Girsanov's Lemma 7 provided that we find for each  $\varepsilon > 0$ , an integer  $N(\varepsilon)$  and a monotone system of open sets  $C_N(t)$  with

$$C_N(t) \in S_t, \quad C_N(s) \subseteq C_N(t) \quad \text{for } t > s$$

and such that:

- (a)  $\{\omega : y(\cdot, \omega) \in C_N(t)\} \in \mathcal{S}_t$ ,
- (b)  $P\{w(\cdot, \omega) \in C_N(1)\} > 1 - \varepsilon$ ,
- (c)  $|g(t, \omega)| < N$  if  $y(\cdot, \omega) \in C_N(t)$ ,
- (d) if  $x(\cdot) \in C_N(s)$ ,  $x(\cdot) \notin C_N(t)$ ,  $t > s$ , then there is a  $\tau$  with  $s < \tau < t$  such that

$$x(\cdot) \in C_N(u) \quad \text{for } u < \tau, \quad x(\cdot) \notin C_N(\tau).$$

Let

$$C_N(t) = \{x \in C : 1 + 2e^2(\kappa + \sup_{s \in [0, t]} |x(s)|^2) < N^2/\kappa\}.$$

It is easy to verify that properties (a), (b) and (d) obtain. To prove (c) consider that the growth condition on  $f$  gives

$$\begin{aligned}
 |w(t)|^2 &\leq 2|y(t)|^2 + 2\left(1 + \int_0^t |w(s)|^2 ds\right), \\
 \sup_{s \in [0,t]} |w(s)|^2 &\leq 2 \sup_{s \in [0,t]} |y(s)|^2 + 2\kappa \left(1 + \int_0^t \sup_{u \in [0,s]} |w(u)|^2 ds\right) \\
 &\leq 2 e^{2\kappa} \left(\kappa + \sup_{s \in [0,t]} |y(s)|^2\right)
 \end{aligned}$$

by Gronwall's inequality. Hence if  $y(\cdot, \omega) \in C_N(t)$ , then

$$\begin{aligned}
 \kappa(1 + |w(t)|^2) &< N^2, \\
 |g(t, \omega)| &< N,
 \end{aligned}$$

which proves (c).

LEMMA 1. *There exists a constant  $\alpha > 1$  such that*

$$\sup_{g \in \mathcal{A}} E e^{\alpha \zeta(g)} < \infty.$$

*Proof.* Consider that

$$\begin{aligned}
 e^{\alpha \zeta(g)} &= \exp \left\{ \zeta(\alpha g) + \frac{\alpha^2 - \alpha}{2} \int_0^1 |g(t, \omega)|^2 dt \right\} \\
 &\leq \exp \left\{ \zeta(\alpha g) + \frac{\alpha^2 - \alpha}{2} \kappa \int_0^1 (1 + |w(t)|^2) dt \right\}.
 \end{aligned}$$

Since  $E e^{\zeta(\alpha g)} = 1$ , the density  $e^{\zeta(\alpha g)}$  makes the function

$$x(t) = w(t) - \alpha \int_0^t g ds$$

a Wiener process. Also

$$|w(t)|^2 \leq 2|x(t)|^2 + 2\alpha^2 \kappa \left(1 + \int_0^t |w(u)|^2 du\right),$$

for which Gronwall's inequality gives

$$|w(t)|^2 \leq 2 \left( \alpha^2 \kappa + \sup_{0 \leq t \leq 1} |x(t)|^2 \right) \exp 2\alpha^2 \kappa.$$

It follows that

$$E e^{\alpha \zeta(g)} \leq e^{(\alpha^2 - \alpha)\kappa/2} E \exp \left\{ \zeta(\alpha g) + \kappa(\alpha^2 - \alpha) \left( \alpha^2 \kappa + \sup_{0 \leq t \leq 1} |x(t)|^2 \right) \exp 2\alpha^2 \kappa \right\}.$$

Since  $x(\cdot)$  is a Wiener process under  $\exp \zeta(\alpha g) dP$ , the expectation on the right is of the form

$$h(\alpha) E \exp \left\{ \kappa(\alpha^2 - \alpha) e^{2\alpha^2 \kappa} \sup_{0 \leq t \leq 1} |w(t)|^2 \right\},$$

where  $h(\cdot)$  is bounded near 1, and so is finite for  $\alpha > 1$  small enough (Doob [loc. cit., p. 392]).

Lemma 1 implies that the set  $\mathcal{D} = \{e^{\zeta(g)}, g \in \mathcal{A}\}$  of attainable densities is uniformly integrable.

LEMMA 2. *If  $g_n, n \geq 1$ , is a sequence of admissible drifts, and*

$$\xi = \text{weak lim}_n \exp \zeta(g_n)$$

*in  $L_1$  or  $L_2$ , then*

$$E\{\xi | \mathcal{S}_t\} > 0 \quad \text{a.s.}, \quad 0 \leq t \leq 1.$$

*Proof.* If  $E\{\xi | \mathcal{S}_t\} = 0$  on a set  $A \in \mathcal{S}_t$ , of positive measure, then

$$\begin{aligned} E\xi(t, \omega)\chi_A &= EE\{\xi | \mathcal{S}_t\}\chi_A \\ &= EE\{\xi\chi_A | \mathcal{S}_t\} \\ &= E\xi\chi_A = 0. \end{aligned}$$

Since  $\xi \geq 0$  a.e. we must have  $\xi = 0$  a.e. on  $A$ . By weak convergence  $\int_A e^{\zeta_0^1(g_n)} dP \rightarrow 0$ . Hence some subsequence, relabeled  $\{e^{\zeta(g_n)}, n \geq 1\}$ , converges to 0 a.e. on  $A$ , and so  $\zeta_0^1(g_n) \rightarrow -\infty$  a.e. on  $A$ . But

$$\zeta_0^1(g_n) = \int_0^1 g_n dw - \frac{1}{2} \int_0^1 |g_n|^2 dt.$$

The second term is bounded by  $\kappa \int_0^1 [1 + |w(s)|^2] ds$ , which is independent of  $n$  and finite a.e., so  $\int_0^1 g_n dw \rightarrow -\infty$  a.e. on  $A$ . Thus

$$P\left\{A \cap \int_0^1 g_n dw > -N\right\} \rightarrow 0 \quad \text{for every } N > 0.$$

Also though,

$$\begin{aligned} P\left\{A \cap \int_0^1 g_n dw \leq -N\right\} &\leq P\left\{\left|\int_0^1 g_n dw\right| \geq N\right\} \\ &\leq N^{-2} E \int_0^1 |g_n|^2 dt, \\ &\leq o(1), \end{aligned}$$

uniformly in  $n$ .

Now observe that

$$P\{A\} = P\left\{A \cap \int_0^1 g_n dw > -N\right\} + P\left\{A \cap \int_0^1 g_n dw \leq -N\right\}.$$

Let  $\varepsilon > 0$  be given. Pick first  $N$  so large that the second term is less than  $\varepsilon/2$  uniformly in  $n$ ; then pick  $n = n(N)$  so large that the first term is less than  $\varepsilon/2$ . Thus  $P\{A\} = 0$ .

LEMMA 3 (H. P. McKean, Jr.). If  $y(t) = \int_0^t \varphi(u, \omega) dw(u)$  and

$$T(t) = \int_0^t |\varphi(u, \omega)|^2 du,$$

$$T^{-1}(x) = \inf t : T(t) = x \quad \text{for } x \leq \int_0^1 |\varphi|^2 du,$$

then

$$y(T^{-1}(t))$$

is a Wiener process terminated at  $T(1)$ .

*Proof.* Let  $\chi_j(t)$  be the indicator of  $T^{-1}(t_j) > t$ , where  $t_1, t_2, \dots, t_n \in [0, \infty)$  and let

$$z(t) = \exp \left\{ i \sum_{j=1}^n \gamma_j \int_0^t \chi_j \varphi dw + \frac{1}{2} \sum_{j,k=1}^n \gamma_j \gamma_k \int_0^t \chi_j \chi_k |\varphi|^2 du \right\}.$$

Note that on  $T(1) > \max t_j$ ,

$$\begin{aligned} \int_0^1 \chi_j \chi_k |\varphi|^2 du &= \int_0^{\min(T^{-1}(t_j), T^{-1}(t_k))} |\varphi|^2 du \\ &= \min(t_j, t_k) \end{aligned}$$

and

$$\int_0^1 \chi_j \varphi dw = y(T^{-1}(t_j)).$$

The stochastic differential of  $z(\cdot)$  is

$$dz(t) = iz(t) \sum_{j=1}^n \gamma_j \chi_j \varphi dw(t)$$

so that

$$z(t) = 1 + i \int_0^t z(u) \sum_{j=1}^n \gamma_j \chi_j(u) \varphi(u) dw(u).$$

Since

$$|z(t)| \leq \exp \frac{1}{2} \sum_{j,k=1}^n |\gamma_j \gamma_k| t_j t_k,$$

we can multiply by  $\chi_{T(1) > \max t_i}$  and take expectations at  $t = 1$  to get

$$Ez(1)\chi_{T(1) > \max t_i} = \Pr \{T(1) > \max t_i\}$$

or

$$E \exp \left\{ i \sum_{j=1}^n \gamma_j y(T^{-1}(t_j)) \mid T(1) > \max t_i \right\} = \exp \frac{1}{2} \sum_{j,k=1}^n \gamma_j \gamma_k \min(t_j, t_k).$$

This says that conditional on its still being defined at the maximum time  $\max t_i$ ,

the process  $y(T^{-1}(t))$  is Gaussian with a covariance appropriate to Brownian motion.

LEMMA 4. For  $w(\cdot)$  a Wiener process, and  $y(\cdot)$  and  $T(\cdot)$  as in Lemma 3,

$$\Pr \left\{ \sup_{0 \leq t \leq 1} |y(t)| > a \right\} \leq \Pr \left\{ \sup_{0 \leq s \leq t} |w(s)| > a \right\} + \Pr \{T(1) > t\}.$$

*Proof.* Let  $t \in [0, 1]$  so  $T(t)$  is defined and  $T^{-1}(T(t)) \leq t$ . Then

$$y(T^{-1}(T(t))) = y(u) \quad \text{for some } u \in [0, t]$$

and so

$$\sup_{0 \leq t \leq 1} |y(t)| > a \Rightarrow \sup_{0 \leq u \leq 1} |y(T^{-1}(T(u)))| > a.$$

Thus, since  $y(T^{-1}(\cdot))$  is a Wiener process as long as it is defined,

$$\begin{aligned} \Pr \left\{ \sup_{0 \leq t \leq 1} |y(t)| > a \right\} &\leq \Pr \left\{ \sup_{0 \leq u \leq 1} |y(T^{-1}(T(u)))| > a \right\} \\ &\leq \Pr \left\{ \sup_{0 \leq s \leq T(1)} |y(T^{-1}(s))| > a \right\} \\ &\leq \Pr \left\{ \sup_{0 \leq s \leq t} |w(s)| > a \right\} + \Pr \{T(1) > t\}. \end{aligned}$$

We shall also need a version [15] of an implicit function lemma due to McShane and Warfield [5]; this version allows simultaneous explicit as well as implicit dependence on the independent variable, provided that this dependence is measurable with respect to the same  $\sigma$ -algebra as is the desired function.

If  $\mathcal{M}$  is a  $\sigma$ -algebra of subsets of a set  $M$ , and  $S$  is a topological space, we say that a function  $g: N \rightarrow S$ ,  $N \in \mathcal{M}$  (defined on  $N$ ) is  $\mathcal{M}$ -measurable if and only if  $g^{-1}(F) \in \mathcal{M}$  for closed  $F \subseteq S$ .

LEMMA 5. Let  $(M, \mathcal{M})$  be a measure space,  $A$  a separable metric space, and  $U$  a compact metric space. Let  $k: M \times U \rightarrow A$  be continuous in its second argument for each value of the first, and  $\mathcal{M}$ -measurable in the first for each value of the second. Let  $y: M \rightarrow A$  be  $\mathcal{M}$ -measurable, with

$$y(x) \in k(x, U), \quad x \in M.$$

Then there exists an  $\mathcal{M}$ -measurable  $u: M \rightarrow U$  such that

$$y(x) = k(x, u(x)).$$

From Lemma 5 we can obtain an implicit function lemma for limits.

LEMMA 6. Let  $(M, \mathcal{M})$ ,  $A$  and  $U$  be as in Lemma 5. Let  $h: M \times U \rightarrow A$  be continuous in its second argument for each value of the first, and let  $u_n: M \rightarrow U$  be a sequence of  $\mathcal{M}$ -measurable functions, and  $z$  a function such that pointwise

$$h(x, u_n(x)) \rightarrow z(x).$$

Then there exists an  $\mathcal{M}$ -measurable function  $u: M \rightarrow U$  such that

$$z(x) = h(x, u(x)).$$

*Proof.* Since  $U$  is compact metric, there is a continuous map  $\psi$  of the Cantor set  $C$  onto  $U$ . Find, by Lemma 5,  $\mathcal{M}$ -measurable functions  $\xi_n: M \rightarrow C$  such that  $u_n(x) = \psi(\xi_n(x))$ ,  $x \in M$ . Define  $\zeta(x) = \limsup_{n \rightarrow \infty} \xi_n(x)$ ,  $u(x) = \psi(\zeta(x))$ . Fix  $x$

and find  $n_i = n_i(x)$ ,  $i = 1, 2, \dots$ , such that

$$\xi_{n_i(x)}(x) \rightarrow \xi(x).$$

Then  $h(x, \psi(\xi_{n_i(x)}(x))) \rightarrow h(x, \psi(\xi(x)))$ . Thus  $z(x) = h(x, \psi(\xi(x)))$ .  $u$  is  $\mathcal{M}$ -measurable because if  $0$  is an open set of  $U$ , then

$$u^{-1}(0) = \{x : \xi(x) \in \psi^{-1}(0)\} \in \mathcal{M},$$

since  $\psi^{-1}(0)$  is open in  $C$ , because  $\psi$  is continuous.

**6. Convexity.** In deterministic control theory there are counterexamples which indicate that restrictions must be placed on the “right-hand side” of the constraining differential equation if optimal control laws are to exist. Conditions of linearity or convexity have been imposed for this purpose. The fundamental paper [16] of Markus and Lee postulated a form linear in the control  $u(\cdot)$ :

$$(8) \quad \dot{x}(t) = f(t, x(t), u(t)) = g(t, x(t)) + H(t, x(t))u(t),$$

with  $H$  a matrix and  $u$  a vector. Roxin [17] suggested that it was enough to assume convexity of  $f(t, y, \Gamma)$  for each  $t, y$ , with  $\Gamma$  the space of control points. These conditions (of linearity or convexity) were used to show that a certain function obtained as a weak limit by a compactness argument was indeed an admissible “right-hand side”, i.e., was  $f(\cdot, x(\cdot), u(\cdot))$  for some measurable  $u(\cdot)$  with values in  $\Gamma$ .

A similar, and in some respects a worse, situation holds for stochastic optimal control. In the deterministic case an admissible control has only to be measurable and to take values in the right set. As has been noted by Fleming [18, p. 79], in the stochastic case the concept of admissibility is much more complicated. This is because now control can depend with advantage on available information. As a result, difficulties arise in showing that functions obtained as weak limits are indeed of the desired form, e.g., admissible drifts ( $g \in \mathcal{A}$ ) or attainable densities ( $\xi \in \mathcal{D}$ ).

In our setup, the important and difficult problem seems to be that there are drawbacks to using each of the natural ways in (or places at) which to take convex combinations in order to turn weak convergence into strong by the Banach-Saks theorem. One can take them in  $\mathcal{A}$ , and try to show under convexity of  $\mathcal{A}$  that from a minimizing sequence of admissible drifts one can obtain a convergent minimizing sequence by convexifying. Or one can take them in  $\mathcal{D}$ , provided  $\mathcal{D}$  is convex, and then prove that  $\mathcal{D}$  is closed. It turns out that  $\mathcal{A}$  need not be convex even under Roxin’s convexity condition, although special Roxin-type conditions suffice for certain forms of  $f$  in (8). And when  $\mathcal{A}$  is convex, the property of being a minimizing sequence is not known to transfer from  $\{g_n\} \subseteq \mathcal{A}$  to a sequence of convex combination of  $\{g_n\}$ . Further, we have managed to prove  $\mathcal{D}$  convex only in the case of complete information about the past ( $G_t = S_t$ ); counterexamples suggest that  $\mathcal{D}$  is rarely convex. (Some of these results and counterexamples are the meat of this section.) The trouble then is that it is difficult to show either that convex combinations of minimizing admissible drifts are minimizing, or that convex combinations of minimizing attainable densities are attainable; the second alternative is feasible when  $G_t = S_t$ , and most of our results concern this case.



We have described the structure of the available information by saying that for each time  $t$  there is a  $\sigma$ -algebra  $G_t$  on  $C$  representing information available at  $t$ , and that the control law at  $t$  must be measurable on  $G_t$ . However, the information in  $G_t$  may be very different from (usually it is much smaller than) that on which the system depends at  $t$ . In other words for some  $u \in \Gamma$ ,  $f(t, \cdot, u)$  may be nowhere near being  $G_t$ -measurable. As a result  $\mathcal{A}$  may not be convex. In particular, Roxin's condition does not ensure that the set of admissible drifts is convex unless the system depends on no more than the controller knows, a situation atypical in practice.

We give two examples of the failure of Roxin's condition to guarantee convexity of  $\mathcal{A}$ . These examples are included mostly because they illustrate how the pattern of available information affects convexity, not because convexity of  $\mathcal{A}$  is useful. Indeed the properties of the functional  $\exp \zeta(\cdot)$  have precluded our finding a role for convexity of  $\mathcal{A}$  so far.

*Example 1.* Let  $\Gamma = [0, 1]$ ,  $f(t, y, u) = \exp \{uy(t)\}$  for  $y(\cdot) \in C$  scalar, and suppose that  $G_t$  is the four element algebra  $\{\emptyset, A_t, A_t^c, C\}$ , where  $A_t$  is some  $S_t$ -measurable set. This situation corresponds to knowing at  $t$  only whether  $y(\cdot) \in A_t$  or not. We have

$$f(t, y, \Gamma) = [1, e^{y(t)}],$$

a closed convex set, so Roxin's condition is satisfied. Let  $u_1$  and  $u_2$  be two admissible control laws. Clearly  $u_1(t, \cdot)$  and  $u_2(t, \cdot)$  are constant on  $A_t$ . Suppose now that  $u_0$  is an admissible control law such that

$$f(t, y, u_0(t, y)) = \lambda f(t, y, u_1(t, y)) + (1 - \lambda)f(t, y, u_2(t, y)).$$

Then, on  $A_t$ , dropping arguments on the  $u$ 's, we have

$$e^{u_0 y(t)} = \lambda e^{u_1 y(t)} + (1 - \lambda) e^{u_2 y(t)},$$

$$u_0 = \frac{1}{y(t)} \log \{ \lambda e^{u_1 y(t)} + (1 - \lambda) e^{u_2 y(t)} \}.$$

But  $u_0$  must be constant on  $A_t$ , while the right-hand side above obviously depends on what function  $y(\cdot)$  from  $A_t$  one has. Thus if  $A_t$  contains two functions assuming different values at time  $t$ , then  $u_0$  cannot be constant on  $A_t$ , and so  $\mathcal{A}$  cannot be convex.

*Example 2.* The drift in Example 1 is not at most linear in growth unless  $u \equiv 0$ . To get one that is, let  $\Gamma$  and  $G_t$  be as before, but take

$$f(t, y, u) = (u + y^2(t))^{1/2}.$$

Then  $f(t, y, \Gamma) = [|y(t)|, (1 + y^2(t))^{1/2}]$ , a convex set, so Roxin's condition is satisfied. Let now  $u_0, u_1$ , and  $u_2$  be admissible control laws such that

$$(u_0 + y^2(t))^{1/2} = \lambda(u_1 + y^2(t))^{1/2} + (1 - \lambda)(u_2 + y^2(t))^{1/2},$$

with arguments on the  $u$ 's omitted. Then

$$u_0 = (\lambda(u_1 + y^2(t))^{1/2} + (1 - \lambda)(u_2 + y^2(t))^{1/2})^2 - y^2(t)$$

must be independent of  $y(\cdot)$  for  $y(\cdot)$  in  $A_t$ . This is impossible unless  $y(t)$  is the same number for all  $y \in A_t$ .

When the system depends on no more than the controller knows, then Roxin's condition implies the convexity of  $\mathcal{A}$ . Such is the content of the following theorem.

**THEOREM 2.** *If for each  $t, u$   $f(t, \cdot, u)$  is  $G_t$ -measurable and if for each  $t, y$   $f(t, y, \Gamma)$  is convex, then  $\mathcal{A}$  is convex.*

*Proof.*  $f(\cdot, \cdot, u)$  is  $G$ -measurable for each  $u$ , and  $f(t, y, \cdot)$  is continuous for each  $t, y$ . By Roxin's condition, with  $u_1$  and  $u_2$  admissible,

$$\lambda f(t, y, u_1(t, y)) + (1 - \lambda)f(t, y, u_2(t, y)) \in f(t, y, \Gamma).$$

Then the result follows from the implicit function Lemma 5.

We turn now to consider the convexity of the set  $\mathcal{D}$  of attainable densities. The principal result is the following theorem.

**THEOREM 3.** *If  $G_t = S_t$ , i.e., if the whole past is known, and if  $f(t, y, \Gamma)$  is convex for  $t, y \in [0, 1] \times C$ , then  $\mathcal{D}$  is convex.*

*Proof.* For  $e^{\zeta(g_i)} \in \mathcal{D}$ ,  $a_i \geq 0$ ,  $i = 1, \dots, n$ , and  $\sum_{i=1}^n a_i = 1$ , we are to find an admissible drift  $g$  such that

$$(9) \quad e^{\zeta(g)} = \sum_{i=1}^n a_i e^{\zeta(g_i)}.$$

Consider the process

$$g(t, \omega) = \frac{\sum_{i=1}^n a_i e^{\zeta_b(g_i)} g_i(t, \omega)}{\sum_{i=1}^n a_i e^{\zeta_b(g_i)}}.$$

We note that  $g(t, \cdot)$  is  $\mathcal{L}_t$ -measurable, and that  $g(\cdot, \omega)$  is Lebesgue measurable for almost all  $\omega$ . Thus  $g(\cdot, \cdot)$  differs from an  $\mathcal{F}$ -measurable function at most on a null set. By Roxin's condition, with  $\lambda =$  Lebesgue measure,

$$(10) \quad g(t, \omega) \in f(t, w(\omega), \Gamma) \quad \text{a.e.} \quad [\lambda \times P].$$

By changing  $g(t, \omega)$  on a set of measure zero, we can induce it to satisfy (10) everywhere. Thus by Lemma 5, there is an  $\mathcal{F}$ -measurable function  $\gamma: [0, 1] \times \Omega \rightarrow \Gamma$  such that

$$g(t, \omega) = f(t, w(\omega), \gamma(t, \omega)),$$

and further [15], there is a  $G$ -measurable function  $u: [0, 1] \times C \rightarrow \Gamma$  such that a.e.

$$u(t, w(\omega)) = \gamma(t, \omega).$$

Hence  $g(t, \omega) = f(t, w(\omega), u(t, w(\omega)))$ , which shows that  $g$  is an admissible drift,  $g \in \mathcal{A}$ . To see that (9) holds we note that the stochastic differential of the ratio

$$e^{\zeta_b(g)} / \sum_{i=1}^n a_i e^{\zeta_b(g_i)}$$

is zero.

**7. When is  $\mathcal{D}$  a subset of  $L_2$ ?** Let  $L_2 = L_2(\Omega, P, \mathcal{S}_1)$  be the real Hilbert space of functions measurable on  $\mathcal{S}_1 = w^{-1}(S_1)$  and square-integrable with respect to  $P$ . When  $\mathcal{D} \subseteq L_2$ ,  $L_2$  provides a particularly convenient topology for  $\mathcal{D}$ . Unfortunately it turns out that  $\mathcal{D}$  is demonstrably a subset of  $L_2$  only when the “right-hand side” function  $f(t, y, u)$  in the equation of interest,

$$dx = f(t, x, u(t, x)) dt + dw,$$

increases with  $y$  either slower than linearly or linearly at a slow enough rate. While these conditions cover many cases of interest, they do not cover the general case of at most linear growth. Further, there are choices of  $f$  that are linear in  $y$  for which  $\mathcal{D} \not\subseteq L_2$ . In this section we prove some of these facts.

**LEMMA 7.** *If  $|f(t, y, u)|^2 \leq \kappa(1 + |y(t)|^{2\alpha})$  for some  $\alpha < 1$ , then  $\mathcal{D}$  is a bounded set of  $L_2$ .*

*Proof.* Take  $g \in \mathcal{A}$ ,  $T = \int_0^1 |g|^2 dt$ . Then since  $2\alpha < 2$ ,

$$\begin{aligned} T > t &\Rightarrow \int_0^1 |w(t)|^{2\alpha} dt > \frac{t}{\kappa} - 1 \\ &\Rightarrow \left( \int_0^1 |w(t)|^2 dt \right)^{1/2\alpha} > \left( \frac{t}{\kappa} - 1 \right)^{1/2\alpha} \\ &\Rightarrow \left( \int_0^1 |w(t)|^2 dt \right)^{1/2} > \left( \frac{t}{\kappa} - 1 \right)^{1/2\alpha}. \end{aligned}$$

Thus for any number  $0 < \lambda < \pi^2/8$  ( $\pi^2/8$  is the abscissa of convergence of  $E \exp \lambda \|w\|^2$ ) [19],

$$\begin{aligned} P\{T > t\} &\leq P\left\{ \|w\|^2 > \left( \frac{t}{\kappa} - 1 \right)^{1/\alpha} \right\} && \left( \|w\|^2 = \int_0^1 |w|^2 dt \right) \\ &\leq e^{-\lambda(t/\kappa - 1)^{1/\alpha}} E \exp \lambda \|w\|^2. \end{aligned}$$

Lemma 4 now gives

$$\begin{aligned} P\{e^{\xi(g)} > a\} &\leq P\left\{ \int_0^1 g dw > \log a \right\} \\ &\leq \frac{1}{\log a} \sqrt{\frac{2t}{\pi}} e^{-(\log^2 a)/2t} + \text{const.} e^{-\lambda(t/\kappa - 1)^{1/\alpha}} \end{aligned}$$

Let  $t = \log a/(4 + 2\varepsilon)$ , so that

$$P\{e^{\xi(g)} > a\} \leq \sqrt{\frac{2 + \varepsilon}{\pi \log a}} e^{-(2 + \varepsilon)\log a} + \text{const.} e^{-(t\lambda/\kappa)(t/\kappa - 1)^{1/\alpha - 1}}.$$

Since  $(t/\kappa - 1)^{1/\alpha - 1}$  is eventually larger than  $\kappa/\lambda$ , the result follows.

To illustrate how  $e^{\xi(g)}$  can fail to belong to  $L_2$  consider the scalar case ( $d = 1$ )  $f(t, y, u) = \kappa y(t) + u$ , take  $\Gamma = [0, 1]$  together with the control law that

is identically 0, so that  $g(t, \omega) = \kappa w(t)$ , and

$$\begin{aligned} \zeta(g) &= \kappa \int_0^1 w(t) dw(t) - \frac{\kappa^2}{2} \int_0^1 |w(t)|^2 dt \\ &= \frac{\kappa}{2} w^2(1) - \frac{\kappa}{2} - \frac{\kappa^2}{2} \int_0^1 |w(t)|^2 dt. \end{aligned}$$

If  $\kappa$  is negative, i.e., if feedback is negative, then  $e^{\zeta(g)}$  is actually bounded by  $\exp -\frac{1}{2}\kappa$ . In general, we can use a result [19] of L. A. Shepp to find for what values of  $\kappa$   $e^{\zeta(g)}$  has what moments. He showed that for  $m(\cdot)$  a measure on  $[0, 1]$ ,

$$E \exp\left(-\frac{1}{2} \int_0^1 w^2(t) dm(t)\right) < \infty$$

if and only if the integral equation

$$g(t) = 1 - \int_t^1 (t - u)g(u) dm(u)$$

has a solution positive in  $[0, 1]$ ; if it does, the value of the expectation is  $g(0)^{-1/2}$ .

Choosing  $dm(u) = \lambda\kappa^2 du - \lambda\kappa\delta(u - 1) du$ , we have

$$\begin{aligned} g(t) &= 1 - \lambda\kappa \int_t^1 (t - u)g(u) du + \lambda\kappa \int_t^1 (t - u)g(u)\delta(u - 1) du \\ &= 1 - \lambda\kappa^2 \int_t^1 (t - u)g(u) du + \lambda\kappa(t - 1)g(1), \\ g'(t) &= -\lambda\kappa^2 \int_t^1 g(u) du + \alpha\kappa g(1), & g'(1) &= \lambda\kappa g(1), \\ g''(t) &= \lambda\kappa^2 g(t), & g(1) &= 1. \end{aligned}$$

We find

$$2g(t) = (1 + \sqrt{\lambda}) e^{\kappa\sqrt{\lambda}(t-1)} + (1 - \sqrt{\lambda}) e^{-\kappa\sqrt{\lambda}(t-1)}, \quad 0 \leq t \leq 1.$$

This has a zero at  $t$  if and only if

$$\begin{aligned} e^{2\kappa\sqrt{\lambda}t} &= \frac{\sqrt{\lambda} - 1}{\sqrt{\lambda} + 1} e^{2\kappa\sqrt{\lambda}}, \\ t &= 1 + \frac{1}{2\kappa\sqrt{\lambda}} \log \left( \frac{\sqrt{\lambda} - 1}{\sqrt{\lambda} + 1} \right). \end{aligned}$$

Thus the function on the right is out of the range  $[0, 1]$  if and only if the expectation exists. In particular, if  $\lambda = 2$  and

$$\kappa > -\frac{1}{2\sqrt{2}} \log \frac{\sqrt{2} - 1}{\sqrt{2} + 1},$$

then  $E e^{2\zeta(g)} = +\infty$ .

**8. Closure and existence.** We come finally to closure and existence theorems similar to those of deterministic control theory. When  $\mathcal{D}$  is a bounded subset of  $L_2$ , these results are proved in a natural way using strong and weak  $L_2$ -topologies. When  $\mathcal{D}$  is known only to be a subset of  $L_p$ , they are proved in a more laborious way by using stopping times  $\tau_N \rightarrow 1$  as  $N \rightarrow \infty$ , for which

$$\exp \zeta_0^{\tau_N}(g)$$

is in  $L_2$  for each  $N$ .

**THEOREM 4.**  $L_2 \cap \mathcal{D}$  is closed in  $L_2$ -norm topology.

*Proof.* Let  $\exp \zeta(g_n)$  converge strongly to  $\xi \in L_2(P)$ . By Itô's representation [20] of square-integrable functionals of Brownian motion we can write

$$\xi(\omega) = 1 + \int_0^1 \varphi(s, \omega) dw(s),$$

where  $\varphi$  is a nonanticipative functional with

$$E \int_0^1 |\varphi(t, \omega)|^2 dt < \infty.$$

Let  $\xi(t) = E\{\xi | \mathcal{S}_t\}$ . Then

$$\begin{aligned} \xi(t) &= 1 + E \left\{ \int_0^t \varphi(s, \omega) dw(s) | \mathcal{S}_t \right\} + E \left\{ \int_t^1 \varphi(s, \omega) dw(s) | \mathcal{S}_t \right\} \\ &= 1 + \int_0^t \varphi(s, \omega) dw(s) \end{aligned}$$

because the first stochastic integral is  $\mathcal{S}_t$ -measurable, and the second conditional expectation is zero.

Jensen's inequality gives

$$\begin{aligned} |E\{\xi - e^{\zeta(g_n)} | \mathcal{S}_t\}| &\leq E\{|\xi - e^{\zeta(g_n)}| | \mathcal{S}_t\}, \\ E|E\{\xi - e^{\zeta(g_n)} | \mathcal{S}_t\}| &\leq E|\xi - e^{\zeta(g_n)}|. \end{aligned}$$

The right-hand side goes to zero, by Schwarz's inequality. This shows that for each  $t$ ,

$$E|e^{\zeta_0^{\delta(g_n)}} - \xi(t)| \leq E|\xi - e^{\zeta(g_n)}|^2 = o(1).$$

Hence  $e^{\zeta_0^{\delta(g_n)}} \rightarrow \xi(\cdot)$  in  $L_1(\lambda \times P)$ .

Now since

$$E \left| \int_0^1 e^{\zeta_0^{\delta(g_n)}} g_n(s) dw(s) - \int_0^1 \varphi(s) dw(s) \right|^2 = E \int_0^1 \left| e^{\zeta_0^{\delta(g_n)}} g_n(s) - \varphi(s) \right|^2 ds,$$

there exists a subsequence, assumed relabeled, such that

$$\begin{aligned} e^{\zeta_0^{\delta(g_n)}} &\rightarrow \xi(\cdot), \\ e^{\zeta_0^{\delta(g_n)}} g_n(\cdot) &\rightarrow \varphi(\cdot) \end{aligned}$$

a.s.  $[\lambda \times P]$  ( $\lambda =$  Lebesgue measure).

Noting that

$$\xi(t)g_n(t) - \varphi(t) = [\xi(t) - e^{\zeta_b(g_n)}]g_n(t) + e^{\zeta_b(g_n)}g_n(t) - \varphi(t)$$

and that  $|g_n(t)|^2 \leq \text{const.} [1 + |w(t, \omega)|^2]$ , we can conclude from Lemma 2 that  $g_n$  converges almost surely to a function  $g$ . By Lemma 6,  $g$  is an admissible control, and

$$\xi(t) = 1 + \int_0^t \xi(s)g(s) dw(s).$$

Hence [14],  $\xi(t) = \exp \zeta_0^t(g)$  with  $g \in \mathcal{A}$ . Thus  $\xi \in \mathcal{D}$ , and  $\mathcal{D}$  is closed in norm topology.

*Remark.* If  $\mathcal{D}$  is a bounded and weakly sequentially closed subset of  $L_2$ , then an optimal control law exists.

*Proof.* This is almost obvious. Let  $g_n$  be a minimizing sequence of admissible drifts.  $\mathcal{D}$  is weakly sequentially compact, being bounded. Thus we may suppose that  $\exp \zeta(g_n)$  converges weakly in  $L_2$  to an  $L_2$ -function  $\xi \in \mathcal{D}$ . Hence there is an admissible drift  $g$  such that  $\xi = \exp \zeta(g)$ . The cost of using drift  $g_n$  is

$$E_{W_0}(1) e^{\zeta(g_n)} \rightarrow E_{W_0}(1) e^{\zeta(g)}.$$

Since  $g_n$  is minimizing,  $g$  is an optimal admissible drift.

Our basic  $L_2$  existence result is the following theorem.

**THEOREM 5.** *If  $G_t = S_t$ , if  $f(t, y, \Gamma)$  is convex, and if  $\mathcal{D}$  is  $L_2$ -bounded, then an optimal control law exists.*

*Proof.* The first two hypotheses imply, by Theorem 3, that  $\mathcal{D}$  is convex. Theorem 4 implies that  $\mathcal{D}$  is strongly closed. Hence it is weakly closed, and existence of an optimal admissible drift follows from the preceding remark.

**LEMMA 8.** *If  $\tau_N = \min \{1, \inf t : |w(t)| = N\}$ , then*

$$\sup_{g \in \mathcal{A}} E \exp 2\zeta_0^{\tau_N}(g) < \infty.$$

*Proof.* If  $\chi_N = \chi_{\{\sup_{0 \leq s \leq t} |w(s)| \leq N\}}$ , we have  $\zeta_0^{\tau_N}(g) = \zeta_0^1(\chi_N g)$ , and with  $T(t) = \int_0^t |g|^2 dt$ ,  $T = T(1)$ ,  $W(t) = \int_0^{T^{-1}(t)} g dw$ ,  $T^{-1}(t) = \inf u : T(u) = t$ , the argument in Lemma 4 gives

$$\chi_{\zeta(\chi_N g) > \log a} \leq \chi_{\{T \leq t, \sup_{0 \leq s \leq t} |W(s)| > \log a\}} + \chi_{\{T > t\}}.$$

Take  $t = \kappa\{1 + N^2\}$  to get

$$\begin{aligned} \chi_{\zeta(\chi_N g) > \log a} &\leq \chi_{\{\sup_{0 \leq s \leq \kappa(1 + N^2)} |W(s)| > \log a\}}, \\ P\{e^{\zeta(\chi_N g)} > a\} &= P\{\zeta(\chi_N g) > \log a\} \\ &\leq P\{\sup_{0 \leq s \leq \kappa(1 + N^2)} |W(s)| > \log a\} \\ &\leq \frac{1}{\log a} \sqrt{\frac{2}{\pi}} \sqrt{\kappa(1 + N^2)} \exp - \frac{\log^2 a}{2\kappa(1 + N^2)} \end{aligned}$$

$$\begin{aligned} &\leq \frac{\text{const.}}{\log a} a^{-(\log a)/2\kappa(1+N^2)} \\ &\leq \text{const. } a^{-2-\varepsilon} \quad \text{for } \frac{\log a}{2\kappa(1+N^2)} > 2 + \varepsilon. \end{aligned}$$

But

$$E \exp(2\zeta_0^{\tau N}(g)) = \int_0^\infty P\{\exp \zeta(\chi_{Ng}) > a\} a \, da$$

< const. depending on  $N$  but not on  $g$ .

LEMMA 9. Let  $\xi$  be a random variable measurable on  $\mathcal{S}_1$ , with  $E|\xi|^p < \infty$  for some  $p > 1$ , and such that

$$\xi(t) = E\{\xi | \mathcal{S}_t\}$$

is a continuous martingale. Let  $\tau$  be a stopping time and let  $\mathcal{S}_\tau$  be the  $\sigma$ -algebra of all  $\mathcal{S}_t$  sets  $A$  such that  $A \cap \{\tau \leq t\} \in \mathcal{S}_t$ . Then a.s.

$$\xi(\tau) = E\{\xi | \mathcal{S}_\tau\}.$$

*Proof.* We shall show that  $\xi(\tau)$  is  $\mathcal{S}_\tau$ -measurable and that  $B \in \mathcal{S}_\tau$  implies

$$\int_B \xi(\tau) \, dP = \int_B \xi \, dP.$$

To show that  $\xi(\tau)$  is  $\mathcal{S}_\tau$ -measurable it is sufficient to show that every set of the form  $\{\omega : \xi(\tau) \in A, \tau \leq t\}$ ,  $A$  closed, belongs to  $\mathcal{S}_t$ . Let  $\pi_m$  be a countable cover of  $[0, t]$  by open sets of diameter  $\leq (\frac{1}{2})^m$ . Then to within a set of measure zero [15],

$$\{\omega : \xi(\tau) \in A, \tau \leq t\} = \bigcap_m \bigcup_{S \in \pi_m} \{\tau \in S \text{ and } \xi^{-1}(A) \cap S \neq \emptyset\}.$$

Also, with  $D$  a countable set dense in  $S$ ,

$$\begin{aligned} \{\omega : \xi^{-1}(A) \cap S = \emptyset\} &= \bigcup_{s \in S} \{\omega : \xi(\omega, s) \in A\} \\ &= \bigcup_{s \in D} \{\omega : \xi(\omega, s) \in A\}. \end{aligned}$$

Thus  $\xi(\tau)$  is  $\mathcal{S}_\tau$ -measurable. Moreover, with  $B \in \mathcal{S}_\tau$ , and  $0 = t_0 < t_1 < \dots < t_n = 1$ ,

$$\begin{aligned} \int_B \xi \, dP &= \sum_{i=1}^n \int_{B \cap \{t_{i-1} < \tau \leq t_i\}} \xi \, dP \\ &= \sum_{i=1}^n \int_{B \cap \{t_{i-1} < \tau \leq t_i\}} \xi(t_i) \, dP \\ &= \int_B \xi(\tau) \, dP + \int_B \sum_{i=1}^n \chi_{\{t_{i-1} < \tau \leq t_i\}} [\xi(t_i) - \xi(\tau)] \, dP. \end{aligned}$$

With  $\eta(\omega) = \xi(t_{i+1})$  on  $\{t_i < \tau \leq t_{i+1}\}$ , we have

$$\int_B |\eta - \xi(\tau)| \, dP \leq \left( \int_B |\eta - \xi(\tau)|^p \, dP \right)^{1/p} P^{p/(1-p)}(B).$$

Clearly  $|\eta - \xi(\tau)|^p \leq 2 \sup_{0 \leq s \leq 1} |\xi(s)|^p$ , and by the martingale property,

$$E \sup_{0 \leq s \leq 1} |\xi(s)|^p \leq \left( \frac{p}{p-1} \right)^p E|\xi(1)|^p.$$

Thus the integrand  $|\eta - \xi(\tau)|^p$  approaches zero a.s. as  $\max_{0 < i \leq n} |t_i - t_{i-1}| \rightarrow 0$  and is dominated by an integrable function; hence the second integral above goes to zero.

LEMMA 10. *If  $\xi$  is an integrable Brownian functional measurable on  $\mathcal{S}_1$ , with  $\xi > 0$  a.s., then there is a functional  $\psi$  such that*

$$\xi(t) = E\{\xi | \mathcal{S}_t\} = e^{\zeta_b(\psi)} E\xi,$$

where  $\int_0^1 |\psi|^2 dt < \infty$  a.s. and  $\psi \cdot, \cdot$  is nonanticipative.

*Proof.* It is easy to see that  $\xi(t) > 0$  a.s. for each  $t$ . Now J. M. C. Clark has shown [21] that if  $\xi$  is an integrable  $\mathcal{S}_1$ -measurable function, then there is a non-anticipative function  $\varphi(\cdot, \cdot)$  such that for  $s < t$ ,

$$E\{\xi | \mathcal{S}_t\} - E\{\xi | \mathcal{S}_s\} = \int_s^t \varphi(u, \omega) dw(u),$$

$$P \left\{ \int_0^1 |\varphi|^2 du < \infty \right\} = 1,$$

$\mathcal{S}_0$  is trivial, so  $\xi(0) = E\xi$ . Let now

$$\psi(t, \omega) = \frac{\varphi(t, \omega)}{E\xi + \int_0^t \varphi dw} = \frac{\varphi(t, \omega)}{\xi(t, \omega)}.$$

For fixed  $\omega$ ,  $\psi(\cdot, \omega)$  is a Lebesgue function of  $t$ ; for fixed  $t$ , it is  $\mathcal{S}_t$ -measurable; so it is nonanticipative. Evidently

$$\xi(t) = E\xi + \int_0^t \xi(s)\psi(s)dw(s).$$

McKean [14] has shown that the only solution to this stochastic equation is

$$\xi(t) = \exp \zeta'_0(\psi) E\xi.$$

It remains to prove that  $\int_0^1 |\psi|^2 dt < \infty$  a.s. This will follow from

$$P \left\{ \inf_{0 \leq t \leq 1} \xi(t) \leq 0 \right\} = 0.$$

Since  $\xi(\cdot)$  is a martingale, Jensen's inequality gives

$$E\{\log \xi(t) | \mathcal{S}_s\} \leq \log E\{\xi(t) | \mathcal{S}_s\} = \log \xi(s),$$



so that  $-\log \zeta(\cdot)$  is a submartingale. Hence

$$\begin{aligned} P\left\{\inf_{0 \leq t \leq 1} \zeta(t) \geq e^{-n}\right\} &= P\left\{\sup_{0 \leq t \leq 1} [-\log \zeta(t)] \leq n\right\} \\ &\geq 1 - n^{-1} \int_{\left\{\sup_{0 \leq t \leq 1} [-\log \zeta(t)] \geq n\right\}} \zeta \, dP, \\ P\left\{\inf_{0 \leq t \leq 1} \zeta(t) < e^{-n}\right\} &\leq n^{-1} E\zeta. \end{aligned}$$

Now

$$\left\{\omega: \inf_{0 \leq t \leq 1} \zeta(t) \leq 0\right\} \subseteq \bigcap_{n \geq 1} \left\{\omega: \inf_{0 \leq t \leq 1} \zeta(t) < e^{-n}\right\}$$

and the probability of the intersection on the right is zero.

**THEOREM 6.** *If  $G_t = S_t$ , and if  $f(t, y, \Gamma)$  is convex for  $t, y \in [0, 1] \times C$  (Roxin's condition), then  $\mathcal{D}$  is weakly sequentially closed in  $L_1$ .*

*Proof.* Let  $e^{\zeta(g_n)} \in \mathcal{D}$  approach  $\zeta \in L_1$  weakly. By Lemma 10, the martingale  $\xi(t) = E\{\xi | \mathcal{S}_t^i\}$  is representable as  $\exp \zeta_0^t(\psi)$  with  $\psi(t, \omega)$  nonanticipating, and square-integrable in  $t$  a.s. Thus  $\xi(\cdot)$  has continuous sample paths almost surely. Since (Lemma 1)  $\mathcal{D}$  is a bounded set in  $L_p$  for some  $p > 1$  we can show that  $\xi$  belongs to the same  $L_p$ . The processes  $\exp \zeta_0^t(g_n) = E\{e^{\zeta(g_n)} | \mathcal{S}_t^i\}$ ,  $0 \leq t \leq 1$ , are also continuous  $L_p$ -martingales. Introduce the stopping times

$$\tau_N = \min \{1, \inf t: |w(t)| = N\}.$$

Lemma 9 implies that a.s.

$$\begin{aligned} \xi(\tau_N) &= E\{\xi | \mathcal{S}_{\tau_N}^i\}, \\ \exp \zeta_0^{\tau_N}(g_n) &= E\{e^{\zeta(g_n)} | \mathcal{S}_{\tau_N}^i\}. \end{aligned}$$

Since  $E\{\cdot | \mathcal{S}_{\tau_N}^i\}$  is "self-adjoint" we have for  $h \in L_\infty$ ,

$$\begin{aligned} EhE\{e^{\zeta(g_n)} | \mathcal{S}_{\tau_N}^i\} &= Ee^{\zeta(g_n)}E\{h | \mathcal{S}_{\tau_N}^i\}, \\ EhE\{\xi | \mathcal{S}_{\tau_N}^i\} &= E\xi E\{h | \mathcal{S}_{\tau_N}^i\}. \end{aligned}$$

It follows that for each  $N$ ,  $\exp \zeta_0^{\tau_N}(g_n)$  approaches  $\xi(\tau_N)$  in weak  $L_1$ . By Lemma 8, for each  $N$ , the sequence  $\exp \zeta_0^{\tau_N}(g_n)$  is bounded in  $L_2$ . Hence the convergence above holds in weak  $L_2$ -topology, and  $\xi(\tau_N) \in L_2$  for each  $N$ .

Let  $\chi_N = \chi_N(t, \omega) = \chi_{i(\tau_N > t)}$ , so that  $\xi(\tau_N) = \exp \zeta(\chi_N \psi)$  and  $\zeta_0^{\tau_N}(g_n) = \zeta_0^1(\chi_N g_n)$ . We have shown that for each  $N$ ,

$$\exp \zeta(\chi_N g_n) \rightarrow \exp \zeta(\chi_N \psi) \quad \text{in weak } L_2.$$

Let  $N$  be fixed, and choose a subsequence of  $\{g_n\}$ , assumed relabeled, so that

$$\frac{1}{n} \sum_{i=1}^n \exp \zeta(\chi_N g_i) \rightarrow \exp \zeta(\chi_N \psi) \quad \text{in strong } L_2.$$

As in the convexity result Theorem 3, the function

$$g_n^*(t, \omega) = \frac{\sum_{i=1}^n e^{\zeta_b^i(g_i)} g_i(t)}{\sum_{i=1}^n e^{\zeta_b^i(g_i)}}$$

is an admissible drift such that

$$\frac{1}{n} \sum_{i=1}^n \exp \zeta(\chi_N g_i) = \exp \zeta(\chi_N g_n^*).$$

The strong  $L_2$ -convergence implies, as in Theorem 4, that  $\chi_N \psi \in \chi_N f(t, w(\omega), \Gamma)$  a.s. It follows that

$$\psi(t, \omega) \in f(t, w(\omega), \Gamma) \quad \text{a.s.}$$

Thus by Lemma 5 and the elementary result in § 3,  $\psi$  is an admissible drift. This proves Theorem 6.

**THEOREM 7.** *If  $\mathcal{D}$  is weakly sequentially closed in  $L_1$ , then an optimal control law exists.*

*Proof.* By Lemma 1,  $\mathcal{D}$  is uniformly integrable. The inequality

$$\int_A e^{\zeta(g)} dP \leq \int_{\zeta(g) > \log N} e^{\zeta(g)} dP + NP(A)$$

then shows that  $\lim \int_A e^{\zeta(g)} dP = 0$  as  $P(A) \rightarrow 0$ , uniformly in  $g \in \mathcal{A}$ .  $\mathcal{D}$  is clearly bounded, and so it is weakly sequentially compact in  $L_1$ . Thus we may assume as given a minimizing sequence  $g_n, n \geq 1$ , of admissible drifts such that  $e^{\zeta(g_n)}$  converges weakly in  $L_1$  to a member  $e^{\zeta(g)}$  of  $\mathcal{D}$ . The cost of using  $g_n$  is

$$Ex_0(1)e^{\zeta(g_n)} \rightarrow \inf_{u \in \mathcal{U}} k(u).$$

We next note that for some  $p > 1, q = p/(p - 1)$ , and

$$\chi_N = \chi_{\{ \sup_{0 \leq s \leq 1} |w(s)|^2 > N \}},$$

we have for  $h \in \mathcal{A}$ ,

$$|E_{W_0}(1)\chi_N e^{\zeta(h)}| \leq E^{1/q} \chi_N |W_0(1)|^q \cdot E^{1/p} e^{p\zeta(h)}.$$

Lemma 1 shows that the second factor is bounded uniformly in  $h \in \mathcal{A}$ ; the first factor vanishes as  $N$  increases. Let  $\varepsilon > 0$  be given, and choose  $N$  so large that

$$|E_{W_0}(1)\chi_N e^{\zeta(h)}| < \frac{\varepsilon}{3}$$

for all  $h \in \mathcal{A}$ . Then by the weak  $L_1$ -convergence choose  $n_0$  so large that  $n > n_0$  implies

$$|E_{W_0}(1)(1 - \chi_N)(e^{\zeta(g)} - e^{\zeta(g_n)})| < \frac{\varepsilon}{3}.$$

It will follow that  $|k(g_n) - k(g)| < \varepsilon$  for  $n > n_0$ . Since  $\varepsilon$  was arbitrary,  $g$  is an optimal admissible drift.

Our final result removes from Theorem 5 the hypothesis that  $\mathcal{D}$  be included in a ball of  $L_2$ .

**THEOREM 8.** *If  $G_t = S_t$ , if  $f(t, y, \Gamma)$  is convex, and if  $|f(t, y, u)|^2 \leq \text{const.}(1 + |y(t)|^2)$ , then an optimal control law exists.*

*Proof.* Theorem 6 implies that  $\mathcal{D}$  is weakly sequentially closed in  $L$ . The result follows from Theorem 7.

**Acknowledgment.** The author is indebted to S. R. S. Varadhan for calling his attention to Girsanov's work and for suggesting a martingale proof of Girsanov's theorem; to H. S. Witsenhausen for two years of suggesting, encouraging, and criticizing; to H. P. McKean, Jr. for the random time change arguments in Lemmas 3 and 4; to T. T. Kadota and L. A. Shepp for helpful discussions; to T. E. Duncan and P. Varaiya for clarifying the ideas underlying Lemma 0 by sending him a preprint [23] of work based on an earlier version of this paper.

## REFERENCES

- [1] I. V. GIRSANOV, *On transforming a certain class of stochastic processes by absolutely continuous substitution of measures*, Theor. Probability Appl., 5 (1960), pp. 285–301.
- [2] W. H. FLEMING AND M. NISIO, *On the existence of optimal stochastic controls*, J. Math. Mech., 15 (1966), pp. 777–794.
- [3] A. V. SKOROKHOD, *Studies in the Theory of Random Processes*, Kiev, 1961; English translation: Addison-Wesley, Reading, Massachusetts, 1965.
- [4] R. H. CAMERON AND W. T. MARTIN, *Transformations of Wiener integrals under translations*, Ann. of Math., 45 (1944), pp. 386–96.
- [5] E. J. MCSHANE AND R. B. WARFIELD, *On Filippov's implicit functions lemma*, Proc. Amer. Math. Soc., 18 (1967), pp. 41–47.
- [6] H. J. KUSHNER, *On the existence of optimal stochastic controls*, this Journal, 3 (1965), pp. 463–474.
- [7] W. H. FLEMING, *Optimal control of partially observable diffusions*, this Journal, 6 (1968), pp. 194–214.
- [8] A. D. VENTCEL, *Additive functionals of multidimensional Wiener processes*, Dokl. Akad. Nauk SSSR, 139 (1961), pp. 13–16.
- [9] R. E. MORTENSEN, *Optimal control of continuous-time stochastic systems*, Rep. ERL66-1, Electronic Research Laboratory, University of California, Berkeley, 1966.
- [10] T. KAILATH, *A general likelihood-ratio formula for random signals in gaussian noise*, IEEE Trans. Information Theory, IT-15 (1969), pp. 350–361.
- [11] T. E. DUNCAN, *Probability densities for diffusion processes with applications to nonlinear filtering theory and detection theory*, Rep. 7001–4 and 7050–12, Center for Systems Research, Stanford University, Stanford, Calif., 1967.
- [12] G. KALLIANPUR AND C. STRIEBEL, *Estimation of stochastic systems; arbitrary system process with additive white noise observations error*, Ann. Math. Statist., 39 (1968), pp. 785–801.
- [13] D. STROOCK AND S. R. S. VARADHAN, *Diffusion processes with continuous coefficients*, Comm. Pure Appl. Math., 22 (1969), pp. 345–400 and 479–530.
- [14] H. P. MCKEAN, JR., *Stochastic Integrals*, Academic Press, New York and London, 1969.
- [15] V. E. BENEŠ, *Existence of optimal strategies based on specified information, for a class of stochastic additive white noise observations error*, Ann. Math. Statist., 39 (1968), pp. 785–801.
- [16] E. B. LEE AND L. MARKUS, *Optimal control for nonlinear processes*, Arch. Rational Mech. Anal., 8 (1961), pp. 36–58.
- [17] E. ROXIN, *The existence of optimal controls*, Michigan Math. J., 9 (1962), pp. 109–119.
- [18] W. H. FLEMING, *Optimal continuous-parameter stochastic control*, SIAM Rev., 11 (1969), pp. 470–509.
- [19] L. A. SHEPP, *Radon–Nikodym derivatives of Gaussian measures*, Ann. Math. Statist. 37 (1966), pp. 321–354.
- [20] K. ITÔ, *Multiple Wiener integral*, J. Math. Soc. Japan, 3 (1951), pp. 158–161.
- [21] J. M. C. CLARK, *The representation of functionals of Brownian motion by stochastic integrals*, Ann. Math. Statist., 41 (1970), pp. 1282–1295.
- [22] J. L. DOOB, *Stochastic Processes*, John Wiley, New York, 1953.
- [23] T. E. DUNCAN AND P. VARAIYA, *On the solutions of a stochastic control system*, to appear.

## STOCHASTIC CONTROL FOR SMALL NOISE INTENSITIES\*

WENDELL H. FLEMING†

**Abstract.** This paper is concerned with the approximate solution of stochastic optimal control problems which arise by perturbing the system equations in the deterministic Pontryagin control model, through an additive white noise term with small coefficient. The system states are assumed completely observable. Mathematically the problem becomes one of singular perturbation of the Hamilton–Jacobi equation by a small second order term. Our main results concern expansions of solutions of the perturbed equation in powers  $\varepsilon, \varepsilon^2, \varepsilon^3, \dots$  of the noise variance coefficient  $\varepsilon$ . The results obtained hold in regions where the corresponding solution of the Hamilton–Jacobi equation is sufficiently well-behaved.

**1. Introduction.** This paper is concerned with the approximate solution of some stochastic optimal control problems. The models of control systems which we consider are perturbations of the deterministic optimal control model of Pontryagin [22]. We suppose that, in the stochastic problem, the system is perturbed by random disturbances which take the form of an additive white noise term in the system equations (1.1). If the noise coefficient  $\sigma$  in (1.1) is small, it is reasonable to seek an approximate solution of the stochastic problem in terms of quantities computable from the Pontryagin problem (with  $\sigma = 0$ ). Some rather weak statements of this kind were proved in [7, pp. 269, 276]. Stronger results were obtained in [9, p. 527] for the randomly perturbed simplest problem in calculus of variations. In the present paper still sharper statements are proved; see the theorems in §§ 6 and 7.

Mathematically, the problem becomes one of singular perturbation of the Hamilton–Jacobi equation (1.8<sup>0</sup>) by a small second order term (equation (1.8<sup>ε</sup>)). Our main results concern the validity of the approximate formulas (1.11)–(1.12), in regions where the solution  $\varphi^0$  of (1.8<sup>0</sup>) is sufficiently well-behaved. Such regions are called *regions of strong regularity* (§ 3). A rather general first order equation  $\varphi_s + F(s, x, \varphi_x) = 0$  can be regarded as the Hamilton–Jacobi equation for some control problem (indeed, for the simplest problem in calculus of variations). Hence our results apply to such equations. See Theorem 7.2.

We consider the same kind of stochastic optimal control model as in [7]. Let

$$\xi(t) = (\xi_1(t), \dots, \xi_n(t))$$

denote the state of the system at time  $t$ . Thus  $\xi(t)$  is a vector in  $n$ -dimensional space  $R^n$ . The state process  $\xi$  evolves according to stochastic differential equations, written in vector-matrix notation as

$$(1.1) \quad d\xi = f(t, \xi(t), u(t)) dt + \sigma dw,$$

where  $w$  is an  $n$ -dimensional Brownian motion with  $w(s) = 0$ . At the initial time  $s$ , an initial state vector

$$(1.2) \quad x = \xi(s)$$

\* Received by the editors May 26, 1970.

† Department of Mathematics, Brown University, Providence, Rhode Island 02912. This work was supported in part by the Air Force Office of Scientific Research under Grant AF-AFOSR 67-0693A and in part by the National Science Foundation under Grant GP6733.

is known. We require that  $(s, x) \in Q$ , where  $Q$  is a given open subset of  $R^{n+1}$ . The control  $u(t)$  applied in (1.1) at time  $t$  is a vector subject to

$$(1.3) \quad u(t) \in K,$$

where the “control set”  $K$  is a closed convex subset of some  $R^k$ . The control  $u$  in (1.3), as well as  $\xi$  and  $w$ , is generally a stochastic process defined on some probability space  $\Omega$ . Following custom, we simplify the notation by writing  $u(t)$ ,  $\xi(t)$ ,  $w(t)$  instead of  $u(t, \omega)$ ,  $\xi(t, \omega)$ ,  $w(t, \omega)$ ,  $\omega \in \Omega$ .

Let  $\tau$  denote the first time  $t \geq s$  when  $(t, \xi(t)) \notin Q$ ;  $\tau$  is called the *exit time* from  $Q$ , and is sometimes denoted by  $\tau_Q$  to indicate its dependence on  $Q$ . The criterion of performance of the system is

$$(1.4) \quad J = E \int_s^\tau L(t, \xi(t), u(t)) dt,$$

where  $E$  denotes expected value.

If  $\sigma = 0$  and  $u$  is a (deterministic) function of  $t$ , then (1.1) is a vector ordinary differential equation. Instead of (1.4) we have the performance criterion

$$(1.4^0) \quad J^0 = \int_s^\tau L(t, \xi(t), u(t)) dt.$$

The problem of finding a control function  $u^0$  such that  $J^0$  is minimum is precisely that of Pontryagin. There is an extensive literature; see for instance [13], [19], [22].

When  $\sigma \neq 0$  we must specify the information available to the controller. Let us assume that the state  $\xi(t)$  is completely observable at each time  $t$ . This assumption is essential if the stochastic control problem is to be reduced to one in partial differential equations via dynamic programming. It is interesting to try to replace it by other observability assumptions. Wonham [26] proved a separation principle reducing certain partially observable problems to completely observable ones. Kushner [16] considered the case of random disturbances entering the system as impulses at discrete time instants, rather than as white noise as in (1.1). See also remarks about the open loop problem at the end of § 6.

Since the states are completely observable the controller may take

$$(1.5) \quad u(t) = Y(t, \xi(t)),$$

where  $Y$  is a *control policy* belonging to a certain class  $\mathcal{Y}$  (§ 4). Sometimes  $Y$  is called a *closed loop*, or *Markov*, control policy, also a program [13]. Given a control policy  $Y$  and initial data  $(s, x)$ , the states  $\xi(t)$  evolve according to a Markov process. The policy  $Y$  controls the local drift coefficient  $f^Y$ , defined by

$$f^Y(s, x) = f(s, x, Y(s, x));$$

and (1.1) can be rewritten in the customary form [5], [12, Chap. 8]:

$$(1.1') \quad d\xi = f^Y dt + \sigma dw.$$

The stochastic optimization problem is to find among all  $Y \in \mathcal{Y}$  one for which  $J$  minimizes (1.4).

For simplicity let the noise coefficient  $\sigma$  be a diagonal matrix, with diagonal elements  $(2\varepsilon)^{1/2}$ . The components of  $\sigma w(t)$  are then independent and have mean 0, variance  $2\varepsilon(t - s)$ . In (1.4), let us write  $J = J^\varepsilon(Y; s, x)$ , and

$$(1.6^\varepsilon) \quad \varphi^\varepsilon(s, x) = \inf_{Y \in \mathcal{Y}} J^\varepsilon(Y; s, x).$$

For each  $(s, x)$  and  $p \in R^n$  let

$$(1.7) \quad H(s, x, p) = \min_{y \in K} [L(s, x, y) + pf(s, x, y)].$$

Let  $\psi_s, \psi_{x_i}$  denote the partial derivatives of a function  $\psi(s, x)$ . Let

$$\psi_x = (\psi_{x_1}, \dots, \psi_{x_n}), \quad \Delta_x \psi = \sum_{i=1}^n \psi_{x_i x_i}$$

denote the gradient vector and Laplacean in the variables  $x$ . If  $Q$  belongs to a class  $\mathcal{S}$  of regions suitable from the viewpoint of parabolic equations (§ 4), if  $K$  is compact, and if we impose suitable assumptions on  $f, L$  (§ 2), then  $\varphi^\varepsilon$  satisfies in  $Q$  the parabolic partial differential equation

$$(1.8^\varepsilon) \quad \varphi_s^\varepsilon + \varepsilon \Delta_x \varphi^\varepsilon + H(s, x, \varphi_x^\varepsilon) = 0$$

together with the data

$$(1.9) \quad \varphi^\varepsilon(s, x) = 0 \quad \text{for } (s, x) \in \partial^* Q.$$

Here  $\partial^* Q$  is the “essential” portion of the boundary  $\partial Q$  (§ 4). The optimal policy  $Y^\varepsilon$  has the property that

$$(1.10) \quad L(s, x, y) + \varphi_x^\varepsilon(s, x)f(s, x, y) = \min \quad \text{on } K$$

when  $y = Y^\varepsilon(s, x)$ .

Thus the completely observable optimal control problem is in principle reduced to solving the boundary problem (1.8 $^\varepsilon$ )–(1.9) for  $\varphi^\varepsilon$ , and then minimizing  $L + \varphi_x^\varepsilon f$  over  $K$  for each  $(s, x) \in Q$ . This is usually difficult to do in practice, although some approximate methods have been proposed. See for instance [15], [25]. In the present paper we seek approximate formulas for  $\varphi^\varepsilon$  and  $\varphi_x^\varepsilon$  for small positive  $\varepsilon$ .

When  $\varepsilon = 0$  we define  $\varphi^0$  by

$$(1.6^0) \quad \varphi^0(s, x) = \inf_{u \in \mathcal{U}_s} J^0(u; s, x),$$

where  $J^0$  is as in (1.4 $^0$ ) and  $\mathcal{U}_s$  is an appropriate class of (open loop) controls (§ 3). It can be shown that also

$$\varphi^0(s, x) = \inf_{Y \in \mathcal{Y}} J^0(Y; s, x).$$

Formally,  $\varphi^0$  satisfies in  $Q$  the Hamilton–Jacobi equation

$$(1.8^0) \quad \varphi_s^0 + H(s, x, \varphi_x^0) = 0$$

with  $\varphi^0 = 0$  on  $\partial^* Q$ ; an optimal policy  $Y^0$  is found formally by minimizing  $L + \varphi_x^0 f$  over  $K$ .

We wish to find  $\varphi^\varepsilon, \varphi_x^\varepsilon$  (and hence also  $Y^\varepsilon$ ) approximately in terms of quantities computable from  $\varphi^0$ . Two apparent difficulties are immediately seen. One is

mathematical, the other practical. First, the function  $\varphi^0$  is not differentiable everywhere; one must usually interpret (1.8<sup>0</sup>) in some generalized sense [7]. The policy  $Y^0$  need not be continuous, and one does not know that the system equations  $d\xi = f^{Y^0} dt$  can be solved for the optimal deterministic trajectories. The apparent practical difficulty is that the boundary problem (1.8<sup>0</sup>)–(1.9) does not seem easier to solve numerically than (1.8 <sup>$\varepsilon$</sup> )–(1.9).

However, let us attempt to solve the stochastic problem, not in all of  $Q$ , but merely near some particular optimal trajectory  $\gamma^0$ . We suppose that  $\gamma^0$  lies in some region  $N$  of strong regularity (§ 3). In  $N$ ,  $\varphi^0$  is rather smooth and satisfies (1.8<sup>0</sup>) in the classical sense. We seek in  $N$  expansions of the type

$$(1.11) \quad \varphi^\varepsilon = \varphi^0 + \varepsilon\theta_1 + \varepsilon^2\theta_2 + \cdots + \varepsilon^l\theta_l + o(\varepsilon^l),$$

$$(1.12) \quad \varphi_x^\varepsilon = \varphi_x^0 + \varepsilon(\theta_1)_x + \varepsilon^2(\theta_2)_x + \cdots + \varepsilon^l(\theta_l)_x + o(\varepsilon^l).$$

The coefficients  $\theta_1, \dots, \theta_l$  satisfy the equations found by differentiating (1.8 <sup>$\varepsilon$</sup> ) repeatedly with respect to  $\varepsilon$  and setting  $\varepsilon = 0$ . These equations involve the partial derivatives of  $H$  of corresponding orders and of  $\varphi^0$ . The validity of such expansions will depend on smoothness properties of  $H$ , discussed in § 2. For certain problems, including the linear time optimal problem with polyhedral  $K$ , we know only (1.11) with  $l = 1$ ; in such problems,  $H$  is not even of class  $C^1$ . In other problems where a strong convexity condition holds we may take  $l = 2$  in (1.11) and  $l = 1$  in (1.12). See § 6. In § 7 we discuss problems without control constraints (thus of the type in calculus of variations). For such problems one can take  $l$  arbitrarily large, but finite. One cannot let  $l \rightarrow \infty$  and obtain  $\varphi^\varepsilon$  as a convergent power series in  $\varepsilon$ ;  $\varphi^\varepsilon$  is at best a  $C^\infty$ , not real analytic, function of  $\varepsilon$  for  $\varepsilon \geq 0$ . It turns out that the optimal trajectory  $\gamma^0$  is a characteristic curve for the equations which  $\theta_j, (\theta_j)_x$  satisfy. Therefore the values  $\theta_j(s, x), (\theta_j)_x(s, x)$  at the left endpoint  $(s, x)$  of  $\gamma^0$  are integrals along  $\gamma^0$  of expressions involving partial derivatives of  $\varphi^0$  and  $H$  (§ 6). These partial derivatives can be calculated from the method of characteristics in a neighborhood of  $\gamma^0$ . See the Appendix. When there are no control constraints (§ 7) one needs for this purpose a solution to the secondary minimum problem in calculus of variations. For  $\gamma^0$  to lie in a region of strong regularity it suffices in this case that:

- (i) the initial point  $(s, x)$  of  $\gamma^0$  be regular (meaning that  $\gamma^0$  is the unique optimal deterministic trajectory with initial point  $(s, x)$ ); and
- (ii) the classical Jacobi condition holds (meaning that  $(s, x)$  is not a conjugate point).

One might expect that the optimal stochastic trajectories, with the same initial point  $(s, x)$  as  $\gamma^0$ , remain in  $N$  with probability very nearly 1. A sharp estimate of this kind is Corollary 5.7.

Our method also answers the following question. Suppose that the controller does not solve the stochastic optimization problem, but instead uses the optimal deterministic policy  $Y^0$ . (More precisely, suppose that he uses a policy  $Y \in \mathcal{Y}$  agreeing with  $Y^0$  in  $N$ .) How close to the optimum is the performance  $J^\varepsilon(Y^0; s, x)$  in the stochastic problem? See the end of § 6.

The program of the paper is as follows. We begin with some properties of  $H$  (§ 2) and a discussion of the Pontryagin problem (§ 3). In §§ 4, 5 we obtain some

preliminary estimates for the convergence of  $\varphi^\varepsilon$  and  $\varphi_x^\varepsilon$  as  $\varepsilon \rightarrow 0$ , which are weaker than those eventually proved in §§ 6, 7. A crucial step is to establish convergence of  $\varphi_x^\varepsilon$  to  $\varphi_x^0$  uniformly on compact subsets of  $N$  (Lemma 5.5).

In §§ 4–7 we impose a strong convexity condition (2.2'). This insures that the optimal policies  $Y^\varepsilon$  for  $\varepsilon > 0$  belong to  $\mathcal{Y}$ . The standard Ito conditions then guarantee that the system equations (1.1) have a solution when  $Y = Y^\varepsilon$ . In § 8 we give some partial results when (2.2') does not hold. The optimal  $Y^\varepsilon$  may then be discontinuous; however, the existence of solutions to (1.1) for  $\varepsilon > 0$  is still guaranteed in the sense of [24]. An interesting open question is to prove precise formulas for the perturbation of switching surfaces for  $Y^\varepsilon$ , for example, in the linear time-optimal problem. Formal studies of such perturbation problems were made in [4] and [23].

In § 9 the autonomous problem is mentioned. Equation (1.8 $^\varepsilon$ ) is then of elliptic type. In § 10 some examples are given, including one studied by E. Hopf [14] in which (1.8 $^\varepsilon$ ) is equivalent to Burgers' equation. These examples show that our results are in many respects best possible.

**2. Assumptions on  $f, L, K$ ; Properties of  $H$ .** If  $\Gamma$  is an open set, we write  $g \in C^l(\Gamma)$  to mean that the function  $g$  together with its partial derivatives of orders  $j = 1, \dots, l$  are continuous on  $\Gamma$ . If  $\Gamma$  is not open, then  $g \in C^l(\Gamma)$  means that  $g$  agrees on  $\Gamma$  with a function  $g' \in C^l(\Gamma')$ , where  $\Gamma'$  is open and  $\Gamma \subset \Gamma'$ . If  $\Gamma$  is clear from the context, we sometimes write merely  $g \in C^l$ .

Throughout, we assume that  $f, L \in C^\infty(R^{n+k+1})$ . However, we shall use their values only on  $[T_0, T] \times R^n \times K$ , where  $[T_0, T]$  is a fixed (finite) time interval and the control set  $K \subset R^k$  is closed and convex. (In § 9 we consider the autonomous case, with  $f, L$  defined on  $R^{n+k}$ .)

The following additional assumptions are made.

(2.1)  $f$  is linear in the control variables  $y = (y_1, \dots, y_k)$ , namely,  $f(s, x, y) = A(s, x) + B(s, x)y$ . The functions  $A, B$  (respectively vector- and  $n \times k$  matrix-valued) are bounded together with their first order partial derivatives.

$L(s, x, y)$  is convex in  $y$ , namely,

(2.2) 
$$vL_{yy}v \geq 0 \quad \text{for all vectors } v \in R^k.$$

(2.3)  $L \geq c_1 > 0$ . Moreover,  $L$  and its gradient  $L_x$  in the variables  $x = (x_1, \dots, x_n)$  are bounded on  $[T_0, T] \times R^n \times K_1$  for any compact  $K_1 \subset K$ .

In (2.2), as throughout the paper, we use vector-matrix notation. Thus  $L_{yy}$  denotes the matrix of second order partial derivatives, and

$$vL_{yy}v = \sum_{i,j=1}^k L_{y_i y_j} v_i v_j.$$

If in (1.4),  $\tau \equiv T$ , which means  $Q = (T_0, T) \times R^n$ , we can replace  $L$  by  $L + \text{const}$ . In that case (2.3) is equivalent to the assumption that  $L$  is bounded below. (The boundedness assumptions on  $A, B, L$  can be replaced by suitable growth conditions as  $|x| \rightarrow \infty$ ; see [10]. However, when we work in some bounded region  $N$ , the behavior of these functions for large  $|x|$  is of interest only in that they insure the



global existence of the solutions  $\varphi^\epsilon$  to (1.8 $^\epsilon$ )–(1.9) and a uniform bound for  $\varphi^\epsilon$  on  $N$ .)

To avoid some analytical difficulties, and to obtain stronger results, we replace (2.2) in §§ 4–7 by the stronger condition

$$(2.2') \quad \nu L_{yy} \nu \geq c(|y|)|\nu|^2, \quad c > 0,$$

for all  $\nu \in R^n$ , where  $c(r)$  is a positive, nonincreasing function of  $r$ . We also suppose in §§ 4–6 that  $K$  is compact. In that case  $c(|y|) \geq c_0 > 0$  for all  $y \in K$ .

**Properties of  $H$ .** If  $K$  is compact, then  $H$  is clearly well-defined by (1.7). If  $K$  is not compact (as in § 7), then we shall always assume (2.2') and the following.

(2.3) For each  $(s, x, p)$ ,  $L + pf$  has a minimum on  $K$  at a unique  $y = V(s, x, p)$ .  
 Moreover,  $|p| \leq \nu$  implies  $|V| \leq R(\nu)$  for some nondecreasing  $R(\nu)$ .

We are interested in continuity and differentiability properties of  $H$ . We begin with a known lemma. For completeness a proof is included.

LEMMA 2.1. Let  $\Gamma \subset R^m$  be open and  $K \subset R^k$  compact. Let  $\Phi$  and its gradient  $\Phi_q$  in the variables  $q$  be continuous on  $\Gamma \times K$ ; and let

$$(2.4) \quad \Psi(q) = \min_{y \in K} \Phi(q, y).$$

Then:

- (a)  $\Psi$  satisfies a Lipschitz condition on any compact subset  $\Gamma_1$  of  $\Gamma$ .
- (b) If  $\Phi(q_0, y)$  is minimum on  $K$  at  $y_0$  and  $\Psi$  is differentiable at  $y_0$ , then

$$(2.5) \quad \Psi_q(q_0) = \Phi_q(q_0, y_0).$$

- (c) Suppose that, for every  $q \in \Gamma$ ,  $\Phi(q, y)$  has a minimum on  $K$  at a unique point  $y = V(q)$ . Then  $\Psi \in C^1(\Gamma)$ .

*Proof.* Since  $\Phi_q$  is bounded on  $\Gamma_2 \times K$ , where  $\Gamma_2$  is a neighborhood of  $\Gamma_1$ , a Lipschitz condition

$$|\Phi(q_1, y) - \Phi(q_0, y)| \leq C|q_1 - q_0|$$

holds for  $q_0, q_1 \in \Gamma_1$  and  $y \in K$ . (The constant  $C$  depends on  $\Gamma_1$ .) Let  $\Phi(q_i, y)$  be minimum on  $K$  at  $y_i, i = 0, 1$ . Then

$$(2.6) \quad \begin{aligned} \Phi(q_1, y_1) - \Phi(q_0, y_1) &\leq \Psi(q_1) - \Psi(q_0) \leq \Phi(q_1, y_0) - \Phi(q_0, y_0), \\ |\Psi(q_1) - \Psi(q_0)| &\leq C|q_1 - q_0|. \end{aligned}$$

This proves (a). To prove (b), we subtract  $(q_1 - q_0)\Phi_q(q_0, y_0)$  from the middle and right-hand terms of (2.6) and divide by  $|q_1 - q_0|$ . This gives

$$(2.7) \quad \limsup_{q_1 \rightarrow q_0} |q_1 - q_0|^{-1} [\Psi(q_1) - \Psi(q_0) - (q_1 - q_0)\Phi_q(q_0, y_0)] \leq 0.$$

Since  $\Psi$  is differentiable at  $q_0$ , this implies that  $\Phi_q = \Psi_q$  (take  $q_1 - q_0 = \pm \rho h$ , where  $h$  is a fixed unit vector,  $\rho \rightarrow 0$ .) This proves (b).

To prove (c) by compactness of  $K$ ,  $V$  is continuous. It suffices to verify that  $\Psi$  is differentiable at each  $q_0 \in \Gamma$  and to use (b) with  $V(q_0) = y_0$ . Let us subtract  $(q_1 - q_0)\Phi_q(q_1, y_1)$  from the middle and left-hand terms of (2.6), where  $y_1 = V(q_1)$ . By continuity of  $V$  and  $\Phi_q$  we get

$$\liminf_{q_1 \rightarrow q_0} |q_1 - q_0|^{-1} [\Psi(q_1) - \Psi(q_0) - (q_1 - q_0)\Phi_q(q_0, y_0)] \geq 0.$$

This together with (2.7) implies that  $\Psi$  is differentiable at  $q_0$ . This proves (c).

*Note.* In the older literature differentiability is called total differentiability. By Rademacher's theorem [6, p. 216],  $\Psi$  is differentiable almost everywhere in  $\Gamma$  if  $\Psi$  is Lipschitz on any compact  $\Gamma_1 \subset \Gamma$ .

Now let  $q = (s, x, p)$ ,  $\Phi = L + pf$ , and  $\Psi = H$ . According to Lemma 2.1, we have the following lemma.

LEMMA 2.2. *Let  $K$  be compact. Then:*

- (a)  *$H$  is Lipschitz on any compact set.*
- (b) *If  $H$  is differentiable at  $(s, x, p)$ , then*

$$(2.8) \quad H_s = L_s + pf_s, \quad H_x = L_x + pf_x, \quad H_p = f.$$

(c) *Let  $\Gamma \subset R^{2n+1}$  be open. If  $L + pf$  is minimum on  $K$  at a unique  $y = V(s, x, p)$  for every  $(s, x, p) \in \Gamma$ , then  $H \in C^1(\Gamma)$  and (2.8) holds there.*

In (2.8),  $L_s, f_s, \dots, f$  are evaluated at a point  $y$ , where  $L + pf$  is minimum on  $K$ . In (c),  $y = V$ . The symbol  $f_x$  denotes the matrix of partial derivatives  $\partial f_i / \partial x_j$ , where  $f = (f_1, \dots, f_n)$ .

*Example 2.1.* Let  $f = y, L = 1, K = \{|y| \leq 1\}$ . Then  $H(p) = 1 - |p|$  is not differentiable at  $p = 0$ .  $V(p) = -|p|^{-1}p, p \neq 0$ . This is a simple special case of the time-optimal problem. The autonomous version of (1.8<sup>0</sup>) is just the eikonal equation  $|\varphi_x^0| = 1$ . See [2, p. 88].

LEMMA 2.3. *Assume (2.2'), and either  $K$  compact or (2.3'). Then:*

- (a)  $H \in C^1$ .
- (b)  $H_s, H_x, H_p, V$  are Lipschitz on compact sets.
- (c) *If  $\Gamma$  is a set such that  $V \in C^l(\Gamma)$ , then  $H \in C^{l+1}(\Gamma)$ .*

*Proof.* First consider the case  $K$  compact. By reasoning in [7, p. 259], noting that (2.1) implies  $\Phi_y = L_y + pB$ , we find that  $V$  is Lipschitz on any compact set. The remaining assertions follow from (2.8), with  $y = V(s, x, p)$ .

If  $K$  is not compact, let  $|p| \leq v$ . Then  $|V(s, x, p)| \leq r$  for  $r = R(v)$ , and we may replace  $K$  by the compact set  $K_r = K \cap \{|y| \leq r\}$ .

*Example 2.2.* Let  $f = y, L = \frac{1}{2}|y|^2, K = \{|y| \leq 1\}$ . Then

$$H(p) = \begin{cases} -\frac{1}{2}|p|^2, & \text{if } |p| \leq 1, \\ \frac{1}{2} - |p|, & \text{if } |p| \geq 1. \end{cases}$$

$H_p$  is continuous, but  $H_{pp}$  is discontinuous when  $|p| = 1$ .  $V(p) = -p$  if  $|p| \leq 1$ ,  $V(p) = -|p|^{-1}p$  if  $|p| \geq 1$ .  $V$  is Lipschitz, but not  $C^1$  where  $|p| = 1$ .

If there are no control constraints, then more can be said.

LEMMA 2.4. *Assume (2.2'), (2.3') and  $K = R^k$ . Then  $H \in C^\infty$  and  $L_y(s, x, V) = -pB(s, x)$ .*

*Proof.* By calculus,

$$0 = \Phi_y = L_y + pB.$$

By (2.2') and the implicit function theorem,  $V \in C^\infty$ . Hence  $H \in C^\infty$ .

*Note.* If  $f = y$  and  $B$  is the identity matrix, the dual formulas  $L_y = -p, H_p = y$  are the classical Legendre transformation and its inverse.

**3. The Pontryagin problem; Regions of strong regularity.** In this section we collect some definitions, and recall the classical formula for solution of a linear first order partial differential equation by characteristics. Then we prove a certain lemma, Lemma 3.1, which is similar to a lemma in [9, p. 523].

Given  $s \in (T_0, T)$  let  $\mathcal{U}_s$  denote the class of all bounded measurable functions  $u$  from  $[s, T]$  into the control set  $K$ . (These  $u$  are sometimes called open loop controls.) As usual  $u$  and  $v$  are regarded as the same if  $u(t) = v(t)$  almost everywhere. Let  $Q$  be open, with  $Q \subset (T_0, T) \times R^n$ . The Pontryagin problem is: given initial data  $(s, x) \in Q$  minimize  $J^0(s, x; u)$  in (1.4<sup>0</sup>) among all  $u \in \mathcal{U}_s$ . Let us suppose that an optimal  $u^0$  exists. (This is true if  $K$  is compact and (2.1)–(2.3) hold. For noncompact  $K$  additional assumptions are needed; a special case is treated in § 7.)

Let  $\gamma^0 = \gamma^0(s, x)$  denote the optimal trajectory corresponding to  $u^0$ :

$$(3.1) \quad \begin{aligned} \gamma^0 &= \{(t, \xi^0(t)) : s \leq t \leq \tau^0\}, \\ \frac{d\xi^0}{dt} &= f(t, \xi^0(t), u^0(t)), \end{aligned}$$

with  $\xi^0(s) = x$  and  $\tau^0$  the exit time of  $\gamma^0$  from  $Q$ . We often write  $z^0 = (\tau^0, \xi^0(\tau^0))$  for the exit place. If (2.2') holds, then any optimal  $u^0$  is continuous on  $[s, \tau^0]$ .

There may be several optimal control functions  $u^0$  for the same initial point  $(s, x)$ . Such points  $(s, x)$  are known to present difficulty from the viewpoint of constructing an optimal control policy  $Y^0$ . We call them *irregular points*. In the situation considered in § 7, the irregular points are exactly those where the function  $\varphi^0$  in (1.6<sup>0</sup>) fails to be differentiable [9, p. 520].

We shall need to assume that  $(s, x)$  belongs to a region  $N$  where stronger properties hold. Let us write  $\bar{D} = D \cup \partial D$  for the closure of a set  $D$ .

**DEFINITION 3.1.** Let  $N = \bar{Q} \cap \mathcal{O}_0$ , where  $\mathcal{O}_0$  is open. Then  $N$  is a *region of strong regularity* if there exist disjoint  $C^\infty$  hypersurfaces  $\Sigma_1, \dots, \Sigma_m$  and disjoint open sets  $N_1, \dots, N_m$  such that:

- (a)  $N - (\Sigma_1 \cup \dots \cup \Sigma_m) = N_1 \cup \dots \cup N_m$ .
- (b)  $N \cap \Sigma_1 = N \cap \partial Q, \partial N_j \subset \Sigma_j \cup \Sigma_{j+1} \cup (\partial N - \partial Q)$  for  $j = 1, \dots, m - 1, \partial N_m \subset \Sigma_m \cup (\partial N - \partial Q)$ .
- (c) For any  $(s, x) \in N$ , there is a unique optimal  $u^0 \in \mathcal{U}_s$ . The corresponding optimal trajectory  $\gamma^0(s, x)$  is contained in  $N$ .
- (d) For  $(s, x) \in N_j, \gamma^0(s, x)$  meets  $\Sigma_i$  nontangentially at a single point  $(s_i, x_i)$  for  $i \leq j$ , and  $\gamma^0(s, x)$  does not meet  $\Sigma_i$  for  $i > j$ .
- (e)  $\varphi^0 \in C^1(\bar{N})$  and  $\varphi^0 \in C^\infty(\bar{N}_j)$  for  $j = 1, \dots, m$ .
- (f) For  $(s, x) \in N_j, j = 1, \dots, m$ , and  $(s, x) \in \Sigma_1, L + \varphi_x^0 f$  has a minimum on  $K$  at a unique point  $y = Y^0(s, x)$ , where  $Y^0 \in C^\infty(\bar{N}_j)$ .
- (g) For  $(s, x) \in N_j, j = 1, \dots, m$ , and  $(s, x) \in \Sigma_1$ , the function  $V$  in Lemma 2.3 belongs to  $C^\infty(\Gamma)$  for some open set  $\Gamma$  containing  $(s, x, \varphi_x^0(s, x))$ .

See Fig. 1.

The classical technique for constructing regions of strong regularity is the method of characteristics, applied to the Hamilton–Jacobi equation. (In calculus of variations this is called constructing fields of extremals (see [13, Chap. 3]).) In the Appendix the method is outlined in the setting of the Pontryagin problem.

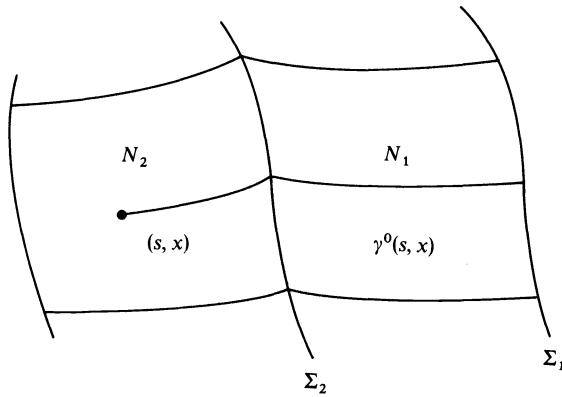


FIG. 1

The functions in (f), (g) are related by

$$Y^0(s, x) = V(s, x, \varphi_x^0(s, x)).$$

Since  $\varphi^0 \in C^1(\bar{N})$  the Hamilton–Jacobi equation (1.8<sup>0</sup>) holds in  $\bar{N}$ . Moreover,

$$(3.2) \quad u^0(t) = Y^0(t, \xi^0(t))$$

except perhaps for  $t = s_j$ . For this reason we call  $Y^0$  the *optimal control policy in N*. We do not require that the  $C^\infty$  extensions of  $Y^0$  from  $N_{j-1}$  and  $N_j$  to  $\bar{N}_{j-1}$  and  $\bar{N}_j$  agree on  $\Sigma_j$ . When  $Y^0$  is discontinuous across  $\Sigma_j$ , we call  $\Sigma_j$  a *switching surface*. (In (d) nontangency means from both sides of  $\Sigma_j$ .) If  $Y^0$  is continuous across  $\Sigma_j$  but its first order partial derivatives are discontinuous there, then we call  $\Sigma_j$  a *transition surface*.

A switching or transition surface  $\Sigma_j$  is usually determined by an equation of the form  $G_j(s, x, \varphi_x^0) = 0$ . The hypersurface  $G_j(s, x, p) = 0$  often separates two regions in  $R^{2n+1}$  such that  $V$  obeys some control constraint in one region but not in the other.

If the strong convexity condition (2.2') holds, then there are no switching surfaces. If besides (2.2') there are no control constraints (§ 7), then we may take  $m = 1$ . There are in that case neither switching nor transition surfaces.

For notational simplicity we set  $f^0 = f^{Y^0}$ . Thus

$$f^0(s, x) = f(s, x, Y^0(s, x)).$$

Consider a linear partial differential equation

$$(3.3) \quad \psi_s + f^0 \psi_x + g(s, x) = 0 \quad \text{in } N$$

with the data  $\psi = 0$  on  $\Sigma_1$ . Suppose that  $g \in C^\infty(\bar{N}_j)$  for  $j = 1, \dots, m$ . Let

$$(3.4) \quad \psi(s, x) = \int_{\gamma^0(s, x)} g = \int_s^{t^0} g(t, \xi^0(t)) dt, \quad (s, x) \in N.$$

Then  $\psi$  is continuous in  $N$ ,  $\psi \in C^\infty(\bar{N}_j \cap N)$  for  $j = 1, \dots, m$ , and  $\psi$  is a solution to (3.3) with the data on  $\Sigma_1$ . This follows by the classical method of characteristics [2]. If  $Y^0$  and  $g$  are continuous across  $\Sigma_j$ , then the first order partial derivatives of  $\psi$  are continuous there. Otherwise, the discontinuities in the first order partials of

$\psi$  across  $\Sigma_j$  can be calculated. The derivatives of  $\psi$  in directions tangent to  $\Sigma_j$  are the same from both sides, while the derivatives tangent to the characteristic curves from (3.3) satisfy  $d\psi/dt = -g$ .

In particular, take

$$\psi = \varphi_{x_i}^0, \quad g = H_{x_i}^0, \quad \text{and} \quad \psi = \varphi_s^0, \quad g = H_s^0,$$

where  $H_{x_i}^0 = H_{x_i}(s, x, \varphi_{x_i}^0)$ , etc. Recall (2.8). We find that the second order partial derivatives of  $\varphi^0$  are continuous across a transition surface  $\Sigma_j$  (however, the third order partial derivatives are generally discontinuous across  $\Sigma_j$ ). Across a switching surface  $\Sigma_j$ , first order, but not second order, partials of  $\varphi_0$  are continuous.

The following notations are used in the next lemma, and frequently in later sections:

(3.5)  $N$  denotes a region of strong regularity;  
 $N'$  denotes a bounded set with  $\bar{N}' \subset N$ .

(3.6)  $d(z) = \text{distance}(z, \partial Q)$  for  $z \in Q$ .

(3.7)  $\|\xi\|_t = \max_{s \leq r \leq t} |\xi(r)|$ .

In case  $t = T$  and the meaning is clear, we write

$$\|\xi\|_T = \|\xi\|.$$

In the lemma it is convenient to denote the solution of the deterministic system equations by  $\eta$  (rather than  $\xi$ ):

(3.8)  $\frac{d\eta}{dt} = f(t, \eta(t), u(t)), \quad s \leq t \leq T,$

with  $\eta(s) = x$ . Note that  $u, \eta$  are defined on  $[s, T]$ , although their values after the exit time  $\tau$  from  $Q$  do not affect  $J^0$ . If  $u = u^0$  is optimal, we write as above  $\eta = \xi^0$  and  $\tau^0, z^0$  for the corresponding exit time and place.

LEMMA 3.1. *Let  $K$  be compact, and  $N', N$  as above. Given  $a > 0$  there exists  $\delta = \delta(a, N') > 0$  with the following property: If  $(s, x) \in N', d(t, \eta(t)) < \delta$ , and*

$$\int_s^t L(r, \eta(r), u(r)) dr < \varphi^0(s, x) + \delta,$$

then

$$|(t, \eta(t)) - z^0| < a, \quad \|\eta - \xi^0\|_t < a, \quad \int_s^t |u(r) - u^0(r)| dr < a.$$

*Proof.* The conclusions can be restated as follows. Suppose that  $(s_m, x_m) \in N', m = 1, 2, \dots, (s_m, x_m) \rightarrow (s, x), (t_m, \eta_m(t_m)) \rightarrow z^1$  as  $m \rightarrow \infty$ , and

$$\liminf_{m \rightarrow \infty} \int_s^{t_m} L(t, \eta_m, u_m) dt \leq \lim_{m \rightarrow \infty} \varphi^0(s_m, x_m) = \varphi^0(s, x).$$

Here  $\eta_m$  satisfies (3.8) with  $u = u_m, \eta_m(s_m) = x_m$ , and  $z^1 \in \partial Q$ . We must show that  $z^1 = z^0$ , and as  $m \rightarrow \infty$ ,

$$\|\eta_m - \xi^0\|_{t_m} \rightarrow 0, \quad \int_s^{t_m} |u_m - u^0| dt \rightarrow 0.$$

Since  $K$  is compact we may assume (by taking subsequences) that  $u_m$  tends weak\* in  $L^\infty[s, T]$  to some  $u^1 \in \mathcal{U}_s$ . Then using (2.1), we have that  $\eta_m$  tends uniformly on  $[s, T]$  to the solution  $\eta^1$  of  $\dot{\eta}^1 = f(t, u^1, \eta^1)$  with  $\eta^1(s) = x$ . By lower semicontinuity of  $J^0$  with respect to weak\* convergence (using (2.2)) and the strict positivity of  $L$  in (2.3), we conclude that  $u^1$  is optimal. By uniqueness of  $u^0$ , (c) of the definition of strong regularity (Definition 3.1),

$$u^1(t) = u^0(t), \quad \eta^1(t) = \xi^0(t) \quad \text{for } s \leq t \leq \tau^0,$$

and  $z^1 = z^0$ . Since  $t_m \rightarrow \tau^0$ ,  $\|\eta_m - \xi^0\|_{t_m} \rightarrow 0$ . (So far we have used only the fact that each point in  $\bar{N}'$  is regular.)

To complete the proof, let

$$\begin{aligned} h(t, y) &= L(t, \xi^0(t), y) + p^0(t)f(t, \xi^0(t), y), \\ p^0(t) &= \varphi_x^0(t, \xi^0(t)). \end{aligned}$$

Then  $h(t, y)$  has a minimum on  $K$  at the unique point  $y = u^0(t)$ , except perhaps for a finite number of  $t = s_j$ ,  $j = 2, \dots, m$ , as in (d) of Definition 3.1. Let  $I$  be any compact subset of  $[s, \tau^0]$  not containing these  $s_j$ . Given  $\alpha > 0$  there exists  $\beta = \beta(\alpha)$  tending to 0 as  $\alpha \rightarrow 0$  such that

$$(3.9) \quad h(t, y) - h(t, u^0(t)) < \alpha \Rightarrow |y - u^0(t)| < \beta,$$

provided  $t \in I$ . By linearity (2.1) of  $f$  in  $y$  and convergence of  $t_m$  to  $\tau^0$ ,

$$\lim_{m \rightarrow \infty} \int_s^{t_m} p^0(t)f(t, \xi^0, u_m) dt = \int_s^{\tau^0} p^0(t)f(t, \xi^0, u^0) dt.$$

By the first part of the proof,

$$\lim_{m \rightarrow \infty} \int_s^{t_m} L(t, \xi^0, u_m) dt = \lim_{m \rightarrow \infty} \int_s^{t_m} L(t, \xi_m, u_m) dt = \int_s^{\tau^0} L(t, \xi^0, u^0) dt.$$

Therefore,

$$\lim_{m \rightarrow \infty} \int_s^{t_m} h(t, u_m(t)) dt = \int_s^{\tau^0} h(t, u^0(t)) dt.$$

Since  $h(t, u_m) \geq h(t, u^0)$  and  $t_m$  tends to  $\tau^0$ ,  $h(t, u_m) \rightarrow h(t, u^0)$  in measure on  $I$ . By (3.9),  $u_m$  tends to  $u^0$  in measure on  $I$ . Since  $u_m(t) \in K$ ,  $K$  is compact, and the measure of  $[s, \tau^0] - I$  can be made arbitrarily small,

$$\lim_{m \rightarrow \infty} \int_s^{t_m} |u_m - u^0| dt = 0.$$

**4. The stochastic problem; Preliminary estimates.** Throughout §§ 4–6 we assume that  $K$  is compact, and that (2.1), (2.2'), (2.3) hold. The main results of the present section are the estimates contained in Lemmas 4.4 and 4.5.

We begin with some notation and preliminary remarks.

(a) *The class  $\mathcal{Y}$  of control policies.* Let  $\mathcal{Y}$  consist of all functions  $Y$  from  $[T_0, T] \times R^n$  into  $K$  with the following property: for each  $T' < T$  there exist

positive constants  $C, \alpha$  (depending on  $Y$  and  $T'$ ) such that

$$|Y(s', x) - Y(s, x)| \leq C|s' - s|^\alpha,$$

$$|Y(s, x') - Y(s, x)| \leq C|x' - x|$$

for every  $x, x' \in R^n$  and  $T_0 \leq s, s' \leq T'$ .

(b) *The region  $Q$ .* To insure that the boundary problem (1.8<sup>ε</sup>)–(1.9) has a well-defined solution, we take  $Q$  an open subset of  $(T_0, T) \times R^n$  of a shape “suitable” from the viewpoint of parabolic equations. Thus we assume that

$$\partial Q = B_{T_0} \cup S \cup B_T,$$

where  $B_{T_0}, B_T$  are open subsets of the hyperplanes  $\{T_0\} \times R^n, \{T\} \times R^n$  respectively, and the “lateral boundary”  $S$  is a compact subset of a  $C^\infty$  manifold  $S_0$ . Moreover, if  $v = (v_0, v_1, \dots, v_n)$  is normal to  $S_0$  at a point of  $S$ , then  $v_i \neq 0$  for some  $i > 0$ . Let  $\mathcal{S}$  denote the class of all such  $Q$ .

The “essential” part of  $\partial Q$  is denoted by

$$\partial^*Q = S \cup B_T.$$

Let  $C^{1,2}(Q)$  denote the class of functions  $\psi$  bounded and continuous on  $\bar{Q}$  with  $\psi_s, \psi_{x_i}, \psi_{x_i x_j}$  continuous on  $Q, i, j = 1, \dots, n$ . If  $Q \in \mathcal{S}$ , then the function  $\varphi^\varepsilon$  in (1.6<sup>ε</sup>) is the unique solution of (1.8<sup>ε</sup>)–(1.9) in  $C^{1,2}(Q)$ . See [7, § 2], [11, Chap. 7] or [18]. Actually,  $\varphi_x^\varepsilon$  is Hölder continuous on  $\bar{Q}$ , and  $\varphi_s^\varepsilon, \varphi_{x_i x_j}^\varepsilon$  are Hölder continuous on any compact subset of  $\bar{Q} - (S \cap B_T), i, j = 1, \dots, n$ . The optimal policy  $Y^\varepsilon$  defined by (1.10) belongs to  $\mathcal{Y}$ ; more precisely,  $Y^\varepsilon$  has an extension for  $(s, x) \notin Q$  which belongs to  $\mathcal{Y}$ . See [7, p. 261].

If  $S$  is empty, then  $\partial^*Q$  is the hyperplane  $\{T\} \times R^n$ . The boundary data (1.9) are then Cauchy data.

We shall write

$$(4.1) \quad \Lambda^\varepsilon = \frac{\partial}{\partial s} + \varepsilon \Delta_x.$$

We also put  $f^\varepsilon = f^{Y^\varepsilon}, L^\varepsilon = L^{Y^\varepsilon}$ , where as in § 1

$$g^Y(s, x) = g(s, x, Y(s, x));$$

and we put  $H^\varepsilon = H(s, x, \varphi_x^\varepsilon)$ . Thus (1.8<sup>ε</sup>) becomes

$$0 = \Lambda^\varepsilon \varphi^\varepsilon + H^\varepsilon = \Lambda^\varepsilon \varphi^\varepsilon + f^\varepsilon \varphi_x^\varepsilon + L^\varepsilon.$$

(c) *Some probabilistic notions.* Let  $\Omega$  be the space of continuous functions  $\omega$  from  $[T_0, T]$  into  $R^n$ . We shall suppose that all processes which appear are defined on this sample space  $\Omega$ . Fix  $s \in [T_0, T]$ . Let  $B_{st}$  denote the least  $\sigma$ -algebra containing all sets  $\{\omega \in \Omega : \omega(r) \in A\}, s \leq r \leq t, A \subset R^n$  open. A process  $\xi$  on the time interval  $[s, T]$  is called *nonanticipative* if  $\xi(t)$  (meaning  $\xi(t, \cdot)$ ) is  $B_{st}$  measurable for  $s \leq t \leq T$ . A random variable  $\tau$  with  $s \leq \tau \leq T$  is a (nonanticipative) *stopping time* if the event  $\tau > t$  is  $B_{st}$  measurable,  $s \leq t \leq T$ . If  $\tau$  and  $\tau'$  are stopping times, then

$$\tau \wedge \tau' = \min(\tau, \tau')$$

is also a stopping time.

Let  $\xi$  be a continuous nonanticipative process on  $[s, T]$ , with values in  $R^n$ . Such a process is specified by a measure  $\pi$  on  $B_{sT}$ , if we set  $\xi(t, \omega) = \omega(t)$ ,  $\omega \in \Omega$ . In particular, let  $w$  be an  $n$ -dimensional Brownian motion on  $[s, T]$ , with  $w(s) = 0$ ; the associated  $\pi$  is a Wiener measure.

Let  $D \subset [T_0, T] \times R^n$ , with interior  $D_0$  not empty. If  $(s, x) \in D_0$ , let

$$\tau_D = \inf \{t : (t, \xi(t)) \notin D_0\}.$$

We call  $\tau_D$  the *exit time*; it is a stopping time. Thus in formula (1.4),  $\tau = \tau_D$ . We also write  $z_D = (\tau_D, \xi(\tau_D))$  for the exit place.

Given a control policy  $Y \in \mathcal{Y}$ , a continuous nonanticipative process  $\xi$  is determined by the system equations (1.1), the initial data (1.2), and (1.5). See [3, Chap. VI, § 3] or [12, Chap. 8]. The process  $\xi$  is in fact Markov and strongly Feller [5, Chaps. 11, 13], also [24].

If  $Y^e$  is the optimal policy, then write  $\xi^e$  for this process:

$$d\xi^e = f^e(t, \xi^e(t)) dt + (2\varepsilon)^{1/2} dw, \quad s \leq t \leq T,$$

with  $\xi^e(s) = x$ . The corresponding exit times and places are denoted by  $\tau_D^e, z_D^e$ .

(d) *Two estimates for maxima.* Suppose that

$$d\xi = f(t, \xi(t), u(t)) dt + (2\varepsilon)^{1/2} dw,$$

$$d\eta = f(t, \eta(t), u(t)) dt,$$

with  $\xi(s) = \eta(s)$ . By a standard estimate for ordinary differential equations,

$$(4.2) \quad \|\xi - \eta\|_t \leq C\varepsilon^{1/2} \|w\|_t, \quad s \leq t \leq T,$$

where  $C$  is a positive constant depending on  $T - T_0$  and a bound for  $|f_x|$ . (This notation  $\|\cdot\|_t$  was defined in § 3.)

Since  $w$  is a Brownian motion with  $w(s) = 0$ , the random variables  $w_1(t), \dots, w_n(t)$  are independent with mean 0 and variance  $t - s$ . For any  $\lambda > 0$ ,

$$\Pr \{ \|w\|_t \geq \lambda \} \leq 4n \Pr \left\{ w_i(t) \geq \frac{\lambda}{n} \right\} \leq \frac{4n^2(t-s)^{1/2}}{(2\pi)^{1/2}\lambda} \exp \left\{ -\frac{\lambda^2}{n^2(t-s)} \right\}.$$

See [3, p. 392]. Therefore for  $\lambda = e\varepsilon^{-1/2}$  we get, after writing  $\|w\|$  for  $\|w\|_T$ ,

$$(4.3) \quad \Pr \{ \varepsilon^{1/2} \|w\| \geq e \} \leq C\varepsilon^{1/2} \exp(-\beta/\varepsilon),$$

for suitable positive  $C, \beta$  depending on  $e$ .

(e) *A stochastic analogue of the fundamental theorem of calculus.* Suppose that  $\xi$  satisfies the stochastic differential equation

$$d\xi = b(t, \xi(t)) dt + (2\varepsilon)^{1/2} dw, \quad s \leq t \leq T,$$

with  $\xi(s) = x$ . Assume that  $b$  is continuous on  $[T_0, T] \times R^n$ , bounded, and satisfies a uniform Lipschitz condition in  $x$  on  $[T_0, T] \times R^n$  if  $T' < T$ . In particular, one can take  $b = f^Y$ , where  $Y \in \mathcal{Y}$ . When  $Y = Y^e$ ,  $b = f^e$ ; another useful choice will be  $b = b^e$ , where  $b^e$  is defined in § 6. Let  $(s, x) \in D$ , with  $D$  open, and  $\tau'$  a stopping time such that  $s \leq \tau' \leq \tau_D$ . If  $\psi \in C^{1,2}(D)$ , with  $\Lambda^e \psi + b\psi_x$  bounded, then

$$(4.4) \quad \psi(s, x) = -E \int_s^{\tau'} (\Lambda^e \psi + b\psi_x) dt + E\psi(z'),$$



where  $z' = (\tau', \xi(\tau'))$ . This is a consequence of the following well-known [12, p. 391] Ito formula for stochastic differentials of composites. Let  $\tilde{\psi} \in C^{1,2}(R^{n+1})$ . Then

$$(4.4') \quad d\tilde{\psi}(t, \xi(t)) = (\Lambda^e \tilde{\psi} + b\tilde{\psi}_x) dt + (2e)^{1/2} \tilde{\psi}_x dw.$$

To get (4.4) from (4.4'), let  $\bar{D}_1 \subset D$ . Take  $\tilde{\psi}$  such that  $\tilde{\psi} = \psi$  on  $\bar{D}_1$ . If  $\tau' \leq \tau_{D_1}$ , then we get (4.4) by integrating (4.4') between  $s$  and  $\tau'$  and taking expected values. From properties of stochastic integrals,

$$E \int_s^{\tau'} \tilde{\psi}_x dw = E \int_s^T \chi \tilde{\psi}_x dw = 0,$$

where  $\chi$  is the characteristic function of the interval  $[s, \tau')$ . Finally,  $D$  is the union of an expanding sequence  $D_1, D_2, \dots$  with  $\bar{D}_j \subset D$ ; and  $\tau'$  is the limit of the monotone sequence of stopping times  $\tau_j = \tau' \wedge \tau_{D_j}$ .

*Note.* Equation (4.4) is actually valid under the weaker assumption that  $\psi$  belongs to the class  $\mathcal{F}(D)$  defined in § 8. This can be shown by an approximation argument, which we shall carry out under the following somewhat stronger assumption than  $\psi \in \mathcal{F}(D)$ . Suppose that  $\psi \in C^1(\bar{D})$  with the second order partials  $\psi_{x_i x_j}$  bounded and continuous except on a finite number of  $C^\infty$  hypersurfaces  $\Sigma_2, \dots, \Sigma_m$ . (In the proofs of Theorems 6.4 and 6.5 we encounter such functions  $\psi$ .) Take  $\tilde{\psi}_1, \tilde{\psi}_2, \dots$  in  $C^{1,2}(R^{n+1})$  such that, for each  $i = 1, 2, \dots, \tilde{\psi}_j$  tends to  $\psi$  uniformly on  $\bar{D}_i$  as  $j \rightarrow \infty$  while  $\Lambda^e \tilde{\psi}_j + b(\tilde{\psi}_j)_x$  is bounded on  $D_i$  and tends to  $\Lambda^e \psi + b\psi_x$  on  $D_i - (\Sigma_2 \cup \dots \cup \Sigma_m)$ . For each  $t > s$  the probability distribution of  $\xi(t)$  is absolutely continuous with respect to Lebesgue measure. In fact, the probability density is a weak solution of the forward equation; see for instance [8]. Therefore,

$$\Pr \{ \xi(t) \in \Sigma_2 \cup \dots \cup \Sigma_m \} = 0$$

for almost all  $t$ , which implies

$$\lim_{j \rightarrow \infty} E \int_s^{\tau_j} [\Lambda^e \tilde{\psi}_j + b(\tilde{\psi}_j)_x] dt = E \int_s^{\tau_i} [\Lambda^e \psi + b\psi_x] dt.$$

From this we get (4.4) as before.

In the next lemma the notation is as follows. Let  $\mathcal{O}$  be an open set and  $D = \bar{\mathcal{Q}} \cap \mathcal{O}$ . We suppose that  $D \subset N$ , where  $N$  is a region of strong regularity; moreover,  $f^0 = f^{Y^0}$  is assumed to be defined outside  $N$  so that  $f^0 \in C^1(\mathcal{O} - D)$ . See (f) of Definition 3.1. (We do not require outside  $N$  that  $L + \varphi_x^0 f$  be minimum at  $Y^0$ .)

Let  $C_0$  be such that

$$|f^0(s, x') - f^0(s, x)| \leq C_0 |x' - x|$$

for  $(s, x), (s, x') \in \mathcal{O}$ . Let

$$\begin{aligned} d\xi^0 &= f^0(t, \xi^0(t)) dt, & s \leq t \leq t_1, \\ \xi^0(s) &= x, & (t, \xi^0(t)) \in D \text{ for } s \leq t \leq \tau^0, \\ & & (t, \xi^0(t)) \in \mathcal{O} - D \text{ for } \tau^0 < t \leq t_1. \end{aligned}$$

Let  $b$  and the process  $\xi$  be as in part (e) of this section,  $\xi(s) = x$ .

LEMMA 4.1. Given  $a > 0$  there exists  $e > 0$  with the following property: If  $|f^0 - b| < e$  in  $\mathcal{O}$  and  $\varepsilon^{1/2}\|w\|_{\tau_D} < e$ , then:

- (i)  $|\tau_D - \tau^0| < a, \quad (\tau_D, \zeta(\tau_D)) \in \partial^*Q.$
- (ii)  $\|\xi - \xi^0\|_{\tau_D} < a.$

*Proof.* Let  $\tilde{w}(t) = w(t \wedge \tau_D),$

$$\tilde{\xi}(t) = x + \int_s^t b(r, \tilde{\xi}(r)) dr + (2\varepsilon)^{1/2}\tilde{w}(t).$$

Then  $\|\tilde{w}\| = \|w\|_{\tau_D}$  and  $\tilde{\xi}(t) = \xi(t)$  for  $s \leq t \leq \tau_D.$  By standard estimates for ordinary differential equations,

$$|\tilde{\xi}(t) - \xi^0(t)| \leq 3e + C_0 \int_s^t |\tilde{\xi}(r) - \xi^0(r)| dr,$$

provided  $\varepsilon^{1/2}\|\tilde{w}\| < e$  and  $s \leq t \leq t_1 \wedge \tau_\emptyset.$  We may assume that  $t_1 \geq \tau^0 + a.$  For sufficiently small  $e,$  the hypotheses of the lemma imply

$$\begin{aligned} \tau_\emptyset > t_1, \quad \|\tilde{\xi} - \xi^0\|_{t_1} < a, \\ (t, \xi(t)) \in D \quad \text{for } s \leq t \leq \tau^0 - a, \\ (t, \tilde{\xi}(t)) \in \mathcal{O} - D \quad \text{for } \tau^0 + a \leq t \leq t_1. \end{aligned}$$

These inequalities imply (i) and (ii).

*Note.* Let  $\gamma^0$  be the trajectory corresponding to  $\xi^0$  on  $[s, \tau^0],$  as in § 3. Then  $e$  depends on  $a, C_0,$  and distance  $(\gamma^0, \partial\mathcal{O})$  but not on  $\varepsilon.$

The following lemma gives an upper estimate for  $\varphi^\varepsilon$  near a boundary point  $z$  in a region of strong regularity. The method is to define locally a solution  $\psi$  of a linear equation (4.5) which serves as a ‘‘barrier’’ at  $z.$  In Lemma 5.1 a corresponding lower estimate for  $\varphi^\varepsilon$  near  $z$  will be proved.

By *relative neighborhood* of  $z$  we mean a set  $M = \bar{Q} \cap \mathcal{O},$  where  $\mathcal{O}$  is open and contains  $z.$

LEMMA 4.2. Let  $N$  be a region of strong regularity and  $z \in N \cap \partial^*Q.$  Then given  $m_1 > 1$  there exist  $\varepsilon_0 > 0$  and a relative neighborhood  $M'$  of  $z$  such that  $\varphi^\varepsilon(s, x) \leq m_1\varphi^0(s, x)$  for all  $(s, x) \in M'$  and  $0 < \varepsilon < \varepsilon_0.$

*Proof.* Let  $M$  be a relative neighborhood of  $z$  such that  $M$  is of class  $\mathcal{S}$  and  $\bar{M} \subset N_1$  (where  $N_1$  is as in § 3). Let  $Y$  be some control policy such that  $Y = Y^0$  in  $M, Y \in \mathcal{Y},$  and  $Y \in C^1(\mathcal{O})$  for some open  $\mathcal{O}$  with  $M = \bar{Q} \cap \mathcal{O}.$  Let  $\psi \in C^{1,2}(M)$  satisfy

$$(4.5) \quad \Lambda^\varepsilon\psi + f^Y\psi_x + L^Y = 0 \quad \text{in } M$$

with  $\psi = 0$  on  $\partial^*M.$  Let  $\xi$  be the solution of (1.1) corresponding to  $Y$  with initial data  $(s, x) \in M.$  Since  $Y = Y^0$  in  $M,$  by (4.4),

$$(4.5') \quad \psi(s, x) = E \int_s^{\tau_M} L^0 dt, \quad (s, x) \in M.$$

Let

$$\tilde{\psi}(s, x) = J^\varepsilon(s, x; Y) = E \int_s^{\tau_\emptyset} L^Y dt.$$

By the definition (1.6<sup>ε</sup>),  $\varphi^\varepsilon \leq \tilde{\psi}$ . By (4.4) and the fact that  $\tilde{\psi}$  also satisfies (4.5),

$$\tilde{\psi}(s, x) = \psi(s, x) + E\tilde{\psi}(z_M), \quad (s, x) \in M.$$

Let  $\mu > 0$  and  $\varphi = (1 + \mu)\varphi^0$ . Then in  $M$ ,

$$(4.6) \quad \Lambda^\varepsilon \varphi + f^0 \varphi_x + L^0 = \varepsilon(1 + \mu)\Delta_x \varphi^0 - \mu L^0.$$

Since  $\Delta_x \varphi^0$  is continuous in  $\bar{M}$  and  $L^0 \geq c_1 > 0$  by (2.3), the right side of (4.6) is negative for small  $\varepsilon$ . Moreover,  $\varphi \geq 0$  and  $\psi = 0$  on  $\partial^*M$ . By the maximum principle for parabolic equations,  $\psi \leq \varphi$  on  $M$ . Moreover,  $\tilde{\psi}$  is bounded on  $Q$  by  $C = (T - T_0) \sup L$ , and  $\tilde{\psi} = 0$  on  $\partial^*Q$ . Thus

$$(4.7) \quad \tilde{\psi}(s, x) \leq \varphi(s, x) + C \Pr \{z_M \in \partial M - \partial^*Q\}.$$

To estimate the probability on the right side, we use Lemma 4.1 with  $b = f^0$ . Let  $a > 0$  and a relative neighborhood  $M'$  of  $z$  be chosen such that  $\bar{M}' \subset M$  and, for all  $(s, x) \in M'$ ,

$$\text{distance}(\gamma^0, \partial\mathcal{O}) \geq a.$$

As usual  $\gamma^0 = \gamma^0(s, x)$ . Thus, for some  $e > 0$ ,

$$\tilde{\psi}(s, x) \leq \varphi(s, x) + C \Pr \{\varepsilon^{1/2} \|w\|_{\tau_M} \geq e\}, \quad (s, x) \in M'.$$

Since  $\tau_M$  is a random (not fixed) time we cannot apply (4.3). However, let

$$\tilde{w}(t) = w(\tau_M \wedge t) = \int_s^t \chi(r) dw(r),$$

where  $\chi$  is the characteristic function of  $[s, \tau_M)$ . By properties of stochastic integrals [12, p. 383],

$$\begin{aligned} \Pr \{\|w\|_{\tau_M} \geq \lambda\} &= \Pr \{\|\tilde{w}\|_T \geq \lambda\} \\ &\leq \lambda^{-2} E \int_s^T \chi^2(r) dr = \lambda^{-2} E(\tau_M - s). \end{aligned}$$

Since  $L^0 \geq c_1 > 0$ , by (4.5')

$$E(\tau_M - s) \leq c_1^{-1} \psi \leq c_1^{-1} \tilde{\psi}.$$

Let  $\lambda = e\varepsilon^{-1/2}$ . For suitable  $C_1$ ,

$$\Pr \{z_M \in \partial M - \partial^*Q\} \leq C_1 \varepsilon e^{-2} \tilde{\psi}(s, x).$$

From this and (4.7),

$$(1 - C_2 \varepsilon) \tilde{\psi} \leq (1 + \mu)\varphi^0.$$

Since  $\varphi^\varepsilon \leq \tilde{\psi}$  we get Lemma 4.2.

*Note.* If  $A \subset N \cap \partial^*Q$  and  $A$  is compact, then the same kind of estimate in Lemma 4.2 holds in a relative neighborhood of  $A$ .

For making certain estimates it is a useful technical device to consider not only control policies, but also nonanticipative control processes. Such a process will be denoted by  $u$ . The system equations (1.1) with the initial data again have a well-defined solution  $\xi$ . The process  $\xi$  is nonanticipative, but not necessarily Markov.

Roughly speaking, a nonanticipative  $u$  may use past data in an arbitrary way, while a control policy  $Y$  uses at time  $t$  past data only through the current state  $\xi(t)$ . One cannot obtain better system performance using nonanticipative controls than with control policies. In fact, let  $Y^\varepsilon$  be the optimal control policy and  $\xi^\varepsilon$  as in part (c) of this section. Then

$$u^\varepsilon(t) = Y^\varepsilon(t, \xi^\varepsilon(t))$$

is an optimal nonanticipative control. See [7, p. 263]. This result is to be expected from the assumption of complete observability.

The following was proved in [7, p. 264].

LEMMA 4.3. *Let  $u$  be any nonanticipative control process and  $\tau'$  any stopping time with  $s \leq \tau' \leq \tau_Q$ . Then*

$$\varphi^\varepsilon(s, x) \leq E \int_s^{\tau'} L(t, \xi(t), u(t)) dt + E\varphi^\varepsilon(\tau', \xi(\tau')).$$

Equality holds if  $u = u^\varepsilon, \xi = \xi^\varepsilon$ .

The proof of the next lemma follows that for a corresponding result [9, Theorem 4a] for the Cauchy problem ( $\tau^0 = \tau^\varepsilon = T$ ). However, we need Lemmas 4.2 and 4.3 to deal with the case when exit may occur through the lateral part  $S$  of  $\partial^*Q$ . For the Cauchy problem the estimate for  $\varphi^\varepsilon - \varphi^0$  in Lemma 4.4 is valid globally, not merely in  $N'$ .

LEMMA 4.4. *Let  $N, N'$  be regions of strong regularity with  $\bar{N}' \subset N$ . There exist positive  $C$  and  $\varepsilon_0$  such that*

$$|\varphi^\varepsilon(s, x) - \varphi^0(s, x)| \leq C\varepsilon^{1/2},$$

$$\Pr \{z_N^\varepsilon \in \partial N - \partial^*Q\} \leq C\varepsilon^{1/2}$$

for all  $(s, x) \in N'$  and  $0 < \varepsilon < \varepsilon_0$ .

*Proof. Part 1.* Let  $(s, x) \in N'$ . The optimal open loop control  $u^0$  is suboptimal in the stochastic problem. The corresponding process  $\xi$  is defined by

$$d\xi = f(t, \xi(t), u^0(t)) dt + (2\varepsilon)^{1/2} dw$$

with  $\xi(s) = x$ . Let  $\tau_N$  be the exit time for that process, and

$$\tau' = \tau^0 \wedge \tau_N, \quad z' = (\tau', \xi(\tau')).$$

By Lemma 4.3,

$$\varphi^\varepsilon(s, x) \leq E \int_s^{\tau'} L(t, \xi, u^0) dt + E\varphi^\varepsilon(z').$$

Since  $L > 0$  and  $\tau' \leq \tau^0$ ,

$$\int_s^{\tau'} L(t, \xi^0, u^0) dt \leq \varphi^0(s, x),$$

$$\varphi^\varepsilon(s, x) \leq \varphi^0(s, x) + E \int_s^{\tau'} L \Big|_{\xi^0, u^0}^{\xi, u^0} dt + E\varphi^\varepsilon(z').$$

From (4.2) with  $u = u^0$ ,  $\eta = \xi^0$ , and boundedness of  $L_x$  the middle term on the right side is in absolute value  $\leq C\varepsilon^{1/2}$  (we use  $C$  for any large enough constant). By the note after Lemma 4.2 there is a relative neighborhood  $M$  of  $A = \bar{N}' \cap \partial^*Q$  such that  $\bar{M} \subset N$  and  $\varphi^\varepsilon \leq m_1\varphi^0$  for all  $(s, x) \in M$  and small  $\varepsilon$ . Since  $N'$  is a region of strong regularity,  $\gamma^0(s, x) \subset N'$ . Hence there exists  $e > 0$  (not depending on the choice of  $(s, x) \in N'$ ) such that  $\varepsilon^{1/2}\|w\| < e$  implies  $z' \in M$ . Define  $d(z)$  by (3.6). If  $z' \in M$ , then either

$$\tau' = \tau_N, \quad z' \in \partial^*Q, \quad d(z') = 0;$$

or else

$$\tau' = \tau^0, \quad d(z') \leq |\xi(\tau^0) - \xi^0(\tau^0)|.$$

In either case,  $d(z') \leq \|\xi - \xi^0\|_{\tau'}$ . From Lemma 4.2 we have for such  $w$ ,

$$\begin{aligned} \varphi^\varepsilon(z') &\leq m_1\varphi^0(z') \leq Cd(z') \leq C\varepsilon^{1/2}\|w\|; \\ E\{\varphi^\varepsilon(z'); \varepsilon^{1/2}\|w\| < e\} &\leq C\varepsilon^{1/2}. \end{aligned}$$

On the other hand,  $\varphi^\varepsilon$  is bounded on  $N$  (by  $(T - T_0) \sup L$ ). Hence from (4.3),

$$E\{\varphi^\varepsilon(z'); \varepsilon^{1/2}\|w\| \geq e\} \leq C\varepsilon^{1/2} \exp(-\beta/\varepsilon).$$

These two inequalities give

$$(4.8) \quad \varphi^\varepsilon(s, x) \leq \varphi^0(s, x) + C\varepsilon^{1/2}.$$

*Part 2.* To get an inequality opposite to (4.8) we use the fact that the optimal stochastic  $u^\varepsilon$  is suboptimal in the deterministic Pontryagin problem with probability 1. Again, let  $(s, x) \in N'$ . From the definitions of  $\xi^\varepsilon, u^\varepsilon$ ,

$$d\xi^\varepsilon = f(t, \xi^\varepsilon, u^\varepsilon) dt + (2\varepsilon)^{1/2} dw.$$

Define  $\eta^\varepsilon$  by

$$d\eta^\varepsilon = f(t, \eta^\varepsilon, u^\varepsilon) dt,$$

with  $\eta^\varepsilon(s) = \xi^\varepsilon(s) = x$ . Let  $\tau^\varepsilon = \tau_Q^\varepsilon, \tau_1^\varepsilon =$  exit time from  $Q$  of  $(t, \eta^\varepsilon(t))$ , and

$$\tau = \tau^\varepsilon \wedge \tau_1^\varepsilon, \quad z = (\tau, \eta^\varepsilon(\tau)).$$

With probability 1,

$$(4.9) \quad d(z) \leq \|\xi^\varepsilon - \eta^\varepsilon\|_\tau \leq C\varepsilon^{1/2}\|w\|,$$

$$(4.10) \quad \varphi^0(s, x) \leq \int_s^\tau L(t, \eta^\varepsilon, u^\varepsilon) dt + \varphi^0(z).$$

Let  $0 < a \leq \frac{1}{2}$  distance  $(N', \partial N - \partial^*Q)$ , and  $\delta$  be as in Lemma 3.1. Consider the three mutually exclusive events

$$\begin{aligned} B_1 &= \{C\varepsilon^{1/2}\|w\| \geq \delta\}, \\ B_2 &= \{C\varepsilon^{1/2}\|w\| < \delta, z_N^\varepsilon \in \partial N - M\}, \\ B_3 &= \{C\varepsilon^{1/2}\|w\| < \delta, z_N^\varepsilon \in M\}, \end{aligned}$$

where  $M$  is as in Part 1,  $C$  as in (4.9). Using (4.9) we also take  $\delta$  small enough that

$\|\eta^\varepsilon - \xi^0\|_\tau > a$  on  $B_2$ . The event  $\{z_N^\varepsilon \in \partial N - \partial^*Q\}$  is contained in  $B_1 \cup B_2$ . Thus

$$(4.11) \quad \Pr \{z_N^\varepsilon \in \partial N - \partial^*Q\} \leq \Pr \{B_1\} + \Pr \{B_2\}.$$

By (4.3) and boundedness of  $\varphi^0$ ,

$$E\{\varphi^0(z); B_1\} \leq C\varepsilon^{1/2} \exp(-\beta/\varepsilon).$$

By Lemma 3.1 with  $t = \tau, \eta = \eta^\varepsilon$ ,

$$(4.12) \quad \varphi^0(s, x) + \delta \leq \int_s^\tau L(t, \eta^\varepsilon, u^\varepsilon) dt \quad \text{on } B_2.$$

By (4.9),

$$\begin{aligned} \varphi^0(z) &\leq C\varepsilon^{1/2}\|w\| \quad \text{on } B_3, \\ E\{\varphi^0(z); B_3\} &\leq C\varepsilon^{1/2}. \end{aligned}$$

We now take expectations in (4.10), to obtain

$$\varphi^0(s, x) + \delta \Pr \{B_2\} \leq E \int_s^\tau L(t, \eta^\varepsilon, u^\varepsilon) dt + C\varepsilon^{1/2}.$$

From  $L > 0, \tau \leq \tau^\varepsilon, L_x$  bounded, and (4.2), we have

$$\begin{aligned} E \int_s^\tau L(t, \eta^\varepsilon, u^\varepsilon) dt &\leq E \int_s^\tau L(t, \xi^\varepsilon, u^\varepsilon) dt + C\varepsilon^{1/2} \\ &\leq E \int_s^{\tau^\varepsilon} L(t, \xi^\varepsilon, u^\varepsilon) dt + C\varepsilon^{1/2}. \end{aligned}$$

The first term on the right side is  $\varphi^\varepsilon(s, x)$ . Thus,

$$(4.13) \quad \varphi^0(s, x) + \delta \Pr \{B_2\} \leq \varphi^\varepsilon(s, x) + C\varepsilon^{1/2}.$$

This together with inequality (4.8) implies  $|\varphi^\varepsilon - \varphi^0| \leq C\varepsilon^{1/2}$  and  $\Pr \{B_2\} \leq C\varepsilon^{1/2}$ . By the (sharper) estimate (4.3) for  $\Pr \{B_1\}$  and (4.11), this completes the proof of Lemma 4.4.

In proving (4.13) we have also established

$$(4.14) \quad E \int_s^{\tau^\varepsilon} L(t, \xi^\varepsilon, u^\varepsilon) dt \leq C\varepsilon^{1/2}.$$

This will be used in proving the next lemma.

LEMMA 4.5 *Let  $N, N'$  be as in Lemma 4.4. Given  $a > 0$  there exist positive  $C$  and  $\varepsilon_0$  such that:*

- (i)  $\Pr \{|z_N^\varepsilon - z^0| \geq a\} \leq C\varepsilon^{1/2},$
- (ii)  $\Pr \{\|\xi^\varepsilon - \xi^0\|_{\tau^\varepsilon} \geq a\} \leq C\varepsilon^{1/2},$
- (iii)  $\Pr \left\{ \int_s^{\tau^\varepsilon} |u^\varepsilon - u^0| dt \geq a \right\} \leq C\varepsilon^{1/2},$

for all  $(s, x) \in N'$  and  $0 < \varepsilon < \varepsilon_0$ . Here  $\tau^\varepsilon = \tau_Q^\varepsilon$ .

*Proof.* We may suppose that  $0 < a < \text{distance}(N', \partial N - \partial^*Q)$ . Define  $B_1, B_2, B_3$  as in Part 2 of the preceding proof, but with  $\delta = \delta(a/2, N')$ . Let

$$B_4 = B_3 \cap \left\{ \int_s^\tau L(t, \eta^\varepsilon, u^\varepsilon) dt < \varphi^0(s, x) + \delta \right\}.$$

The same proof as for (4.13) shows that  $P\{B_3 - B_4\} \leq C\varepsilon^{1/2}$ . Since  $P\{B_1\} + P\{B_2\} \leq C\varepsilon^{1/2}$ , it suffices to consider the case when  $B_4$  holds. By Lemma 3.1, on  $B_4$ ,

$$|z - z^0| < \frac{a}{2}, \quad \|\eta^\varepsilon - \xi^0\|_\tau < \frac{a}{2}, \quad \int_s^\tau |u^\varepsilon - u^0| dt < \frac{a}{2}.$$

From the definitions of  $z_N^\varepsilon, z$ ,

$$\begin{aligned} |z_N^\varepsilon - z^0| &\leq |z_N^\varepsilon - z| + |z - z^0|, \\ |z_N^\varepsilon - z| &\leq |\tau_N^\varepsilon - \tau| + |\zeta^\varepsilon(\tau_N^\varepsilon) - \eta^\varepsilon(\tau)|. \end{aligned}$$

Moreover,  $\tau_N^\varepsilon = \tau^\varepsilon$  on  $B_3$  (hence on  $B_4$ ) and  $\tau \leq \tau^\varepsilon$ . Hence

$$|z_N^\varepsilon - z| \leq \tau^\varepsilon - \tau + \|\zeta^\varepsilon - \eta^\varepsilon\| + |\eta^\varepsilon(\tau^\varepsilon) - \eta^\varepsilon(\tau)|.$$

Since  $d\eta^\varepsilon/dt$  is bounded, we get on  $B_4$

$$|z_N^\varepsilon - z^0| \leq C(\tau^\varepsilon - \tau) + C\varepsilon^{1/2}\|w\| + \frac{1}{2}a.$$

Similarly, on  $B_4$ ,

$$\begin{aligned} \|\zeta^\varepsilon - \xi^0\|_{\tau^\varepsilon} &\leq C(\tau^\varepsilon - \tau) + C\varepsilon^{1/2}\|w\| + \frac{1}{2}a, \\ \int_s^{\tau^\varepsilon} |u^\varepsilon - u^0| dt &\leq C(\tau^\varepsilon - \tau) + \frac{1}{2}a \end{aligned}$$

for some positive  $C$ . Let  $\mu < a/(2C) - \delta/C$ . From (4.14) and (2.3),

$$\Pr \{ \tau^\varepsilon - \tau \geq \mu \} \leq (\mu c_1)^{-1} E \int_\tau^{\tau^\varepsilon} L(t, \zeta^\varepsilon, \mu^\varepsilon) dt \leq C\varepsilon^{1/2}.$$

This together with the estimate (4.3) proves Lemma 4.5.

**5. Convergence of  $\varphi_x^\varepsilon$  to  $\varphi_x^0$ .** This is proved uniformly on compact subsets of a region  $N$  of strong regularity (Lemma 5.5). As a consequence of this and assumption (2.2') the optimal control policy  $Y^\varepsilon$  tends in  $N$  to  $Y^0$  (Corollary 5.6). Like the estimates in Lemma 4.4, Lemma 5.5 is a preliminary step toward more precise results in § 6.

We begin with a lower estimate for  $\varphi^\varepsilon$  near boundary points to complement the upper estimate in Lemma 4.2.

**LEMMA 5.1.** *Under the hypotheses of Lemma 4.2, given  $m_2 < 1$  there exist  $\varepsilon_0 > 0$  and a relative neighborhood  $M'$  of  $z$  such that  $\varphi^\varepsilon(s, x) \geq m_2\varphi^0(s, x)$  for all  $(s, x) \in M'$  and  $0 < \varepsilon < \varepsilon_0$ .*

*Proof.* From equations (1.8 $^\varepsilon$ ), (1.8 $^0$ ), and (1.7),

$$\begin{aligned} \Lambda^\varepsilon \varphi^\varepsilon + f^\varepsilon \varphi_x^\varepsilon + L^\varepsilon &= 0 \quad \text{in } Q, \\ \Lambda^\varepsilon \varphi^0 + f^\varepsilon \varphi_x^0 + L^\varepsilon &\geq \varepsilon \Delta_x \varphi^0 \quad \text{in } N. \end{aligned}$$

(Recall § 4(b), (c) for notation.) Let  $\psi = \varphi^\varepsilon - \varphi^0$ . Then

$$\Lambda^\varepsilon \psi + f^\varepsilon \psi_x \leq -\varepsilon \Delta_x \varphi^0 \quad \text{in } N.$$

Let  $M$  be a relative neighborhood of  $z$  with  $\bar{M} \subset N$ . Set  $\tau' = \tau_M^\varepsilon, z' = z_M^\varepsilon$ . By (4.4) with  $\xi = \xi^\varepsilon$ ,

$$\psi(s, x) \geq \varepsilon E \int_s^{\tau'} \Delta_x \varphi^0 dt + E\psi(z'), \quad (s, x) \in M.$$

From Lemma 4.4,  $|\psi(z')| \leq C\varepsilon^{1/2}$ ; while by (1.9),  $\psi(z') = 0$  if  $z' \in \partial^*Q$ . Thus

$$(5.1) \quad \psi(s, x) \geq \varepsilon E \int_s^{\tau'} \Delta_x \varphi^0 dt - C\varepsilon^{1/2} \Pr \{z' \in \partial M - \partial^*Q\}.$$

Since  $\tau' \leq \tau_Q^\varepsilon$  and  $L \geq c_1$ ,

$$(5.2) \quad E \int_s^{\tau'} \Delta \varphi^0 dt \leq CE(\tau_Q^\varepsilon - s) \leq Cc_1^{-1} \varphi^\varepsilon(s, x).$$

From (1.1) we have the elementary estimate

$$\|\xi^\varepsilon - x\|_t \leq c(t - s) + (2\varepsilon)^{1/2} \|w\|_t, \quad s \leq t \leq T,$$

where  $|f| \leq c$ . Take  $a > 0$  and a relative neighborhood  $M'$  of  $z$  such that

$$\text{distance}((s, x), \partial M - \partial^*Q) \geq 2a \quad \text{if } (s, x) \in M'.$$

Then  $z' \in \partial M - \partial^*Q$  implies either  $\tau' - s \geq ac^{-1}$  or  $(2\varepsilon)^{1/2} \|w\|_{\tau'} \geq a$ . But

$$\Pr \{\tau' - s \geq ac^{-1}\} \leq a^{-1} c E\{\tau' - s\} \leq (ac_1)^{-1} c \varphi^\varepsilon(s, x),$$

$$\Pr \{(2\varepsilon)^{1/2} \|w\|_{\tau'} \geq a\} \leq C\varepsilon \varphi^\varepsilon(s, x)$$

by reasoning used to prove Lemma 4.2. Thus

$$(5.3) \quad \Pr \{z' \in \partial M - \partial^*Q\} \leq C\varphi^\varepsilon(s, x).$$

From (5.1), (5.2), (5.3) and Lemma 4.2,

$$\psi(s, x) \geq -C\varepsilon^{1/2} \varphi^0(s, x), \quad (s, x) \in M'.$$

This proves Lemma 5.1.

LEMMA 5.2. *Let  $N$  be a region of strong regularity. Then  $\varphi_x^\varepsilon \rightarrow \varphi_x^0, \varphi_s^\varepsilon \rightarrow \varphi_s^0$  as  $\varepsilon \rightarrow 0$ , uniformly on any compact set  $A \subset N \cap \partial^*Q$ .*

*Proof.* Since  $\varphi^\varepsilon = \varphi^0 = 0$  on  $N \cap \partial^*Q$ , their derivatives in directions tangent to  $\partial^*Q$  are 0 there. It suffices to show uniform convergence of the normal derivatives on  $A$ . Let  $v$  be the interior unit normal at  $z \in A$ . Then

$$\frac{\partial \varphi^\varepsilon}{\partial v} - \frac{\partial \varphi^0}{\partial v} = \lim_{(s,x) \rightarrow z} \frac{\varphi^\varepsilon(s, x) - \varphi^0(s, x)}{d(s, x)} = \left( \frac{\partial \varphi^0}{\partial v} \right)^{-1} \lim_{(s,x) \rightarrow z} \frac{\varphi^\varepsilon(s, x) - \varphi^0(s, x)}{\varphi^0(s, x)}.$$

According to Lemmas 4.2 and 5.1 the latter limit approaches 0 as  $\varepsilon \rightarrow 0$ , uniformly on  $A$ . This proves Lemma 5.2.

To prove convergence of  $\varphi_x^\varepsilon$  to  $\varphi_x^0$  at points of  $N - \partial^*Q$  we need a formula expressing  $\varphi_x^\varepsilon$  in terms of the processes  $u^\varepsilon$  and  $\xi^\varepsilon$ . Let us write  $f_x^\varepsilon(t, \xi) = f_x(t, \xi, u^\varepsilon)$ .



Given  $s$ , define the fundamental matrices  $W^\varepsilon = (W_{ij}^\varepsilon)$ ,  $i, j = 1, \dots, n$ , for each  $\varepsilon \geq 0$  by

$$(5.4) \quad dW^\varepsilon = f_x^\varepsilon(t, \xi^\varepsilon(t))W^\varepsilon dt, \quad W^\varepsilon(s) = \text{identity}.$$

Suppose that  $\tilde{\psi}_i \in C^{1,2}(R^{n+1})$  for  $i = 1, \dots, n$ , and let  $\zeta_i(t) = \tilde{\psi}_i(t, \xi^\varepsilon(t))$ . The collection of processes  $\xi^\varepsilon, \zeta, W^\varepsilon$  satisfies the stochastic differential equations (1.1) with  $u = u^\varepsilon$ , (5.4), and

$$(5.5) \quad d\zeta_i = (\Lambda^\varepsilon \tilde{\psi}_i + f^\varepsilon(\tilde{\psi}_i)_x) dt + (2\varepsilon)^{1/2}(\tilde{\psi}_i)_x dw.$$

Hence, we may compute the stochastic differential  $d(\zeta_i W_{ij}^\varepsilon)$  by the Ito rule, obtaining

$$(5.6) \quad d(\zeta_i W_{ij}^\varepsilon) = (d\zeta_i)W_{ij}^\varepsilon + \zeta_i dW_{ij}^\varepsilon.$$

The additional term which normally appears in the product rule for stochastic differentials is missing since no noise term appears in (5.4).

LEMMA 5.3. *Let  $\tau'$  be any stopping time,  $s \leq \tau' \leq \tau_Q^\varepsilon$ , and let  $z' = (\tau', \xi^\varepsilon(\tau'))$ . Then for  $\varepsilon > 0$  and  $(s, x) \in Q$ ,*

$$(5.7) \quad \varphi_x^\varepsilon(s, x) = E \int_s^{\tau'} L_x^\varepsilon W^\varepsilon dt + E\{\varphi_x^\varepsilon(z')W^\varepsilon(\tau')\},$$

where  $L_x^\varepsilon = L_x(t, \xi^\varepsilon(t), u^\varepsilon(t))$ .

*Proof.* Let  $P_i = \varphi_{x_i}^\varepsilon$  and differentiate (1.8 $^\varepsilon$ ) with respect to  $x_i$ , obtaining

$$(5.8) \quad \Lambda^\varepsilon P_i + H_{x_i}^\varepsilon + (P_i)_x H_p^\varepsilon = 0, \quad i = 1, \dots, n,$$

where  $H_x^\varepsilon, H_p^\varepsilon$  are evaluated at  $(s, x, P)$ . As noted in § 4(b),  $P_i$  is continuous on  $\bar{Q}$ . Using Lemma 2.3(b), we have that  $H_{x_i}^\varepsilon$  and  $(P_i)_x H_p^\varepsilon$  are Hölder continuous on any compact subset of  $Q$ ; hence so is  $\Lambda^\varepsilon P_i$ . This implies that  $P_i \in C^{1,2}(Q)$ . See [11] or [18]. Using (2.8) we rewrite (5.8) as

$$(5.8') \quad \Lambda^\varepsilon P_i + f^\varepsilon(P_i)_x = -(L_{x_i}^\varepsilon + P f_{x_i}^\varepsilon).$$

For the moment suppose that  $\tau' \leq \tau_D$ , where  $\bar{D} \subset Q$ . In (5.6) take  $\tilde{\psi}_i \in C^{1,2}(R^{n+1})$  such that  $\tilde{\psi}_i = P_i$  on  $\bar{D}$ . Then, for  $(s, x) \in D$ ,

$$P_i(s, x)W_{ij}^\varepsilon(s) = -E \int_s^{\tau'} \{[\Lambda^\varepsilon P_i + f^\varepsilon(P_i)_x]W_{ij}^\varepsilon dt + P_i dW_{ij}^\varepsilon\} + E\{P_i(z')W_{ij}^\varepsilon(\tau')\}.$$

Let us sum on  $i$ , and use (5.4), (5.8'),  $W(s) = \text{identity}$ . We get (5.7), in case  $\tau' \leq \tau_D$ . By writing  $Q$  as the union of an expanding sequence  $D_1, D_2, \dots$  and  $\tau'$  as the limit of the monotone sequence  $\tau' \wedge \tau_{D_j}$ , we then get Lemma 5.3.

LEMMA 5.4. *Let  $N$  be a region of strong regularity, and  $\bar{N}' \subset N$ . Then there exist positive  $C$  and  $\varepsilon_0$  such that  $|\varphi_{x_i}^\varepsilon| \leq C$  for  $(s, x) \in N'$  and  $0 < \varepsilon < \varepsilon_0$ .*

*Proof.* First of all, let  $1 < \rho < \frac{3}{2}$ . The following crude bound holds:

$$(5.9) \quad |\varphi_x^\varepsilon| \leq C\varepsilon^{-\rho}, \quad (s, x) \in \partial^*Q.$$

(This is rather similar to an estimate of Oleinik [21].) Clearly (5.9) holds on  $B_T$  since  $\varphi^\varepsilon = 0$  there; see § 4(b) for notation. Suppose that  $z_0 \in S$ . If  $v = (v_0, v_1, \dots, v_n)$

is the interior unit normal to  $S$  at  $z_0$ , then  $v_i \neq 0$  for some  $i > 0$ . Let  $z_1 = z_0 - rv$ , where  $r$  is chosen (simultaneously for all  $z_0$ ) small enough that the sphere  $|z - z_1| \leq r$  meets  $\bar{Q}$  only at  $z_0$ . Consider the function

$$\begin{aligned} \psi^\varepsilon(s, x) &= 1 - \exp[-\varepsilon^{-\rho}(|z - z_1|^2 - r^2)], \\ \mathcal{O} &= \{-\frac{1}{2}\varepsilon^\rho < |z - z_1|^2 - r^2 < \frac{1}{2}\varepsilon^\rho\}, \quad M = \mathcal{O} \cap \bar{Q}, \end{aligned}$$

where  $\varepsilon > 0$  is chosen small enough that  $x \neq x_1$  if  $z = (s, x) \in \mathcal{O}$ ,  $z_1 = (s_1, x_1)$ . An easy calculation shows that for small  $\varepsilon$ ,

$$\begin{aligned} \Lambda^\varepsilon \psi^\varepsilon + f^\varepsilon \psi^\varepsilon_x + L^\varepsilon &< 0 \quad \text{in } \mathcal{O}, \\ \psi^\varepsilon &= 1 - e^{-1/2} \quad \text{on } (\partial\mathcal{O}) \cap \bar{Q}. \end{aligned}$$

Since

$$\begin{aligned} \Lambda^\varepsilon \varphi^\varepsilon + f^\varepsilon \varphi^\varepsilon_x + L^\varepsilon &= 0 \quad \text{in } Q, \\ \Lambda^\varepsilon(\alpha\psi^\varepsilon - \varphi^\varepsilon) + f^\varepsilon(\alpha\psi^\varepsilon - \varphi^\varepsilon)_x &< 0 \quad \text{in } M, \end{aligned}$$

for any  $\alpha \geq 1$ . On  $\partial^*Q$ ,  $\varphi^\varepsilon = 0$ ,  $\psi^\varepsilon \geq 0$ ; while on  $\partial M - \partial^*Q$ ,  $\alpha\psi^\varepsilon - \varphi^\varepsilon \geq 0$  for suitable  $\alpha$ . By the maximum principle,  $\varphi^\varepsilon \leq \alpha\psi^\varepsilon$  in  $M$  for such  $\alpha$ . In  $Q$ ,  $\varphi^\varepsilon \geq 0$ ; while  $\varphi^\varepsilon = \psi^\varepsilon = 0$  at  $z_0$ . Hence

$$0 \leq \frac{\partial\varphi^\varepsilon}{\partial v} \leq \alpha \frac{\partial\psi^\varepsilon}{\partial v} \leq C\varepsilon^{-\rho}$$

at  $z_0$ . Since the derivatives of  $\varphi^\varepsilon$  in directions tangent to  $S$  are 0 at  $z_0$ , this implies (5.9).

Now let  $D_1, D_2, D_3$  be regions of strong regularity with

$$N' \subset D_1, \quad \bar{D}_1 \subset D_2, \quad \bar{D}_2 \subset D_3, \quad \bar{D}_3 \subset N.$$

Since  $z_Q^\varepsilon \in \partial^*Q - \bar{D}_3$  implies  $\tau_{D_3}^\varepsilon < \tau_Q^\varepsilon$ , we have by Lemma 4.4,

$$\Pr\{z_Q^\varepsilon \in \partial^*Q - \bar{D}_3\} \leq \Pr\{z_{D_3}^\varepsilon \in \partial D_3 - \partial^*Q\} \leq C_1\varepsilon^{1/2}$$

for  $(s, x) \in D_2$ . By Lemma 4.2,  $\varphi_x^\varepsilon$  is bounded on  $(\partial^*Q) \cap \bar{D}_3$ . Moreover,  $L_x^\varepsilon$  and  $W^\varepsilon$  are bounded. By Lemma 5.3 with  $\tau' = \tau_Q^\varepsilon$ ,

$$\begin{aligned} |\varphi_x^\varepsilon| &\leq C_2 + C_3 E\{|\varphi_x^\varepsilon(z_Q^\varepsilon)|; z_Q^\varepsilon \in \partial^*Q - \bar{D}_3\}, \\ |\varphi_x^\varepsilon| &\leq C_4\varepsilon^{-\rho+1/2} \end{aligned}$$

for  $(s, x) \in \bar{D}_2$  and sufficiently small  $\varepsilon$ . We now apply Lemmas 4.2 and 5.3 again, with  $\tau' = \tau_{D_2}^\varepsilon$ , and use the estimate just obtained when  $z' \in \partial D_2 - \partial^*Q$ . Thus  $|\varphi_x^\varepsilon| \leq C\varepsilon^{-\rho+1}$  when  $(s, x) \in \bar{D}_1$ . Since  $\rho < \frac{3}{2}$  another application of these lemmas gives  $|\varphi_x^\varepsilon| \leq C$ ,  $(s, x) \in N'$ .

LEMMA 5.5. Let  $N, N'$  be as in Lemma 5.4. Then  $\varphi_x^\varepsilon \rightarrow \varphi_x^0$  as  $\varepsilon \rightarrow 0$ , uniformly on  $N'$ .

Proof. Let  $(s, x) \in N'$ . For  $\varepsilon = 0$  formula (5.7) becomes

$$(5.7^0) \quad \varphi_x^0(s, x) = \int_s^{\tau^0} L_x^0 W^0 dt + \varphi_x^0(z^0)W^0(\tau^0).$$

Let  $D$  be a region of strong regularity with  $\bar{N}' \subset D$  and  $\bar{D} \subset N$ . Set  $\tau' = \tau_D^\varepsilon$ ,  $z' = z_D^\varepsilon$ . From Lemma 4.5,  $\|W^\varepsilon - W^0\|_{\tau'}$ ,  $L_x^\varepsilon - L_x^0$ , and  $z' - z^0$  tend to 0 in probability as  $\varepsilon \rightarrow 0$ , uniformly with respect to  $(s, x) \in N'$ . Moreover,  $L_x^\varepsilon, W^\varepsilon$  are bounded, while  $\varphi_x^\varepsilon(z')$  is bounded by Lemma 5.4. By Lemma 5.2,

$$\varphi_x^\varepsilon(z') - \varphi_x^0(z^0) = \varphi_x^\varepsilon(z') - \varphi_x^0(z') + \varphi_x^0(z') - \varphi_x^0(z^0)$$

tends uniformly to 0 in probability. Similarly,

$$|W^\varepsilon(\tau') - W^0(\tau^0)| \leq \|W^\varepsilon - W^0\|_{\tau'} + |W^0(\tau') - W^0(\tau^0)|$$

tends uniformly to 0 in probability. Using (5.7),  $\varphi_x^\varepsilon$  tends uniformly on  $N'$  to the right side of (5.7<sup>0</sup>) as  $\varepsilon \rightarrow 0$ . This proves Lemma 5.5.

The optimal control policy  $Y^\varepsilon$  is given by

$$(5.10) \quad Y^\varepsilon(s, x) = V(s, x, \varphi_x^\varepsilon(s, x)).$$

By Lemma 2.3,  $V$  is locally Lipschitz; thus Lemma 5.5 has the following corollary.

**COROLLARY 5.6.**  $Y^\varepsilon$  tends to  $Y^0$  uniformly on  $N'$ .

In § 8 we shall see that Lemma 5.5 remains true in cases where  $Y^\varepsilon$  may have switching surfaces. However, Corollary 5.6 is not true in such cases.

From Corollary 5.6 there is a much sharper estimate than the one in Lemma 4.4 for the probability that the optimal stochastic trajectory leaves  $N$  before time  $\tau_Q^\varepsilon$ .

**COROLLARY 5.7.** *There exist positive  $C, \beta, \varepsilon_0$  such that*

$$\Pr \{z_N^\varepsilon \in \partial N - \partial^*Q\} \leq C\varepsilon^{1/2} \exp(-\beta/\varepsilon)$$

for  $(s, x) \in N'$  and  $0 < \varepsilon < \varepsilon_0$ .

*Proof.* Let  $D$  be as in the proof of Lemma 5.5, and  $\mathcal{O}$  open with  $D = \mathcal{O} \cap \bar{Q}$ . By Corollary 5.6,  $Y^\varepsilon$  tends to  $Y^0$  uniformly on  $\bar{D}$ . We can define  $Y^\varepsilon$  outside  $Q$  so that  $Y^\varepsilon \in \mathcal{Y}$  and  $Y^\varepsilon$  tends to  $Y^0$  uniformly on  $\bar{\mathcal{O}}$ . In Lemma 4.1, let  $b = f^\varepsilon$ . There exists  $e > 0$  such that  $(s, x) \in N', |f^0 - f^\varepsilon| < e$  in  $\mathcal{O}$ , and  $\varepsilon^{1/2}\|w\| < e$  imply

$$z_D^\varepsilon = z_N^\varepsilon = z_Q^\varepsilon \in \partial^*Q.$$

We apply (4.3).

**6. Asymptotic formulas for  $\varphi^\varepsilon, \varphi_x^\varepsilon$ .** We are now ready to state the main results, Theorems 6.2, 6.4, 6.5. At the end of the section we indicate how the same methods tell the goodness of the policy  $Y^0$  in the stochastic problem.

As in §§ 4, 5 the control set  $K$  is assumed compact. We again assume (2.1), (2.3). Theorem 6.2 is true assuming the convexity condition (2.2). However, in Theorems 6.4, 6.5 we need the stronger condition (2.2'). The following observation will be used several times.

**LEMMA 6.1.** *Let  $G(s, x, p)$  be a real-valued function on  $N \times R^n$ , satisfying for any  $v$  a Lipschitz condition on  $N \times \{|p| \leq v\}$ . Let  $P, P^0$  be bounded, measurable on  $N$ . Then:*

(a) *There exists a bounded Borel measurable  $b$  such that*

$$G(s, x, P) - G(s, x, P^0) = b(s, x) \cdot (P - P^0).$$

(b) *If  $G$  is  $C^1$  in  $p$ , then we can take*

$$b(s, x) = \int_0^1 G_p(s, x, P^0 + \lambda(P - P^0)) d\lambda.$$

*Proof.* In (a) let

$$b(s, x) = \begin{cases} \frac{G(s, x, P) - G(s, x, P^0)}{|P - P^0|^2} (P - P^0) & \text{if } P \neq P^0, \\ 0 & \text{if } P = P^0. \end{cases}$$

Part (b) is just the integral form of the mean value theorem.

We now consider formula (1.11) with  $l = 1$ . For brevity we set  $\theta_1 = \theta$ .

**THEOREM 6.2.** *Let  $N$  be a region of strong regularity, and  $N'$  such that  $\bar{N}' \subset N$ . Let  $\theta$  be defined in  $N$  by*

$$(6.1) \quad \theta_s + f^0 \theta_x + \Delta_x \varphi^0 = 0$$

with  $\theta = 0$  on  $N \cap \partial^* Q$ . Then

$$(6.2) \quad \varphi^\varepsilon = \varphi^0 + \varepsilon \theta + o(\varepsilon),$$

where  $\varepsilon^{-1} o(\varepsilon) \rightarrow 0$  as  $\varepsilon \rightarrow 0$ , uniformly for  $(s, x) \in N'$ .

*Proof.* Let us assume (2.2'); changes in the proof needed when merely (2.2) holds will be indicated in § 8. By (3.4) with  $g = \Delta_x \varphi^0$ ,

$$(6.3) \quad \theta(s, x) = \int_s^{\tau^0} \Delta_x \varphi^0(t, \zeta^0(t)) dt.$$

Since (2.2') holds,  $\varphi^0 \in C^2(\bar{N})$  and  $\theta \in C^1(\bar{N})$ . Both  $\varphi^0$  and  $\theta$  are in  $C^\infty(\bar{N}_j)$  for  $j = 1, \dots, m$ . For  $\varepsilon > 0$  let

$$\theta^\varepsilon = \varepsilon^{-1}(\varphi^\varepsilon - \varphi^0).$$

We need to show that  $\theta^\varepsilon \rightarrow \theta$  uniformly on  $N'$  as  $\varepsilon \rightarrow 0$ . Define  $b^\varepsilon$  by Lemma 6.1(b) with  $G = H, P = \varphi_x^\varepsilon, P^0 = \varphi_x^0$ . Now  $P^0 \in C^1(\bar{N})$ , while  $P, P_x$  are continuous on  $\bar{N}$  and  $H_p$  is Lipschitz on compact sets by Lemma 2.3(b). Hence  $b^\varepsilon$  is continuous on  $\bar{N}$  and satisfies a uniform Lipschitz condition in  $x$ . We extend  $b^\varepsilon$  outside  $\bar{N}$  to have these same properties. Then the Ito conditions hold, insuring that the process  $\tilde{\zeta}^\varepsilon$  satisfying  $\tilde{\zeta}^\varepsilon(s) = x$  and

$$d\tilde{\zeta}^\varepsilon = b^\varepsilon(t, \tilde{\zeta}^\varepsilon(t)) dt + (2\varepsilon)^{1/2} dw, \quad s \leq t \leq T,$$

is well-defined.

Let  $\mathcal{O}, D$  be as in the proof of Corollary 5.7. By Lemma 5.5,  $b^\varepsilon \rightarrow f^0$  uniformly on  $\bar{D}$  as  $\varepsilon \rightarrow 0$ . We may arrange the above extension of  $b^\varepsilon$  outside  $\bar{N}$  such that also  $b^\varepsilon \rightarrow f^0$  uniformly on  $\bar{\mathcal{O}}$ . By definition of  $\theta^\varepsilon$  and (1.8 $^\varepsilon$ ), (1.8 $^0$ ),

$$(6.1^\varepsilon) \quad \Lambda^\varepsilon \theta^\varepsilon + b^\varepsilon \theta_x^\varepsilon + \Delta_x \varphi^0 = 0 \quad \text{in } N.$$

From (4.4) we get, writing  $\tau_D, z_D$  for the exit time and place for  $(t, \tilde{\zeta}^\varepsilon(t))$ ,

$$(6.3^\varepsilon) \quad \theta^\varepsilon(s, x) = E \int_s^{\tau_D} \Delta_x \varphi^0(t, \tilde{\zeta}^\varepsilon(t)) dt + E \theta^\varepsilon(z_D).$$

We apply Lemma 4.1 with  $b = b^\varepsilon$ . As  $\varepsilon \rightarrow 0, \|\tilde{\zeta}^\varepsilon - \zeta^0\|_{\tau_D}$  and  $\tau_D - \tau^0$  tend to 0 in probability, uniformly with respect to  $(s, x) \in N'$ . Moreover (as in the proof of Corollary 5.7),

$$\Pr \{z_D \in \partial D - \partial^* Q\} \leq C\varepsilon^{1/2} \exp(-\beta/\varepsilon)$$

for small  $\varepsilon$  and  $(s, x) \in N'$ . By (1.9),  $\theta^\varepsilon(z_D) = 0$  if  $z_D \in \partial^*Q$ . Since  $\varphi^\varepsilon, \varphi^0$  are bounded,

$$|\theta^\varepsilon(s, x)| \leq C\varepsilon^{-1}, \quad (s, x) \in Q.$$

As  $\varepsilon \rightarrow 0$ ,  $E\theta^\varepsilon(z_D) \rightarrow 0$ ; while the middle term in (6.3 $^\varepsilon$ ) tends to the right side of (6.3) uniformly on  $N'$ . This proves Theorem 6.2.

LEMMA 6.3. *Using the above notation  $\theta_s^\varepsilon \rightarrow \theta_s^0, \theta_x^\varepsilon \rightarrow \theta_x^0$  as  $\varepsilon \rightarrow 0$  uniformly on any compact set  $A \subset N \cap \partial^*Q$ .*

*Proof.* As in the proof of Lemma 5.2 it suffices to show that  $(\varphi^0)^{-1}(\theta^\varepsilon - \theta^0)$  is uniformly small on a relative neighborhood  $M'$  of any  $z \in A$  for small  $\varepsilon$ . Choose  $M, M'$  as in the proof of Lemma 4.2. By Lemma 4.1 there exists  $e > 0$  such that  $z_M \in \partial M - \partial^*Q$  implies  $\varepsilon^{1/2} \|\omega\|_{\tau_M} \geq e$  whenever  $(s, x) \in M'$  and  $|b^\varepsilon - f^0| < e$  on  $\mathcal{O}$ . Since  $b^\varepsilon$  tends to  $f^0$  uniformly on  $\mathcal{O}$ , we have (see proof of Lemma 4.2)

$$(6.4) \quad \Pr \{z_M \in \partial M - \partial^*Q\} \leq CE\{\tau_M - s\}, \quad (s, x) \in M'.$$

Let  $\omega^\varepsilon = \theta^\varepsilon - \theta$ . By (6.1), (6.1 $^\varepsilon$ ),

$$(6.5) \quad \begin{aligned} \Lambda^\varepsilon \omega^\varepsilon + b^\varepsilon \omega_x^\varepsilon + (b^\varepsilon - f^0)\theta_x + \varepsilon \Delta_x \theta &= 0 \quad \text{in } N, \\ \omega^\varepsilon(s, x) &= E \int_s^{\tau_M} [(b^\varepsilon - f^0) \cdot \theta_x + \varepsilon \Delta_x \theta] dt + E\omega^\varepsilon(z_M), \end{aligned}$$

the integrand being evaluated at  $(t, \tilde{z}^\varepsilon(t))$ . On  $N \cap \partial^*Q$ ,  $\omega^\varepsilon = 0$ . By Theorem 6.2,  $\omega^\varepsilon(z_M)$  is uniformly small for small  $\varepsilon$ ; moreover,  $\omega^\varepsilon(z_M) = 0$  if  $z_M \in \partial^*Q$ . Since  $\theta \in C^\infty(\bar{M})$ ,  $\theta_x$  and  $\Delta_x \theta$  are bounded there. By (6.4) and (6.5) it suffices to obtain for small  $\varepsilon$  an estimate

$$(6.6) \quad E\{\tau_M - s\} \leq C\varphi^0(s, x), \quad (s, x) \in M.$$

By (4.4),  $E\{\tau_M - s\} = \psi(s, x)$ , where  $\psi$  is the unique bounded solution in  $C^{1,2}(M)$  of

$$\Lambda^\varepsilon \psi + b^\varepsilon \psi_x + 1 = 0$$

with  $\psi = 0$  on  $\partial^*M$ . Let  $\varphi = C\varphi^0$  with  $C > c_1^{-1}$ ,  $c_1$  as in (2.3). Then

$$\Lambda^\varepsilon(\varphi - \psi) + b^\varepsilon(\varphi - \psi)_x + (f^0 - b^\varepsilon)\varphi_x - \varepsilon \Delta_x \varphi \leq -c < 0,$$

where  $c = Cc_1 - 1$ . For small  $\varepsilon$ ,  $|(f^0 - b^\varepsilon)\varphi_x - \varepsilon \Delta_x \varphi| < c$  in  $M$ . Since  $\varphi - \psi \geq 0$  on  $\partial^*M$ , the maximum principle implies  $\varphi \geq \psi$  in  $M$  for such  $\varepsilon$ . This establishes (6.6) and hence the lemma.

THEOREM 6.4. *Let  $N, N'$  be as in Theorem 6.2. Assume (2.2'). Then*

$$(6.7) \quad \varphi_x^\varepsilon = \varphi_x^0 + \varepsilon \theta_x + o(\varepsilon),$$

where  $\varepsilon^{-1}o(\varepsilon) \rightarrow 0$  as  $\varepsilon \rightarrow 0$ , uniformly for  $(s, x) \in N'$ .

*Proof.* In the notation above, (6.7) is equivalent to the statement  $\theta_x^\varepsilon \rightarrow \theta_x^0$  as  $\varepsilon \rightarrow 0$ . For notational simplicity let

$$\pi_i = \theta_{x_i}^\varepsilon = \varepsilon^{-1}(\varphi_{x_i}^\varepsilon - \varphi_{x_i}^0).$$

Recall (see proof of Lemma 5.3) that  $\varphi_{x_i}^\varepsilon \in C^{1,2}(\bar{N})$ ; while from § 3,  $\varphi_{x_i}^0 \in C^1(\bar{N})$  and  $\varphi_{x_i}^0 \in C^\infty(\bar{N}_j)$  for  $j = 1, \dots, m$ . Using (5.8) we get

$$(6.8^\varepsilon) \quad \Lambda^\varepsilon \pi_i + H_p^\varepsilon(\pi_i)_x + \varepsilon^{-1}(H_{x_i}^\varepsilon - H_{x_i}^0) + \varepsilon^{-1}(H_p^\varepsilon - H_p^0)\varphi_{xx_i}^0 + (\Delta_x \varphi^0)_{x_i} = 0,$$

where  $H_x^\varepsilon = H_x(s, x, \varphi_x^\varepsilon)$ , etc.

From Lemmas 6.1(a) and 2.3(b),

$$\varepsilon^{-1}(H_x^\varepsilon - H_x^0) = \pi B_1^\varepsilon, \quad \varepsilon^{-1}(H_p^\varepsilon - H_p^0) = \pi B_2^\varepsilon,$$

where the matrices  $B_i^\varepsilon(s, x)$  are bounded. Moreover, by Lemmas 6.1(b) and 2.3(c),  $B_1^\varepsilon \rightarrow H_{px}^0$ ,  $B_2^\varepsilon \rightarrow H_{pp}^0$  uniformly on any compact subset of  $N - (\Sigma_2 \cup \dots \cup \Sigma_m)$ . Define fundamental matrices  $X^\varepsilon$  by

$$dX^\varepsilon = (B_1^\varepsilon + B_2^\varepsilon \varphi_{xx}^0) X^\varepsilon dt, \quad X^\varepsilon(s) = \text{identity},$$

where  $B_i^\varepsilon, \varphi_{xx}^0$  are evaluated at  $(t, \xi^\varepsilon(t))$ . Let  $D$  be as in the proof of Theorem 6.2 and  $\tau' = \tau_D^\varepsilon, z' = z_D^\varepsilon$ . By Lemma 2.3,  $H_p^\varepsilon = f^\varepsilon$ . As in the proof of Lemma 5.3 we get, from (6.8<sup>ε</sup>),

$$\theta_x^\varepsilon(s, x) = E \int_s^{\tau'} (\Delta_x \varphi^0)_x X^\varepsilon dt + E\{\theta_x^\varepsilon(z') X^\varepsilon(\tau')\}$$

with  $(\Delta_x \varphi^0)_x$  evaluated at  $(t, \xi^\varepsilon(t))$ .

From previous estimates,

$$\Pr \{z' \in \partial D - \partial^* Q\} \leq C\varepsilon^{1/2} \exp(-\beta/\varepsilon)$$

for  $(s, x) \in N'$ . Since  $\varphi_x^\varepsilon, \varphi_x^0$  are bounded on  $D, |\theta_x^\varepsilon(z')| \leq C\varepsilon^{-1}$ . Using Lemma 4.1, with  $b = f^\varepsilon$ , Lemma 6.3, and previous reasoning we find that uniformly on  $N'$ ,

$$(6.9) \quad \lim_{\varepsilon \rightarrow 0} \theta_x^\varepsilon(s, x) = \int_s^{\tau^0} (\Delta_x \varphi^0)_x X^0 dt + \theta_x(z^0) X^0(\tau^0),$$

where  $(\Delta_x \varphi^0)_x$  is now evaluated at  $(t, \xi^0(t))$  and

$$\frac{dX^0}{dt} = (H_{px}^0 + H_{pp}^0 \varphi_{xx}^0) X^0, \quad t \geq s,$$

$$X^0(s) = \text{identity}.$$

Let us differentiate (6.1) with respect to  $x_i, i = 1, \dots, n$ , recalling that  $f^0 = H_p^0$ . We get the linear system of partial differential equations

$$(\theta_{x_i})_s + f^0(\theta_{x_i})_x + \theta_x(H_{px_i}^0 + H_{pp}^0 \varphi_{xx_i}^0) + (\Delta_x \varphi^0)_{x_i} = 0, \quad i = 1, \dots, n.$$

By the method of characteristics,  $\theta_x(s, x)$  equals the right side of (6.9). This proves Theorem 6.4.

Using Theorem 6.4 we can sharpen Theorem 6.2 by including the second order term in the expansion (1.11). For brevity we write  $\theta_2 = \chi$ .

**THEOREM 6.5.** *Keep the same hypotheses as Theorem 6.4. Let  $\chi$  be defined in  $N$  by*

$$(6.10) \quad \chi_s + f^0 \chi_x + \frac{1}{2} \theta_x H_{pp}^0 \theta_x + \Delta_x \theta = 0$$

with  $\chi = 0$  on  $N \cap \partial^* Q$ . Then

$$(6.11) \quad \varphi^\varepsilon = \varphi^0 + \varepsilon \theta + \varepsilon^2 \chi + o(\varepsilon^2),$$

where  $\varepsilon^{-2} o(\varepsilon^2) \rightarrow 0$  as  $\varepsilon \rightarrow 0$ , uniformly for  $(s, x) \in N'$ .

*Proof.* Let

$$\chi^\varepsilon = \varepsilon^{-1}(\theta^\varepsilon - \theta) = \varepsilon^{-2}(\varphi^\varepsilon - \varphi^0 - \varepsilon \theta).$$

Then (6.11) states that  $\chi^\epsilon \rightarrow \chi$  as  $\epsilon \rightarrow 0$ . Now

$$\Lambda^\epsilon \chi^\epsilon + b^\epsilon \chi_x^\epsilon + \epsilon^{-1}(b^\epsilon - f^0)\theta_x + \Delta_x \theta = 0.$$

Moreover,

$$\epsilon^{-1}(b^\epsilon - f^0) = \epsilon^{-1} \int_0^1 (H'_p - H_p^0) d\lambda,$$

$$H'_p = H_p(s, x, \varphi_x^0 + \lambda(\varphi_x^\epsilon - \varphi_x^0)).$$

Then  $\epsilon^{-1}(b^\epsilon - f^0)$  is bounded on  $D$  (using  $D$  as above), and by Theorem 6.4,

$$\lim_{\epsilon \rightarrow 0} \epsilon^{-1}(b^\epsilon - f^0) = \frac{1}{2} \theta_x H_{pp}^0$$

uniformly on any compact subset of  $N - (\Sigma_1 \cup \dots \cup \Sigma_m)$ . By the same proof as for Theorem 6.2 we find that

$$(6.12) \quad \lim_{\epsilon \rightarrow 0} \chi^\epsilon(s, x) = \int_s^{\tau^0} (\frac{1}{2} \theta_x H_{pp}^0 \theta_x + \Delta_x \theta) dt$$

uniformly on  $N'$ . By (6.10) and (3.4) the right side is  $\chi(s, x)$ . This proves Theorem 6.5.

One may ask about the validity of formula (1.11) for larger values of  $l$ . Consider first the case when merely (2.2) holds, as in Theorem 6.2, and take  $l = 2$ . Then  $Y^0$  may have switching surfaces. Across a switching surface,  $\theta_x$  and  $H_p^0 = f^0$  are generally discontinuous. It would seem that at a point of  $\gamma^0$  on a switching surface,  $\frac{1}{2} \theta_x H_{pp}^0 \theta_x + \Delta_x \theta$  should be interpreted in (6.12) as contributing a kind of Dirac delta function. (The author does not know which combination of right- and left-hand values is correct.) It, therefore, appears that  $\chi$  will be discontinuous across a switching surface; see Example 10.4, where this occurs. Since  $\chi^\epsilon$  is continuous, it cannot happen that  $\chi^\epsilon$  tends to  $\chi$  uniformly on  $N'$ . Thus Theorem 6.5 is not true without condition (2.2').

Similarly, when (2.2') holds consider  $l = 3$ . Then (1.11) is not correct in the sense of uniform convergence on  $N'$ . Let  $\zeta = \theta_3$ . The equation for  $\zeta$  corresponding to (6.1), (6.10) is

$$(6.13) \quad \zeta_s + f^0 \zeta_x + \frac{1}{2} \theta_x H_{pp}^0 \zeta_x + \frac{1}{6} H_{ppp}^0 \theta_x^3 + \Delta_x \zeta = 0,$$

where

$$H_{ppp}^0 \theta_x^3 = \sum_{i,j,r=1}^n H_{p_i p_j p_r} \theta_{x_i} \theta_{x_j} \theta_{x_r}.$$

The difficulty now occurs across transition surfaces, where  $\chi_x, H_{pp}^0$  are discontinuous.

In § 7 we deal with a case where these difficulties do not arise.

**Goodness of  $Y^0$  in the stochastic problem.** Suppose that, instead of seeking the optimal policy  $Y^\epsilon$  for the stochastic problem, one simply uses the optimal deterministic policy  $Y^0$ . It is plausible that, for small  $\epsilon$ ,  $Y^0$  should give approximately the optimum in the stochastic problem. The theorems just obtained, and their method of proof, put this rough statement on a quantitative basis.

Since  $Y^0$  has been defined only in  $\bar{N}$  (§ 3) we first extend  $Y^0$  outside  $N$  so that  $Y^0 \in \mathcal{Y}$ . However, the values  $Y^0(s, x)$  for  $(s, x) \notin N$  have no effect on the coefficients  $\theta, \tilde{\chi}$  in formula (6.14) below. Let

$$\Phi^\varepsilon(s, x) = J^\varepsilon(Y^0; s, x), \quad (s, x) \in N.$$

In particular,  $\Phi^0 = \varphi^0$ . For  $\varepsilon > 0$ ,  $\Phi^\varepsilon(s, x) - \varphi^\varepsilon(s, x)$  represents how much  $Y^0$  fails to be optimal in the stochastic problem. The function  $\Phi^\varepsilon$  satisfies the linear parabolic equation

$$(6.14^\varepsilon) \quad \Lambda^\varepsilon \Phi^\varepsilon + f^0 \Phi^\varepsilon_x + L^0 = 0 \quad \text{in } Q$$

with  $\Phi^\varepsilon = 0$  on  $\partial^*Q$ . The problem of expanding solutions of (6.14 $^\varepsilon$ ) in powers of  $\varepsilon$  is easier than for the nonlinear equation (1.8 $^\varepsilon$ ). Let us make the same assumptions as in Theorem 6.5. Then

$$(6.15) \quad \Phi^\varepsilon = \varphi^0 + \varepsilon\theta + \varepsilon^2\tilde{\chi} + o(\varepsilon^2)$$

uniformly for  $(s, x) \in N'$ . The first order coefficient  $\theta$  is the same as in (6.1). However,  $\tilde{\chi}$  satisfies, instead of (6.10), the equation

$$(6.16) \quad \tilde{\chi}_s + f^0 \tilde{\chi}_x + \Delta_x \theta = 0 \quad \text{in } N$$

with  $\tilde{\chi} = 0$  on  $N \cap \partial^*Q$ .

Let us indicate a proof of (6.15). Let

$$\tilde{\theta}^\varepsilon = \varepsilon^{-1}(\Phi^\varepsilon - \varphi^0), \quad \tilde{\chi}^\varepsilon = \varepsilon^{-2}(\Phi^\varepsilon - \varphi^0 - \varepsilon\theta).$$

From (1.8 $^0$ ), (6.1), (6.14 $^\varepsilon$ ),

$$\Lambda^\varepsilon \tilde{\theta}^\varepsilon + f^0 \tilde{\theta}^\varepsilon_x + \Delta_x \varphi^0 = 0, \quad \Lambda^\varepsilon \tilde{\chi}^\varepsilon + f^0 \tilde{\chi}^\varepsilon_x + \Delta_x \theta = 0.$$

Given  $(s, x) \in N'$  consider the process defined by

$$d\tilde{\zeta} = f^0(t, \tilde{\zeta}(t)) dt + (2\varepsilon)^{1/2} dw, \quad s \leq t \leq T,$$

with  $\tilde{\zeta}(s) = x$ . Corresponding to (6.3 $^\varepsilon$ ) we have

$$\tilde{\theta}^\varepsilon(s, x) = E \int_s^{\tau_D} \Delta_x \varphi^0(t, \tilde{\zeta}(t)) dt + E\tilde{\theta}^\varepsilon(z_D),$$

$$\tilde{\chi}^\varepsilon(s, x) = E \int_s^{\tau_D} \Delta_x \theta(t, \tilde{\zeta}(t)) dt + E\tilde{\chi}^\varepsilon(z_D).$$

By a simplification of the arguments above (§§ 4–6) it can be shown that  $\tilde{\theta}^\varepsilon \rightarrow \theta$ ,  $\tilde{\chi}^\varepsilon \rightarrow \tilde{\chi}$  as  $\varepsilon \rightarrow 0$ , uniformly on  $N'$ . This gives (6.15).

By comparing (6.11) and (6.15) we find that

$$(6.17) \quad \Phi^\varepsilon(s, x) - \varphi^\varepsilon(s, x) = -\frac{\varepsilon^2}{2} \int_s^{\tau_D} \theta_x H_{pp}^0 \theta_x dt + o(\varepsilon^2).$$

The integral is taken along the optimal trajectory  $\gamma^0(s, x)$  with initial endpoint  $(s, x)$ . Note that, from its definition (1.7),  $H$  is a concave function of  $p$ . Hence the integrand is nonpositive. Formula (6.17) shows that  $Y^0$  gives *within order the square of the noise variance  $2\varepsilon$  of the optimum if the initial data  $(s, x)$  belong to a region of strong regularity.*



If instead of (2.2') we assume the weaker condition (2.2), then instead of (6.17) we can show the weaker estimate  $0 \leq \Phi^\epsilon - \varphi^\epsilon \leq o(\epsilon)$ . This follows from (6.2) and the corresponding formula  $\Phi^\epsilon = \varphi^0 + \epsilon\theta + o(\epsilon)$ .

An even simpler control choice than  $Y^0$  in the stochastic problem would be the optimal open loop  $u^0$  corresponding to particular initial data  $(s_0, x_0)$ . Since  $u^0$  ignores the observations of the states  $\xi(t)$  obtained during system operation,  $u^0$  may be expected to behave more poorly than  $Y^0$  in the stochastic problem. To put this in quantitative terms, let

$$\Psi^\epsilon(s, x) = J^\epsilon(u^0; s, x);$$

then  $\Psi^0(s_0, x_0) = \varphi^0(s_0, x_0)$ , although generally  $\Psi^0(s, x) > \varphi^0(s, x)$  when  $(s, x)$  is not on the optimal trajectory  $\gamma^0(s_0, x_0)$ . By the same proof as for (6.15) we find that

$$(6.18) \quad \Psi^\epsilon = \Psi^0 + \epsilon\Theta + o(\epsilon)$$

uniformly for  $(s, x)$  in some neighborhood of  $\gamma^0(s_0, x_0)$ , provided  $(s_0, x_0) \in N$ . The function  $\Theta$  is defined by

$$\Theta_s + f(s, x, u^0(s))\Theta_x + \Delta_x\Psi^0 = 0$$

with  $\Theta = 0$  on  $N \cap \partial^*Q$ . (The second order coefficient in (6.18) can also be calculated in the same way as  $\tilde{\chi}$ , but is of no interest here.) Instead of (6.17) we get

$$(6.19) \quad \Psi^\epsilon(s_0, x_0) - \varphi^\epsilon(s_0, x_0) = \epsilon \int_{s_0}^{\tau_0} \Delta_x(\Psi^0 - \varphi^0) dt + o(\epsilon).$$

The integral is taken along  $\gamma^0(s_0, x_0)$ . Thus  $u^0$  gives within order  $\epsilon$  (rather than order  $\epsilon^2$ ) of the optimum in the stochastic problem.

We have considered the completely observable stochastic optimal control problem. It is an interesting open question to obtain approximate formulas for the optimal performance and control also for partially observable problems with small noise. In particular, one might study the open loop stochastic problem and the sampled data problem. For these problems necessary conditions for a minimum are known. See [17], [8].

**7. Higher order approximations; The Cauchy problem.** Let us now consider a situation where the approximations (1.11), (1.12) for  $\varphi^\epsilon, \varphi_x^\epsilon$  are valid for any finite  $l$  (Theorem 7.1). Later in the section this is interpreted as a result about the Cauchy problem in the theory of first order partial differential equations (Theorem 7.2).

The following assumptions are made:

$$(7.1) \quad K = R^k \quad (\text{there are no control constraints});$$

$$(7.2) \quad Q = (T_0, T) \times R^n \quad (\text{fixed stopping time } \tau = T);$$

$$(2.1) \quad f(s, x, y) = A(s, x) + B(s)y \quad \text{with } A_x \text{ bounded};$$

also (2.2') and in addition to (2.3'):

$$L \geq -c_0,$$

$$(2.3'') \quad L \text{ is bounded on } Q \times \{|y| \leq r\} \quad \text{for any } r;$$

$$|L_x| \leq c'L + c'' \quad \text{for suitable } c', c''.$$

By (7.1) the Pontryagin problem ( $\varepsilon = 0$ ) is now of Lagrange type in calculus of variations. Condition (7.2) means that the parabolic equation (1.8<sup>ε</sup>) is to be considered in the strip  $Q$  with Cauchy data  $\varphi^\varepsilon(T, x) = 0$ . Since the limits  $s$  and  $T$  of integration in (1.4) are now fixed, the problem is unchanged if  $L$  is replaced by  $L + c_0$ . Thus we may assume that  $L \geq 0$ .

By Lemma 2.4,  $H \in C^\infty$ . We also have, if  $\varepsilon > 0$  and  $N$  is a region of strong regularity,

$$\varphi^\varepsilon \in C^\infty(\bar{Q}), \quad Y^\varepsilon \in C^\infty(\bar{Q}), \quad \varphi^0 \in C^\infty(\bar{N}), \quad Y^0 \in C^\infty(\bar{N}).$$

(In the definition of § 3 we now take  $m = 1$ .)

Let us first verify the a priori bounds on  $Q$ :

$$(7.3) \quad |\varphi^\varepsilon| \leq C, \quad |\varphi_x^\varepsilon| \leq C,$$

for suitable  $C$ . Let  $Y_1$  be some bounded control policy (for instance,  $Y_1(s, x) \equiv 0$ ). Since  $L \geq 0$  we have from (1.6<sup>ε</sup>),

$$0 \leq \varphi^\varepsilon(s, x) \leq J^\varepsilon(s, x; Y_1),$$

and the right side is bounded. By Lemma 5.3 with  $\tau' = T$ ,

$$\varphi_x^\varepsilon(s, x) = E \int_s^T L_x^\varepsilon W^\varepsilon dt,$$

where  $\dot{W}^\varepsilon = A_x^\varepsilon W^\varepsilon$  by (2.1'). Hence  $W^\varepsilon(t)$  is bounded. From the bound on  $L_x$  in (2.3'') and on  $\varphi^\varepsilon$  we get a bound on  $\varphi_x^\varepsilon$ . Thus (7.3) holds. (These estimates also appear in [7, § 7] and [9].) Now

$$(7.4) \quad Y^\varepsilon(s, x) = V(s, x, \varphi_x^\varepsilon(s, x)).$$

From (2.3') and (7.3) this implies a bound  $|Y^\varepsilon| \leq r$  for some  $r$ . We get the same minimum for  $J^\varepsilon(s, x; Y)$  using the compact control set  $K = \{|y| \leq r\}$  instead of  $R^k$ .

**THEOREM 7.1.** *Let  $N$  be a region of strong regularity, and  $N'$  with  $\bar{N}' \subset N$ . Assume (7.1), (7.2), (2.1')–(2.3'), (2.3''). Then (1.11), (1.12) are valid for any  $l$ , uniformly for  $(s, x) \in N'$ . The coefficients  $\theta_l, l = 1, 2, \dots$ , in (1.11) satisfy linear first order partial differential equations, obtained by differentiating (1.8<sup>ε</sup>)  $l$  times with respect to  $\varepsilon$  and formally setting  $\varepsilon = 0$ . The Cauchy data are  $\theta_l(T, x) = 0$ .*

This theorem can be proved by induction on  $l$ . If we let

$$\theta_{l+1}^\varepsilon = \varepsilon^{-1}(\theta_l^\varepsilon - \theta_l) = \varepsilon^{-l-1}(\varphi^\varepsilon - \varphi^0 - \varepsilon\theta_1 - \dots - \varepsilon^l\theta_l),$$

then one shows first that  $\theta_{l+1}^\varepsilon \rightarrow \theta_{l+1}$  and then  $(\theta_{l+1}^\varepsilon)_x \rightarrow (\theta_{l+1})_x$ , in each case uniformly on  $N'$ . The reasoning is like that in § 6. The details are rather tedious, and we omit them. Since we are dealing with the Cauchy problem, the analogue of Lemma 6.3 is unnecessary. For all  $l$ ,

$$\theta_l^\varepsilon = (\theta_l^\varepsilon)_x = \theta_l = (\theta_l)_x = 0 \quad \text{when } s = T.$$

Formula (1.11) gives the Taylor expansion with remainder of  $\varphi^\varepsilon$  about  $\varepsilon = 0$ . One cannot let  $l \rightarrow \infty$  to obtain  $\varphi^\varepsilon$  as the sum of the corresponding power series in  $\varepsilon$ . The coefficients  $\theta_l(s, x)$  are obtained by integrating along the optimal trajectory  $\gamma^0(s, x)$  expressions computable from  $H$  and its partial derivatives. Thus  $\theta_l(s, x)$  depends only on  $H(s', x', p)$  for  $(s', x')$  in a neighborhood of  $\gamma^0(s, x)$ . However,

the solution  $\varphi^\varepsilon$  of the parabolic equation (1.8 $^\varepsilon$ ) depends on  $H(s', x', p)$  for  $(s', x')$  arbitrarily far from  $\gamma^0(s, x)$ .

From (1.12) we can get an approximate formula for the optimal control policy:

$$(7.5) \quad Y^\varepsilon = Y^0 + \varepsilon Z_1 + \varepsilon^2 Z_2 + \dots + \varepsilon^l Z_l + o(\varepsilon^l)$$

uniformly on  $N'$ . The coefficients  $Z_l$  are obtained by repeatedly differentiating (7.4) with respect to  $\varepsilon$  and then setting  $\varepsilon = 0$ . For instance,

$$(7.6) \quad Z_1 = -L_{yy}^0 \theta_x B.$$

*Example 7.1 (linear regulator).* Let  $f = A(s)x + B(s)y$ , and let  $L$  be quadratic in  $(x, y)$  with  $L > 0$  unless  $x = y = 0$ . Then  $\varphi^0$  is quadratic in  $x$ ,  $\theta_1 = \theta_1(s)$ ,  $\theta_l = 0$  for  $l \geq 2$ . The optimal policy  $Y^\varepsilon$  is a function linear in  $x$ , independent of  $\varepsilon$  (in particular,  $Y^\varepsilon = Y^0$ ).

In the more general problem considered in this section, the optimal policy  $Y^0$  can be found approximately in a neighborhood of an optimal trajectory  $\gamma^0(s, x) \subset N$  by solving a linear regulator-type problem. This is done by linearizing  $f$  and expanding  $L$  to quadratic terms about  $(t, \xi^0(t), u^0(t))$ . A solution to the secondary minimum problem is needed for the computations. See [1].

In the linear regulator example the assumption  $L$  bounded on  $Q \times \{|y| \leq r\}$  in (2.3'') is violated. The optimal policies  $Y^0(s, x)$  are linear in  $x$ , hence not bounded. However, these minor difficulties are inessential. It can be shown that Theorem 7.1 remains true if  $N'$  is assumed bounded and the boundedness assumption on  $L$  in question is replaced by certain polynomial-like growth assumptions [10, §§ 2, 5]. Under these growth assumptions, which hold in particular for the linear regulator,  $Y^\varepsilon$  satisfies a linear growth condition  $|Y^\varepsilon(s, x)| \leq C(1 + |x|)$ .

We conjecture that Theorem 7.1 is true without assumption (7.2). A crucial step seems to be necessary to establish an a priori bound  $|\varphi_x^\varepsilon| \leq C$  in  $N$ , like the one in Lemma 5.4 for the case of compact control set  $K$ .

**A result about partial differential equations.** Consider a Cauchy problem of the type

$$(7.7^0) \quad \varphi_s^0 + F(s, x, \varphi_x^0) = 0, \quad T_0 < s < T,$$

$$(7.8) \quad \varphi^0(T, x) = 0,$$

and the corresponding problem when  $\varepsilon > 0$ ,

$$(7.7^\varepsilon) \quad \varphi_s^\varepsilon + \varepsilon \Delta_x \varphi^\varepsilon + F(s, x, \varphi_x^\varepsilon) = 0, \quad T_0 < s < T,$$

$$(7.8^\varepsilon) \quad \varphi^\varepsilon(T, x) = 0.$$

There is a considerable literature dealing with generalized solutions  $\varphi^0$  of (7.7 $^0$ )–(7.8) and with the convergence of  $\varphi^\varepsilon$  to  $\varphi^0$  as  $\varepsilon \rightarrow 0$ . See references in [20], [9]. However, to the author's knowledge, an expansion of the type (1.11) has not been proved for nonlinear  $F$ . We now obtain such an expansion by inventing a control

problem for which (7.7<sup>0</sup>) is the Hamilton–Jacobi equation (1.8<sup>0</sup>) and then applying Theorem 7.1. The control problem will in fact be a simple problem in calculus of variations.

Let us make the following assumptions about  $F$  (see [9, § 2]):

$$(7.9a) \quad F \in C^\infty, \quad F(s, x, p) \leq C, \quad |F(s, x, 0)| \leq C_0,$$

$$\lim_{|p| \rightarrow \infty} \frac{F(s, x, p)}{|p|} = -\infty;$$

$$(7.9b) \quad |F_x| \leq c'(F - pF_p) + c'';$$

$$(7.9c) \quad |F_p| \leq R(|p|). \text{ For each } r > 0 \text{ there exists } v_r \text{ such that } |F_p| \leq r \text{ implies } |p| \leq v_r.$$

$$(7.9d) \quad C_1(|p|)|v|^2 \leq -vF_{pp}v \leq C_2(|p|)|v|^2 \quad \text{for all } v \in R^n,$$

where  $C_1(v), C_2(v)$  are positive and respectively nonincreasing, nondecreasing in  $v$ .

For a control problem we take the simplest dynamical equations and no control constraints:

$$\frac{d\xi}{dt} = u(t), \quad K = R^n.$$

Let  $L$  be dual to  $F$  in the sense of the classical canonical transformation:

$$(7.10) \quad L(s, x, y) = \max_{p \in R^n} (F - yp).$$

Then

$$(7.11) \quad F(s, x, p) = \min_{y \in R^n} (L + yp)$$

and (7.7<sup>ε</sup>) becomes (1.8<sup>ε</sup>) with  $H = F$ , for  $\varepsilon \geq 0$ .

Assumption (2.1') holds with  $A = 0, B = \text{identity}$ . Now  $L_{yy} = -(F_{pp})^{-1}$ , where  $y$  and  $p$  are related by the (classical) formulas  $y = F_p, p = -L_y$  according to (7.10) and (7.11). By assumptions (c) and (d), (2.2') is satisfied with  $c(r) = C_2(v_r)^{-1}$ . Finally, assumptions (a), (b), (c) imply (2.3'), (2.3'').

We now take

$$(7.12) \quad \varphi^0(s, x) = \min \int_s^T L(t, \xi(t), \dot{\xi}(t)) dt,$$

the minimum being taken among all Lipschitz  $\xi$  on  $[s, T]$  with  $\xi(s) = x$ . It is known that the unique bounded solution  $\varphi^\varepsilon \in C^{1,2}(\bar{Q})$  of (7.7<sup>ε</sup>)–(7.8<sup>ε</sup>) tends uniformly to  $\varphi^0$  as  $\varepsilon \rightarrow 0$  [9, p. 527]. Moreover,  $\varphi^0$  satisfies (7.7<sup>0</sup>) almost everywhere [7, p. 271].

Consider  $N$  such that:

$$(7.13a) \quad N = \mathcal{O}_0 \cap \bar{Q}, \quad \mathcal{O}_0 \text{ open};$$

$$(7.13b) \quad \varphi^0 \in C^\infty(\bar{N});$$

$$(7.13c) \quad (s, x) \in N \Rightarrow \gamma^0(s, x) \subset N.$$

Then  $N$  is a region of strong regularity. The uniqueness of the optimal trajectory  $\gamma^0(s, x)$  was proved in [9, p. 520] using an argument of Kuznetsov and Shishkin. The remaining conditions in the definition of § 3 automatically hold with  $m = 1, N_1 = N, \Sigma_1$  the hypersurface  $\{T\} \times R^n$ .

A given optimal  $\gamma^0(s, x)$  is contained in some such  $N$  provided  $(s, x)$  is neither a conjugate point nor an irregular point [9, p. 521]; also see the Appendix.

From Theorem 7.1 we then have the following theorem.

**THEOREM 7.2.** *Let  $F$  satisfy (7.9),  $N$  satisfy (7.13), and  $\bar{N}' \subset N$ . Then the conclusions of Theorem 7.1 hold (with  $H = F$ ).*

In (7.13) we assumed that  $\gamma^0(s, x)$  gives an absolute minimum, not merely a relative minimum. The method of characteristics is a classical method for constructing a smooth solution  $\varphi$  of (7.7<sup>0</sup>) in some  $N$ . (The characteristic ground curves are solutions of Euler's equation.) In this construction the characteristic ground curve  $\gamma(s, x)$  with initial point  $(s, x)$  gives a minimum among curves lying in  $N$ , but not always an absolute minimum. If the minimum is not absolute, then  $\varphi^\varepsilon$  tends in  $N$  not to  $\varphi$  but rather to  $\varphi^0$ .

**8. Discontinuous control policies.** Let us now impose the convexity condition (2.2) on  $L$  instead of the stronger condition (2.2'). For instance, (2.2) but not (2.2') is satisfied in the time-optimal problem ( $L \equiv 1$ ). One can no longer expect an optimal policy in the class  $\mathcal{Y}$  defined in § 4. Simple examples show that the optimal  $Y^\varepsilon$  may have discontinuities.

The purpose of this section is to give a correct formulation of the control problem in this case, and to outline the modifications needed in the proof of Theorem 6.2. (As noted at the end of the proof of Theorem 6.5, that theorem is incorrect without assumption (2.2').) As before we assume that the control set  $K$  is compact. Let  $\mathcal{Y}'$  consist of all Borel measurable  $Y$  from  $[T_0, T] \times R^n$  into  $K$ . For such  $Y$ , the standard (Ito) theory of stochastic differential equations no longer guarantees a solution  $\xi$  to the system equations

$$(1.1) \quad d\xi = f^Y dt + (2\varepsilon)^{1/2} dw,$$

since the drift coefficient  $f^Y$  is merely bounded and Borel measurable. However, for  $\varepsilon > 0$  the noise coefficient matrix in (1.1) is nonsingular; and we may use instead of the standard theory the work of Stroock and Varadhan [24].

Let  $b$  be bounded, Borel measurable, with values in  $R^n$ ; and take  $\Omega$  as in § 4. According to [24], a continuous nonanticipative process  $\xi$  is a solution of the stochastic differential equations

$$(8.1) \quad d\xi = b(t, \xi(t)) dt + (2\varepsilon)^{1/2} dw, \quad s \leq t \leq T,$$

with  $\xi(s) = x$  if, for each  $\alpha \in R^n$ ,

$$\exp \left[ \alpha \left( \xi(t) - x - \int_s^t b(r, \xi(r)) dr \right) - |\alpha|^2 \varepsilon (t - s) \right]$$

is a martingale with respect to the probability measure  $\Pi$  on  $\Omega$  associated with the process  $\xi$ . For any  $Y \in \mathcal{Y}'$ , this process  $\xi$  exists and is unique in the sense that  $\Pi$  is unique [24, § 6]. Moreover,  $\xi$  is a Markov process with the strong Feller property.

The process

$$w(t) = (2\varepsilon)^{-1/2} \left[ \xi(t) - x - \int_s^t b(r, \xi(r)) dr \right]$$

is a Brownian motion [24, Corollary 3.2]. Thus  $\xi$  is a solution of (1.1) in the ordinary sense. In particular, we may take  $b = f^Y, f \in \mathcal{Y}'$ .

Any  $Y \in \mathcal{Y}'$  can be approximated almost everywhere by a sequence  $Y_1, Y_2, \dots$  in  $\mathcal{Y}$ . It follows using facts about parabolic partial differential equations that

$$\lim_{j \rightarrow \infty} J^\varepsilon(s, x; Y_j) = J^\varepsilon(s, x; Y).$$

See [8, esp. Appendix]. Hence the  $\varphi^\varepsilon$  in (1.6 $^\varepsilon$ ) also satisfies

$$\varphi^\varepsilon(s, x) = \min_{Y \in \mathcal{Y}'} J^\varepsilon(s, x; Y).$$

As before  $\varphi^\varepsilon$  is a solution in  $C^{1,2}(Q)$  of (1.8 $^\varepsilon$ )–(1.9). An optimal  $Y^\varepsilon \in \mathcal{Y}'$  is determined, for almost all  $(s, x) \in Q$ , by (1.10).

The following modifications in the proof of Theorem 6.2 are needed. First of all, the only results in §§4 and 5 which must be restated are Lemma 4.1 and Corollary 5.6. For the latter, we can only assert that  $Y^\varepsilon$  tends to  $Y^0$  uniformly on any compact subsets of  $N - \Sigma$ , where  $\Sigma$  is the union of those  $\Sigma_j$  which are switching surfaces (§3). (It is an interesting open problem to prove statements about the perturbation of the optimal switching surfaces.) In applying Lemma 4.1 to prove Lemma 4.2,  $f^0 \in C^\infty(\mathcal{O})$ . However, in the proof of Theorem 6.2,  $f^0$  is discontinuous at points of  $\Sigma$ . Lemma 4.1 must hence be modified to require  $|f^0 - b| < \varepsilon$  except at those  $(s, x)$  with distance  $((s, x), \Sigma) < \varepsilon$ .

Let  $\mathcal{F}(D)$  denote the class of all  $\psi$  such that  $\psi, \psi_x$  are continuous and bounded on  $\bar{D}$  while  $\psi_s, \psi_{x_i x_j}, i, j = 1, \dots, n$ , are square integrable on compact subsets of  $D$ . We need formula (4.4) when  $b$  is bounded, Borel measurable,  $\psi \in \mathcal{F}(D)$  and  $\Lambda^\varepsilon \psi + b\psi_x$  is bounded on  $D$ . As in §4(e), this can be reduced to (4.4') with  $\tilde{\psi} \in \mathcal{F}(R^{n+1})$  and  $G = \Lambda^\varepsilon \tilde{\psi} + b\tilde{\psi}_x$  bounded. By standard smoothing methods there exist  $\tilde{\psi}^j \in C^{1,2}$  such that  $\tilde{\psi}^j, \tilde{\psi}_x^j$  tend uniformly to  $\tilde{\psi}, \tilde{\psi}_x$  as  $j \rightarrow \infty$ , while

$$\Lambda^\varepsilon \tilde{\psi}^j + b\tilde{\psi}_x^j = G^j$$

is uniformly bounded and tends almost everywhere to  $\Lambda^\varepsilon \tilde{\psi} + b\tilde{\psi}_x = G$ . By the Ito differential rule, (4.4') is correct for each  $\tilde{\psi}^j$  [12, p. 390]. From the fact that the random variable  $\xi(r)$  has probability density absolutely continuous with respect to Lebesgue measure [24],  $G^j(r, \xi(r))$  tends to  $G(r, \xi(r))$  with probability 1 for almost all  $r \in [s, T]$ . It follows that (4.4') holds for  $\tilde{\psi}$ .

We applied (4.4') in proving Lemma 5.1, also in Lemma 5.3 ( $\psi = P_i$ ) and Theorem 6.2 ( $\psi = \theta^\varepsilon$ ). By differentiating (1.8 $^\varepsilon$ ) with respect to  $x_i$  and using Lemma 2.2(a) together with results about parabolic equations [8, Appendix] or [18],  $P_i \in \mathcal{F}(Q)$ . Moreover,  $\varphi^0 \in \mathcal{F}(N)$  since  $N$  is a region of strong regularity. Hence each of these  $\psi \in \mathcal{F}(N)$ . In the first two instances,  $b = f^\varepsilon$ ; in the third,  $b = b^\varepsilon$  is defined using Lemma 6.1(a) rather than Lemma 6.1(b).

**9. Autonomous problems.** Consider the autonomous system

$$(9.1) \quad d\xi = f(\xi, u) dt + (2\varepsilon)^{1/2} dw, \quad t \geq 0,$$

with initial data  $\xi(0) = x$ . We suppose that  $x \in B$ , where  $B \subset R^n$  is open, bounded, and  $\partial B$  is a  $C^\infty$  manifold. Let  $\mathcal{Y}_a$  be the class of Lipschitz functions  $Y$  from  $R^n$  into the control set  $K$ . Such  $Y$  are called autonomous control policies. Let

$$(9.2) \quad \begin{aligned} J^\varepsilon(Y; x) &= E \int_0^{\tau_B} L[\xi(t), u(t)] dt, \\ u(t) &= Y[\xi(t)], \\ \varphi^\varepsilon(x) &= \inf_{Y \in \mathcal{Y}_a} J^\varepsilon(Y; x). \end{aligned}$$

The autonomous version of (1.8<sup>ε</sup>) is the elliptic equation

$$(9.3^\varepsilon) \quad \Delta \varphi^\varepsilon + H(x, \varphi_x^\varepsilon) = 0 \quad \text{in } B$$

with  $\varphi^\varepsilon(x) = 0$  on  $\partial B$ .

We are concerned with the analogues of the theorems in § 6. By an *autonomous region of strong regularity* we mean a set  $N = \bar{B} \cap \mathcal{O}$ , where  $\mathcal{O}$  is open,  $\mathcal{O} \subset R^n$ , and the autonomous versions of (a)–(g) in Definition 3.1 hold. Now  $\Sigma_1, \dots, \Sigma_m$  are hypersurfaces in  $R^n$ ,  $\gamma^0(x) = \{\xi^0(t) : 0 \leq t \leq \tau^0(x)\}$  is the unique optimal trajectory with initial point  $x$  for the problem  $J^0(u; x) = \min$ . Here  $\tau^0 = \tau^0(x)$ ,  $\xi^0(\tau^0)$  are the exit time and place from  $B$  for  $\gamma^0(x)$ .

Let  $Y \in \mathcal{Y}_a$  be a policy which agrees on  $N$  with  $Y^0$ , and let  $\tilde{\psi}(x) = J^\varepsilon(Y; x)$ . We assume that  $Y$  can be chosen so that  $\tilde{\psi}(x) \leq C\varepsilon^{-1}$  for all  $x \in B$ . Then the autonomous versions of the theorems in § 6 hold. For instance, let  $N$  be an autonomous region of strong regularity, and assume (2.1), (2.2'), (2.3). Then

$$(9.4) \quad \varphi^\varepsilon(x) = \varphi^0(x) + \varepsilon\theta(x) + \varepsilon^2\chi(x) + o(\varepsilon^2)$$

uniformly on  $N'$ , if  $\bar{N}' \subset N$ , where  $\theta, \chi$  are defined by the autonomous versions of (6.1), (6.10). If we assume only (2.2), then the autonomous version of Theorem 6.2 holds. As in § 8 the optimal policy  $Y^\varepsilon$  in that case belongs to the class of bounded measurable functions from  $R^n$  into  $K$ .

In an autonomous problem no bound is imposed on the exit time  $\tau_B$ . This necessitates a few changes in the proofs. It may not be true that  $\tilde{\psi}$  is bounded on  $B$  independent of  $\varepsilon$ . However, let  $\bar{D} \subset N$ . Then there is such a bound on  $\bar{D}$  as follows. Since  $\varphi^0$  is continuous on the compact set  $\bar{N}$ ,  $\varphi^0$  is bounded there. By (2.3),  $\tau^0(x) \leq c_1^{-1}\varphi^0(x)$ . Take  $T$  large enough that

$$T > \tau^0(x) + 1 \quad \text{for all } x \in N.$$

Let  $\tau' = \tau_N \wedge T$ , where  $\tau_N$  is the exit time for  $\xi$  satisfying  $d\xi = f^Y dt + (2\varepsilon)^{1/2} dw$ ,  $\xi(0) = x$ . Then

$$\tilde{\psi}(x) = E \int_0^{\tau'} L^0 dt + E\tilde{\psi}[\xi(\tau')], \quad x \in N.$$

The first term on the right side is bounded. By assumption,  $\tilde{\psi} \leq C\varepsilon^{-1}$  in  $B$ . Since  $\tilde{\psi} = 0$  on  $\partial B$ ,

$$E\tilde{\psi}[\xi(\tau')] \leq C\varepsilon^{-1} \Pr \{ \xi(\tau') \notin N \cap \partial B \}.$$

There exists  $e > 0$  such that  $x \in \bar{D}$  and  $\varepsilon^{1/2} \|w\|_T < e$  imply  $\tau' = \tau_N < T$  and  $\xi(\tau') \in N \cap \partial B$ . From (4.3),  $E\tilde{\psi}[\xi(\tau')]$  is then bounded, which implies  $\tilde{\psi} \leq C$  on  $\bar{D}$ . Since  $0 \leq \varphi^\varepsilon \leq \tilde{\psi}$ , this also bounds  $\varphi^\varepsilon$  on  $\bar{D}$ .

After that few further changes are needed. For  $Q$  we take the cylinder

$$Q = (0, T) \times B.$$

In Part 2 of the proof of Lemma 4.4 we now have

$$E \int_0^{\tau_D} L[\xi^\varepsilon(t), u^\varepsilon(t)] dt \leq \varphi^\varepsilon(x)$$

rather than equality. In all proofs in §§ 5, 6 where some exit time  $\tau_D$  appears we now take  $\tau' = \tau_D \wedge T$  rather than  $\tau' = \tau_D$ . Thus all limiting arguments as  $\varepsilon \rightarrow 0$  are carried out for  $\tau'$  and  $t$  in the fixed interval  $[0, T]$ .

**10. Examples.** We have found rather good approximate formulas for  $\varphi^\varepsilon(s, x)$  in terms of integrals along the optimal trajectory  $\gamma^0(s, x)$ , provided  $\gamma^0(s, x)$  lies in some region of strong regularity  $N$ . It is an essentially classical result (stated in the Appendix) that such an  $N$  exists if  $(s, x)$  is regular and nonconjugate. Let us now show by examples that the approximate formulas (1.11), (1.12) are not generally correct near points which are irregular or conjugate. From formula (A.9) in the Appendix it is clear that  $\Delta_x \varphi^0$  misbehaves at any conjugate point where  $\partial p / \partial \alpha$  is nonsingular. We consider first an example of the simplest type in calculus of variations.

*Example 10.1.* Let

$$J = \int_s^0 \frac{1}{2} \dot{\xi}(t)^2 dt + \Phi[\xi(0)],$$

$$\Phi(x) = -(x^2 + 1)^{1/2}.$$

This can be put in the form in § 7 with  $n = k = 1, T = 0, f(y) = y, L = \frac{1}{2}y^2 + \Phi'(x)y$ . The extremals are linear functions. Let

$$\xi(t, \alpha) = \alpha - t\Phi'(\alpha).$$

Let  $s(\alpha) = -(\alpha^2 + 1)^{1/2}$ , which implies  $\xi(s(\alpha), \alpha) = 0$ . The line segment

$$\gamma(\alpha) = \{(t, \xi(t, \alpha)) : s(\alpha) \leq t \leq 0\}$$

is an optimal trajectory. For any  $\delta > 0$  the set

$$N = \{(s, x) : s \leq 0, |x| > \delta \text{ or } -1 + \delta < s\}$$

is a region of strong regularity.

The point  $(-1, 0)$  is the initial endpoint of the horizontal segment  $\gamma(0)$ . Since  $\gamma(0)$  is the unique optimal trajectory with this initial endpoint,  $(-1, 0)$  is a regular point. Since  $\partial \xi / \partial \alpha = 0$  when  $\alpha = 0, t = -1, (-1, 0)$  is a conjugate point. (For  $\alpha \neq 0$  the solution  $t(\alpha)$  of the equation  $\partial \xi / \partial \alpha = 0$  also gives in the classical terminology a conjugate point. However,  $t(\alpha) < s(\alpha)$  for  $\alpha \neq 0$ . Since we use only the part  $\gamma(\alpha)$  of the line  $\xi = \xi(t, \alpha)$  which minimizes (globally among all curves with the same initial point) these other conjugate points are of no interest here.) Let

$$p(t, \alpha) = \Phi'(\alpha) - \Phi'(\xi(t, \alpha)).$$



From (A.9), for  $-1 < t \leq 0$ ,

$$\begin{aligned} \varphi_{xx}^0(t, \xi(t, \alpha)) &= \Phi''(\alpha) \left( \frac{\partial \xi}{\partial \alpha} \right)^{-1} - \Phi''(\xi(t, \alpha)), \\ \varphi_{xx}^0(t, 0) &= -(1+t)^{-1} + 1, \\ \theta(s, 0) &= \int_s^0 [-(1+t)^{-1} + 1] dt = \log(1+s) - s, \\ (10.1) \quad \varphi^\varepsilon(s, 0) &= \varphi^0(s, 0) + \varepsilon\theta(s, 0) + o(\varepsilon), \quad -1 < s \leq 0, \end{aligned}$$

the last equation being just (6.3). By (10.1) no estimate of the type  $|\varphi^\varepsilon - \varphi^0| \leq C\varepsilon$  can hold uniformly in a neighborhood of  $(-1, 0)$ .

For  $s < -1$ , two minimizing trajectories  $\gamma(\alpha), \gamma(-\alpha)$  issue from the point  $(s, 0)$ . The partial derivative  $\varphi_x^0$  has a jump discontinuity across the half-line consisting of such  $(s, 0)$ . This introduces difficulties with (1.11) and (1.12) which Example 10.2 will illustrate.

Example 10.1 is of a type studied by E. Hopf [14] in a pioneering paper on the small viscosity method for studying nonlinear conservation laws. The substitution  $v^\varepsilon = \varphi_x^\varepsilon$  in the present instance turns (7.7<sup>ε</sup>) into Burgers' equation with Cauchy data  $v^\varepsilon = \Phi'$  when  $s = 0$ . The same substitution generally converts (7.7<sup>ε</sup>) when  $n = 1$  into

$$v_s^\varepsilon + \varepsilon v_{xx}^\varepsilon + F(s, x, v^\varepsilon)_x = 0.$$

For various results on the behavior of  $v^\varepsilon$  as  $\varepsilon \rightarrow 0$  see [20]. Formula (1.12) gives an expansion of  $v^\varepsilon$  valid in any  $N$  such that the limit  $v^0$  of  $v^\varepsilon$  as  $\varepsilon \rightarrow 0$  belongs to  $C^\infty(\bar{N})$  and

$$v^0(s, x) = \int_{\gamma^0(s, x)} F_x dt + v^0(T, \xi^0(T)), \quad (s, x) \in N,$$

where the characteristic  $\gamma^0(s, x) = \{(t, \xi^0(t)): s \leq t \leq T\}$  is contained in  $N$ .

*Example 10.2.* Consider the following very simple autonomous, 1-dimensional, minimum time problem. Let

$$d\xi = u(t) dt + (2\varepsilon)^{1/2} dw,$$

with  $\xi(0) = x$ ,

$$-1 \leq x \leq 1, \quad -1 \leq u(t) \leq 1.$$

By straightforward calculations using (9.3<sup>ε</sup>), the minimum expected time  $\varphi^\varepsilon(x)$  to reach  $-1$  or  $1$  from  $x$  is

$$\begin{aligned} \varphi^\varepsilon(x) &= \varphi^0(x) + \varepsilon[\exp(-|x|/\varepsilon) - \exp(-1/\varepsilon)], \\ \varphi^{\varepsilon'}(x) &= \varphi^{0'}(x) - \operatorname{sgn} x \exp(-|x|/\varepsilon), \end{aligned}$$

where  $\varphi^0(x) = 1 - |x|$  is the minimum time for the deterministic problem. All coefficients  $\theta_j(x)$  in (1.11) are 0 for  $x \neq 0$ . Formally,  $-(\varphi^0)'$  is a Dirac delta function, which according to (6.3) suggests the value  $\theta_1(0) = 1$ . This makes the formula  $\varphi^{\varepsilon'}(0) = \varphi^{0'}(0) + \varepsilon\theta_1(0) + o(\varepsilon)$  correct. Of course, no formula  $\varphi^\varepsilon(x) = \varphi^0(x) + \varepsilon\theta_1(x) + o(\varepsilon)$  holds uniformly in a neighborhood of 0.

*Example 10.2'.* Let us extend Example 10.2 by taking  $x, \zeta(t), u(t)$  vectors in  $R^n, n \geq 2$ , with  $|\zeta(t)| \leq 1, |u(t)| \leq 1$ . Then  $B$  is the spherical ball  $|x| \leq 1$  in  $R^n$ . It is easily seen that

$$Y^\varepsilon(x) = |x|^{-1}x \quad \text{if } x \neq 0,$$

$$\varphi^\varepsilon(x) = g^\varepsilon(|x|), \quad \text{all } x \in B,$$

where  $g^\varepsilon(r)$  is decreasing on  $0 \leq r \leq 1, g^\varepsilon(1) = 0$ , and by (9.3 $\varepsilon$ ),

$$(10.2) \quad \varepsilon[(g^\varepsilon)'' + \frac{n-1}{r}(g^\varepsilon)'] + (g^\varepsilon)' + 1 = 0.$$

Since  $\varphi_x^\varepsilon$  is continuous at 0,  $(g^\varepsilon)'(0) = 0$ . For  $0 < r_0 < 1$ , the set

$$N = \{x : r_0 < |x| \leq 1\}$$

is a region of strong regularity. The point 0 is both conjugate and irregular. In  $N$  the expansion (1.11) holds with

$$\varphi^0(x) = g^0(|x|) = 1 - |x|, \quad \theta_j(x) = h_j(|x|),$$

$$h'_j + h''_{j-1} + \frac{n-1}{r}h'_{j-1} = 0, \quad h_j(1) = 0,$$

$$h_0(r) = g^0(r) = 1 - r,$$

$$h_1(r) = (n-1) \log r, \quad h_2(r) = (n-1)(n-2)(r^{-1} - 1).$$

For  $r = 0$  we proceed differently. Using the integrating factor  $r^{n-1}e^{-r/\varepsilon}$  to solve (10.2), we get

$$g^\varepsilon(0) = \int_0^1 r^{1-n}e^{-r/\varepsilon} dr \int_0^r \rho^{n-1}e^{\rho/\varepsilon} d\rho.$$

By using integrations by parts on the inner integral we get the approximation

$$g^\varepsilon(0) = g^0(0) + (n-1)\varepsilon \log \varepsilon + c\varepsilon + o(\varepsilon)$$

for suitable  $c$ .

*Example 10.3.* In this example Theorem 6.4 gives some information about transition boundaries for the optimal control policy  $Y^\varepsilon$ . Let  $n = k = 1$  (scalar state and control),  $K = [-1, 1]$ ,

$$f = Ax + By, \quad B > 0, \quad L = \frac{1}{2}Mx^2 + \frac{1}{2}y^2.$$

Take  $Q = (T_0, T) \times R^n$ , the Cauchy problem. Except for the control constraint  $|u(t)| \leq 1$  this would be a linear regulator problem. For  $\varepsilon = 0$  the solution is as follows. There are transition curves  $x = \pm G^0(s)$  such that:

$$Y^0(s, x) = \begin{cases} -1 & \text{for } x > G^0(s), \\ -g(s)x & \text{for } -G^0(s) < x < G^0(s), \\ 1 & \text{for } x < -G^0(s). \end{cases}$$

In the central region  $-G^0(s) < x < G^0(s)$ ,  $\varphi^0(s, x) = \frac{1}{2}g(s)x^2$  with  $g(s) > 0$  for  $s < T$ . The solution coincides there with that of the linear regulator. Obviously  $G^0(s) = g(s)^{-1}$ . The second order, but not third order, partial derivatives of  $\varphi^0$  are continuous across the transition curves. There are no conjugate or irregular points. The formulas (6.7), (6.11) hold everywhere in  $\bar{Q}$ .

In the stochastic problem,

$$Y^\varepsilon(s, x) = -\text{sat} [B\varphi_x^\varepsilon(s, x)],$$

where  $\text{sat } a = a$  if  $|a| \leq 1$ ,  $\text{sat } a = a|a|^{-1}$  if  $|a| > 1$ . From linearity of  $f$ , strict convexity of  $L$  in  $x, y$  and the optimality of  $Y^\varepsilon$  among nonanticipative controls, it can be seen that  $\varphi^\varepsilon$  is a strictly convex function of  $x$ . Thus  $\varphi_x^\varepsilon$  is a strictly increasing function of  $x$ ; and  $\varphi_x^\varepsilon$  tends uniformly to  $\varphi_x^0$  as  $\varepsilon \rightarrow 0$ . Therefore

$$\begin{aligned} Y^\varepsilon(s, x) &= -1 && \text{for } x > G^\varepsilon(s), \\ -1 < Y^\varepsilon(s, x) < 1 && \text{for } -G^\varepsilon(s) < x < G^\varepsilon(s), \\ Y^\varepsilon(s, x) &= 1 && \text{for } x < -G^\varepsilon(s), \end{aligned}$$

where  $G^\varepsilon(s) \rightarrow G^0(s)$  as  $\varepsilon \rightarrow 0$ . Since  $\varphi_x^\varepsilon(s, G^\varepsilon(s)) = \varphi_x^0(s, G^0(s)) = B^{-1}$ , we have

$$0 = \varepsilon^{-1}[\varphi_x^\varepsilon(s, G^\varepsilon) - \varphi_x^0(s, G^\varepsilon)] + \varepsilon^{-1}[\varphi_x^0(s, G^\varepsilon) - \varphi_x^0(s, G^0)].$$

Using Theorem 6.4 we get as  $\varepsilon \rightarrow 0$ ,

$$\theta_x + \varphi_{xx}^0 \frac{\partial G^\varepsilon}{\partial \varepsilon} \Big|_{\varepsilon=0} = 0$$

at  $(s, G^0(s))$ . Since  $\varphi_{xx}^0(s, G^0(s)) = g(s) > 0$ , the derivative  $\partial G^\varepsilon / \partial \varepsilon$  exists at  $\varepsilon = 0$ . However,  $\theta = \theta(s)$  in the central region (see linear regulator example in § 7). Since  $\theta$  is  $C^1$  across the transition curves,  $\theta_x(s, G^0(s)) = 0$ . We have shown that in this example,

$$\frac{\partial G^\varepsilon}{\partial \varepsilon} = 0 \quad \text{when } \varepsilon = 0.$$

A multidimensional steady state version of this example was studied in [25], where a computational method for generating reasonably good suboptimal controls was found.

We do not know how to apply the method of Example 10.3 to study perturbation of optimal switching surfaces. Both  $\theta_x$  and  $\varphi_{xx}^0$  are discontinuous across a switching surface for  $Y^0$ . It is not clear what combination of right- and left-hand values is needed. A related problem is to determine when the switching surface for  $Y^\varepsilon$  perturbs toward the terminal surface  $\Sigma_1$  and when away from  $\Sigma_1$  (notation in § 3).

In the following one-dimensional autonomous example, the optimal switching point can be found by direct calculation. Some two-dimensional examples have been studied in [4], [23].

*Example 10.4.* Consider the autonomous system

$$d\xi = u(t) dt + (2\varepsilon)^{1/2} dw, \quad \xi(0) = x,$$

with  $0 < a \leq u(t) \leq 1$ ,  $-b < x < 0$ . The problem is to maximize (rather than minimize)

$$J = E \int_0^\tau [\xi(t) + l(u(t))] dt,$$

where  $\tau$  is the exit time from  $B = (-b, 0)$  for  $\xi(t)$ . Let  $l_a = l(a)$ ,  $l_1 = l(1)$ . We choose  $a, b, l$  such that  $l(u) = c_1 u + c_2$  is linear and

$$0 < al_1 < l_a, \quad -(1 - a)b < al_1 - l_a.$$

Equation (9.3<sup>e</sup>) becomes

$$\varepsilon(\varphi^\varepsilon)'' + \max_{a \leq u \leq 1} [l(u) + (\varphi^\varepsilon)'u] + x = 0$$

with  $\varphi^\varepsilon(-b) = \varphi^\varepsilon(0) = 0$ . (For  $\varepsilon = 0$  no boundary condition is imposed at  $-b$ .) Define  $x^0$  by

$$(1 - a)x^0 = al_1 - l_a.$$

Then  $x^0$  is the optimal switching point in the deterministic problem:

$$Y^0(x) = \begin{cases} 1 & \text{if } x < x^0, \\ a & \text{if } x^0 < x \leq 0. \end{cases}$$

Any interval  $(x_1, 0]$  with  $-b < x_1$  is a region of strong regularity. When  $\varepsilon > 0$  the switching point  $x^\varepsilon$  is determined from the equation

$$l_a + (\varphi^\varepsilon)'a = l_1 + (\varphi^\varepsilon)'.$$

By straightforward calculations we find that

$$x^\varepsilon = x^0 + \varepsilon + O(\exp(-\beta/\varepsilon))$$

for some  $\beta > 0$ . In this example  $\Sigma_1 = \{0\}$ ,  $\Sigma_2 = \{x^0\}$ ; and  $x^\varepsilon$  perturbs toward  $\Sigma_1$ .

**Appendix.** We outline the method of characteristics for the Hamilton–Jacobi equation (1.8<sup>o</sup>), and how quantities needed in the asymptotic formulas (1.11), (1.12) can be computed from it.

The characteristic equations for (1.8<sup>o</sup>) form the following system of  $2n$  ordinary differential equations for functions  $\xi(t), p(t)$ :

$$\frac{d\xi}{dt} = H_p, \quad \frac{dp}{dt} = -H_x.$$

See [2, Chap. II]. Define the control function  $u(\cdot)$  by

$$(A.1) \quad L(t, \xi(t), y) + p(t)f(t, \xi(t), y) = \min \text{ on } K \quad \text{when } y = u(t).$$

By (2.8) the characteristic equations become

$$(A.2) \quad \frac{d\xi}{dt} = f, \quad \frac{dp}{dt} = -pf_x - L_x.$$

Equations (A.2) are the state and “costate” equations, and (A.1) the minimum condition in Pontryagin’s principle [13], [22]. Suppose that (A.1), (A.2) hold for  $\bar{\tau} \leq t \leq \tau$  with  $z = (\tau, \xi(\tau)) \in \Sigma_1$ , whereas in § 3,  $\Sigma_1$  is a  $C^\infty$  hypersurface,  $\Sigma_1 \subset \partial Q$ .

We impose the transversality condition:

$$(A.3) \quad \text{For } t = \tau, (L + pf, -p) \text{ is a normal vector to } \Sigma_1 \text{ at } z.$$

The trajectory  $\gamma = \{(t, \xi(t)) : \bar{\tau} \leq t \leq \tau\}$  is called a *characteristic ground curve*.

Suppose now that a triple  $\xi(t, \alpha), p(t, \alpha), u(t, \alpha)$  has been constructed for each vector  $\alpha = (\alpha_1, \dots, \alpha_n)$  in the closure  $\bar{\mathbf{U}}$  of some open set  $\mathbf{U}$  and for  $\bar{\tau}(\alpha) \leq t \leq \tau(\alpha)$ . Moreover, suppose functions  $\tau_j(\alpha), \xi_j(\alpha)$  have been found for  $j = 2, \dots, m$ , such that the following hold:

$$A1. \quad \begin{aligned} \tau_1(\alpha) &= \tau(\alpha), & \tau_{m+1}(\alpha) &= \bar{\tau}(\alpha), \\ \tau_j(\alpha) &> \tau_{j+1}(\alpha), & j &= 1, \dots, m, \quad \alpha \in \bar{\mathbf{U}}. \end{aligned}$$

A2. Let  $\mathcal{B}_j = \{(t, \alpha) : \tau_{j+1}(\alpha) < t < \tau_j(\alpha), \alpha \in \mathbf{U}\}$ . Then  $u \in C^\infty(\bar{\mathcal{B}}_j), j = 1, \dots, m$ . Moreover, (A.1)–(A.3) hold.

A3. If  $\tau_{m+1}(\alpha) < s < \tau_1(\alpha), \alpha \in \mathbf{U}$ , then  $u(\cdot, \alpha)$  is the unique optimal control for the initial data  $(s, \xi(s, \alpha))$ .

A4. There exist  $C^\infty$  hypersurfaces  $\Sigma_2, \dots, \Sigma_m$  such that the characteristic ground curve  $\gamma(\alpha)$  meets  $\Sigma_j$  nontangentially at the single point  $(\tau_j(\alpha), \xi_j(\alpha)), j = 1, \dots, m, \alpha \in \bar{\mathbf{U}}$ .

A5. The mapping  $(t, \alpha) \rightarrow (t, \xi(t, \alpha))$  is one-to-one and the matrices  $\partial \xi / \partial \alpha$  are nonsingular for  $\tau_{m+1}(\alpha) \leq t \leq \tau_1(\alpha), \alpha \in \bar{\mathbf{U}}$ .

A6. For  $t \neq \tau_j(\alpha), j = 2, \dots, m$ :

(a) the minimum in (A.1) is attained only for  $y = u(t, \alpha)$ ;

(b) the function  $V$  in Lemma 2.3 belongs to  $C^\infty(\Gamma)$  for some open set  $\Gamma$  containing  $(t, \xi(t, \alpha), p(t, \alpha))$ .

We allow  $u$  (and hence  $\dot{\xi} = \partial \xi / \partial t$ ) to be discontinuous when  $t = \tau_j(\alpha), j = 2, \dots, m$ . However,  $\xi$  and  $p$  are continuous there.

By the transversality condition (A.3) and (2.3), condition A4 holds when  $t = \tau_1$ . If  $\alpha$  is a set of local coordinates for  $\Sigma_1$ , then condition A5 holds when  $t = \tau_1$ . Let

$$\begin{aligned} N &= \{(t, \xi(t, \alpha)) : \tau_{m+1}(\alpha) < t \leq \tau_1(\alpha), \alpha \in \mathbf{U}\}, \\ N_j &= \{(t, \xi(t, \alpha)) : (t, \alpha) \in \mathcal{B}_j\}, \quad j = 1, \dots, m. \end{aligned}$$

Then  $N$  is a region of strong regularity (Definition 3.1). We have

$$(A.4) \quad \varphi^0(t, \xi(t, \alpha)) = \int_t^{\tau_1(\alpha)} L(r, \xi(r, \alpha), u(r, \alpha)) dr,$$

$$(A.5) \quad \varphi_x^0(t, \xi(t, \alpha)) = p(t, \alpha),$$

$$(A.6) \quad V(t, \xi(t, \alpha), p(t, \alpha)) = u(t, \alpha).$$

The following result gives a sufficient condition that a given optimal trajectory lie in some region of strong regularity  $N$ . Suppose that  $u^0$  is optimal for initial data  $(s, x)$  with corresponding  $\xi^0, p^0, \gamma^0 = \gamma^0(s, x)$  such that the necessary conditions (A.1)–(A.3) for a minimum hold. Suppose also that a triple  $\xi, p, u$  satisfying A1, A2, A4, A6 has been constructed, coinciding with  $\xi^0, p^0, u^0$  for  $\alpha = \alpha^0$  and  $s \leq t \leq \tau^0 = \tau_1(\alpha^0)$ . Here we suppose that  $\bar{\tau}(\alpha^0) < s < \tau_1(\alpha^0)$ . Following classical terminology we say that  $(t, \xi^0(t))$  is a *conjugate point* if the matrix  $\partial \xi / \partial \alpha$  is singular at  $(t, \alpha^0)$ .

Besides the general assumptions in § 2, let us assume (as in previous sections) either that  $K$  is compact or the conditions in § 7 hold.

**THEOREM A.1.** *If  $(s, x)$  is regular and not conjugate, then there exists a region of strong regularity  $N$  containing  $\gamma^0(s, x)$ .*

The proof uses classical reasoning in calculus of variations together with Lemma 3.1. See [9, Theorem 2]. In fact, one can take

$$N = \{(t, \xi(t, \alpha)) : \tau < t \leq \tau_1(\alpha), \alpha \in \mathfrak{U}_1\},$$

where  $\mathfrak{U}_1$  is some neighborhood of  $\alpha^0$  and  $\bar{\tau}(\alpha^0) < \tau < s$ ,  $\tau$  sufficiently near  $s$ .

Frequently, the required triple  $\xi, p, u$  is constructed in a neighborhood of  $\alpha^0$  as follows. Let  $\alpha$  be a set of local coordinates for  $\Sigma_1$  near  $z^0 = (\tau^0, \xi^0(\tau^0))$ ; points of  $\Sigma_1$  near  $z^0$  are of the form  $(\tau_1(\alpha), \xi_1(\alpha))$ . By Lemma 2.2,  $(L_u + pf_u)V_p = 0$ . Using this fact and the implicit function theorem,  $u(\tau_1(\alpha), \alpha), p(\tau_1(\alpha), \alpha)$  are determined near  $\alpha^0$  by (A.3), (A.6) and  $u^0(\tau^0)$ . Equations (A.2), (A.6) then determine  $\xi(t, \alpha), p(t, \alpha), u(t, \alpha)$ . Let  $\Delta$  denote a set of  $(s, x, p)$  such that  $V \in C^\infty(R^{2n+1} - \Delta)$ ; recall the examples in §§ 2, 10. Let us suppose that  $(t, \xi^0(t), p^0(t)) \in \Delta$  for  $t = s_j, j = 2, \dots, m$ . Moreover, suppose that, in a neighborhood of  $(s_j, \xi^0(s_j), p^0(s_j))$ ,  $\Delta$  is described by an equation  $U_j(s, x, p) = 0$  with  $U_j \in C^\infty$  and

$$\frac{d}{dt}U_j(t, \xi^0(t), p^0(t)) \neq 0 \quad \text{for } t = s_j.$$

Then  $\tau_j(\alpha)$  is determined for  $\alpha$  near  $\alpha^0$  by

$$U_j(t, \xi(t, \alpha), p(t, \alpha)) = 0 \quad \text{for } t = \tau_j(\alpha),$$

with  $s_j = \tau_j(\alpha^0)$ . Let

$$(A.7) \quad \xi_j(\alpha) = \xi(\tau_j(\alpha), \alpha), \quad j = 2, \dots, m.$$

Assume that  $\partial \xi / \partial \alpha$  is nonsingular for  $\alpha = \alpha^0, s \leq t \leq \tau^0$  (no conjugate points). Then

$$(A.8) \quad \frac{\partial \xi_j}{\partial \alpha} = \xi^0 \frac{\partial \tau_j}{\partial \alpha} + \frac{\partial \xi^0}{\partial \alpha}.$$

This implies that the vectors  $(\partial \tau_j / \partial \alpha_\lambda, \partial \xi_j / \partial \alpha_\lambda), \lambda = 1, \dots, n$ , are linearly independent when  $\alpha = \alpha^0$ . Hence,  $(\tau_j(\alpha), \xi_j(\alpha))$  lies in a  $C^\infty$  manifold  $\Sigma_j$  of dimension  $n$ , for  $\alpha$  near  $\alpha^0$ ; and the trajectory  $\gamma(\alpha)$  is not tangent to  $\Sigma_j$  at  $(\tau_j(\alpha), \xi_j(\alpha))$ .

According to (6.3) to find the coefficient  $\theta(s, x)$  in the approximation (6.2) we need to calculate  $\Delta_x \varphi^0$  along  $\gamma^0(s, x)$ . Now  $\Delta_x \varphi^0$  is the trace of the matrix  $\varphi_{xx}^0$ . By differentiating (A.5) with respect to  $\alpha$  we get the matrix equation

$$(A.9) \quad \varphi_{xx}^0 = \frac{\partial p}{\partial \alpha} \left( \frac{\partial \xi}{\partial \alpha} \right)^{-1}.$$

To find  $\theta_x$ , which appears in the approximation (6.7), we can use (6.9). By differentiating (A.5) again with respect to  $\alpha$  we get

$$(A.10) \quad (\varphi_{x_i x_x}^0)_{xx} \frac{\partial \xi}{\partial \alpha_\lambda} \frac{\partial \xi}{\partial \alpha_\mu} + (\varphi_{x_i x_x}^0)_x \frac{\partial^2 \xi}{\partial \alpha_\lambda \partial \alpha_\mu} = \frac{\partial^2 p_i}{\partial \alpha_\lambda \partial \alpha_\mu}, \quad i, \lambda, \mu = 1, \dots, n,$$

from which  $(\Delta_x \varphi^0)_x$  is expressed using (A.9) in terms of first and second order  $\alpha$  partial derivatives of  $\xi$  and  $p$ . These derivatives obey ordinary differential equations obtained by taking  $\partial/\partial\alpha_\lambda$ ,  $\partial^2/\partial\alpha_\lambda\partial\alpha_\mu$  in (A.2) and in the formula (A.6). From (6.3), with  $x = \xi(s, \alpha)$  and  $\tau^0 = \tau_1(\alpha)$ , we have the data

$$\theta_x = (\Delta_x \varphi^0) \left( \frac{\partial \tau_1}{\partial \alpha} \right) \left( \frac{\partial \xi}{\partial \alpha} \right)^{-1} \quad \text{at } t = \tau_1(\alpha).$$

Let  $p_j(\alpha) = p(\tau_j(\alpha), \alpha)$ . Then

$$\frac{\partial p_j}{\partial \alpha} = \dot{p} \frac{\partial \tau_j}{\partial \alpha} + \frac{\partial p}{\partial \alpha}.$$

This, together with (A.2) and (A.8), gives a relation between right- and left-hand values of  $\partial \xi / \partial \alpha$ ,  $\partial p / \partial \alpha$  at  $\tau_j(\alpha)$ ,  $j = 2, \dots, m$ , if  $\Sigma_j$  is a switching surface. If  $\Sigma_j$  is a transition surface, then  $\xi$ ,  $\dot{p}$ ,  $\partial \xi / \partial \alpha$ ,  $\partial p / \partial \alpha$  are continuous at  $\tau_j(\alpha)$ . Another differentiation in  $\alpha$  gives a relation between right- and left-hand values of  $\partial^2 \xi / \partial \alpha_\lambda \partial \alpha_\mu$ ,  $\partial^2 p / \partial \alpha_\lambda \partial \alpha_\mu$  at  $\tau_j(\alpha)$ . If  $\alpha$  is taken as a set of local coordinates for  $\Sigma_1$ , then data for all of these quantities at time  $\tau_1(\alpha)$  are determined from (A.3) and (A.6).

To calculate the coefficient  $\chi$  in (6.11) from (6.12), one needs  $\Delta_x \theta$  as well as  $\theta_x$  along  $\gamma^0(s, x)$ . However,

$$\theta_{xx} = \frac{\partial q}{\partial \alpha} \left( \frac{\partial \xi}{\partial \alpha} \right)^{-1},$$

where

$$q(t, \alpha) = \theta_x(t, \xi(t, \alpha)),$$

and  $\partial q / \partial \alpha$  can be found by using (6.9) and similar calculations to those above.

#### REFERENCES

- [1] J. V. BREAKWELL, J. L. SPEYER AND A. E. BRYSON, *Optimization and control of nonlinear systems using the second variation*, this Journal, 1 (1963), pp. 193–223.
- [2] R. COURANT AND D. HILBERT, *Methods of Mathematical Physics*, vol. 2, Interscience, New York, 1961.
- [3] J. L. DOOB, *Stochastic Processes*, John Wiley, New York, 1953.
- [4] P. DORATO, CHANG-MING HSIEH AND P. N. ROBINSON, *Optimal bang-bang control of linear stochastic systems with a small noise parameter*, IEEE Trans. Automatic Control, AC-12 (1967), pp. 682–689.
- [5] E. B. DYNKIN, *Markov Processes*, Springer-Verlag, Berlin, 1965.
- [6] H. FEDERER, *Geometric Measure Theory*, Springer-Verlag, Berlin, 1969.
- [7] W. H. FLEMING, *Duality and a priori estimates in Markovian optimization problems*, J. Math. Anal. Appl., 16 (1966), pp. 254–279; Erratum, 19 (1969), p. 204.
- [8] ———, *Optimal control of partially observable diffusions*, this Journal, 6 (1968), pp. 194–214.
- [9] ———, *The Cauchy problem for a nonlinear first-order partial differential equation*, J. Differential Equations, 5 (1969), pp. 515–530.
- [10] ———, *Controlled diffusions under polynomial growth conditions*, Control Theory and the Calculus of Variations, A. V. Balakrishnan, ed., Academic Press, 1969, pp. 209–234.
- [11] A. FRIEDMAN, *Partial Differential Equations of Parabolic Type*, Prentice-Hall, Englewood Cliffs, N.J., 1964.
- [12] I. I. GIKHMAN AND A. V. SKOROKHOD, *Introduction to the Theory of Random Processes*, W. B. Saunders, Philadelphia, 1969.
- [13] M. R. HESTENES, *Calculus of Variations and Optimal Control Theory*, John Wiley, New York, 1966.

- [14] E. HOPF, *The partial differential equation  $u_t + u_{xx} = \mu u_{xx}$* , Comm. Pure Appl. Math., 3 (1950), pp. 201–230.
- [15] H. J. KUSHNER AND A. J. KLEINMAN, *Numerical methods for the solution of degenerate nonlinear elliptic equations arising in optimal stochastic control theory*, IEEE Trans. Automatic Control, AC-13 (1968), pp. 344–353.
- [16] H. J. KUSHNER, *Near optimal control in the presence of small stochastic perturbations*, Trans. ASME Ser. D, J. Basic Engrg., 87 (1965), pp. 103–108.
- [17] ———, *On the stochastic maximum principle: fixed time of control*, J. Math. Anal. Appl., 11 (1965), pp. 78–92.
- [18] O. A. LADYZHENSKAYA, V. A. SOLONNIKOV AND N. N. URAL'SEVA, *Linear and Quasilinear Equations of Parabolic Type*, American Mathematical Society, Providence, R.I., 1968.
- [19] J. P. LASALLE AND H. HERMES, *Functional Analysis and Time Optimal Control*, Academic Press, New York, 1969.
- [20] O. A. OLEINIK, *Discontinuous solutions of nonlinear differential equations*, Uspekhi Mat. Nauk, 12 (1957), no. 3, pp. 3–73; English transl., Amer. Math. Soc. Transl. (2), no. 26, pp. 95–172.
- [21] ———, *On a problem of Fichera*, Dokl. Akad. Nauk SSSR, 157 (1964), pp. 1297–1301 = Soviet Math. Dokl., 5 (1964), pp. 1129–1133.
- [22] L. S. PONTRYAGIN, V. G. BOLTYANSKII, R. V. GAMKRELIDZE AND E. F. MISHCKENKO, *The Mathematical Theory of Optimal Processes*, John Wiley, New York, 1962.
- [23] R. L. STRATONOVICH, *On the theory of optimal control. An asymptotic method for solving the diffusive alternative equation*, Automat. Remote Control, 23 (1962), pp. 1352–1360.
- [24] D. W. STROOCK AND S. R. S. VARADHAN, *Diffusion processes with continuous coefficients*, Comm. Pure Appl. Math., 22 (1969), pp. 345–400, 479–530.
- [25] W. M. WONHAM AND W. F. CASHMAN, *A computational approach to optimal control of stochastic saturating systems*, Internat. J. Control, 10 (1969), pp. 77–98.
- [26] W. M. WONHAM, *On the separation theorem of stochastic control*, this Journal, 6 (1968), pp. 312–326.



## WEAK SOLUTIONS OF A PARTIAL DIFFERENTIAL EQUATION OF DYNAMIC PROGRAMMING\*

RAYMOND W. RISHEL†

**Abstract.** A formal dynamic programming argument relates the "value function of a stochastic optimal control problem" with the solution of a nonlinear parabolic partial differential equation. In cases in which the partial differential equation is degenerate, it may not have a classical solution but may have a weak solution in the sense of the theory of Schwartz distributions. It has been an open question as to whether a weak solution of the partial differential equation does equal the value function of the stochastic optimal control problem.

This paper shows that, roughly, whenever an associated uncontrolled system has an appropriately behaved density function, the equality holds. It is also shown that an optimal control law may be determined by minimizing a Hamiltonian formed in terms of the partial derivatives of the weak solution of the partial differential equation.

**1. Introduction.** In [3] Fleming studied the partial differential equation of dynamic programming connected with a stochastic optimal control problem. He showed, when the partial differential equation was uniformly parabolic, that there was a smooth solution and this solution was the value function for the stochastic optimal control problem. He called the case in which the partial differential equation was not uniformly parabolic, degenerate. Under mild assumptions he showed in the degenerate case that the partial differential equation had a weak solution. It was left as an open conjecture whether this weak solution equaled the value function of the stochastic control problem. Since in most applications the partial differential equation is degenerate, the degenerate case is especially important. For instance, for any stochastic control system governed by a single  $n$ th order differential equation with white noise input, the corresponding partial differential equation is degenerate.

This paper studies the relationship between a weak solution of the degenerate partial differential equation and the value function of the optimal control problem. Theorem 3 shows, under assumptions similar to those of [3], that if the partial derivatives of the weak solution are functions in appropriate  $L_p$ -spaces and solutions of an associated uncontrolled stochastic differential equation have a density which is in an appropriate  $L_p$ -space, then the weak solution equals the value function of the stochastic control problem. It also follows from Theorem 3 that there is a control function which minimizes a Hamiltonian function formed in terms of the partial derivatives of the weak solution, and this control is an optimal control for the stochastic control problem.

The question of obtaining precise conditions under which solutions of stochastic differential equations have densities is an interesting open question. It is known [6, Lemma 8.6] that solutions of nondegenerate stochastic differential equations have densities. Recent papers of Kushner [5] and Zakai [8] consider special cases of degenerate stochastic differential equations whose solutions have densities.

The optimal control problem is formulated for a class of controls which are merely bounded measurable functionals on the past of the processes involved.

---

\* Received by the editors August 4, 1970, and in revised form March 19, 1971.

† Bell Telephone Laboratories, Whippany, New Jersey 07981.

A technique of Girsanov [4] is used to assure there are processes corresponding to this very general type of control. The author would like to express his appreciation to Dr. V. E. Beneš for pointing out to him the importance of Girsanov's technique. Beneš used this technique in [2] to obtain existence theorems for optimal stochastic controls.

## 2. The stochastic control problem.

**2.1. Definition of notation.** Let  $f(t, x, y)$  be a bounded continuous function and  $g(t, x, y, u)$  a bounded Borel measurable function:

$$(1) \quad f: E^1 \times E^n \times E^m \rightarrow E^n, \quad g: E^1 \times E^n \times E^m \times E^k \rightarrow E^m.$$

Suppose that  $f$  satisfies

$$(2) \quad |f(t, x, y) - f(t, x', y')| \leq K_1[|x - x'| + |y - y'|]$$

and that  $g(t, x, y, u)$  is continuous as a function of  $u$  for fixed  $(t, x, y)$ . Let  $M_m$  be the space of nonsingular  $m \times m$  matrices with norm defined so that the square of the norm equals the sum of the squares of elements of the matrix. Let  $\sigma(t, y)$  be a bounded continuous function,

$$(3) \quad \sigma: E^1 \times E^m \rightarrow M_m,$$

such that

$$(4) \quad |\sigma(t, y) - \sigma(t, y')| \leq K_2|y - y'|.$$

Suppose, in addition, that the inverse matrix  $\sigma(t, y)^{-1}$  of  $\sigma(t, y)$  is a bounded continuous function. Let  $C(t, x, y, u)$ ,

$$(5) \quad C: E^1 \times E^n \times E^m \times E^k \rightarrow E^1,$$

be a continuous real-valued function. Let  $(\Omega, \beta, P)$  denote a triple of a probability space  $\Omega$ , Borel field  $\beta$ , and probability measure  $P$ . Let  $\zeta(t)$  denote an  $m$ -dimensional Brownian motion process, with respect to  $(\Omega, \beta, P)$ , whose covariance matrix is given by

$$(6) \quad E\{\zeta(t) \cdot \zeta(s)\} = \min(t, s)I,$$

where  $I$  is the identity matrix.

Let  $F[t_0, T]$  denote the space of all functions on the interval  $[t_0, T]$  to  $E^{m+n}$ . Let  $\mathcal{M}$  denote the Borel field on  $F[t_0, T]$  generated by the cylinder sets. A function  $\eta$ ,

$$(7) \quad \eta: E^1 \times F[t_0, T] \rightarrow E^k,$$

will be said to be nonanticipative if for each  $t$  and  $x_1(s), x_2(s) \in F[t_0, T]$  such that  $x_1(s) = x_2(s)$  if  $t_0 \leq s \leq t$ , then  $\eta(t, x_1(\cdot)) = \eta(t, x_2(\cdot))$ .

**2.2. The optimal control problem.** Let  $U$  be a compact subset of  $E^k$ . A function

$$(8) \quad \eta: E^1 \times F[t_0, T] \rightarrow U$$

will be called an admissible control if it is nonanticipative and Borel measurable.

Let  $B$  be a bounded closed region in  $E^{n+m}$  whose boundary is a twice continuously differentiable manifold. Let  $t_0$  be an initial time and  $(x, y)$  an initial state satisfying respectively  $0 \leq t_0 \leq T$ ;  $(x, y) \in B$ . Let trajectories of the control system be defined by the equations

$$(9) \quad x(t) = x + \int_{t_0}^t f(s, x(s), y(s)) ds,$$

$$(10) \quad y(t) = y + \int_{t_0}^t g(s, x(s), y(s), u(s)) ds + \int_{t_0}^t \sigma(s, y(s)) d\xi(s),$$

$$(11) \quad u(s) = \eta(s, x(\cdot), y(\cdot)).$$

The last integral in (10) is as a stochastic Ito integral. Let  $\tau$  denote the stopping time for the process which is the first time  $(x(t), y(t))$  hits the boundary of  $B$ . Let  $\tau_1 = \min [\tau, T]$ .

The optimal control problem is to find a control  $\eta$  in the class of admissible controls so that

$$(12) \quad E \left\{ \int_{t_0}^{\tau_1} C(s, x(s), y(s), u(s)) ds \right\}$$

is a minimum.

**3. Comments on the statement of the optimal control problem.** Since the controls are only assumed to be Borel measurable functionals on the past of the process and the function  $g(t, x, y, u)$  is only assumed to be Borel measurable, the question of whether solutions of (9)–(11) exist arises. The following procedure of I. V. Girsanov [4] can be used to show that there are solutions of (9)–(11). Under the assumptions above, classical theorems [7, p. 47] show that there is a unique solution to the equations

$$(13) \quad x(t) = x + \int_{t_0}^t f(s, x(s), y(s)) ds,$$

$$(14) \quad y(t) = y + \int_{t_0}^t \sigma(s, y(s)) d\xi(s).$$

Define  $\zeta_{t_0}^t(g)$  by

$$(15) \quad \begin{aligned} \zeta_{t_0}^t(g) = & \int_{t_0}^t g(s, x(s), y(s), u(s))^T [\sigma(s, y(s))^T]^{-1} d\xi \\ & - \frac{1}{2} \int_{t_0}^t g(s, x(s), y(s), u(s))^T [\sigma(s, y(s))^T]^{-1} \\ & \cdot [\sigma(s, y(s))]^{-1} g(s, x(s), y(s), u(s)) ds, \end{aligned}$$

where  $u(s) = \eta(s, x(\cdot), y(\cdot))$  and  $x(t), y(t)$  are solutions of (13)–(14).

**THEOREM 1** (Girsanov [4, p. 287 and p. 296]). *Define a measure  $\tilde{P}$  by*

$$(16) \quad \tilde{P}(d\omega) = \exp [\zeta_{t_0}^T(g)] P(d\omega)$$

*and a stochastic process  $\tilde{\xi}(t)$  for  $t_0 \leq t \leq T$  by*

$$(17) \quad \tilde{\xi}(t) = \xi(t) - \int_{t_0}^t \sigma(s, y(s))^{-1} g(s, x(s), y(s), u(s)) ds.$$

Then  $\tilde{P}$  is a probability measure on  $(\Omega, \beta)$ .<sup>1</sup> The process  $\tilde{\xi}(t, \omega)$  is a Brownian motion process on  $(\Omega, \beta, \tilde{P})$ . The solution of (13)–(14),  $x(t, \omega), y(t, \omega)$ , when considered with respect to the probability triple  $(\Omega, \beta, \tilde{P})$ , is a solution of (9)–(11) with the Brownian motion  $\tilde{\xi}(t)$  replacing  $\xi(t)$ .

It will always be assumed that the solutions of (9)–(11) which are involved in the optimal control problem are constructed from solutions of (13)–(14) in the manner stated in the theorem. The question of whether (9)–(11) has a unique solution appears to be still open. In [6] Stroock and Varadhan prove in a more restricted case analogous to (9)–(11) that solutions of the equations always induce a unique probability measure on the space of continuous functions.

**4. Intermediate results.** The following lemma shows how densities of solutions of (13)–(14) are related to densities of solutions of (9)–(11) constructed according to Theorem 1. The proof of this lemma is nearly the same as that of Lemma 8.6 of [6].

LEMMA 1. Let  $t_0, T, h$  and  $\alpha$  be numbers such that  $t_0 < t_0 + h < T$  and  $\alpha > 1$ . If for a solution  $x(t), y(t)$  of (13)–(14), the random vector  $x(t), y(t)$  has a probability density  $p(t, x, y)$  such that

$$(18) \quad \int_B \int_{t_0+h}^T p(t, x, y)^\alpha dt dx dy < \infty,$$

then for the solution  $x(t), y(t)$  of (9)–(11) constructed by Theorem 1, the random vector  $x(t), y(t)$  has a density  $q(t, x, y)$  such that for  $1 < \alpha' < \alpha$ ,

$$(19) \quad \int_B \int_{t_0+h}^T q(t, x, y)^{\alpha'} dt dx dy < \infty.$$

Since the solution of (9)–(11) will depend on which control  $\eta(t, x(\cdot), y(\cdot))$  is used, the density  $q(t, x, y)$  also will depend on which control  $\eta$  is being used. To save burdening the notation, this dependence will not be indicated in the formulas. However, the reader is cautioned to recall that these quantities will always depend on the control being considered.

The proof of Lemma 1 requires the following lemma of [4].

LEMMA 2 (Girsanov [4, p. 291]). Let  $\mathcal{F}_t$  denote the Borel field generated by the random vectors  $x(s), y(s)$  for  $t_0 \leq s \leq t$ . Let  $\zeta_{t_0}^t(g)$  be defined by (15). Then for any  $F_t$ -measurable function  $\eta(\omega)$ ,

$$(20) \quad \int_\Omega \eta(\omega) \tilde{P}(d\omega) = \int_\Omega \eta(\omega) \exp [\zeta_{t_0}^t(g)] P(d\omega).$$

*Proof of Lemma 1.* For a real number  $\rho > 1$  let

$$(21) \quad \int_\Omega \exp [\rho \zeta_{t_0}^t(g)] P(d\omega) = N_\rho.$$

The previous assumptions imply that there is a number  $k$  such that

$$(22) \quad |g(s, x(s), y(s), u(s))^T [\sigma(s, y(s))^T]^{-1} [\sigma(s, y(s))]^{-1} g(s, x(s), y(s), u(s))| < k.$$

<sup>1</sup> This first conclusion of Theorem 1 follows from the remark on p. 296 of [4]. The remaining conclusions are those of [4, Theorem 1, p. 287].

Using an abbreviated notation, we have

$$\begin{aligned}
 N_\rho &= \int_\Omega \exp \left[ \rho \int_{t_0}^t g^T \sigma^{-1T} \sigma^{-1} dy - \frac{1}{2} \rho \int_{t_0}^t g^T \sigma^{-1T} \sigma^{-1} g ds \right] P(d\omega) \\
 (23) \quad &\leq \exp \left[ \frac{1}{2} k(t - t_0) \rho (\rho - 1) \right] \int_\Omega \exp \left[ \int_{t_0}^t \rho g^T \sigma^{T-1} \sigma^{-1} dy \right. \\
 &\quad \left. - \frac{1}{2} \int_{t_0}^t \rho g^T \sigma^{T-1} \sigma^{-1} \rho g ds \right] P(d\omega) \\
 &\leq \exp \frac{1}{2} k(t - t_0) \rho (\rho - 1).
 \end{aligned}$$

The last step follows because, by applying Lemma 2 with  $g$  replaced by  $\rho g$ , the integral can be shown to equal one. By Lemma 2, if  $A$  is a set measurable with respect to the Borel field generated by  $x(t), y(t)$ ,

$$(24) \quad \int_A \tilde{P}(d\omega) = \int_A E\{\exp [\zeta_{t_0}^t(g)] | x(t), y(t)\} P(d\omega);$$

hence solutions of (9)–(11) have a density  $q(t, x, y)$  given by

$$q(t, x, y) = E\{\exp [\zeta_{t_0}^t(g)] | x(t) = x, y(t) = y\} p(t, x, y).$$

Denote  $E\{\exp [\zeta_{t_0}^t(g)] | x(t) = x, y(t) = y\}$  more briefly by  $r(t, x, y)$ . Since conditional expectation does not increase the  $L_\rho$ -norm, definition (21) implies

$$(25) \quad \int_{E^{m+n}} |r(t, x, y)|^\rho p(t, x, y) dx dy \leq N_\rho.$$

Consequently,

$$\begin{aligned}
 (26) \quad &\int_{t_0+h}^T \int_{E^{m+n}} |r(t, x, y)|^\rho p(t, x, y) dx dy dt \\
 &\leq (T - t_0) \exp \left[ \frac{1}{2} k(T - t_0) \rho (\rho - 1) \right].
 \end{aligned}$$

This implies

$$\begin{aligned}
 (27) \quad &\int_{t_0+h}^T \int_{E^{m+n}} |q(t, x, y)|^\rho |p(t, x, y)|^{1-\rho} dx dy dt \\
 &\leq (T - t_0) \exp \left[ \frac{1}{2} k(T - t_0) \rho (\rho - 1) \right].
 \end{aligned}$$

If  $1 < \alpha' < \alpha$  and  $h, k, \rho$  and  $\beta$  are chosen so that  $h > (\alpha - 1)(\alpha - \alpha')^{-1}$ ,  $h^{-1} + k^{-1} = 1$ ,  $\rho = h\alpha'$ , and  $\beta = h^{-1}(\rho - 1)$ , then an elementary argument shows that  $1 < \beta k < \alpha$ . Therefore, writing  $q^{\alpha'} = q^\alpha p^{-\beta} p^\beta$ , applying Hölder's inequality with parameters  $h$  and  $k$ , and using (27) and (18) gives, if  $1 < \alpha' < \alpha$ , that

$$(28) \quad \int_{t_0+h}^T \int_B |q(t, x, y)|^{\alpha'} dt dx dy < \infty,$$

which is the conclusion of Lemma 1.

Define the matrix  $a(t, y)$  by  $a(t, y) = \sigma(t, y)\sigma(t, y)^T$ . For the stochastic optimal control problem define a value function as follows. Consider the optimal control

problem starting from a variable time  $t$  and point  $(x, y)$ , where  $0 \leq t \leq T$  and  $(x, y) \in B$ . Let  $V(t, x, y)$  be the value of the infimum of the performance index over the class of admissible controls, given that the processes started at time  $t$  from the point  $(x, y)$ . That is,

$$(29) \quad V(t, x, y) = \inf_{\eta} E \left\{ \int_t^{\tau_1} C(s, x(s), y(s), u(s)) ds \mid x(t) = x, y(t) = y \right\},$$

where  $u(s) = \eta(s, x(\cdot), y(\cdot))$ .

The method of dynamic programming gives a formal argument which connects the value function  $V(t, x, y)$  with the solution  $\phi(t, x, y)$  of

$$(30) \quad \phi_t + \sum_{i=1}^m a_{ij}(t, y)\phi_{y_i y_j} + \sum_{i=1}^n f_i(t, x, y)\phi_{x_i} + \min_{u \in U} \left\{ \sum_{i=1}^m g_i(t, x, y, u)\phi_{y_i} + C(t, x, y, u) \right\} = 0,$$

which satisfies the boundary condition

$$(31) \quad \phi(t, x, y) = 0 \quad \text{if } (x, y) \in \text{boundary } B \quad \text{or} \quad t = T.$$

We shall say for  $\beta > 1$  that a function  $\phi$  is in class  $F_\beta$  if  $\phi(t, x, y)$  is continuous on  $[0, T] \times B$ , and the Schwartz distribution theory second partial derivatives with respect to  $y$  variables and first partial derivatives with respect to  $t$  and  $x$  variables are functions which are in  $L_\beta([0, T] \times B)$ .

Define a norm on the functions of class  $F_\beta$  by

$$(32) \quad |\phi| = \left[ \sup_{0 \leq t \leq T, (x, y) \in B} \phi(t, x, y) \right] + \sum_{z=t, x_i, y_i} \left[ \int_0^T \int_B |\phi_z|^\beta dx dy dt \right]^{1/\beta} + \sum_{i, j=1}^m \left[ \int_0^T \int_B |\phi_{y_i y_j}|^\beta dx dy dt \right]^{1/\beta}.$$

A function  $\phi$  in  $F_\beta$  will be said to be a weak solution of (30) if (30) is satisfied almost everywhere with respect to  $(m + n + 1)$ -dimensional Lebesgue measure on  $[0, T] \times B$ .

Denote by  $\Lambda(t, x, y, u)$  the partial differential operator for which

$$(33) \quad \Lambda(t, x, y, u)[\phi] = \phi_t + \sum_{i=1}^m a_{ij}(t, y)\phi_{y_i y_j} + \sum_{i=1}^n f_i(t, x, y)\phi_{x_i} + \sum_{i=1}^m g_i(t, x, y, u)\phi_{y_i}.$$

In cases in which the dependence of  $\Lambda$  on the variables and functions involved in its coefficients can be inferred from the context, the shortened notation  $\Lambda[\phi]$  will be used for (33).

**THEOREM 2.** *Let the hypothesis of Lemma 1 be satisfied. Let  $\phi(t, x, y)$  be a solution of (30) in  $F_\beta$ , where  $\beta = \alpha'(\alpha' - 1)^{-1}$ . Let  $x(t), y(t)$  be solutions of (9)–(11) corresponding to the control  $\eta(t, x(\cdot), y(\cdot))$ . Let  $\tau$  denote the first time  $x(t), y(t)$  hits boundary  $B$  and  $\tau_1 = \min(\tau, T)$ ,  $\tau_h = \min(\tau, t_0 + h)$ . Let  $\chi(t)$  be the characteristic*

function of the set  $\{\omega : \tau(\omega) \geq t\}$ . Then

$$\begin{aligned}
 & E\{\phi(\tau_1, x(\tau_1), y(\tau_1))\} - E\{\phi(\tau_h, x(\tau_h), y(\tau_h))\} \\
 (34) \quad & = \int_{t_0+h}^T \int_B E\{\chi(s)\Lambda(s, x(s), y(s), \eta(s, x(\cdot), y(\cdot)))[\phi]|x(s) = x, y(s) = y\} \\
 & \quad \cdot q(s, x, y) dx dy ds.
 \end{aligned}$$

*Proof of Theorem 2.* Let  $\phi^n(t, x, y)$  be a sequence of twice continuously differentiable functions such that  $\phi^n$  converges to  $\phi$  in  $F_\beta$ . Since  $\phi^n$  is twice continuously differentiable, Ito's formula [7] gives

$$\begin{aligned}
 & E\{\phi^n(\tau_1, x(\tau_1), y(\tau_1))\} - E\{\phi^n(\tau_h, x(\tau_h), y(\tau_h))\} \\
 (35) \quad & = E\left\{ \int_{t_0+h}^T \chi(s)\Lambda(s, x(s), y(s), \eta(s, x(\cdot), y(\cdot)))[\phi^n] \right\} ds.
 \end{aligned}$$

Now by Fubini's theorem,

$$\begin{aligned}
 & E \int_{t_0+h}^T \chi(s)\Lambda[\phi^n] ds = \int_{t_0+h}^T E\{\chi(s)\Lambda\phi^n\} ds \\
 (36) \quad & = \int_{t_0+h}^T \int_B E\{\chi(s)\Lambda\phi^n|x(s) = x, y(s) = y\} q(s, x, y) dx dy dt.
 \end{aligned}$$

Looking at definition (33) we see that the equality

$$\begin{aligned}
 & E\{\chi(s)g_i(s, y(s), \eta(s, x(\cdot), y(\cdot)))\phi^n_{y_i}(s, x(s), y(s))|x(s) = x, y(s) = y\} \\
 (37) \quad & = \phi^n_{y_i}(s, x, y)E\{\chi(s)g_i(s, x(s), y(s), \eta(s, x(\cdot), y(\cdot)))|x(s) = x, y(s) = y\}
 \end{aligned}$$

is typical for the terms of  $E\{\chi(s)\Lambda(s)[\phi^n]|x(s) = x, y(s) = y\}$ . A similar result holds when  $\phi$  replaces  $\phi^n$  in (37). Therefore, the expression

$$(38) \quad E\{\chi(s)(\Lambda(s)[\phi^n] - \Lambda(s)[\phi])|x(s) = x, y(s) = y\}$$

can be written as a sum of bounded terms times terms which are converging to zero in  $L_\beta([t_0, T] \times B)$ . Therefore, we conclude that (38) converges to zero in  $L_\beta([t_0, T] \times B)$ .

By applying Hölder's inequality, involving (19) and (38), we see that the right side of (36) converges to the right side of (34). Since  $F_\beta$  convergence implies uniform convergence of  $\phi_n$  to  $\phi$ , the left side of (35) converges to the left side of (34), which completes the proof of Theorem 2.

If  $\phi(t, x, y)$  is a weak solution of (30), an implicit function theorem of Beneš [1] asserts there is a measurable function  $\gamma(t, x, y)$ ,

$$(39) \quad \gamma : [t_0, T] \times B \rightarrow U,$$

such that

$$\begin{aligned}
 & \sum_{i=1}^m g_i(t, x, y, \gamma(t, x, y))\phi_{y_i} + C(t, x, y, \gamma(t, x, y)) \\
 (40) \quad & = \min_{u \in U} \left\{ \sum_{i=1}^m g_i(t, x, y, u)\phi_{y_i} + C(t, x, y, u) \right\}
 \end{aligned}$$

almost everywhere with respect to Lebesgue measure on  $[t_0, T] \times B$ . Extend  $\gamma$  in an arbitrary way to a Borel measurable function on  $[t_0, T] \times E^{m+n}$  with values in  $U$ .

This function induces a measurable function on  $[t_0, T] \times F[t_0, T]$  by the formula

$$(41) \quad \gamma(t, x(\cdot), y(\cdot)) = \gamma(t, x(t), y(t)).$$

**5. Fundamental theorem.**

**THEOREM 3.** *Let the solution of (13)–(14) have a density  $p(t, x, y)$  such that for some  $\alpha, \alpha > 1$ ,*

$$(42) \quad \int_{t_0+h}^T p(t, x, y)^\alpha dx dy dt < \infty.$$

*Let  $1 < \alpha' < \alpha$ , and let  $\phi(t, x, y)$  be a solution of (30) which is in  $F_\beta$  for  $\beta = \alpha'(\alpha' - 1)^{-1}$  which satisfies the boundary condition (31). Let  $\gamma(t, x(\cdot), y(\cdot))$  be the control defined in (40)–(41).*

*Then if  $x(t), y(t)$  is a solution of (9)–(11) corresponding to the control  $u(t) = \gamma(t, x(\cdot), y(\cdot))$ ,*

$$(43) \quad \phi(t_0, x, y) = E \left\{ \int_{t_0}^{t_1} C(s, x(s), y(s), u(s)) ds \right\}.$$

*If  $x(t), y(t)$  is a solution of (9)–(11) corresponding to an arbitrary control  $u(t) = \eta(t, x(\cdot), y(\cdot))$ ,*

$$(44) \quad \phi(t_0, x, y) \leq E \left\{ \int_{t_0}^{t_1} C(s, x(s), y(s), u(s)) ds \right\}.$$

Notice that  $t_0, x, y$  could be any variables for which  $0 < t_0 < T$  and  $(x, y) \in B$ . Thus Theorem 3 implies that  $\phi(t, x, y) = V(t, x, y)$  on  $[0, T] \times B$ . That is, if there is a solution of (30) with boundary condition (31), it agrees with the value function of the stochastic control problem. Notice also that (43) and (44) imply that the control  $\gamma$  constructed from (40)–(41) is an optimal control law for the stochastic control problem.

*Proof of Theorem 3.* Since  $\phi$  is a weak solution of (30), the definition of  $\gamma$  implies that there is a set  $N$  of  $(m + n + 1)$ -dimensional Lebesgue measure zero such that

$$(45) \quad -\Lambda(t, x, y, \gamma(t, x, y))[\phi] = C(t, x, y, \gamma(t, x, y))$$

on  $[t_0, T] \times B - N$ .

Let  $u(t) = \gamma(t, x(t), y(t))$  and  $x(t), y(t)$  be corresponding solutions of (8)–(11). Since  $x(t), y(t)$  has a density  $q(t, x, y)$ , (45) and Fubini's theorem imply that for almost every  $t$ ,

$$(46) \quad -\chi(t)\Lambda(t, x(t), y(t), u(t))[\phi] = c(t, x(t), y(t), u(t))$$

with probability one and



$$(47) \quad \begin{aligned} E\{-\chi(t)\Lambda(t, x(t), y(t), u(t))[\phi] | x(t) = x, y(t) = y\} \\ = -\Lambda(t, x, y, \gamma(t, x, y))[\phi] = C(t, x, y, \gamma(t, x, y)) \end{aligned}$$

except at a set of measure zero on  $B$ .

Substituting (47) in (34) gives

$$(48) \quad \begin{aligned} E\{\phi(\tau_h, x(\tau_h), y(\tau_h))\} - E\{\phi(\tau_1, x(\tau_1), y(\tau_1))\} \\ = \int_{t_0+h}^T E\{\chi(t)C(t, x(t), y(t), u(t))\} dt. \end{aligned}$$

Since  $\phi$  satisfies the boundary condition (31),  $\phi(\tau_1, x(\tau_1), y(\tau_1)) = 0$ . Since  $\phi$  is bounded and continuous and  $x(t), y(t)$  has continuous paths, using Lebesgue's convergence theorem and passing to the limit as  $h$  approaches zero gives

$$(49) \quad \begin{aligned} \phi(t_0, x, y) &= \lim_{h \rightarrow 0} E\{\phi(\tau_h, x(\tau_h), y(\tau_h))\} \\ &= E \int_{t_0}^{\tau_1} C(t, x(t), y(t), u(t)) dt. \end{aligned}$$

To establish (44) notice that since  $\phi$  is a weak solution of (30) that there is a set  $N$  of  $(m + n + 1)$ -dimensional Lebesgue measure zero, which does not depend on  $u$ , such that for all  $u \in U$ ,

$$(50) \quad -\Lambda(t, x, y, u)[\phi] \leq C(t, x, y, u)$$

on  $[t_0, T] \times B - N$ .

For any control  $\eta(t, x(\cdot), y(\cdot))$ , let  $x(t), y(t)$  be a solution of (8)–(11) corresponding to  $\eta$ . Since  $x(t), y(t)$  has a density  $q(t, x, y)$ , (50) implies

$$(51) \quad -\chi(t, \omega)\Lambda(t, x(t, \omega), y(t, \omega), u)[\phi] \leq \chi(t, \omega)C(t, x(t, \omega), y(t, \omega), u)$$

for all  $u \in U$  and  $dt \times P$  almost every  $(t, \omega)$  in  $[t_0, T] \times \Omega$ . Letting  $u(t, \omega) = \eta(t, x(\cdot, \omega), y(\cdot, \omega))$ , we have that inequality (51) implies

$$(52) \quad \chi(t, \omega)\Lambda(t, x(t, \omega), y(t, \omega), u(t, \omega))[\phi] \leq \chi(t, \omega)C(t, x(t, \omega), y(t, \omega), u(t, \omega)),$$

$dt \times P$  almost everywhere on  $[t_0, T] \times \Omega$ .

By substituting (52) in (34) and proceeding in a manner similar to that of (48)–(49), it follows that (44) holds.

**6. Concluding remarks.** By combining Theorem 3 with results of Fleming [3], a fairly complete theory is available for the type of stochastic control problem discussed in this paper. Fleming [3, Theorem 5.1, p. 269] proves that there exists a weak solution of the partial differential equation (30) with boundary condition (31) which is in  $F_2$ . Then under the combined hypothesis of [3] and Theorem 3, where the exponent  $\alpha$  in Theorem 3 is larger than two, there exists an optimal control, this optimal control minimizes the Hamiltonian as in (40), and the optimal value of the performance starting at time  $t$  from the point  $x, y$  is given by the solution of (30). It also follows from Theorem 3, under the hypothesis of the theorem, that a weak solution of the partial differential equation of dynamic programming (30) in  $F_\beta$  with boundary condition (31) is unique.

## REFERENCES

- [1] V. E. BENEŠ, *The existence of optimal strategies based on specified information for a class of stochastic decision problems*, this Journal, 8 (1970), pp. 179–188.
- [2] ———, *The existence of optimal stochastic control laws*, this Journal, 9 (1971), pp. 446–472.
- [3] W. H. FLEMING, *Duality and a priori estimates in Markovian optimization problems*, J. Math. Anal. Appl., 16 (1966), pp. 254–279.
- [4] I. V. GIRSANOV, *On transforming a certain class of stochastic processes by absolutely continuous substitution of measures*, Theor. Probability Appl., 5 (1960), pp. 285–301.
- [5] H. J. KUSHNER, *The Cauchy problem for a class of degenerate parabolic equations and asymptotic properties of related diffusion processes*, J. Differential Equations, 6 (1969), pp. 209–231.
- [6] D. W. STROOCK AND S. R. S. VARADHAN, *Diffusion processes with continuous coefficients. II*, Comm. Pure Appl. Math., 22 (1969), pp. 479–530.
- [7] A. V. SKOROKHOD, *Studies in the Theory of Random Processes*, Addison-Wesley, Reading, Mass., 1965.
- [8] M. ZAKAI, *A Lyapunov criterion for the existence of stationary probability distributions for systems perturbed by noise*, this Journal, 7 (1969), pp. 390–397.

## GENERALIZED CUTTING PLANE ALGORITHMS\*

B. CURTIS EAVES† AND W. I. ZANGWILL‡

**Abstract.** This paper introduces a master cutting plane algorithm for nonlinear programming that isolates the points it generates from one another until a solution is achieved. The master algorithm provides a foundation for the study of cutting plane algorithms and directs the way for development of procedures which permit deletion of old cuts.

**Introduction.** Since Kelley [5], Cheney and Goldstein [2] introduced the first cutting plane method for nonlinear programming over ten years ago, and through the developments of the Veinott [8] supporting hyperplane method and the Dantzig and Wolfe [3] dual cutting plane method, the following question has persisted. Is it possible to drop previously generated cutting planes and still guarantee convergence? All of the aforementioned methods require that every previous cut be retained. Then, since the linear programming subproblem must on every iteration increase in size, the subproblem soon becomes unwieldy and difficult, if not impossible, to solve. To circumvent this complication in practical application, one would simply drop old cuts and solve the resulting smaller subproblems. It was not known whether or not this procedure would converge.

One of the authors recently suggested a unifying theory for cutting plane algorithms [9, Chap. 14]. Although this theory did bring together the previous methods, it still did not permit old cuts to be dropped. By extending and generalizing this theory, this paper presents precise methods on how to drop cuts and still guarantee convergence.

The paper is divided into three parts. Section 1 introduces the master cutting plane algorithm (MCPA) and states important background material for convergence of sequences. In § 2 we derive four methods from the MCPA. One of these four will in § 3 further specialize into the Kelley, Cheney, Goldstein (KCG), the Veinott (VN), and the Dantzig-Wolfe methods (DW), while the other three will illustrate not only different approaches to cutting plane methods, but also how to drop old cuts.

The authors have recently been informed that Algorithm IV for the unique optimum case is similar to an approach developed independently and simultaneously by Professor Donald M. Topkis (see [6], [7]). Professor Topkis examined the uniformly strictly concave case and used an entirely different method of proof.

**1. Fundamentals.** Following the approach in [9, Chap. 14]—cutting plane methods generate a sequence of points  $z^k$ ,  $k = 1, 2, 3, \dots$ , in an effort to calculate a point in a set  $G$ , termed the goal set. The set  $G$  will be closed in a metric space and

---

\* Received by the editors August 21, 1969, and in final revised form February 16, 1971. This work was supported in part by the U.S. Office of Naval Research under Grant Nonr-047-069, the U.S. Army Research Office—Durham under Contract DAHC-04-67-C-0028, and the National Science Foundation under Grant GK-5695.

† Department of Operations Research, Stanford University, Stanford, California 94305.

‡ Center for Research in Management Science, University of California, Berkeley, California. Now at U.S. Office of Education, Washington, D.C. 20202.

specially selected so that when the MCPA calculates a point in it, the nonlinear programming problem will be solved.

The points generated by the MCPA will be in a closed set  $X$  and satisfy the special property that if  $z^k \notin G$ , then  $z^{k+1}$  is separated from the points  $z^1, z^2, \dots, z^k$  by at least a certain distance. This distance will be determined by a separator function that we now define.

DEFINITION. Given the closed sets  $\emptyset \neq G \subset X$  in a metric space, an extended real-valued function  $\delta$  on  $X \sim G$  is a *separator*

- (a) if it is nonnegative

and

- (b) if  $z^k \rightarrow z^\infty$  and  $\delta(z^k) \rightarrow 0$ , imply  $z^\infty \in G$ .

An immediate consequence is that a separator is positive on  $X \sim G$ . For, suppose  $z^* \in X \sim G$  and  $\delta(z^*) = 0$ . Then, letting  $z^k = z^*$  for all  $k$ , we have

$$z^k \rightarrow z^*, \quad \delta(z^k) \rightarrow 0$$

and by (b),

$$z^* \in G.$$

The contradiction is clear and  $\delta$  must be positive on  $X \sim G$ .

The next lemma also follows immediately from the definition.

LEMMA 1. *Let  $\delta$  be a real-valued function on  $X \sim G$ . If  $\delta$  is positive and lower semicontinuous, it is a separator.*

An example of a separator is the distance  $d(z, G)$  between a point  $z \in X$  and the set  $G$ . It is easy to show that  $\delta$  is a separator if and only if  $\delta$  has a positive lower bound on each compact set in  $X \sim G$ . The minimum of a finite collection of separators is also a separator as is any positive multiple of a separator. Moreover, if  $\delta$  is a separator and  $\delta' \geq \delta$  then so is  $\delta'$ . A condition which is frequently imposed on a separator is  $\delta \leq d(\cdot, G)$ , or in other words, that  $\delta(z)$  not exceed the distance  $d(z, G)$  from  $z$  to  $G$  for all  $z \in X \sim G$ .

We now pose the MCPA, a simple procedure that underlies the cutting plane methods.

MASTER CUTTING PLANE ALGORITHM. Let  $\delta \leq d(\cdot, G)$  be a separator on  $X \sim G$  and suppose the points  $z^1, \dots, z^k$  in  $X$  have been generated. If  $z^k \in G$ , then terminate; otherwise calculate  $z^{k+1} \in X$  such that

$$z^{k+1} \notin B(z^i, \delta(z^i)), \quad i = 1, \dots, k,$$

where  $B(z, \delta(z))$  is the open ball of radius  $\delta(z)$  at  $z$ .

THEOREM 2. *Given the closed sets  $\emptyset \neq G \subset X$ , let the sequence  $z^1, z^2, \dots$  in  $X$  be generated by the MCPA. Then any cluster point is in  $G$ .*

*Proof.* Suppose for some subsequence

$$z^k \rightarrow z^\infty.$$

By construction of the algorithm, this could occur only if

$$\delta(z^k) \rightarrow 0.$$

Then from the definition of a separator,  $z^\infty \in G$ .

*Spaces of sets.* To apply the MCPA we must first explore the suggestion in [9] that cutting plane algorithms operate by generating closed sets  $Z^k$  contained in  $X$ . Specifically, we must explore the space of all such sets  $Z^k$  in  $X$ .

Assume that  $X$  is a compact metric space, that  $G \subset X$ , and that  $G$  is compact. Then define  $\mathcal{Z}$  to be the collection of all closed sets  $Z$  for which

$$G \subset Z \subset X.$$

$\mathcal{Z}$  can be topologized [4, pp. 166–172] to form a metric space. With this topology,  $\lim Z^k = Z^\infty$  or equivalently  $Z^k \rightarrow Z^\infty$  if and only if both (a) and (b) hold.

(a) If  $z^k \in Z^k$  for  $k = 1, 2, \dots$  and  $z^k \rightarrow z^\infty$ , then  $z^\infty \in Z^\infty$ .

(b) If  $z \in Z^\infty$ , then there exists  $z^k \in Z^k$  for  $k = 1, 2, \dots$  such that  $z^k \rightarrow z$ .

Intuitively, if  $Z^k \rightarrow Z^\infty$ , then as  $k \rightarrow \infty$  the set  $Z^k$  becomes more and more like  $Z^\infty$ . We summarize a key lemma [4] about this set space.

LEMMA 3. *Suppose  $X$  is a compact metric space and  $G \subset X$  is a compact subset.*

Let

$$\mathcal{Z} = \{Z | G \subset Z \subset X, Z \text{ closed}\}.$$

Then  $\mathcal{Z}$  is compact (in the set topology).

The set space permits us to define cuts.

*Cuts and maps.* We know that cutting plane algorithms operate by taking cuts, yet the MCPA was not stated in terms of a cut but was phrased in terms of a separator. Now we must link the two concepts of a cut and a separator by defining a cut map. Recall that a point-to-set map from  $X \sim G$  to  $\mathcal{Z}$  takes points of  $X \sim G$  into subsets of  $\mathcal{Z}$ .

DEFINITION. A point-to-set map  $\varphi : X \sim G \rightarrow \mathcal{Z}$  is a *cut map* if there is a separator such that for each  $x \in X \sim G$  and each  $Z \in \varphi(x)$ ,

$$Z \cap B(x, \delta(x)) = \emptyset.$$

Note that any set  $Z \in \mathcal{Z}$  can be a cut, but we shall only call  $Z$  a cut if  $Z \in \bigcup \varphi(X \sim G)$ , that is, if  $Z \in \varphi(z)$  for some  $z \in X \sim G$ . Also notice that

$$G = \bigcap \bigcup \varphi(X \sim G)$$

which is to say that  $G$  is the intersection of all cuts. In particular, one can observe that if all cuts are convex, then  $G$  must also be convex. (In the computational procedures of § 3, cuts will be of the form  $H \cap X$ , where  $H$  is a half-space of  $E^n$ ,  $n$ -dimensional Euclidean space, and  $X$  is a compact polyhedral convex set in  $E^n$ .)

In practice we may know that a map  $\varphi$  is a cut map without knowing an underlying separator. However, note that we can easily generate other separators from any cut map. Simply define  $\delta$  on  $X \sim G$  by setting  $\delta(z) = d(z, Z)$  for some  $Z \in \varphi(z)$ .

LEMMA 4. *Let  $\varphi$  be a cut map,  $z^k \rightarrow z^\infty \notin G$ ,  $C^k \in \varphi(z^k)$ , and  $C^k \rightarrow C^\infty$ . Then  $z^\infty \notin C^\infty$ .*

*Proof.* Let  $\delta$  be a separator underlying  $\varphi$  and consider

$$\limsup \delta(z^k) \leq \lim d(z^k, C^k) = d(z^\infty, C^\infty).$$

If  $d(z^\infty, C^\infty) = 0$ , then  $\delta(z^\infty) \rightarrow 0$  and by definition  $z^\infty \in G$ .

*Closed maps.* In several connections we shall use the notion of a closed map. Let  $U$  and  $V$  be metric spaces and let  $\beta:U \rightarrow V$  be a point-to-set map ( $\beta(u) \subset V$ ). We say that  $\beta$  is a closed map if

$$u^i \rightarrow u^\infty, \quad v^i \rightarrow v^\infty, \quad v^i \in \beta(u^i) \quad \text{imply} \quad v^\infty \in \beta(u^\infty)$$

for all such sequences and limit points.

The next lemma indicates a relationship between closed and cut maps.

LEMMA 5. Let  $\emptyset \neq G \subset X$  be compact sets and let  $\varphi: X \sim G \rightarrow \mathcal{Z}$  be a closed map such that  $z \notin Z$  for  $Z \in \varphi(z)$ . Then  $\varphi$  is a cut map.

*Proof.* The result follows from the fact that

$$\delta(z) = \inf \{d(z, Z) | Z \in \varphi(z)\}$$

is a separator. Assume  $z^k \rightarrow z^\infty$  and  $\delta(z^k) \rightarrow 0$ . Using closedness, choose  $Z^k \in \varphi(z^k)$  such that  $\delta(z^k) = d(z^k, Z^k)$ . There is a subsequence for which  $Z^k \rightarrow Z^\infty$  and we have  $z^\infty \in Z^\infty$ . If  $z^\infty \notin G$ , by closedness  $Z^\infty \notin \varphi(z^\infty)$  and from the hypothesis  $z^\infty \notin Z^\infty$ , which is a contradiction.

*Objective function.* Our abstract algorithms of § 2 solve the following non-linear programming problem:

$$\text{maximize } r(y), \\ y \in G$$

where  $r$  is a real-valued continuous function on  $X$ .

What is interesting is that the objective function  $r$  and a separator  $\delta$  induce another separator  $\delta'$ .

LEMMA 6. Given a compact set  $X$  in a metric space, let  $\delta$  be a separator on  $X \sim G$  and  $r$  be a continuous function on  $X$ . Then the function  $\delta'$  defined on  $X \sim G$  by

$$\delta'(z) = \inf \{d(z, x) | r(x) \leq r(z) - \delta(z), x \in X\}$$

is a separator.

*Proof.* (We adopt the convention that  $\inf \emptyset = +\infty$ , and hence,  $\delta'(z) = +\infty$  may occur.) If  $\delta'(z) < +\infty$ , then by compactness and continuity there is an  $x \in X$  such that  $\delta'(z) = d(z, x)$  and  $r(x) \leq r(z) - \delta(z)$ .

Now suppose

$$z^k \rightarrow z^\infty \quad \text{and} \quad \delta'(z^k) \rightarrow 0.$$

Let

$$\delta'(z^k) = d(z^k, x^k)$$

and

$$r(x^k) \leq r(z^k) - \delta(z^k).$$

Then  $d(z^k, x^k) \rightarrow 0$ , and it follows that  $\delta(z^k) \rightarrow 0$ . As  $\delta$  is a separator,  $z^\infty \in G$ . This completes the proof.

Note that if  $z^k, k = 1, 2, \dots$ , is a sequence in  $X \sim G$  such that  $z^k \rightarrow z^\infty$  and  $r(z^{k+1}) \leq r(z^k) - \delta(z^k)$ , then  $z^\infty \in G$ . This fact follows from Lemma 6, but it also has a trivial proof, namely,  $z^k \rightarrow z^\infty$  implies  $\delta(z^k) \leq r(z^k) - r(z^{k+1}) \rightarrow 0$  which implies  $z^\infty \in G$ .

We now adopt a function and a map that depend upon the nonlinear programming problem which facilitates further discussion.

Let  $r^*$  be the real-valued function on  $\mathcal{Z}$  defined by

$$r^*(Z) = \max_{z \in Z} r(z).$$

In words,  $r^*$  is the optimal objective value as a function of the feasible region  $Z$ .

Next we define an “arg max” map on  $\mathcal{Z}$  by

$$\Gamma(Z) = \{x \in Z | r(x) = r^*(Z)\}.$$

The map  $\Gamma$  specifies all the optimal solution points to the nonlinear programming problem.

The following lemma is proved in [1].

LEMMA 7. *Let  $r$  be a continuous real-valued function on  $X$ , a compact set in a metric space. Then the function  $r^* : \mathcal{Z} \rightarrow E^1$  is continuous and the point-to-set map  $\Gamma : \mathcal{Z} \rightarrow X$  is a closed map.*

**2. Abstract cutting plane methods.** In this section, we describe four abstract cutting plane methods, Algorithms I–IV, for solving the problem

$$(1) \quad \underset{y \in G}{\text{maximize}} \quad r(y).$$

Each of the methods generates a sequence of sets  $Z^k$  and points  $z^k \in \Gamma(Z^k)$ , such that cluster points of the  $z^k$  sequence will be in  $\Gamma(G)$ . Consequently, any cluster point will solve the programming problem (1).

These four algorithms certainly do not exhaust the possibilities; nevertheless, an effort has been made to pick a representative selection from those of which we are aware.

The algorithms, I–IV, have been ordered so that, roughly, each succeeding one relies more upon improvement in objective value. Algorithm I retains all generated cuts. Algorithms II and III provide for the deletion of cuts, but just how many cuts will be dropped is problem dependent and could vary from none to numerous to almost all. Algorithm IV requires more stringent conditions than I–III, but it achieves the ideal with regard to the deletion of cuts; namely, all inactive cuts can be dropped ( $Z$  is inactive for  $z$  if  $z$  is interior to  $Z$  with respect to the space  $X$ ). Algorithm II will probably prove to be the most useful.

We shall always assume that  $\emptyset \neq G \subset X$  are compact sets in a metric space. Also

$$\mathcal{Z} = \{Z | G \subset Z \subset X, Z \text{ closed}\}$$

and  $r$  is continuous.

Algorithm 0 is the general algorithm that serves as a basis for the others. It assumes that a separator is given either implicitly or explicitly.

ALGORITHM 0. Given a separator  $\delta \leq d(\cdot, G)$ , let the set  $Z^1 \in \mathcal{Z}$  be selected arbitrarily and choose  $z^1 \in \Gamma(Z^1)$ . Assume that  $Z^i$  and  $z^i, i = 1, \dots, k$ , have been generated. If  $z^k \in G$ , terminate; otherwise select  $Z^{k+1}$  and  $z^{k+1} \in \Gamma(Z^{k+1})$  such that

$$z^{k+1} \notin B(z^i, \delta(z^i)) \quad \text{for } i = 1, \dots, k.$$

Clearly, Algorithm 0 operates as the MCPA so that any cluster point  $z^\infty$  of the  $z^k$  sequence is in  $G$ . Moreover, by specification of  $\mathcal{Z}$ ,  $Z^k \supset G$ , and consequently  $r(z^k) \geq r^*(G)$ . From the continuity of  $r$ ,  $r(z^\infty) \geq r^*(G)$ , and we conclude that  $z^\infty \in \Gamma(G)$ . In other words,  $z^\infty$  solves (1).

ALGORITHM I. This algorithm, as will be seen in § 3 of the paper, is a generalization of the KCG, VN, and DW methods. It retains all the old cuts (except perhaps for the redundant ones). Given a cut map  $\varphi: X \sim G \rightarrow \mathcal{Z}$ , the algorithm is as follows.

Assume that  $Z^1, \dots, Z^k$  and  $z^1, \dots, z^k$  have been generated. If  $z^k \in G$ , terminate; otherwise select  $C^k \in \varphi(z^k)$ , form  $C^k \cap Z^k$ , and let  $Z^{k+1} \in \mathcal{Z}$  be any subset of  $C^k \cap Z^k$ . Then calculate the point  $z^{k+1} \in \Gamma(Z^{k+1})$ .

To show that this operates as Algorithm 0, note that  $z^{k+1} \in Z^{k+1}$  since  $z^{k+1} \in \Gamma(Z^{k+1})$ . This algorithm produces

$$Z^{k+1} \subset Z^k \cap C^k \subset Z^k \subset \dots \subset Z^i \cap C^i \quad \text{for } 1 \leq i \leq k,$$

where  $C^i \in \varphi(z^i)$ . Thus

$$z^{k+1} \in C^i,$$

and by definition of the cut map, a separator exists such that

$$z^{k+1} \notin B(z^i, \delta(z^i)), \quad i = 1, \dots, k.$$

Consequently, the algorithm is a special case of Algorithm 0, and any cluster point  $z^\infty$  solves (1).

ALGORITHM II. This algorithm drops cuts after  $r$  has made sufficient progress as judged by a separator  $\delta$ . Briefly, if cut  $C^i$  is introduced to cut off the point  $z^i$ , then  $C^i$  can be dropped later on iteration  $k + 1$  if

$$r(z^{k+1}) \leq r(z^i) - \delta(z^i)$$

and if the cuts which remain “maintain objective value.” Let  $\varphi$  be a cut map and let  $\delta$  be a separator on  $X \sim G$ . At each iteration  $k$  an index set  $I^k$  will specify the cuts that have not yet been dropped. The algorithm is as follows.

Initiate the algorithm with  $I^1 = \emptyset$  and  $z^1 \in \Gamma(X)$ . Assume  $I^i \subset \{1, \dots, i - 1\}$  and  $z^i$  for  $i = 1, \dots, k$  and  $C^i$  for  $i = 1, \dots, k - 1$  have been generated. If  $z^k \in G$ , terminate; otherwise select  $C^k \in \varphi(z^k)$ , let

$$Y^{k+1} = \left( \bigcap_{i \in I^k} C^i \right) \cap C^k$$

and let  $y^{k+1} \in \Gamma(Y^{k+1})$ . Then let

$$Z^{k+1} = \bigcap_{i \in I^{k+1}} C^i$$

and  $z^{k+1} \in \Gamma(Z^{k+1})$ , where  $I^{k+1}$  is chosen so that:

- (a)  $I^{k+1} \subset I^k \cup \{k\}$ ,
- (b)  $\{i = 1, \dots, k \mid r(y^{k+1}) > r(z^i) - \delta(z^i)\} \subset I^{k+1}$ ,
- (c)  $r(z^{k+1}) = r(y^{k+1})$ .

This specifies the algorithm.



Let us clarify the selection of  $I^{k+1}$ . From (a) the indices in  $I^k$  and the index  $k$  are the only candidates. We may drop such a candidate index if from (b) there is sufficient decrease in objective value, i.e.,

$$r(y^{k+1}) \leq r(z^i) - \delta(z^i),$$

and if from (c), dropping it does not cause the objective value to increase over its value at  $y^{k+1}$ . Note (b) and (c) must both be satisfied before  $i$  can be eliminated.

To prove convergence, let  $\delta'$  be the separator corresponding to  $\varphi$  and let  $\delta''$  be the separator induced from  $\delta$  and  $r$  as in Lemma 6. We shall use the separator

$$\delta^* = \min [\delta', \delta'']$$

to show that

$$z^{k+1} \notin B(z^i, \delta^*(z^i)), \quad i = 1, \dots, k.$$

Convergence will then follow from Algorithm 0.

If  $i \notin I^{k+1}$ , then by the definition of  $\delta''$  in Lemma 6,

$$z^{k+1} \notin B(z^i, \delta''(z^i)).$$

If  $i \in I^{k+1}$ , then since

$$z^{k+1} \in Z^{k+1} \subset C^i,$$

where  $C^i \in \varphi(z^i)$ ,

$$z^{k+1} \notin B(z^i, \delta'(z^i)).$$

Therefore,  $z^{k+1} \notin B(z^i, \delta^*(z^i))$  and convergence is proved.

ALGORITHM III. This algorithm uses features of both I and II and drops cuts en masse if sufficient improvement has been made in objective value. Let  $\delta$  be a separator and  $\varphi$  a cut map, both on  $X \sim G$ .

To begin iteration  $(k, i)$  we have  $z^1, \dots, z^k = y_1^k, \dots, y_i^k$ , and  $Z_i^k$  where  $y_i^k \in \Gamma(Z_i^k)$ . If  $y_i^k \in G$ , terminate; otherwise:

(a) let  $C \in \varphi(y_i^k)$ ;

(b) let  $y \in \Gamma(Z_i^k \cap C)$ :

(i) if  $r(y) \leq r(z^k) - \delta(z^k)$ , then choose  $Z_1^{k+1}$  such that  $r^*(Z_1^{k+1}) \leq r(y)$  and let  $y_1^{k+1} = z^{k+1} \in \Gamma(Z_1^{k+1})$ ;

(ii) otherwise, let  $Z_{i+1}^k = Z_i^k \cap C$  and  $y_{i+1}^k = y$ .

Iteration  $(k, i)$  is complete.

If the sequence  $z^1, z^2, \dots$  generated is infinite, then

$$r(z^k) \leq r(z^i) - \delta(z^i)$$

for  $i < k$ . Hence  $z^k \notin B(z^i, \delta'(z^i))$ , where  $\delta'$  is the separator of Lemma 6 induced by  $r$  and  $\delta$ , and convergence follows from Algorithm 0. If the generated sequence  $y_1^k, y_2^k, \dots$  is infinite, then the algorithm relapses into Algorithm I and any cluster point of  $y_1^k, y_2^k, \dots$  solves (1).

ALGORITHM IV. Let  $\mathcal{Z}_c$  be a compact subset of  $\mathcal{Z}$  which is closed under intersection and such that if  $Z \in \mathcal{Z}_c$ , then

$$(2) \quad z \in \Gamma(Z), \quad z \notin G \quad \text{imply} \quad \Gamma(Z) = \{z\}.$$

Let  $\varphi: X \sim G \rightarrow \mathcal{L}_c$  be a cut map. The algorithm is:

Assume  $Z^i \in \mathcal{L}_c$  and  $z^i \in \Gamma(Z^i)$  for  $i = 1, \dots, k$  have been generated. Let  $C^k \in \varphi(z^k)$  and select  $Z^{k+1} \in \mathcal{L}_c$  such that

$$r^*(Z^{k+1}) \leq r^*(Z^k \cap C^k).$$

Let  $z^{k+1} \in \Gamma(Z^{k+1})$ .

To prove convergence, suppose  $z^k \rightarrow z^\infty \notin G$  on some subsequence  $K$ . Without loss of generality and using Lemma 3, we can assume that  $Z^k \rightarrow Z^\infty$  and  $C^k \rightarrow C^\infty$  on the same subsequence  $K$ . By Lemma 7 and (2), we have  $\{z^\infty\} = \Gamma(Z^\infty)$ . By Lemma 4,  $z^\infty \notin C^\infty$  and we have  $r^*(Z^\infty \cap C^\infty) < r(z^\infty)$ . But consider  $r(z^\infty) = \lim r(z^k) = \lim r(z^{k+1}) \leq \lim r^*(Z^k \cap C^k) \leq r^*(Z^\infty \cap C^\infty) < r(z^\infty)$ . This is a contradiction.

Observe that the bonafide example which meets the conditions of Algorithm IV is where  $G, X$ , and all cuts are convex and where  $r$  is strictly quasi-concave.

**3. Computational procedures.** The methods, Algorithms I–IV, of § 2 will now be further specialized into the specific computational algorithms KCG, VN, and DW. In each of these procedures,  $X$  will be a compact polyhedral convex subset of  $E^n$ ,  $G$  will be a compact convex subset of  $X$ , and cuts will be of the form  $X \cap H \supset G$ , where  $H$  is a half-space of  $E^n$ .

The cuts for the procedures are generated from a common form. In detail, let us define the half-space

$$H(v) = \{x|a(v) + b(v)x \geq 0\} \subset E^n$$

for given functions  $a: Y \rightarrow E^1$ , and  $b: Y \rightarrow E^n$  where  $Y \subset E^m$ . Then using the point to set map  $\alpha: X \sim G \rightarrow Y$ , specify the point to set map  $\varphi: X \sim G \rightarrow \mathcal{L}$  by mapping  $z$  to

$$\{H(v) \cap X|v \in \alpha(z)\}.$$

From a point  $z$ , after calculating  $v \in \alpha(z)$ , we shall construct the cutting half space  $H(v)$ . The next proposition demonstrates that  $\varphi$  will be a cut map if  $a$  and  $b$  are continuous functions, the map  $\alpha$  is closed, and if the half spaces “cut” off the point  $z$ .

LEMMA 8. Let  $\emptyset \neq G \subset X$  be compact convex sets in  $E^n$ ,  $Y \subset E^m$  be compact, and define the map  $\varphi: X \sim G \rightarrow \mathcal{L}$  by

$$\varphi(z) = \{H(v) \cap X|v \in \alpha(z)\}$$

where

$$H(v) = \{x|a(v) + b(v)x \geq 0\}.$$

Suppose:

- (i) the functions  $a: Y \rightarrow E^1$  and  $b: Y \rightarrow E^n$  are continuous and the map  $\alpha: X \sim G \rightarrow Y$  is closed;
- (ii)  $z \in X \sim G$  and  $H \cap X \in \varphi(z)$  imply  $z \notin H$ .

Then  $\varphi$  is a cut map.

Proof. Define  $\delta(z)$  for  $z \in X \sim G$  to be

$$\inf \{d(z, y)|y \in H(v) \cap X, v \in \alpha(z)\}.$$

We shall show  $\delta$  to be a separator and establish the result; assume that  $z^k \rightarrow z^\infty$  and  $\delta(z^k) \rightarrow 0$ . By compactness of  $X$  and  $Y$  and by closedness of  $\alpha$  we can select  $v^k \in \alpha(z^k)$  and  $y^k \in H(v^k) \cap X$  so that  $\delta(z^k) = d(z^k, y^k)$ . Then, again, by compactness and closedness, for appropriate subsequences  $v^k \rightarrow v^\infty \in \alpha(z^\infty)$  and  $y^k \rightarrow y^\infty$ . By continuity of  $a$  and  $b$ ,  $y^\infty \in H(v^\infty) \cap X$ . Since  $d(z^k, y^k) \rightarrow 0$ ,  $y^\infty = z^\infty \in H(v^\infty) \cap X$ . But if  $z^\infty \notin G$ , then by hypothesis,  $z^\infty \notin H(v^\infty)$ , which is a contradiction, thus completing the proof.

The map  $\alpha$  and functions  $a$  and  $b$  will differ for the various procedures KCG, VN, and DW. The sets  $Z^k$  will always have the form

$$X \cap \left( \bigcap_{i \in I} H^i \right),$$

where  $H^i \cap X$  is a cut. The objective function  $r$  will be linear for Algorithms I, II, and III, and consequently, finding  $z^k \in \Gamma(Z^k)$  will be a linear program. In Algorithm IV we shall let  $\mathcal{L}_c$  be the set of convex sets in  $\mathcal{L}$ ; however, a linear function would not necessarily meet the condition (2). Here, in general,  $r$  will not be linear.

To specify KCG, VN, and DW and their variants in terms of Algorithms I–IV, we need to define the goal set  $G$ , the polyhedral set  $X$  which encloses  $G$ , a cut map  $\varphi$ , an objective function  $r$ , and a separator  $\delta$ . Algorithm I for KCG, VN, and DW corresponds to the original statements of these algorithms. Algorithms II and III are variants which permit deletion of cuts, but again the number of cuts which will be dropped is problem dependent, and at this time, is essentially unpredictable. These procedures, II and III, will probably be very useful for some problems and worthless for others, but in any case the effort to implement them seems comparatively small. Algorithm IV requires more conditions than the others, namely (2), but offers the advantage that all inactive cuts can be dropped.

*The Kelley, Cheney, Goldstein concave cutting plane method (KCG).* KCG is designed to solve the program

$$(3) \quad \begin{aligned} &\text{maximize} && qx = \sum q_i x_i \\ &\text{subject to} && g(x) \geq 0, \end{aligned}$$

where  $g$  is real-valued, continuous, and concave. That the problem is equivalent to the general concave programming problem is discussed in [9]. We define  $G$  as the feasible set

$$G = \{x | g(x) \geq 0\}$$

and assume that  $G$  is contained in the compact set  $X$ , where

$$X = \{x | Ax \geq b\}.$$

Let  $U(z)$  be the set of vectors  $u$  such that

$$g(x) \leq g(z) + u(x - z)$$

for all  $x \in X$ . We assume that  $U(z) \neq \emptyset$  for  $z \in X$  and that

$$\bigcup_{z \in X} U(z)$$

is bounded. To define the cut map  $\varphi$  let  $\alpha$  be the closed map

$$\alpha(z) = \{(z, u) \mid u \in U(z)\}$$

and let

$$\varphi(z) = \{H(z, u) \cap X \mid (z, u) \in \alpha(z)\},$$

where

$$H(z, u) = \{x \mid g(z) + u(x - z) \geq 0\}.$$

Since  $\alpha$  is closed, using Lemma 8 we see that  $\varphi$  is a cut map.

For Algorithms I, II, and III we let  $r(x) = qx$ . For Algorithm IV we assume that we have a continuous function  $p$  on  $X$  such that  $r(x) = qx + p(x)$  satisfies (2) (where  $\mathcal{L}_c$  is the collection of convex sets in  $\mathcal{L}$ ) and that  $r(x) = qx$  for  $x \in G$  (for the procedure to be practicable, we would probably also need  $r$  to be quasi-concave on  $X$ ). Possible candidates for  $p$  would be functions of type  $\min(0, g(z))$  or  $-(\min(0, g(z)))^2$ .

For separators for Algorithms II and III one could use, for example,  $\delta(z) = g^2(z)$  or  $-g(z)$ .

Each method, of course, generates sequences  $z^k$  whose cluster points solve (3). We now discuss the individual algorithms in detail.

*Algorithm I for the KCG cut map.* Algorithm I is the original KCG method. Let  $Z^1 = X$ ; then given  $Z^k$ ,

$$Z^{k+1} = Z^k \cap H^k,$$

where  $H^k \in \varphi(z^k)$  and  $z^k \in \Gamma(Z^k)$ .

Explicitly,  $Z^k$  will have the form

$$Z^k = \{x \mid Ax \geq b, g(z^i) + u^i(x - z^i) \geq 0, i = 1, \dots, k - 1\}.$$

Then, given  $Z^k$ , solve the problem

$$\text{maximize } qx_{x \in Z^k}$$

for the optimal point  $z^k$ . If  $g(z^k) \geq 0$ , so that  $z^k \in G$ , terminate. Otherwise, obtain  $u^k \in U(z^k)$ , set

$$H^k = \{x \in X \mid g(z^k) + u^k(x - z^k) \geq 0\},$$

and let

$$Z^{k+1} = Z^k \cap H^k.$$

Observe that the dimension of the linear programming subproblem increases by one each iteration.

*Algorithm II with the KCG cut map.* Use any separator  $\delta$ . The set  $Z^k$  will be of the form

$$Z^k = X \cap \left( \bigcap_{i \in I^k} H^i \right) = \left\{ x \mid \begin{array}{l} Ax \geq b \\ g(z^i) + u^i(x - z^i) \geq 0, i \in I^k \end{array} \right\},$$

and  $z^k \in \Gamma(Z^k)$  will be given. Calculate

$$z^{k+1} \in \Gamma(Z^k \cap H^k),$$

where  $H^k \in \varphi(z^k)$ . Define  $J \subset \{k\} \cap I^k$  to be the set of indices  $i$  such that both

$$g(z^i) + u^i(z^{k+1} - z^i) > 0$$

and

$$r(z^{k+1}) \leq r(z^i) - \delta(z^i).$$

Then let

$$I^{k+1} = I^k \cup \{k\} \sim J.$$

*Algorithm III with the KCG cut map.* Given  $z^1, \dots, z^k = y_1^k, \dots, y_i^k$ , and  $Z_i^k$ , we iterate.

- (a) Let  $H \in \varphi(y_i^k)$ .
- (b) Let  $y \in \Gamma(Z_i^k \cap H)$ .
- (i) If

$$r(y) \leq r(z^k) - \delta(z^k),$$

set  $y_1^{k+1} = z^{k+1} = y$  and let  $Z^{k+1}$  be the intersection of  $X$  with the active cuts at  $z^{k+1}$ .

- (ii) Otherwise, let  $Z_{i+1}^k = Z_i^k \cap H$  and  $y_{i+1}^k = y$ .

The iteration is complete and the next step is (a).

Briefly, at  $Z^k$  add cuts until the optimal objective function value of the sub-problem drops by  $\delta(z^k)$ . Determine the active constraints at the corresponding optimal point and intersect these with  $X$  to form  $Z^{k+1}$ .

*Algorithm IV and KCG.* Assuming that  $r(z) = qz + p(z)$  satisfies (2) and that  $r(z) = qz$  for  $z \in G$ , the specifics of Algorithm IV are straightforward. Given  $Z^k$  and  $z^k \in \Gamma(Z^k)$ , compute  $z^{k+1} \in \Gamma(Z^k \cap H^k)$ , where  $H^k \in \varphi(z^k)$ . Then drop the inactive cuts at  $z^{k+1}$  so that  $Z^{k+1}$  will be  $X$  intersected with the active half-spaces at  $z^{k+1}$ .

Observe that at most  $n$  (the dimension of the space) cuts need to be retained (although it might be necessary to drop some active but redundant cuts). Of course, this procedure could also be applied to solve

$$\begin{aligned} &\text{maximize } r(x) \\ &\text{subject to } g(x) \geq 0, \end{aligned}$$

where  $r$  is strictly quasi-concave on  $X$ .

*The Veinott supporting hyperplane method (VN).* VN is designed to solve

$$(4) \quad \begin{aligned} &\text{maximize } qx \\ &\text{subject to } g_i(x) \geq 0, \quad i = 1, \dots, m, \end{aligned}$$

where the  $g_i$  are quasi-concave and continuously differentiable. As in KCG, we let

$$G = \{x | g_i(x) \geq 0, i = 1, \dots, m\}$$

and assume that

$$X = \{x | Ax \geq b\}$$

is given and that  $G \subset X$  are compact. It is assumed that  $g_i(a) > 0$  for  $i = 1, \dots, m$  for some  $a$  and that  $g_i(u) = 0$  implies  $\nabla g_i(u) \neq 0$  for  $i = 1, \dots, m$  for each boundary point  $u$  of  $G$ .

Given any point  $z \in X \sim G$ , let  $u(z)$  be the point on the boundary of  $G$  that is on the line segment between  $z$  and  $a$ . Note that  $u$  is continuous on  $X \sim G$ . Let  $\alpha(z)$  be the set

$$\{(u, j) | g_j(u) = 0, j = 1, \dots, m\}.$$

We define the cut map  $\varphi$  by assigning to  $z$  the cuts  $H \cap X$ , where

$$H = \{x | \nabla g_j(u)'(x - u) \geq 0\},$$

where  $(u, j) \in \alpha(z)$ ; see Lemma 8.

The objective  $r$  for Algorithms I, II, III, and IV are chosen exactly as in KCG. Here a candidate for the function  $p$  is

$$- \sum_1^m (\min(0, g_i(z)))^2.$$

For a separator, Algorithms II and III can use

$$\delta(z) = \sum_1^m (g_i(z))^2.$$

Algorithm I for these  $G, X, \varphi, r$  is the original VN. Again, Algorithms II and III provide a vehicle for dropping cuts. Algorithm IV, if  $p$  is available, only requires that at most  $n$  constraints be retained; the subproblems, nevertheless, are nonlinear. Each of these procedures solves (4).

The Algorithm IV version of VN can also be used to solve the problem

$$\begin{aligned} &\text{maximize } r(x) \\ &\text{subject to } g_i(x) \geq 0, \quad i = 1, \dots, m, \end{aligned}$$

where  $r$  is strictly quasi-concave.

*Dantzig-Wolfe method (DW)*. The DW method solves the concave programming problem

$$(5) \quad \begin{aligned} &\text{maximize } f(t) \\ &\text{subject to } g(t) \geq 0, \quad t \in T, \end{aligned}$$

where  $g = (g_1, \dots, g_m)$ , the functions  $f$  and  $g_i$  are concave and continuous, and the set  $T$  is convex and compact. It is further assumed that there is a given point  $\bar{t} \in T$  for which  $g(\bar{t}) > 0$ .

As in [9] we shall deviate slightly from the Dantzig and Wolfe attitude in [3] by treating the dual of the problem, which is

$$(6) \quad \begin{aligned} &\text{minimize } u \\ &\text{subject to } (\lambda, u) \in H(t) \text{ for all } t \in T \end{aligned}$$

and

$$\lambda \geq 0,$$

where  $\lambda = (\lambda_1, \dots, \lambda_m) \in E^m$ ,  $u \in E^1$ , and  $H(t) = \{(\lambda, u) | f(t) + \lambda g(t) \leq u\}$ . Here the constraints range over the entire set  $T$ .

Observe that we can adjoin the additional constraint

$$u \leq \bar{u},$$

where

$$\bar{u} \geq \max \{f(t) | t \in T\}$$

without altering the solution set. Moreover, because  $g(\bar{t}) > 0$ , this additional constraint produces a compact feasible region, as is shown in [3] or [9].

We identify the dual problem with (1) by letting

$$X = \{(\lambda, u) | f(\bar{t}) + \lambda g(\bar{t}) \leq u, u \leq \bar{u}, \lambda \geq 0\},$$

$$G = \{(\lambda, u) | (\lambda, u) \in H(t) \text{ for all } t \in T\} \cap X$$

and  $r(\lambda, u) = -u$ . Consequently, determining  $z^* \in \Gamma(G)$  will solve the dual problem.

For each of the algorithms, the  $Z^k$  sets will be of the form  $X \cap (\bigcap_{i \in I^k} H(t^i))$ , where  $I^k$  is an appropriate index set. Thus, if we let  $t^0 = \bar{t}$  and always maintain  $0 \in I^k$  for  $k = 0, 1, 2, \dots$ , then the subproblems,  $\max \{r(z) : z \in Z^k\}$ , are the linear programs

$$\text{maximize } -u$$

$$\text{subject to } f(t^i) + \lambda g(t^i) \leq u, \quad i \in I^k, \quad \lambda \geq 0.$$

Note we have dropped the constraint  $u \leq \bar{u}$  as it clearly does not alter the solution.

Given  $(\lambda, u) \in X \sim G$  the map  $\alpha(\lambda, u)$  determines the set of all  $t^* \in T$  that solve

$$\text{maximize}_{t \in T} f(t) + \lambda g(t).$$

Then for  $t^* \in \alpha(\lambda, u)$ , we specify  $H$ , where  $H \cap X \in \varphi(\lambda, u)$ , by

$$H = \{(\lambda, u) | f(t^*) + \lambda g(t^*) \leq u\}.$$

This calculation of  $H$  determines the cut map  $\varphi$  for DW.

It follows easily that the procedures converge to a solution of the dual problem (6).

For the separator of Algorithms II and III use the one induced by the cut map; that is, for each  $(\lambda, u) \in X \sim G$  let

$$\delta(\lambda, u) = \varepsilon(f(t^*) - \lambda g(t^*) - u)$$

for some  $t^* \in \alpha(\lambda, u)$ , where  $0 < \varepsilon < 1$  is fixed (here a small  $\varepsilon$  seems best). Specifications of Algorithms I, II and III are precisely analogous to the KCG method except for the different cut map  $\varphi$  and separator  $\delta$ . Algorithm I is the original DW method.

Although Algorithms I, II and III solve the dual problems, these procedures simultaneously solve the primal problem (5) as observed in [3] or more directly in [9].

*Remark.* This paper has developed cutting plane methods from the underlying principle of separators. However, it is evident that this principle is applicable not just to cutting plane methods but to nonlinear programming in general.

## REFERENCES

- [1] C. BERGE, *Topological Spaces*, Oliver and Boyd, Edinburgh and London, 1963.
- [2] E. W. CHENEY AND A. A. GOLDSTEIN, *Newton's method of convex programming and Tchebycheff approximation*, Numer. Math., 1 (1959), pp. 253–268.
- [3] G. B. DANTZIG, *Linear Programming and Extensions*, Princeton University Press, Princeton, N.J., 1963.
- [4] F. HAUSDORFF, *Set Theory*, 2nd ed., Chelsea, New York, 1962.
- [5] J. E. KELLEY, *The cutting-plane method for solving convex programs*, J. Soc. Indust. Appl. Math., 8 (1960), pp. 703–712.
- [6] D. M. TOPKIS, *Cutting plane methods without nested constraint sets*, Operations Res., 18 (1970), pp. 404–413.
- [7] ———, *A note on cutting plane methods without nested constraint sets*, Ibid., to appear.
- [8] A. F. VEINOTT, JR., *The supporting hyperplane method for unimodal programming*, Ibid., 15 (1967), pp. 147–152.
- [9] W. I. ZANGWILL, *Nonlinear Programming: A Unified Approach*, Prentice-Hall, Englewood Cliffs, N.J., 1969.



## A NOTE ON COMPLETE CONTROLLABILITY\*

INGE TROCH†

**Abstract.** For linear autonomous multivariable control systems the conditions are studied under which the system is always—that is, for every control matrix of appropriate rank—completely controllable, which means that for every initial state one can construct a control which steers the system to any given final state in finite time. The resultant necessary and sufficient condition says that the number of independent controls either must equal the system's dimension or must be smaller by one. In the latter case no eigenvalue of the system matrix may be real. These results are valid for continuous systems as well as for sampled-data systems and are derived by means of well-known criteria for complete controllability and some theorems of matrix theory. Analogous results which are valid for the question of complete observability are formulated too. As a special case, previous results of Hsin Chu are covered.

**1. Introduction.** In a previous work [1] Hsin Chu has shown that linear continuous multivariable control systems, which can be steered by a single control, are always completely controllable if and only if the dimension of the control system equals two and the eigenvalues of the system matrix are not real. In this work, for control systems of dimension  $n$  with  $m$  independent controls, conditions will be studied which guarantee that the control system is always completely controllable. By the principle of duality similar theorems for complete observability can be formulated. For the concepts mentioned above and terminology the reader is referred to [5] or [7].

**2. Notation and formulation of the problem.** In the following, linear multivariable control systems will be studied, which can be described by

$$(1) \quad \dot{x} = Ax + Bu, \quad y = Cx + Du,$$

where  $x$  denotes the  $n$ -dimensional state vector ( $x \in R^n$ ),  $u \in R^m$  the control, and  $y \in R^s$  the output vector;  $A, B, C, D$  are constant matrices of appropriate dimension. Moreover,  $\lambda$  denotes an eigenvalue of the matrix  $A$ , that is, a solution of the equation ( $I = \text{diag}(1, 1, \dots, 1)$ , unity matrix)

$$(2) \quad \det(A - \lambda I) = 0;$$

$(\cdot, \cdot)$  denotes the scalar product of two vectors;  $A^T$  denotes the transpose of the matrix  $A$ ; and  $a_{ik}$  denotes the elements of the matrix  $A$ . As has been said, the system has  $m$  independent controls, which means that the matrix  $B$  is of rank  $m$ . Without loss of generality we may restrict ourselves to matrices  $B$  with  $m$  linearly independent column vectors  $b_1, b_2, \dots, b_m$ . There is now the problem of deriving conditions on  $A$  and  $m$  so that for every matrix  $B$  with

$$(3) \quad \text{rank } B = m,$$

the system (1) is completely controllable in the sense introduced by Kalman [5]. We note that, speaking intuitively, the system (1) is called completely controllable if, knowing the matrices  $A$  and  $B$  and the initial state, one can construct a control  $u$  which will bring the state to the zero state in finite time.

\* Received by the editors November 3, 1970, and in revised form March 15, 1971.

† Institut für Mathematik der Technischen Hochschule in Wien, Vienna, Austria.

**3. The condition on  $A$ .** As is well known [7] a necessary and sufficient condition for complete controllability of (1) is

$$(4) \quad \text{rank } Q := \text{rank } (B, AB, \dots, A^{n-1}B) = n.$$

As the formulated problem is quite trivial for  $m = n$ , this case will be omitted at first.

Assume now, that  $A$  has at least one real eigenvalue  $\lambda_1$ . Then there exists a matrix  $T$  over the real field, so that

$$A^{(1)} = TAT^{-1}$$

has only one nonzero element in its last line:

$$(5) \quad \begin{aligned} a_{nk}^{(1)} &= 0, & k &= 1, 2, \dots, n-1, \\ a_{nm}^{(1)} &= \lambda_1. \end{aligned}$$

Now let  $B$  be chosen in such a way that for the elements of the corresponding matrix  $B^{(1)} = TB$  the equations

$$b_{nk}^{(1)} = 0 \quad \text{for } k = 1, 2, \dots, m$$

hold. Then the last line of  $Q$  (built with  $A^{(1)}$  and  $B^{(1)}$ ) has no nonzero element, and therefore  $Q$  cannot be of rank  $n$ . We have the first necessary condition.

FIRST NECESSARY CONDITION. *The system (1) can be completely controllable for every matrix  $B$  of rank  $m < n$  only if no eigenvalue of the matrix  $A$  is real.*

**4. The condition on the rank  $B$ .** Note first that there exists a regular matrix  $T_1$  over the real field such that  $A$  is similar to a generalized diagonal matrix

$$A^{(2)} = T_1AT_1^{-1} = \text{diag } (P_1, \dots, P_k),$$

where  $P_j, j = 1, 2, \dots, k$ , corresponds to the  $j$ th elementary divisor (over the real field!) of  $A$ . If  $\mu_j \pm iv_j (v_j \neq 0$  because of § 3) denotes the corresponding eigenvalues, then  $P_j$  is of the form

$$(6) \quad \begin{pmatrix} D_j & I & 0 & \dots & 0 \\ & & & \dots & \\ & 0 & \dots & D_j & I \\ & 0 & & & D_j \end{pmatrix}$$

with

$$D_j = \begin{pmatrix} \mu_j & v_j \\ -v_j & \mu_j \end{pmatrix}, \quad j = 1, \dots, k.$$

For the proof consult [2, p. 106].

Assume now that  $m \leq n - 2$ , and take a matrix  $B$  of rank  $m$  with  $(B^{(2)} = T_1B)$

$$(7) \quad b_{n-1,j}^{(2)} = b_{n,j}^{(2)} = 0 \quad \text{for } j = 1, 2, \dots, m.$$

It then follows that the last two lines of the matrix  $A^{(2)}B^{(2)}$ —and also of  $A^{(2)2}B^{(2)}, \dots, A^{(2)m-1}B^{(2)}$ —have no nonzero elements, and therefore the rank of  $Q$  (built with  $A^{(2)}$  and  $B^{(2)}$ ) is less than or equal to  $n - 2$ . This gives the second necessary condition.

SECOND NECESSARY CONDITION. *The system (1) can be completely controllable for every matrix  $B$  of rank  $m$  only if  $m \geq n - 1$ .*

**5. Necessary and sufficient condition.** It shall now be shown that the two necessary conditions together are sufficient too, that is, the validity of the following theorem is shown.

**THEOREM 1.** *For every matrix  $B$  of rank  $m$  system (1) is completely controllable if and only if  $(\alpha)$  or  $(\beta)$  holds:*

$(\alpha)$   $m = n$ ;

$(\beta)$   $m = n - 1$ , and no eigenvalue of  $A$  is real.

*Proof.* Necessity follows from §§ 3 and 4.

The proof of sufficiency is still left: Condition (4) is trivially fulfilled if  $(\alpha)$  holds. Now let condition  $(\beta)$  hold and assume that there exists a matrix  $B$  of rank  $n - 1$  such that

$$(8) \quad \text{rank } Q = n - 1.$$

Note first that there exists a vector  $c \in R^n$ , which is orthogonal to every vector  $b_j$  of  $B$ :

$$(9) \quad (c, b_j) = 0 \quad \text{for } j = 1, 2, \dots, n - 1,$$

and which is unique except for multiplication by a scalar. From (8) it follows that every vector  $Ab_j$ ,  $j = 1, 2, \dots, n - 1$ , is linearly dependent on the vectors  $b_1, \dots, b_{n-1}$ . But from this, (9) and the properties of the scalar product, we have also

$$(10) \quad (c, Ab_j) = 0 \quad \text{for every } j = 1, 2, \dots, n - 1.$$

Now let  $S_{n-1}$  be the  $(n - 1)$ -dimensional subspace of  $R^n$  spanned by the vectors  $b_1, \dots, b_{n-1}$ . Then  $R^n$  can be written as the direct sum (see e.g. [4])

$$R^n = S_{n-1} \oplus c.$$

Conditions (9) and (10) can be written as

$$(9') \quad B^T c = 0$$

and

$$(10') \quad (AB)^T c = 0,$$

and the last equation is equivalent to

$$(11) \quad (B^T A^T) c = B^T (A^T c) = 0.$$

From (9'), (11) and the uniqueness of  $c$  (in the defined sense) it follows that

$$A^T c = \lambda c \quad \text{for some } \lambda \in R,$$

which means that  $A$  possesses  $c$  as a left eigenvector belonging to the real eigenvalue  $\lambda$ , in contradiction to assumption  $(\beta)$ . This completes the proof.

**6. The dual statement.** We shall now make use of the second equation in (1) and state a theorem—dual to Theorem 1—which gives necessary and sufficient

conditions for the system (1) to be completely observable for every matrix  $C$  with appropriate rank. We just remind the reader that complete observability means that the knowledge of (1) and the output  $y(t)$  over a finite interval is sufficient to determine uniquely the initial state of the system for every initial state. The precise definition can again be found for example in [7]. By the principle of duality we have the following theorem.

**THEOREM 2.** *System (1) is completely observable for every matrix  $C$  of rank  $s$  if and only if  $(\gamma)$  or  $(\delta)$  holds:*

$(\gamma)$   $s = n$ ;

$(\delta)$   $s = n - 1$ , and no eigenvalue of  $A$  is real.

**7. Sampled-data systems.** As is well known, for the discrete sampled-data system

$$(12) \quad \begin{aligned} x((k+1)T) &= Ax(kT) + Bu(kT), \\ y(kT) &= Cx(kT) + Du(kT), \quad k = 0, 1, 2, \dots, \end{aligned}$$

with regular matrix  $A$ , equation (4) is a necessary and sufficient condition for complete controllability too. This means that Theorems 1 and 2 remain valid and unchanged for the multivariable control system described by (12).

#### REFERENCES

- [1] HSIN CHU, *A remark on complete controllability*, this Journal, 3 (1966), pp. 439–442.
- [2] E. A. CODDINGTON AND N. LEVINSON, *Theory of Ordinary Differential Equations*, McGraw-Hill, New York, 1955.
- [3] F. R. GANTMACHER, *Matrizenrechnung*, Bd.I, Deutscher Verlag der Wissenschaften, Berlin, 1958.
- [4] W. GRÖBNER, *Matrizenrechnung*, BI-Hochschultaschenbücher, Bd. 103/103a, Mannheim, 1966.
- [5] R. E. KALMAN, *On the general theory of control systems*, Proc. First International Congress on Automatic Control, Eyre and Spottiswoode, London, 1961.
- [6] R. E. KALMAN, Y. C. HO AND K. S. NARENDRA, *Controllability of linear dynamical systems*, Contributions to Differential Equations, vol. I, John Wiley, New York, 1961, pp. 189–213.
- [7] L. A. ZADEH AND C. A. DESOER, *Linear System Theory*, McGraw-Hill, New York, 1963.

## ABSTRACT MODELS FOR THE SYNTHESIS OF OPTIMIZATION ALGORITHMS\*

GERARD G. L. MEYER AND E. POLAK†

**Abstract.** In this paper, we present a systematic approach to the problem of synthesis of optimization algorithms. First, we develop abstract models for algorithms. These models guide the inventive process toward "conceptual" algorithms, namely, algorithms which may consist of operations that are inadmissible in a practical method (for example, a conceptual algorithm may require us to find the limit of an infinite sequence at each iteration). Once the abstract models are established, we present a set of methods for converting "conceptual" algorithms falling into the class defined by the abstract models, into "implementable" iterative procedures.

**Introduction.** The convergence of optimization algorithms has been studied extensively in recent years (see [4], [5], [6], [7], [8], [9], [10], [11], [12], [13]). The approach generally adopted in these studies consisted of defining a class of algorithms and then giving convergence theorems which applied to every algorithm in this class. This approach has resulted in the development of general procedures which have considerably simplified the task of establishing whether an algorithm is convergent. However, the emphasis so far has been on *analysis*. Very few of the existing results provide guidelines for the *synthesis* of algorithms.

In this paper we present a systematic approach to the problem of synthesis of optimization algorithms. The development of an algorithm usually evolves through three phases. The first is a heuristic, or invention phase in which intuition plays an extremely important part. In the second phase one transforms one's intuitive ideas into a "conceptual" algorithm, i.e., an algorithm which may consist of operations that are inadmissible in a practical method. (For example, a conceptual algorithm may require us to find the limit of an infinite sequence at each iteration.) The last phase consists of converting the "conceptual" algorithm into an "implementable" algorithm. Our approach to the problem of synthesis consists of two parts. First we develop abstract models for algorithms. These models guide the inventive process towards "conceptual" algorithms which will later be easily made implementable. Once the abstract models are established, we present a set of methods for converting "conceptual" algorithms, falling into the class defined by the abstract models, into "implementable" iterative procedures.

One of the most frequently occurring difficulties in the implementation of a conceptual algorithm is the requirement that an implicit relation be solved at each iteration, e.g., minimize a function along a line, maximize a linear function in a convex set, etc. Generally, we only have methods for constructing a sequence whose limit point satisfies such an implicit relation. Now it is well known that no limit point of an infinite sequence can be determined on a digital computer in a finite time. Consequently, the task most frequently encountered in the design of a transition from a "conceptual" algorithm to an "implementable" one is that of finding methods for avoiding the need to construct limit points.

\* Received by the editors January 19, 1970, and in final revised form February 12, 1971.

† Department of Electrical Engineering and Computer Sciences, Electronics Research Laboratory, University of California, Berkeley, California 94720. The first author is now at Department of Electrical Engineering, University of Southern California, Los Angeles, California 90007. This work was sponsored by the National Aeronautics and Space Administration under Grant NGL-05-003-016(Sup.6).

In this paper, we propose two methods for obviating the need for constructing infinite sequences in subprocedures. The first is a truncation procedure and is presented in § 2, the second is an  $\varepsilon$ -approximation procedure which is presented in § 4.

The scope, to which a paper must be held, does not permit us to copiously illustrate the applicability of the ideas presented. It is our hope, however, that the two examples given in § 3 will convince the reader of the great usefulness of the approach described.

**1. Abstract models for a class of iterative procedures.** Throughout this paper, we shall assume that we are given a closed and bounded subset  $T$  of  $R^n$  in which we wish to find points with a specific property  $\pi$ . We shall call points in  $T$  with the property  $\pi$  *desirable*.

The simplest algorithms for finding desirable points in  $T$  are composed of a map  $\xi(\cdot)$  from  $T$  into  $R^1$  and of a map  $A(\cdot)$  from  $T$  into all the subsets of  $T$ , having the following form.

ALGORITHM 1.

Step 0. Compute a point  $z_0$  in  $T$  and set  $i = 0$ .

Step 1. Compute a point  $z_{i+1}$  in  $A(z_i)$ .

Step 2. If  $\xi(z_{i+1}) < \xi(z_i)$ , set  $i = i + 1$  and go to Step 1; otherwise stop.

DEFINITION 1. We shall say that an iterative procedure of the form of Algorithm 1 is *convergent* if any sequence of points  $\{z_i\}$ , generated by it, satisfies one of the following conditions:

(i) If the sequence  $\{z_i\}$  is finite, i.e.,  $\{z_i\} = \{z_0, z_1, \dots, z_k\}$ , then  $z_{k-1}$  is desirable.

(ii) If the sequence  $\{z_i\}$  is infinite, then any of its cluster points is desirable.

It is quite easy to show that Algorithm 1 is convergent under the following assumption (see [11]).

HYPOTHESIS 1 [Zangwill].

(i)  $z$  in  $T$  is desirable if there exists at least one point  $a$  in  $A(z)$  such that

$$\xi(a) - \xi(z) \geq 0;^1$$

(ii)  $\xi(\cdot)$  is continuous on  $T$ ;

(iii)  $A(\cdot)$  is upper semicontinuous on  $T$ .<sup>2</sup>

Hypothesis 1 is not the only assumption which ensures that Algorithm 1 is convergent. We now state another assumption which ensures that Algorithm 1 is convergent (see [9]).

HYPOTHESIS 2 [Polak].

(i)  $z$  in  $T$  is desirable if there exists at least one point  $a$  in  $A(z)$  such that  $\xi(a) - \xi(z) \geq 0$ .

(ii)  $\xi(\cdot)$  is bounded from below on  $T$ .

(iii) If  $z$  in  $T$  satisfies  $\xi(a) - \xi(z) < 0$  for all  $a$  in  $A(z)$ , then there exist an  $\varepsilon > 0$  and a  $\delta > 0$  such that  $\xi(a') - \xi(z') \leq -\delta < 0$  for all  $z'$  in  $T$  such that  $\|z' - z\| \leq \varepsilon$ , for all  $a'$  in  $A(z')$ .

<sup>1</sup> Note: Part (i) of Hypothesis 1 simply states that the set of points  $z$  in  $T$ , for which there exists at least one point  $a$  in  $A(z)$  satisfying  $\xi(a) - \xi(z) \geq 0$ , is a subset of the set of desirable points.

<sup>2</sup> Note: Let  $A(\cdot)$  be a map from  $T$  into all the subsets of  $T$ . If for any sequence  $\{y_i\}$  converging to  $y^*$  and for any sequence  $\{a_i\}$  converging to  $a^*$ , with  $a_i$  in  $A(y_i)$ ,  $a^*$  belongs to  $A(y^*)$ , then we say that the map  $A(\cdot)$  is upper semicontinuous on  $T$ .

*Remark.* It can be shown that Hypothesis 2 is weaker than Hypothesis 1 (see [1]).

In this paper we differentiate between explicit and implicit algorithms. This differentiation is largely *heuristic* but is extremely important in the construction of computationally efficient algorithms.

By an *explicit* algorithm we shall mean an algorithm of the form of Algorithm 1 in which the computation of a point  $y$  in  $A(z)$  for  $z$  in  $T$  can be carried out in a reasonably straightforward manner.

Explicit algorithms do not lead to computational difficulties and therefore we shall say no more about them. Implicit algorithms on the other hand cannot be readily implemented on a digital computer, as we shall shortly show, and must therefore be regarded as “conceptual” rather than as “practical” algorithms. The following sections of this paper will be devoted to developing methods for modifying convergent implicit algorithms in such a way as to produce convergent algorithms in explicit form.

We shall consider two abstract models of implicit algorithms. The first one, which is defined below, uses a map  $A(\cdot)$  such that to compute a point  $z$  in  $A(z)$ , we must solve an implicit equation. To obtain some motivation for the specific decomposition of  $A(\cdot)$  in Definition 2, below, the reader should digress for a moment and examine Problem 1, the steepest descent Algorithm 8 and Definition 5, in § 3.1.

**DEFINITION 2.** Let  $U(\cdot)$  be a map from  $T$  into all the subsets of  $T$  and let  $\alpha(\cdot)$  and  $\gamma(\cdot)$  be maps from  $T$  into  $R^1$ . Then, for every  $z$  in  $T$ , we define the set  $A(z)$  as consisting of all  $y$  in  $U(z)$  such that  $\gamma(y) = \alpha(z)$ .

The following assumption ensures that the set  $A(z)$  is not empty.

**HYPOTHESIS 3.** Given any point  $z$  in  $T$ , there exists a point  $y$  in  $U(z)$  such that  $\alpha(z) = \gamma(y)$ .

In this case, Algorithm 1 takes on the following expanded form.

**ALGORITHM 2.**

*Step 0.* Compute a point  $z_0$  in  $T$  and set  $i = 0$ .

*Step 1.* Compute a point  $z_{i+1}$  in  $U(z_i)$ , satisfying  $\gamma(z_{i+1}) = \alpha(z_i)$ .

*Step 2.* If  $\xi(z_{i+1}) < \xi(z_i)$ , set  $i = i + 1$  and go to Step 1; otherwise stop.

**HYPOTHESIS 4.**

(i) A point  $z$  in  $T$  is desirable if there exists at least one point  $y$  in  $U(z)$  such that  $\alpha(z) = \gamma(y)$  and  $\xi(y) - \xi(z) \geq 0$ .

(ii) The maps  $\alpha(\cdot)$ ,  $\gamma(\cdot)$  and  $\xi(\cdot)$  are continuous on  $T$ .

(iii) The map  $U(\cdot)$  is upper semicontinuous on  $T$ .

The proof of the following proposition is easy and has been omitted.

**PROPOSITION 1.** *If the maps  $\alpha(\cdot)$ ,  $\gamma(\cdot)$ ,  $\xi(\cdot)$ , and  $U(\cdot)$  satisfy Hypotheses 3 and 4, then the map  $A(\cdot)$  given by Definition 2 is upper semicontinuous on  $T$  and Algorithm 2 is convergent.*

We now present the second specific form of the map  $A(\cdot)$  that we wish to consider. This form is characterized by the fact that, to find points in  $A(z)$ , we must compute intermediate points. To obtain some motivation for the specific decomposition of the map  $A(\cdot)$  in Definition 3, below, the reader should digress for a moment and examine Problem 2, the Frank and Wolfe algorithm (Algorithm 11) and Definition 7, in § 3.2.

DEFINITION 3. Let  $U$  be a closed and bounded subset of  $R^n$ , let  $\beta(\cdot, \cdot)$  be a map from  $T \times U$  into  $R$ , let  $b(\cdot, \cdot)$  be a map from  $T \times U$  into  $T$  and let  $\alpha(\cdot)$  be a map from  $T$  into  $R$ . Then for every  $z$  in  $T$ , we define the set  $A(z)$  as follows:

$$A(z) = \{y | y = b(z, w), w \in U \text{ such that } \beta(z, w) = \alpha(z)\}.$$

The following assumption ensures that for every  $z$  in  $T$  the set  $A(z)$  is non-empty.

HYPOTHESIS 5. Given any  $z$  in  $T$ , there exists a point  $w$  in  $U$  such that  $\beta(z, w) = \alpha(z)$ .

With the map  $A(\cdot)$  defined as above, Algorithm 1 expands as follows.

ALGORITHM 3.

Step 0. Compute a point  $z_0$  in  $T$  and set  $i = 0$ .

Step 1. Compute a point  $w_i$  in  $U$ , satisfying  $\alpha(z_i) = \beta(z_i, w_i)$ .

Step 2. Set  $z_{i+1} = b(z_i, w_i)$ .

Step 3. If  $\xi(z_{i+1}) < \xi(z_i)$ , set  $i = i + 1$  and go to Step 1; otherwise stop.

HYPOTHESIS 6.

(i) A point  $z$  in  $T$  is desirable if there exists at least one point  $w$  in  $U$  satisfying  $\alpha(z) = \beta(z, w)$  such that  $\xi(b(z, w)) - \xi(z) \geq 0$ .

(ii) The maps  $\alpha(\cdot)$  and  $\xi(\cdot)$  are continuous on  $T$ .

(iii) The maps  $\beta(\cdot, \cdot)$  and  $b(\cdot, \cdot)$  are jointly continuous on  $T \times U$ .

PROPOSITION 2. If the maps  $\alpha(\cdot)$ ,  $\beta(\cdot, \cdot)$ ,  $b(\cdot, \cdot)$  and  $\xi(\cdot)$  satisfy Hypotheses 5 and 6, then the map  $A(\cdot)$  given by Definition 3 is upper semicontinuous on  $T$  and Algorithm 3 is convergent.

The transformation of an implicit algorithm into an explicit one usually depends upon the specific conceptual method one has in mind for finding points in the set  $A(z)$ . We shall now introduce two infinite subprocedures for calculating points in  $A(z)$ , and in the next section we shall show how these subprocedures can be truncated to produce convergent explicit algorithms.

We begin by considering a subprocedure for computing points in the set  $A(z)$  when the map  $A(\cdot)$  is defined as in Definition 2 (we suppose, of course, that we are unaware of a more straightforward method for calculating points in  $A(z)$ ). Let us denote by  $N$  the set of positive integers and suppose that we have a mapping  $m(\cdot, \cdot, \cdot)$  from  $T \times T \times N$  into  $T$  which satisfies the following assumption.

HYPOTHESIS 7.

(i)  $m(z, y, j)$  belongs to  $U(z)$  for all  $z$  in  $T$ ,  $y$  in  $U(z)$  and  $j$  in  $N$ .

(ii) The sequence  $\{y(m(z, y, j))\}_{j=0}^\infty$  converges to  $\alpha(z)$  for all  $z$  in  $T$  and  $y$  in  $U(z)$ .<sup>3</sup>

Clearly, for any  $z$  in  $T$  and  $y$  in  $U(z)$ , every cluster point of the sequence  $\{m(z, y, j)\}_{j=0}^\infty$  is in  $A(z)$  given by Definition 2. When the map  $m(\cdot, \cdot, \cdot)$  is introduced, Algorithm 2 assumes the following specific form.

ALGORITHM 4.

Step 0. Compute a point  $z_0$  in  $T$  and set  $i = 0$ .

Step 1. Compute a point  $y_i$  in  $U(z_i)$  and let  $z_{i+1}$  be any cluster point of the sequence  $\{m(z_i, y_i, j)\}_{j=0}^\infty$ .

<sup>3</sup> Note that part (ii) of Hypothesis 7 does not imply that the sequence  $m(z, y, j)$  converges.



*Step 2.* If  $\xi(z_{i+1}) < \xi(z_i)$ , set  $i = i + 1$  and go to Step 1; otherwise stop.

In view of Proposition 1, the following proposition is obvious.

**PROPOSITION 3.** *If the maps  $\alpha(\cdot)$ ,  $\gamma(\cdot)$ ,  $\xi(\cdot)$ ,  $U(\cdot)$  satisfy Hypotheses 3 and 4, and the map  $m(\cdot, \cdot, \cdot)$  satisfies Hypothesis 7, then Algorithm 4 is convergent.*

*Remark.* Obviously, Algorithm 4 cannot be implemented on a digital computer since it would inevitably jam up in Step 1.

We now consider a subprocedure for computing points in  $A(z)$  when the map  $A(\cdot)$  is defined as in Definition 3. Thus, suppose that we have a map  $m(\cdot, \cdot, \cdot)$  from  $T \times U \times N$  into  $T$  which satisfies the following assumption.

**HYPOTHESIS 8.**

(i)  $m(z, y, j)$  belongs to  $U$  for all  $z$  in  $T$ ,  $y$  in  $U$  and  $j$  in  $N$ .

(ii) The sequence  $\{\beta(z, m(z, y, j))\}_{j=0}^{\infty}$  converges to  $\alpha(z)$  for all  $z$  in  $T$  and  $y$  in  $U$ .

When such a map is introduced, Algorithm 3 assumes the following specific form.

**ALGORITHM 5.**

*Step 0.* Compute a point  $z_0$  in  $T$  and set  $i = 0$ .

*Step 1.* Compute a point  $y_i$  in  $U$  and let  $w_i$  be any cluster point of the sequence  $\{m(z_i, y_i, j)\}_{j=0}^{\infty}$ .

*Step 2.* Set  $z_{i+1} = b(z_i, w_i)$ .

*Step 3.* If  $\xi(z_{i+1}) < \xi(z_i)$ , set  $i = i + 1$  and go to Step 1; otherwise stop.

In view of Proposition 2, the following proposition is obvious.

**PROPOSITION 4.** *If the maps  $\alpha(\cdot)$ ,  $\beta(\cdot, \cdot)$ ,  $\xi(\cdot)$  and  $b(\cdot, \cdot)$  satisfy Hypotheses 5 and 6 and the map  $m(\cdot, \cdot, \cdot)$  satisfies Hypothesis 8, then Algorithm 5 is convergent.*

Again it is clear that if implemented on a digital computer, Algorithm 5 would inevitably jam up in Step 1.

**2. Truncation methods.** As we have just seen, Algorithms 4 and 5 will inevitably jam up in Step 1 since it is impossible to compute cluster points of infinite sequences in a finite time by means of a digital computer. Even if one relies on the finite word length of a digital computer to stop calculations after a finite time, this finite time will usually be prohibitively long. Consequently, some sort of truncation procedure must be used in converting these algorithms into a more realistic form.

We begin by defining a class of maps from  $N$  into  $N$  (the set of all positive integers) which will be called truncation functions.

**DEFINITION 4.** We shall say that a map  $\ell(\cdot)$ , from  $N$  into  $N$ , is a *truncation function* if  $\ell(i) \geq 1$  for all  $i$  in  $N$ , and given any  $m$  in  $N$ , there exists a  $k$  in  $N$  such that  $\ell(i) \geq m$  for all  $i \geq k$ ,  $i$  in  $N$ .

We first use truncation functions in Algorithm 4 which then takes on the following form.

**ALGORITHM 6.** Let  $\ell(\cdot)$  be a given truncation function.

*Step 0.* Compute a point  $z_0$  in  $T$  and set  $i = 0$ .

*Step 1.* Compute a point  $y_i$  in  $U(z_i)$  and set  $z_{i+1} = m(z_i, y_i, \ell(i))$ .

*Step 2.* If  $\xi(z_{i+1}) < \xi(z_i)$ , set  $i = i + 1$  and go to Step 1; otherwise stop.

*Remark.* In Step 1 of Algorithm 6, it is required that a point  $y_i$  be found in  $U(z_i)$ . We suppose that this task is easy. In fact, in almost all applications,  $z$  belongs to  $U(z)$ , and a natural choice for  $y_i$  in  $U(z_i)$  consists in letting  $y_i = z_i$ . In order to

ensure that cluster points of infinite sequences generated by Algorithm 6 are desirable the following assumption must be made.

**HYPOTHESIS 9.**

(i) The sequence  $\{\gamma(m(z, y, j))\}_{j=0}^\infty$  is monotonically decreasing for all  $z$  in  $T$  and  $y$  in  $U(z)$ .

(ii) Given any  $z$  in  $T$ ,  $y$  in  $U(z)$ , and  $\delta > 0$ , there exists an  $\varepsilon > 0$  and a  $k$  in  $N$ , depending on  $z, y$  and  $\delta$ , such that  $\gamma(m(z', y', j)) - \gamma(m(z, y, j)) \leq \delta$  for all  $z'$  in  $T$  such that  $\|z' - z\| \leq \varepsilon$ , for all  $y'$  in  $U(z')$  such that  $\|y' - y\| \leq \varepsilon$  and for all  $j \geq k$ .

(iii) If  $\xi(m(z, y, j)) \geq \xi(z)$  for some  $y$  in  $U(z)$  and for some  $j \geq 1$ , then  $z$  is desirable.

**PROPOSITION 5.** *If the maps  $\alpha(\cdot), \gamma(\cdot), \xi(\cdot), U(\cdot)$  satisfy Hypotheses 3 and 4, the map  $m(\cdot, \cdot, \cdot)$  satisfies Hypotheses 7 and 9 and the map  $\ell(\cdot)$  is a truncation function, then Algorithm 6 is convergent.*

*Proof.* Because of part (iii) of Hypothesis 9, the case of finite sequences is trivial.

Now consider an infinite sequence  $\{z_i\}$  generated by Algorithm 6, and let  $z^*$  be a cluster point of this sequence, i.e., let  $K_2$  be a subset of the integers such that the subsequence  $\{z_i\}_{K_2}$  converges to  $z^*$ . Consider the sequences  $\{z_{i+1}\}_{K_2}$  and  $\{y_i\}_{K_2}$  in  $T$ . The compactness of  $T$  implies that there exists  $K_1$ , a subset of  $K_2$ , such that the subsequence  $\{z_{i+1}\}_{K_1}$  converges to  $z^{**}$  in  $T$ , and the subsequence  $\{y_i\}_{K_1}$  converges to  $y^*$  in  $T$ . The property of convergent sequences implies that the subsequence  $\{z_i\}_{K_1}$  also converges to  $z^*$ .

In order to show that  $z^*$  is desirable it is enough to establish the three following facts:

- (i)  $z^{**}$  belongs to  $U(z^*)$ ;
- (ii)  $\alpha(z^*) = \gamma(z^{**})$ ;
- (iii)  $\xi(z^{**}) - \xi(z^*) \geq 0$ .

By construction,  $z_{i+1} = m(z_i, y_i, \ell(i))$  and  $y_i$  is in  $U(z_i)$ . It now follows from (i) of Hypothesis 7 that  $z_{i+1}$  is in  $U(z_i)$ . Next, since the map  $U(\cdot)$  is upper semi-continuous, we must have  $z^{**}$  and  $y^*$  in  $U(z^*)$ .

Let  $\delta > 0$ ; then part (ii) of Hypothesis 7 implies that there exists  $k_4$  in  $N$  such that  $\gamma(m(z^*, y^*, j)) - \delta/2 \leq \alpha(z^*)$  for all  $j \geq k_4$ . From Hypothesis 9 it follows that there exist an  $\varepsilon > 0$  and  $k_3$  such that  $\gamma(m(z', y', j)) - \gamma(m(z^*, y^*, j)) \leq \delta/2$  for all  $z'$  in  $T$  such that  $\|z' - z^*\| \leq \varepsilon$ , for all  $y'$  in  $U(z')$  such that  $\|y' - y^*\| \leq \varepsilon$  and for all  $j \geq k_3$ . Let  $k_2 = \max(k_3, k_4)$ ; then there exists  $k_1$  in  $N$  such that  $\gamma(m(z_i, y_i, k_2)) - \delta \leq \alpha(z^*)$  for all  $i \geq k_1, i$  in  $K_1$ . Since the map  $\ell(\cdot)$  is a truncation function, there exists a  $k_0$  in  $N$  such that  $\ell(i) \geq k_2$  for all  $i \geq k_0, i$  in  $K_1$ . It follows from part (i) of Hypothesis 9 that  $\gamma(m(z_i, y_i, \ell(i))) - \delta \leq \alpha(z^*)$  for all  $i \geq k, i$  in  $K_1$ , where  $k = \max(k_0, k_1)$ . But  $\gamma(\cdot)$  is continuous and therefore  $\gamma(z^{**}) - \delta \leq \alpha(z^*)$ . Since this is true for any  $\delta > 0$ , it follows that  $\gamma(z^{**}) \leq \alpha(z^*)$ .

Now, given any  $z_i$  and  $y_i$  in  $U(z_i)$ , part (i) of Hypothesis 9 implies that  $\gamma(m(z_i, y_i, j)) \geq \alpha(z_i)$  for all  $j$  in  $N$ , i.e.,  $\gamma(m(z_i, y_i, \ell(i))) \geq \alpha(z_i)$ . It follows that  $\gamma(z^{**}) \geq \alpha(z^*)$ , and therefore,  $\gamma(z^{**}) = \alpha(z^*)$ .

Now suppose that  $\xi(z^{**}) - \xi(z^*) < 0$ . The continuity of  $\xi(\cdot)$  implies that there exist a  $\delta > 0$  and an  $\varepsilon > 0$  such that

$$\xi(z'') - \xi(z') \leq -\delta < 0$$

for all  $z''$  in  $T$  such that  $\|z'' - z^{**}\| \leq \varepsilon$ , for all  $z'$  in  $T$  such that  $\|z' - z^*\| \leq \varepsilon$ . It follows that there exists a  $k$  such that  $\xi(z_{i+1}) - \xi(z_i) \leq -\delta$  for all  $i \geq k, i$  in  $K_1$ . The sequence  $\{\xi(z_i)\}$  is monotonically decreasing, the set  $T$  is bounded, and therefore  $\{\xi(z_i)\}$  converges to  $\xi^*$ , which contradicts the fact that  $\xi(z_{i+1}) - \xi(z_i) \leq -\delta$  for all  $i \geq k, i$  in  $K_1$ . It follows that  $\xi(z^{**}) - \xi(z^*) \geq 0$  and  $z^*$  is desirable.

We now use truncation functions in Algorithm 5 to construct the following form.

ALGORITHM 7. Let  $\ell(\cdot)$  be a given truncation function.

Step 0. Compute a point  $z_0$  in  $T$  and set  $i = 0$ .

Step 1. Compute a point  $y_i$  in  $U$  and set  $w_i = m(z_i, y_i, \ell(i))$ .

Step 2. Set  $z_{i+1} = b(z_i, w_i)$ .

Step 3. If  $\xi(z_{i+1}) < \xi(z_i)$ , set  $i = i + 1$  and go to Step 1; otherwise stop.

Remark. In Step 1 of Algorithm 7, it is required that a point  $y_i$  be found in the set  $U$ . We suppose that this task is easy. In fact, in almost all applications,  $z$  belongs to  $U$ , and a natural choice for  $y_i$  consists in letting  $y_i = z_i$ .

In order to ensure that cluster points of infinite sequences generated by Algorithm 7 are desirable, the following assumption must be made.

HYPOTHESIS 10.

(i) The sequence  $\{\beta(z, m(z, y, j))\}_{j=0}^\infty$  is monotonically decreasing for all  $z$  in  $T$  and  $y$  in  $U$ .

(ii) Given any  $z$  in  $T$ ,  $y$  in  $U$  and  $\delta > 0$ , there exist an  $\varepsilon > 0$  and a  $k$  in  $N$ , possibly depending on  $z, y$ , and  $\delta$ , such that

$$\beta(z', m(z', y', j)) - \beta(z, m(z, y, j)) \leq \delta$$

for all  $z'$  in  $T$  such that  $\|z' - z\| \leq \varepsilon$ , for all  $y'$  in  $U$  such that  $\|y' - y\| \leq \varepsilon$  and for all  $j \geq k$ .

(iii) If  $\xi(b(z, m(z, y, j))) \geq \xi(z)$  for some  $y$  in  $U$  and for some  $j \geq 1$ , then  $z$  is desirable.

The following proposition can be proved easily, following the same type of argument as in the proof of Proposition 5.

PROPOSITION 6. *If the maps  $\alpha(\cdot), \beta(\cdot, \cdot), \xi(\cdot), b(\cdot, \cdot)$  satisfy Hypotheses 5 and 6, the map  $m(\cdot, \cdot, \cdot)$  satisfies Hypotheses 8 and 10 and the map  $\ell(\cdot)$  is a truncation function, then Algorithm 7 is convergent.*

**3. Applications.** In order to clarify the concepts and methods exposed in the preceding sections we are now going to examine two specific problems and the algorithms usually used to solve them.

**3.1. Unconstrained minimization problems.** In this subsection we shall examine the following classical problem.

PROBLEM 1. Find a  $\hat{z}$  in  $R^n$  such that

$$f(\hat{z}) \leq f(z) \quad \text{for all } z \text{ in } R^n,$$

where  $f(\cdot)$ , a convex map from  $R^n$  into  $R^1$ , is continuously differentiable, with the property that the set  $\{z | f(z) \leq \alpha\}$  is bounded for every  $\alpha$  in  $R^1$ .

Suppose that we have a point  $z_0$  in  $R^n$ . We define  $T$  as

$$T = \{y \in R^n | f(y) \leq f(z_0)\}.$$

In this case, a point  $\hat{z}$  in  $T$  is desirable if and only if it minimizes  $f(z)$  over  $T$ , i.e.,  $f(\hat{z}) \leq f(z)$  for all  $z$  in  $T$ . Since  $f(\cdot)$  is convex and continuously differentiable, we recognize  $\hat{z}$  to be desirable if and only if  $\nabla f(\hat{z}) = 0$ .

We propose to use the steepest descent method for solving Problem 1, i.e., the following algorithm.

ALGORITHM 8.

Step 0. Let  $z_0$  be a point in  $R^n$  and set  $i = 0$ .

Step 1. Let  $z_{i+1}$  be any point satisfying:

- (i)  $z_{i+1}$  belongs to  $U(z_i)$ ;
- (ii)

$$f(z_{i+1}) = \min \{f(y) | y \in U(z_i)\},$$

where

$$U(z_i) = \{y \in R^n | y = z_i + v\nabla f(z_i), v \in [-1, 0]\}.$$

Step 2. If  $f(z_{i+1}) < f(z_i)$ , set  $i = i + 1$  and go to Step 1; otherwise stop.

Algorithm 8 can be seen to be of the form of Algorithm 2, with the maps  $\alpha(\cdot)$ ,  $\gamma(\cdot)$ ,  $\xi(\cdot)$ ,  $U(\cdot)$  in Algorithm 2 defined as follows.

DEFINITION 5. Let the maps  $\alpha(\cdot)$ ,  $\gamma(\cdot)$ ,  $\xi(\cdot)$  from  $T$  into  $R^1$  and the map  $U(\cdot)$  from  $T$  into all the subsets of  $T$  be defined as:

- (i)  $U(z) = \{y \in T | y = z + v\nabla f(z), v \in [-1, 0]\}$ ;
- (ii)  $\alpha(z) = \min \{f(y) | y \in U(z)\}$ ;
- (iii)  $\gamma(z) = f(z)$ ;
- (iv)  $\xi(z) = f(z)$ .

Since by inspection the maps  $\alpha(\cdot)$ ,  $\gamma(\cdot)$ ,  $\xi(\cdot)$  and  $U(\cdot)$  satisfy Hypotheses 3 and 4, Algorithm 8 is convergent when applied to Problem 1.

At this point, we see the first advantage of using abstract models of algorithms as a tool for proving convergence of algorithms. On the one hand they provide us with a pattern to follow, while on the other, we find that the proof of convergence of a specific algorithm becomes decomposed into fairly simple and independent parts.

In Algorithm 8, the computation of  $z_{i+1}$  from  $z_i$  is not explicit. We therefore proceed as indicated in § 2 in order to produce truncations in the calculation of  $z_{i+1}$  from  $z_i$ .

Making use of Proposition 3, we now define a subalgorithm which we shall use in order to modify Algorithm 8.

For every  $z$  in  $R^n$ ,  $y$  in  $R^n$  and  $j$  a positive integer, consider the following subprocedure.

ALGORITHM 9. Let  $\bar{\alpha} \in (0, 1)$  be given.

Step 0. Set  $y_0 = y$ ,  $i = 0$  and  $v = 0$ .

Step 1. If  $\langle \nabla f(z), \nabla f(y_i) \rangle = 0$ , go to Step 6; otherwise set

$$v = -\text{sgn} \langle \nabla f(z), \nabla f(y_i) \rangle$$

and go to Step 2.

Step 2. Set  $\tilde{y} = y_i + v\nabla f(z)$ .

Step 3. Compute  $\theta$  defined by  $\theta = f(\tilde{y}) - f(y_i) - \bar{\alpha}v \langle \nabla f(z), \nabla f(y_i) \rangle$ .

Step 4. If  $\theta \leq 0$ , set  $y_{i+1} = \tilde{y}$  and go to Step 5; otherwise set  $v = v/2$  and go to Step 2.

Step 5. If  $i < j$ , set  $i = i + 1$  and go to Step 1; otherwise go to Step 6.

Step 6. Set  $\tilde{v} = \text{sat } v$ ,<sup>4</sup> set  $y_j = y_0 + \tilde{v}\nabla f(z)$  and stop.

Remark. We choose this subalgorithm because it has a nontrivial amount of structure rather than because it is the best computationally. The structure of this subalgorithm should serve the purpose of illustrating the complexity that can be found in a subalgorithm for computing the values of the map  $m(\cdot, \cdot, \cdot)$ .

DEFINITION 6. Let  $m(\cdot, \cdot, \cdot)$  be the map from  $T \times T \times N$  into  $T$  defined by  $m(z, y, j) = y_j$ , where  $y_j$  is given by Algorithm 9.

It can be verified (see Appendix), that the map  $m(\cdot, \cdot, \cdot)$ , given by Definition 6, satisfies Hypotheses 7 and 9. Using the map  $m(\cdot, \cdot, \cdot)$  and a truncation function  $\ell(\cdot)$ , we obtain from Algorithm 8 the following "explicit" algorithm.

ALGORITHM 10. Let  $\ell(\cdot)$  be a given truncation function.

Step 0. Compute a  $z_0$  in  $R^n$  and set  $i = 0$ .

Step 1. Set  $z_{i+1} = m(z_i, z_i, \ell(i))$ .

Step 2. If  $f(z_{i+1}) < f(z_i)$ , set  $i = i + 1$  and go to Step 1; otherwise stop.

In view of Proposition 5, the following is obvious.

PROPOSITION 7. Algorithm 10 is convergent for Problem 1.

**3.2. Constrained minimization problem.** In this subsection we shall examine the following classical problem.

PROBLEM 2. Given  $T$  a closed, bounded convex subset of  $R^n$ , and  $t$ , a point in  $R^n$  but not in  $T$ , find  $\hat{z}$  in  $T$  such that

$$\|\hat{z} - t\| \leq \|z - t\| \quad \text{for all } z \text{ in } T.$$

We shall suppose that  $T$  is of the form

$$T = \{z \in R^n \mid f^i(z) \leq 0, i = 1, 2, \dots, m\},$$

where the maps  $f^i(\cdot)$  from  $R^n$  into  $R^1$  are continuously differentiable.

Suppose that we try to solve Problem 2 by means of the Frank-Wolfe algorithm.

ALGORITHM 11 (Frank-Wolfe).

Step 0. Compute a point  $z_0$  in  $T$  and set  $i = 0$ .

Step 1. Compute a point  $w_i$  in  $T$  satisfying  $\langle z_i - t, w_i \rangle \leq \langle z_i - t, w \rangle$  for all  $w$  in  $T$ .

Step 2. Let  $z_{i+1}$  be the point in  $[z_i, w_i]$  satisfying  $\|z_{i+1} - t\| \leq \|z - t\|$  for all  $z$  in  $[z_i, w_i]$ .<sup>5</sup>

Step 3. If  $\|z_{i+1} - t\| < \|z_i - t\|$ , set  $i = i + 1$  and go to Step 1; otherwise stop.

<sup>4</sup> Note: The function  $\text{sat}(\cdot): R^1 \rightarrow R^1$  is defined by

$$\text{sat}(v) = \begin{cases} v & \text{if } |v| \leq 1, \\ 1 & \text{if } v > 1, \\ -1 & \text{if } v < -1. \end{cases}$$

<sup>5</sup> Note: Given two points  $x_1$  and  $x_2$  in  $R^n$ , the set  $\{y \in R^n \mid y = vx_1 + (1 - v)x_2, 0 \leq v \leq 1\}$  is denoted by  $[x_1, x_2]$ .

Algorithm 11 can be seen to be of the form of Algorithm 3 by defining the following maps.

DEFINITION 7. Let the set  $U = T$  and the maps  $\alpha(\cdot)$ ,  $\zeta(\cdot)$  from  $T$  into  $R^1$ ,  $\beta(\cdot, \cdot)$  from  $T \times U$  into  $R^1$ ,  $b(\cdot, \cdot)$  from  $T \times U$  into  $T$ , be defined as follows:

- (i)  $\alpha(z) = \min \{ \alpha | \alpha = \langle z - t, w \rangle, w \in T \}$ ;
- (ii)  $\zeta(z) = \|z - t\|$ ;
- (iii)  $\beta(z, w) = \langle z - t, w \rangle$ ;
- (iv)  $b(z, w)$  is defined by:
  - (a)  $b(z, w)$  belongs to  $[z, w]$ ;
  - (b)  $\|b(z, w) - t\| \leq \|y - t\|$  for all  $y$  in  $[z, w]$ .

In order to show that Algorithm 11 is convergent for Problem 2, it suffices to show that the mappings  $\alpha(\cdot)$ ,  $\zeta(\cdot)$ ,  $\beta(\cdot, \cdot)$  and  $b(\cdot, \cdot)$  satisfy Hypotheses 5 and 6. It is easily verified that this is indeed so and we therefore conclude that Algorithm 11 is convergent for Problem 2.

Once again we see the advantage of using an abstract model in proving the convergence of a specific algorithm.

In Algorithm 11, both the point  $w_i$  and the point  $b(z_i, w_i)$  are defined by implicit relations. However, since the computation of  $b(z_i, w_i)$  from  $z_i$  and  $w_i$  is extremely simple, we shall consider that  $z_{i+1} = b(z_i, w_i)$  is an "explicit" function of  $z_i$  and  $w_i$ . Thus we shall consider that the only real difficulty lies in the computation of  $w_i$ .

To obtain from Algorithm 11 an "explicit" algorithm of the form of Algorithm 7, we must introduce a map  $m(\cdot, \cdot, \cdot)$  satisfying Hypotheses 8 and 10. For example, one can use a method of feasible directions to define such a map.

We now state a method of feasible directions in the required truncated form.

For every  $z$  in  $T$ ,  $y$  in  $T$  and positive integer  $j$ , consider the following algorithm.

ALGORITHM 12. Suppose that  $S$ , a compact neighborhood of the origin in  $R^n$ , and  $\varepsilon$ , a positive scalar, are given. Let  $f^0(\cdot)$  be the map from  $R^n$  into  $R^1$  defined by  $f^0(x) = \langle z - t, x \rangle$  for all  $x$  in  $R^n$ .

Step 0. Set  $y_0 = y$  and  $i = 0$ .

Step 1. Set  $\varepsilon_i = \varepsilon$ .

Step 2. Compute  $\phi_{\varepsilon_i}$  and  $h_{\varepsilon_i}$  by solving the following:

$$\phi_{\varepsilon_i} = \min_{h \in S} \max_{j \in J_{\varepsilon_i}} \langle \nabla f^j(y_i), h \rangle.$$

$h_{\varepsilon_i}$  is any vector in  $S$  such that

$$\phi_{\varepsilon_i} = \max_{j \in J_{\varepsilon_i}} \langle \nabla f^j(y_i), h_{\varepsilon_i} \rangle,$$

where

$$J_{\varepsilon_i} = \{ j \in \{1, 2, \dots, m\} | f^j(y_i) + \varepsilon_i \geq 0 \} \cup \{0\}.$$

Step 3. If  $\phi_{\varepsilon_i} < -\varepsilon_i$ , set  $h_i = h_{\varepsilon_i}$  and go to Step 5; otherwise, compute  $\phi_0$  defined by

$$\phi_0 = \min_{h \in S} \max_{j \in J_0} \langle \nabla f^j(y_i), h \rangle,$$

where

$$J_0 = \{ j \in \{1, 2, \dots, m\} | f^j(y_i) = 0 \} \cup \{0\},$$

and go to Step 4.

Step 4. If  $\phi_0 = 0$ , set  $y_{i+1} = y_i$  and go to Step 7; otherwise, set  $\varepsilon_i = \varepsilon_i/2$  and go to Step 2.

Step 5. Compute  $\lambda_i \geq 0$  such that

$$\lambda_i = \max \{ \lambda \mid f^j(y_i + \lambda h_i) \leq 0 \text{ for all } j = 1, 2, \dots, m \}.$$

Step 6. Set  $y_{i+1} = y_i + \lambda_i h_i$ .

Step 7. If  $i < j$ , set  $i = i + 1$  and go to Step 1; otherwise stop.

DEFINITION 8. Let  $m(\cdot, \cdot, \cdot)$  be the map from  $T \times T \times N$  into  $T$  defined by  $m(z, y, j) = y_j$ , where  $y_j$  is given by Algorithm 12.

It can be verified that the map  $m(\cdot, \cdot, \cdot)$  given by Definition 8 satisfies Hypotheses 8 and 10 (see [9]). Using the map  $m(\cdot, \cdot, \cdot)$  and a truncation function  $\ell(\cdot)$  in Algorithm 11, we obtain the following "explicit" method.

ALGORITHM 13. Let  $\ell(\cdot)$  be a given truncation function.

Step 0. Compute a point  $z_0$  in  $T$ , and set  $i = 0$ .

Step 1. Set  $w_i = m(z_i, z_i, \ell(i))$ .

Step 2. Let  $z_{i+1}$  be the point in  $[z_i, w_i]$  satisfying

$$\|z_{i+1} - t\| \leq \|z - t\| \quad \text{for all } z \text{ in } [z_i, w_i].$$

Step 3. If  $\|z_{i+1} - t\| < \|z_i - t\|$ , set  $i = i + 1$ , and go to Step 1; otherwise stop.

In view of Proposition 6, the following is obvious.

PROPOSITION 8. *Algorithm 13 is convergent for Problem 2.*

The use of the approach defined in this paper may show relations between different well-known algorithms. For example, if a function  $\ell(\cdot)$  from  $N$  into  $N$  defined by  $\ell(i) = 1$  for all  $i$  is used in Algorithm 13 instead of a truncation function, then Algorithm 13 is a method of feasible directions. On the other hand, if a "function"  $\ell(\cdot)$  from  $N$  into  $N$  "defined" by  $\ell(i) = \infty$  for all  $i$  is used in Algorithm 13 instead of a truncation function, then Algorithm 13 is the Frank–Wolfe algorithm. The use of a truncation function in Algorithm 13 thus produces algorithms which are "between" a method of feasible directions and the Frank–Wolfe algorithm.

**4.  $\varepsilon$ -approximations for a class of iterative procedures.** The approach to the synthesis of algorithms described in §§ 1 and 2 is by no means the only one possible. In this section we shall show that an alternative approach exists and gives extremely interesting results.

Throughout this section we shall consider maps  $A(\cdot)$  and  $\zeta(\cdot)$  satisfying Hypothesis 1 in § 1 and we shall suppose that it is impossible to use them in iterative procedures of the form of Algorithm 1 due to the fact that the computation of points  $z_{i+1}$  in  $A(z_i)$  is impossible (or that we are not aware of ways of doing it).

The main idea developed in this section consists in using a map  $B(\cdot, \cdot)$  from  $R_+^1 \times T$  into all the subsets of  $T$  such that the set  $B(0, z)$  is identical to the set  $A(z)$  for all  $z$ .

We shall suppose that the task of finding a  $y$  in  $B(\varepsilon, z)$  for  $\varepsilon > 0$  and  $z$  in  $T$  is relatively easy.

Consider the following algorithm.

ALGORITHM 14. Let  $\bar{\varepsilon} > 0$  be given.

Step 0. Compute  $z_0$  in  $T$  and set  $i = 0$ .

Step 1. Set  $\varepsilon = \bar{\varepsilon}$ .

Step 2. Find a point  $y$  in  $B(\varepsilon, z_i)$ .

Step 3. If  $\xi(y) - \xi(z_i) \leq -\varepsilon$ , set  $\varepsilon_i = \varepsilon$ ,  $z_{i+1} = y$ ,  $i = i + 1$  and go to Step 1; otherwise set  $\varepsilon = \varepsilon/2$  and go to Step 2.

We remark that if we let  $\bar{\varepsilon} = 0$ , then Algorithm 14 is almost of the form of Algorithm 1. The following assumption on  $B(\cdot, \cdot)$  ensures that cluster points of infinite sequences generated by Algorithm 14 are desirable.

HYPOTHESIS 11.

(i)  $B(0, z) = A(z)$  for all  $z$  in  $T$ .

(ii)  $B(\cdot, \cdot)$  is jointly upper semicontinuous on  $T$ ; i.e., for any sequence  $\{\varepsilon_i\}$  converging to  $\varepsilon^*$ , for any sequence  $\{z_i\}$  converging to  $z^*$ , for any sequence  $\{y_i\}$  converging to  $y^*$ , with  $y_i$  in  $B(\varepsilon_i, z_i)$ ,  $y^*$  belongs to  $\beta(\varepsilon^*, z^*)$ .

THEOREM 1. *If the maps  $A(\cdot)$  and  $\xi(\cdot)$  satisfy parts (i) and (ii) of Hypothesis 1, and the map  $B(\cdot, \cdot)$  satisfies Hypothesis 11, then every cluster point of an infinite sequence generated by Algorithm 14 is desirable.*

*Proof.* Let  $\{z_i\}$  be an infinite sequence generated by Algorithm 14, and let  $z^*$  be a cluster point of this sequence. Thus, for some subset  $K_1$  of the integers, the subsequence  $\{z_i\}_{i \in K_1}$  converges to  $z^*$ . Consider the sequence  $\{z_{i+1}\}_{i+1 \in K_1}$  in  $T$  and the sequence  $\{\varepsilon_i\}_{i \in K_1}$  in  $[0, \bar{\varepsilon}]$ . The boundedness of  $T$  and of the interval  $[0, \bar{\varepsilon}]$  ensures that there exists  $K$ , an infinite subset of  $K_1$ , such that the subsequences  $\{z_{i+1}\}_{i+1 \in K}$  and  $\{\varepsilon_i\}_K$  converge to points  $z^{**}$  and  $\varepsilon^*$  in  $T$  and  $[0, \bar{\varepsilon}]$ , respectively. The properties of convergent sequences ensure that the subsequence  $\{z_i\}_K$  also converges to  $z$ .

Now suppose that  $\varepsilon^* > 0$ ; then the form of Algorithm 14 implies that there exists an integer  $k$  such that  $\xi(z_{i+1}) - \xi(z_i) \leq \varepsilon^*/2$  for all  $i \geq k$ ,  $i$  in  $K$ , and this contradicts part (ii) of Hypotheses 1. Consequently  $\varepsilon^* = 0$ .

The map  $B(\cdot, \cdot)$  is jointly upper semicontinuous on  $T$ , and  $z_{i+1}$  belongs to  $B(\varepsilon_i, z_i)$  for all  $i$ . It follows that  $z^{**}$  is in  $B(\varepsilon^*, z^*)$ , i.e., in  $A(z^*)$ . Consequently  $z^*$  is desirable.

HYPOTHESIS 12.

(i)  $B(0, z) = A(z)$  for all  $z$  in  $T$ .

(ii) If  $\xi(a) - \xi(z) < 0$  for all  $a$  in  $A(z)$ , then there exist  $\varepsilon > 0$ ,  $\delta > 0$  and  $\gamma > 0$ , possibly depending on  $z$ , such that

$$\xi(b') - \xi(z') \leq -\gamma < 0$$

for all  $z'$  in  $T$  such that  $\|z' - z\| \leq \varepsilon$  and for all  $b'$  in  $B(v, z')$ ,  $0 \leq v \leq \delta$ .

The proof of the following theorem can be carried out by using the same types of arguments as were used to prove Theorem 1, and it is therefore omitted.

THEOREM 2. *If the maps  $A(\cdot)$  and  $\xi(\cdot)$  satisfy parts (i) and (ii) of Hypothesis 2, and the map  $B(\cdot, \cdot)$  satisfies Hypothesis 12, then every cluster point of an infinite sequence generated by Algorithm 14 is desirable.*

*Remark.* It can be shown that if maps  $A(\cdot)$ ,  $\xi(\cdot)$  and  $B(\cdot, \cdot)$  satisfy Hypotheses 1 and 11, they satisfy Hypotheses 2 and 12 (see [1]).

For examples of how the  $\varepsilon$ -procedure is used in the synthesis of algorithms, see E. Polak [14].



**5. Conclusion.** To conclude, we would like to highlight the two most important aspects of the theory we have presented in this paper. The first is that by using models one can separate out the essential properties of an algorithm from the nonessential ones. Thus, for example, in a gradient method one need not specify in advance exactly which procedure one will use to search along the direction of steepest descent, one only has to specify that the search procedure will have certain properties. The second point that we wish to emphasize is that given a “conceptual” algorithm in which one has to perform in sequence several operations each of which requires an infinite number of iterations, one can obtain an “implementable” algorithm by “shuttling” between these infinite operations, combining them into a single infinite operation. Thus, our method of obtaining an “implementable” algorithm from a plurally infinitely iterative “conceptual” algorithm consists in “parallelizing” the infinite operations of the conceptual algorithm.

Obviously, this paper does not exhaust the study of models for computational methods or of the possibilities of constructing “implementable” algorithms from “conceptual” ones. We hope that this paper will lead to further work, in particular in the study of algorithms with infinite memory.

**Appendix.** We shall assume, throughout this Appendix, that the maps  $f(\cdot)$ ,  $\alpha(\cdot)$ ,  $\gamma(\cdot)$ ,  $\zeta(\cdot)$ ,  $m(\cdot, \cdot, \cdot)$  and  $U(\cdot)$ , and the set  $T$ , are as defined in § 3.1. The fact that the map  $m(\cdot, \cdot, \cdot)$  satisfies part (i) of Hypothesis 7 and part (i) of Hypothesis 9 is an obvious consequence of its definition, and the fact that the map  $m(\cdot, \cdot, \cdot)$  satisfies part (ii) of Hypothesis 7 and part (ii) of Hypothesis 9 is a direct consequence of the following theorem.

**THEOREM A.** *Given any  $z$  in  $T$  and  $\delta > 0$ , there exist  $\varepsilon > 0$  and  $k$  in  $N$ , depending on  $z$  and  $\delta$ , such that  $f(m(z', y', j)) \leq \alpha(z') + \delta$ , for all  $z'$  in  $T$ , satisfying  $\|z' - z\| \leq \varepsilon$ , for all  $y'$  in  $U(z')$ , for all  $j \geq k$ .*

The proof of Theorem A consists of several steps which we shall state as lemmas. These lemmas being quasi-trivial, their proofs have been deleted.

**LEMMA A.1.** *Given any  $z$  in  $T$ , satisfying  $\|\nabla f(z)\| > 0$ , and  $\delta > 0$ , there exist  $\varepsilon > 0$  and  $\rho > 0$ , depending on  $z$  and  $\delta$ , such that  $|\langle \nabla f(z'), \nabla f(y') \rangle| \geq \rho > 0$ , for all  $z'$  in  $T$  such that  $\|z' - z\| \leq \varepsilon$ , for all  $y'$  in  $U(z')$  such that  $f(y') \geq \alpha(z') + \delta$ .*

**LEMMA A.2.** *Given any  $z$  in  $T$ ,  $y$  in  $U(z)$ ,  $j$  in  $N$  and  $\rho > 0$  satisfying  $|\nabla f(m(z, y, j))$ ,  $\nabla f(z)| \geq \rho > 0$ , then there exists  $v > 0$ , depending on  $\rho$  only, such that  $f(m(z, y, j + 1)) - f(m(z, y, j)) \leq -v$ .*

**LEMMA A.3.** *Given any  $z$  in  $T$ , satisfying  $\|\nabla f(z)\| > 0$ , and  $\delta > 0$ , there exist  $\varepsilon > 0$ ,  $\rho > 0$  and  $v > 0$ , depending on  $z$  and  $\delta$ , such that  $f(m(z', y', j + 1)) - f(m(z', y', j)) \leq -v$ , for all  $z'$  in  $T$ , such that  $\|z' - z\| \leq \varepsilon$ , for all  $y'$  in  $U(z')$  such that  $f(m(z', y', j)) \geq \alpha(z') + \delta$ .*

**LEMMA A.4.** *Given any  $z$  in  $T$ , satisfying  $\|\nabla f(z)\| > 0$ , and  $\delta > 0$ , there exist  $\varepsilon > 0$  and  $k$  in  $N$ , depending on  $z$  and  $\delta$ , such that  $f(m(z', y', j)) \leq \alpha(z') + \delta$ , for all  $z'$  in  $T$ , such that  $\|z' - z\| \leq \varepsilon$ , for all  $y'$  in  $U(z')$ .*

**LEMMA A.5.** *Given any  $z$  in  $T$ , satisfying  $\|\nabla f(z)\| = 0$ , and  $\delta > 0$ , there exists  $\varepsilon > 0$ , depending on  $\delta$  only, such that  $f(m(z', y', j)) \leq \alpha(z') + \delta$ , for all  $z'$  in  $T$  satisfying  $\|z' - z\| \leq \varepsilon$ , for all  $y'$  in  $U(z')$ , for all  $j$  in  $N$ .*

## REFERENCES

- [1] E. MICHAEL, *Topologies on spaces of subsets*, Trans. Amer. Math. Soc., 71 (1951), pp. 152–182.
- [2] M. FRANK AND P. WOLFE, *An algorithm for quadratic programming*, U.S. Nav. Res. Logist. Quart., 3 (1956), pp. 95–110.
- [3] G. ZOUTENDIJK, *Methods of Feasible Directions*, Elsevier, Amsterdam, 1960.
- [4] B. T. POLYAK, *Gradient methods for the minimization of functionals*, U.S.S.R. Comput. Math. and Math. Phys., 3 (1963), no. 4, pp. 864–878 = Zh. Vychisl. Mat. i Mat. Fiz., 3 (1963), no. 4, pp. 643–653.
- [5] E. S. LEVITIN AND B. T. POLYAK, *Constrained minimization methods*, U.S.S.R. Comput. Math. and Math. Phys., 6 (1966), no. 5, pp. 1–50 = Zh. Vychisl. Mat. i Mat. Fiz., 6 (1966), no. 5, pp. 787–823.
- [6] W. I. ZANGWILL, *Convergence conditions for nonlinear programming algorithms*, Working Paper 197, Center for Research in Management Science, Univ. of California, Berkeley, 1966.
- [7] D. M. TOPKIS AND A. VEINOTT, *On the convergence of some feasible directions algorithms for nonlinear programming*, this Journal, 5 (1967), pp. 268–279.
- [8] E. POLAK, *On primal and dual methods for solving discrete optimal control problems*, Proc. Second International Conference on Computing Methods in Optimization Problems (San Remo, Italy, 1968), Academic Press, New York, 1969.
- [9] ———, *On the convergence of optimization algorithms*, Rev. Francaise Informatique et Recherche Operationelle, Ser. Rouge, 16 (1969), pp. 17–34.
- [10] ———, *Computational methods in discrete optimal control and nonlinear programming: A unified approach*, Memo. ERL-M261, Electronics Research Laboratory, College of Engineering, Univ. of California, Berkeley, 1969.
- [11] W. I. ZANGWILL, *Nonlinear Programming: A Unified Approach*, Prentice-Hall, Englewood Cliffs, N.J., 1969.
- [12] G. MEYER AND E. POLAK, *Abstract models for the synthesis of optimization algorithms*, Memo. ERL-268, Electronics Research Laboratory, College of Engineering, Univ. of California, Berkeley, 1969.
- [13] R. MEYER, *The validity of a family of optimization methods*, this Journal, 8 (1970), pp. 41–54.
- [14] E. POLAK, *Computational Methods in Optimization: A Unified Approach*, Academic Press, New York, 1970.

## A GENERALIZED LAGRANGE MULTIPLIER RULE FOR EQUALITY CONSTRAINTS IN NORMED LINEAR SPACES\*

D. O. NORRIS†

**Abstract.** It is shown that a necessary and sufficient condition for a differentiable function to have a critical point on an affine constraint set is that a Lagrange multiplier rule hold. The rule applies to problems which are not covered by the standard multiplier theorem. The results are extended to the determination of necessary conditions for a convex function to have a minimum on a nonaffine constraint set.

**1. Introduction.** Let  $X$  and  $Y$  denote  $B$ -spaces and let  $f: X \rightarrow R$  and  $g: X \rightarrow Y$  denote Frechét differentiable functions. Consider the problem of minimizing  $f$  on the set of points for which  $g(x) = 0$ . The problem was apparently first treated in this form by H. H. Goldstine [5] in 1940. He obtained the classical multiplier rule under the hypothesis that the range of  $Dg(x_0)$ , the derivative of  $g$  at  $x_0$ , be closed,  $f$  and  $g$  have continuous derivatives, and any one of four other conditions be satisfied. A proof of the Lagrange multiplier rule which requires that the range of  $Dg(x_0)$  be closed and that  $f$  and  $g$  have continuous derivatives may be found in [9, p. 380] or [8, pp. 243–244]. Specifically, we have the following theorem.

**THEOREM 1.1.** *Let  $X$  and  $Y$  denote  $B$ -spaces. Let  $f: X \rightarrow R$  and  $g: X \rightarrow Y$  each have a continuous Frechét derivative in a neighborhood of  $x_0$ . Let  $f$  have a local minimum on the set  $N(g) = \{x: g(x) = 0\}$  at  $x_0$  and suppose that the range of  $Dg(x_0)$  is closed. Then there is a nonzero  $(\eta^*, y^*) \in R \times Y^*$  such that*

$$\eta^* Df(x_0)h + y^* Dg(x_0)h = 0$$

for all  $h \in X$ .

If  $Dg(x_0)$  is surjective, then both  $\eta^*$  and  $y^*$  may be chosen to be nonzero, whereas if  $Dg(x_0)$  is not surjective, then  $\eta^* = 0^*$  and  $y^*$  nonzero (for some  $y^*$ ) is always a possible solution. The latter, however, is often unsatisfactory since a solution for which both  $\eta^*$  and  $y^*$  are nonzero may exist. In fact, such will be the case if  $Df(x_0) \in R(Dg(x_0)^*)$ , the range of the adjoint of  $Dg(x_0)$ .

In this paper the question of relaxing the requirement that the range of  $Dg(x_0)$  be closed will be considered. It will be shown that for certain classes of problems, necessary and sufficient conditions for the existence of a Lagrange multiplier can be obtained without assuming that the range of  $Dg(x_0)$  is closed. Before proceeding with the development of the theory we present an example for which Theorem 1.1 does not apply, yet a nonzero Lagrange multiplier exists.

The problem of minimizing a function subject to constraints (set constraints, inequality constraints, and equality constraints) has been treated extensively and the question arises as to how some of these other treatments relate to the work considered here. Varaiya [13] and Guignard [6] treat the problem: minimize  $f: X \rightarrow R$  on  $A \cap g^{-1}(B)$ , where  $g: X \rightarrow Y$ ,  $A \subset X$  and  $B \subset Y$ . If  $B = \{0\}$  and  $A = X$ , then their problem reduces to the one considered here. If Varaiya's weak constraint qualification is substituted for the requirement that  $f$  and  $g$  have

\* Received by the editors September 15, 1970, and in final revised form March 25, 1971.

† Department of Mathematics, Ohio University, Athens, Ohio 45701.

continuous derivatives, then his major result reduces to Theorem 1.1 for the problem under discussion. However, the weak constraint qualification may be difficult to verify, and in order to get the multiplier rule in terms of  $Df(x_0)$  and  $Dg(x_0)$  he must assume  $Dg(x_0)$  has a closed range. Similarly, Guignard's major result for the problem treated here requires that  $Dg(x_0)$  have a closed range. Russell [12], Pshenichniy [11], Dubovitskii and Milyutin [1], and Neustadt [10] have treated problems where the inequality constraints and equality constraints are more explicitly apparent. For example, Russell treats the problem of minimizing a function  $f$  subject to an arbitrary number of inequalities and equality constraints given by real-valued functions. His definition of "linearly approximable" is a regularity condition (see Definition 4.3 of this paper) and his multiplier rule is given in set theoretic terms rather than the differential form given in Theorem 1.1. For the problem treated in this paper, Russell's necessary condition for a minimum is that  $Df(x_0)$  be a member of the weak\* closure of the range of  $Dg(x_0)^*$ . Thus, unless  $Df(x_0)$  is actually a member of the range of  $Dg(x_0)^*$ , Russell's result does not apply to the question of solving the equation  $Df(x_0) + Dg(x_0)^*y^* = 0^*$ , and therefore characterization of  $x_0$  may be difficult. Neustadt has developed a very general theory of extremals which accounts for equality constraints, inequality constraints, and set constraints. For the problem discussed here (i.e., equality constraints only, together with differentiability), it is usually possible to recast Neustadt's canonical optimization problem into a problem in which the function  $g$  defining the equality constraints has its range in a finite-dimensional space. In this case, Theorem 1.1 always holds since the range of  $Dg(x_0)$  is closed, being a finite-dimensional subspace. The work of Pshenichniy and Dubovitskii and Milyutin bear a good deal of similarity. Dubovitskii and Milyutin base their treatment of optimization problems on a result they refer to as Euler's equation (Theorem 2.1). For the problem treated in this paper, in which there are no inequality constraints, their result may be stated as follows:

Let  $\Omega_0 \subset X$  denote an open convex cone with vertex at zero and let  $\Omega \subset X$  denote a closed convex cone with vertex at zero.  $\Omega \cap \Omega_0 = \emptyset$  if and only if there are  $x_0^* \in \Omega_0^*$  (the dual cone of  $\Omega$ ) and  $x^* \in \Omega^*$  such that  $x_0^* + x^* = 0$ .

For a differentiable  $f: X \rightarrow R$ ,  $\Omega_0 = \{h: Df(x_0)h < 0\}$ , where  $x_0$  is the desired minimum. If  $g: X \rightarrow Y$  is affine or satisfies a regularity condition,  $\Omega = \{h: Dg(x_0)h = 0\}$ . For this situation Euler's equation may be expressed as  $-Df(x_0) + x^* = 0$  for some  $x^* \in \Omega^*$ . For the problems treated by Dubovitskii and Milyutin it is usually possible to show that the range of  $Dg(x_0)^*$  is closed in which case  $x^*$  can be chosen in the range of  $Dg(x_0)^*$ ; i.e., Theorem 1.1 holds.

Thus, for the papers discussed above and except for [13] and [6], the multiplier rules are not given in the form of Theorem 1.1 and the question of characterizing the range of  $Dg(x_0)$  does not arise until applications are made.

**2. Example.** Represent an element  $x \in l_p$  as  $x = (\xi_1, \xi_2, \dots)$ . Define  $f: l_2 \rightarrow R$  by  $f(x) = \|x\|^2$  and let  $g: l_2 \rightarrow l_1$  such that  $g(x) = (2^{-1}\xi_2 - 2^{-2}, 3^{-1}\xi_3 - 3^{-2}, \dots)$ . Clearly,  $f$  and  $g$  have continuous derivatives. Denote the range and null set of  $g$  by  $R(g)$  and  $N(g)$ , respectively. (Henceforth, this notation will be used for all functions.) We assert the following facts:

- (i)  $x \in N(g)$  if and only if  $x = (\xi_1, \frac{1}{2}, \frac{1}{3}, \dots)$ , where  $\xi_1$  is arbitrary.
- (ii)  $Df(x) = 2x$ .

(iii)  $Dg(x)$  may be represented by

$$Dg(x) = \begin{bmatrix} 0 & \frac{1}{2} & & 0 \\ & 0 & \frac{1}{3} & \\ & & 0 & \frac{1}{4} \\ 0 & & & \ddots \end{bmatrix}.$$

(iv)  $R(Dg(x_0)) \neq l_1$  since, for example, there is no  $h \in l_2$  such that  $Dg(x_0)h = (2^{-3/2}, 3^{-3/2}, \dots) \in l_1$ .

(v)  $R(Dg(x_0))$  is not closed and  $\overline{R(Dg(x_0))} = l_1$ . For, given  $(\eta_1, \eta_2, \dots) \in l_1$  define  $\xi_{n+1} = (n + 1)\eta_n$ ,  $h_n = (0, \xi_2, \dots, \xi_{n-1}, 0, \dots)$ . Thus,  $\lim_{n \rightarrow \infty} Dg(x_0)h_n = (\eta_1, \eta_2, \dots)$ .

(vi) All the hypotheses of Theorem 1.1 are satisfied except that  $R(Dg(x_0))$  is not closed. Hence, the theorem does not apply. However, it is easy to check that with  $\eta^* = 1$  and  $y^*$  represented by  $(-2, -2, \dots) \in l_\infty$ ,  $(\eta^*, y^*)$  is the desired nonzero Lagrange multiplier. It is worth noting that had the range of  $g$  been chosen as  $l_2$ , then since  $(-2, -2, \dots) \notin l_2$  the existence of a continuous linear functional as a multiplier could not have been asserted. However, the multiplier exists as a linear functional (not continuous).

**3. Problems with affine constraints.** We now propose to present results which make it possible to assert the existence of Lagrange multipliers in cases illustrated by the example in § 2. That is, the problem of minimizing a differentiable function subject to an affine constraint will be treated.

First, however, a characterization of the range of the adjoint of a linear operator will be given. This result is needed in the sequel and is less restrictive than the well-known result that if  $R(T)$  is closed, where  $T$  is a linear operator, then  $R(T^*) = \{x^* : N(T) \subset N(x^*)\}$  (e.g., see [2, p. 487]).

**THEOREM 3.1.** *Let  $X$  and  $Y$  denote normed linear spaces and let  $T : X \rightarrow Y$  denote a continuous linear function. For each nonzero  $b^* \in X^*$ , where  $X^*$  is the topological dual of  $X$ , the equation  $T^*y^* = b^*$  has a solution in  $Y^*$  if and only if there is an  $h_0 \in X$  such that*

- (i)  $h_0 \notin \overline{N(b^*)}$ ,
- (ii)  $Th_0 \notin \overline{T(N(b^*))}$ .

*Proof.* To show the necessity, assume that  $T^*y^* = b^*$  has a solution  $y_0^* \in Y^*$ .  $b^* \neq 0^*$ , so  $N(b^*)$  has codimension equal to one and consequently there is an  $h_0 \notin N(b^*)$  such that  $X = \overline{\text{span}\{h_0, N(b^*)\}}$ . Now suppose for every  $h_0 \notin N(b^*)$  it happens that  $Th_0 \in \overline{T(N(b^*))}$ . Since  $X = \overline{\text{span}\{h_0, N(b^*)\}}$ , every  $h \in X$  may be written as  $h = \alpha h_0 + n$  for some real  $\alpha$  and some  $n \in N(b^*)$ . By hypothesis,  $b^*h - y_0^*Th = 0$  for all  $h \in X$ , so  $b^*\alpha h_0 - y_0^*T(\alpha h_0 + n) = 0$  for all real  $\alpha$  and all  $n \in N(b^*)$ . For  $\alpha = 1$ , it follows that  $b^*h_0 - y_0^*Th_0 = \frac{y_0^*Tn}{\alpha}$  for all  $n \in N(b^*)$ . Thus,  $y_0^*Tn = 0$  for all  $n \in N(b^*)$ . Now, note that if  $Th_0 \in \overline{T(N(b^*))}$ , then there is a sequence  $\{T(x_n)\} \subset T(N(b^*))$  such that  $Tx_n \rightarrow Th_0$ .  $y_0^*$  is continuous so  $y_0^*Th_0 = 0$ . Therefore,  $b^*h_0 = 0$  contrary to the assumption that  $h_0 \notin N(b^*)$ .

Conversely, if (i) and (ii) hold, then  $Th_0 \neq 0$ , and by the Hahn-Banach theorem there is a nonzero  $y_0^* \in Y^*$  such that  $y_0^*y = 0$  for all  $y \in \overline{T(N(b^*))}$  and  $y_0^*Th_0 = b^*h_0$ . Then  $b^*h = b^*(\alpha h_0 + n) = \alpha y_0^*Th_0 = y_0^*(T(\alpha h_0 + n)) = y_0^*Th$  for all  $h \in X$ , and thus  $T^*y_0^* = b^*$ .

**COROLLARY 3.1.1.** *Let  $X$  and  $Y$  denote linear spaces and let  $T: X \rightarrow Y$  denote a linear function. For each nonzero  $b' \in X'$ , where  $X'$  is the algebraic dual of  $X$ , the equation  $T'y' = b'$  has a solution in  $Y'$  if and only if there is an  $h_0 \in X$  such that*

- (i)  $h_0 \notin N(b')$ ,
- (ii)  $Th_0 \notin T(N(b'))$ .

*Proof.* The proof is almost the same as that given in Theorem 3.1.

Theorem 3.1 and its corollary will now be connected with the existence of a critical point of  $f$  on an affine subspace. It is a well-known result that if  $f$  has an extremum at  $x_0$  on an affine subspace  $A$ , then  $Df(x_0)h = 0$  for all  $h \in A - x_0$ , i.e.,  $f$  has a critical point on  $A$  at  $x_0$ . To be precise we have the following definition.

**DEFINITION 3.2.** Let  $X$  denote a normed linear space and suppose that  $A$  is an affine subspace in  $X$  (i.e., the translate of a subspace). Let  $f: X \rightarrow R$  denote a differentiable function.  $f$  has a *critical point* at  $x_0 \in A$  if and only if  $Df(x_0)h = 0$  for all  $h \in A - x_0$ .

**THEOREM 3.3.** *Let  $g: X \rightarrow Y$  denote a continuous affine transformation,  $g(x) = T(x - x_1)$ , where  $T$  is linear and continuous. Let  $f$  have a nonzero derivative at  $x_0$ .*

- (a) *If there is an  $h_0 \in X$  such that*

- (i)  $h_0 \notin N(Df(x_0))$  and
- (ii)  $Th_0 \notin T(N(Df(x_0)))$ ,

*then  $f$  has a critical point on  $N(g)$  at  $x_0$ .*

- (b) *Conversely, if  $f$  has a critical point on  $N(g)$  at  $x_0$ , then there is an  $h_0 \in X$  such that*

- (iii)  $h_0 \notin N(Df(x_0))$  and
- (iv)  $Th_0 \notin T(N(Df(x_0)))$ .

*Proof.* Part (a). If  $x_0$  is not a critical point, then there is an  $h_1 \in N(T)$  such that  $Df(x_0)h_1 \neq 0$ . Then  $Df(x_0)h_1 + y^*Th_1 \neq 0$  for all  $y^* \in Y^*$ . In view of conditions (i) and (ii), this contradicts Theorem 3.1.

Part (b). If for every  $h_0 \notin N(Df(x_0))$  it happens that  $Th_0 \in T(N(Df(x_0)))$ , then  $T(X) \subset T(N(Df(x_0)))$  and we have  $T(X) = T(N(Df(x_0)))$ . Since  $Df(x_0) \neq 0^*$ , there is an  $h_2 \notin N(Df(x_0))$  such that  $X = \text{span}\{h_2, N(Df(x_0))\}$ . But,  $T(X) = T(N(Df(x_0)))$  so that there is an  $h_3 \in N(Df(x_0))$  such that  $Th_3 = Th_2$ . Then,  $T(h_3 - h_2) = 0$  and since  $x_0$  is a critical point it follows that  $Df(x_0)(h_3 - h_2) = 0$ . But,  $Df(x_0)h_3 = 0$ , so  $Df(x_0)h_2 = 0$  contrary to the assumption.

The connection with Lagrange multipliers is now easily established.

**THEOREM 3.4.** *Let  $g: X \rightarrow Y$  denote a continuous affine transformation,  $g(x) = T(x - x_1)$ , where  $T$  is linear and continuous. Let  $f$  have a nonzero derivative at  $x_0$ .*

- (a) *If  $f$  has a critical point on  $N(g)$  at  $x_0$ , then there is a nonzero  $y' \in Y'$  such that*

$$Df(x_0)h + y'Dg(x_0)h = 0$$

*for all  $h \in X$ .*

- (b) *If there is a nonzero  $y^* \in Y^*$  such that*

$$Df(x_0)h + y^*Dg(x_0)h = 0$$

*for all  $h \in X$ , then  $f$  has a critical point on  $N(g)$  at  $x_0$ .*

*Proof. Part (a).* If  $f$  has a critical point on  $N(g)$  at  $x_0$ , then the conclusion of Theorem 3.3(b) holds. But  $Dg(x_0) = T$  so from Corollary 3.1.1 it can be concluded that there is a nonzero  $y' \in Y'$  such that  $Dg(x_0)'y' = -Df(x_0)$  which is equivalent to the desired result.

Part (b) follows by a similar argument with the aid of Theorem 3.3(a) and Theorem 3.1.

It can readily be seen that the example of § 2 is covered by Theorem 3.4. Furthermore, the conclusion of part (a) cannot be strengthened to make  $y'$  continuous without the addition of more hypotheses. This is illustrated by the example of § 2 when the range of  $g$  is taken to be  $l_2$  instead of  $l_1$  (see (vi) in § 2).

Theorem 3.4 may be applied to linear distributed parameter control problems. For example, suppose a system is described by the differential equation

$$x'(t) = Ax(t) + u(t), \quad x(0) = x_0,$$

where  $A$  is the infinitesimal generator of a strongly continuous semigroup of operators,  $T(t)$ . Suppose it is desired to find  $u$  which transfers  $x_0$  to 0 at time  $t_1$  such that  $f:U \rightarrow R$  is minimized, where  $f$  is a differentiable function. If  $u(t) \in X$ ,  $X$  a  $B$ -space, and if  $I = [0, t_1]$ , then it may be desirable to choose  $U$  as  $L_2(I, X)$ , where

$$\|u\|_U = \left[ \int_0^{t_1} \|u(t)\|_X^2 dt \right]^{1/2}.$$

Define  $g:U \rightarrow X$  such that

$$g(u) = T(t_1)x_0 + \int_0^{t_1} T(t_1 - s)u(s) ds.$$

The problem may now be reformulated as follows:

$$\text{minimize } f$$

subject to the affine equality constraint

$$g(u) = 0.$$

Theorem 3.4 clearly applies. Furthermore, in distributed parameter control problems,  $X$  is not finite-dimensional, and consequently  $g$  may not have a closed range. In such cases, the classical Lagrange multiplier theorem (Theorem 1.1) will not apply.

**4. Convex functions and nonaffine constraints.** We now propose to extend some of the results of the last section to cover the problem of minimizing a convex, differentiable function subject to a nonaffine equality constraint.

**DEFINITION 4.1.** Let  $X$  denote a normed linear space and suppose  $A \subset X$ . Let  $x_0 \in A$ . The *tangent cone*,  $TC(A, x_0)$ , for  $A$  at  $x_0$  is the set of vectors  $x \in X$  such that:

- (i) there exists a sequence  $\{y_n\} \subset A$  for which  $y_n \rightarrow x_0$ , and
- (ii) there exists a sequence  $\{\lambda_n : \lambda_n \geq 0\}$  for which  $\lambda_n(y_n - x_0) \rightarrow x$ .

Varaiya [13] gives an equivalent definition and calls  $TC(A, x_0)$  the local closed cone of  $A$  at  $x_0$ . Other authors (e.g., see [3], [7]) define the unit vectors of

the tangent cone and require, accordingly, that  $\lambda_n = \|y_n - x_0\|^{-1}$ . The tangent cone is the cone generated by these unit vectors. It is easy to see that these definitions are equivalent.

In the next theorem it is established that when  $f$  is a continuous, convex function and has a minimum on  $A \subset X$  at  $x_0$ , then  $f$  has a minimum on  $x_0 + \text{TC}(A, x_0)$  at  $x_0$ . A similar result where  $A = N(g)$  has been proved by Flett [4, Theorem 1]. If, in addition,  $x_0$  satisfies a regularity condition such as  $N(Dg(x_0)) = \text{TC}(N(g), x_0)$ , then the problem of minimizing the convex function  $f$  subject to a differentiable equality constraint can be reduced to the problem of minimizing  $f$  subject to an affine constraint. When  $f$  is differentiable the results of the last section apply.

**THEOREM 4.2.** *Let  $X$  denote a normed linear space. Let  $A \subset X$  and suppose that  $f : X \rightarrow \mathbb{R}$  is a convex, continuous function. If  $f$  has a minimum on  $A$  at  $x_0$ , then  $f$  has a minimum on  $x_0 + \text{TC}(A, x_0)$  at  $x_0$ .*

*Proof.* If  $\text{TC}(A, x_0) = \{0\}$ , the result is trivially true. Thus, assume that  $\text{TC}(A, x_0) \neq \{0\}$ . Suppose there is an element  $x_1 \in x_0 + \text{TC}(A, x_0)$  such that  $f(x_1) < f(x_0)$ . Then there is a sequence  $\{y_n\} \subset A$  and a sequence  $\{\lambda_n : \lambda_n \geq 0\}$  such that  $y_n \rightarrow x_0$  and  $\lambda_n(y_n - x_0) \rightarrow x_1 - x_0$ .  $f$  is continuous, so for  $n$  sufficiently large

$$(1) \quad f(\lambda_n(y_n - x_0) + x_0) < f(x_0).$$

Now  $\lambda_n \rightarrow \infty$  as  $n \rightarrow \infty$  so for  $\lambda_n > 1$  note that

$$y_n = (1 - \lambda_n^{-1})x_0 + \lambda_n^{-1}(x_0 + \lambda_n(y_n - x_0)).$$

Use (1) and the convexity of  $f$  to conclude that

$$\begin{aligned} f(y_n) &\leq (1 - \lambda_n^{-1})f(x_0) + \lambda_n^{-1}f(x_0 + \lambda_n(y_n - x_0)) \\ &< (1 - \lambda_n^{-1})f(x_0) + \lambda_n^{-1}f(x_0) = f(x_0). \end{aligned}$$

This contradicts the assumption that  $f$  achieves its minimum on  $A$  to  $x_0$ .

**COROLLARY 4.2.1.** *Under the hypothesis of Theorem 4.2, let  $Y$  denote a vector space,  $g : X \rightarrow Y$ , and  $A = N(g)$ . Then  $f$  has a minimum on  $x_0 + \text{TC}(N(g), x_0)$ .*

**DEFINITION 4.3.** Let  $X$  and  $Y$  denote normed linear spaces and let  $g : X \rightarrow Y$  denote a differentiable function.  $x_0 \in N(g)$  is a *regular point* if and only if  $\text{TC}(N(g), x_0) = N(Dg(x_0))$ .

Some authors (e.g., see [5], [8], [9]) define regularity by the requirement that  $Dg(x_0)$  be surjective. However, Flett [3, Theorem 3] has shown that if  $g$  has a continuous derivative and  $Dg(x_0)$  is surjective, then  $\text{TC}(N(g), x_0) = N(Dg(x_0))$ . Definition 4.3 has been given by Hestenes [7, p. 29] in the case that  $g$  has a finite-dimensional range.

**THEOREM 4.4.** *Let  $X$  and  $Y$  denote normed linear spaces. Suppose  $f : X \rightarrow \mathbb{R}$  is a convex differentiable function and  $g : X \rightarrow Y$  is differentiable. If  $f$  has a minimum on  $N(g)$  at  $x_0$  and  $x_0$  is a regular point, then there is a nonzero  $y' \in Y'$  such that*

$$Df(x_0)h + y'Dg(x_0)h = 0$$

for all  $h \in X$ .

*Proof.* By Corollary 4.2.1 it follows that  $f$  has a minimum on  $x_0 + \text{TC}(N(g), x_0)$  at  $x_0$ , and from regularity it follows that  $f$  has a minimum on  $x_0 + N(Dg(x_0))$  at



$x_0$ , i.e., an affine subspace. But, then  $f$  has a critical point on  $x_0 + N(Dg(x_0))$  at  $x_0$ , and by Theorem 3.4(a) the desired conclusion follows.

## REFERENCES

- [1] A. YA. DUBOVITSKII AND A. A. MILYUTIN, *Extremum problems in the presence of restrictions*, U.S.S.R. Comput. Math. and Math. Phys., 5 (1965), no. 3, pp. 1–80.
- [2] N. DUNFORD AND J. T. SCHWARTZ, *Linear Operators, Part I*, Interscience, New York, 1964.
- [3] T. M. FLETT, *On differentiation in normed vector spaces*, J. London Math. Soc., 42 (1967), pp. 523–533.
- [4] ———, *Points of minimum or maximum norm on smooth surfaces in Banach spaces*, Ibid., 44 (1969), pp. 583–586.
- [5] H. H. GOLDSTINE, *Minimum problems in the functional calculus*, Bull. Amer. Math. Soc., 46 (1940), pp. 142–149.
- [6] M. GUIGNARD, *Generalized Kuhn–Tucker conditions for mathematical programming problems in a Banach space*, this Journal, 7 (1969), pp. 232–241.
- [7] M. R. HESTENES, *Calculus of Variations and Optimal Control*, John Wiley, New York, 1966.
- [8] D. G. LUENBERGER, *Optimization by Vector Space Methods*, John Wiley, New York, 1969.
- [9] L. A. LUSTERNIK AND V. J. SOBOLEV, *Elements of Functional Analysis*, Frederick Ungar, New York, 1961.
- [10] L. W. NEUSTADT, *A general theory of extremals*, J. Comput. System Sci., 3 (1969), pp. 57–92.
- [11] B. N. PSHENICHNIY, *Convex programming in a normed space*, Kibernetika, 5 (1965), pp. 46–54.
- [12] D. L. RUSSELL, *The Kuhn–Tucker conditions in Banach space with an application to control theory*, J. Math. Anal. Appl., 15 (1966), pp. 200–212.
- [13] P. P. VARAIYA, *Nonlinear programming in Banach space*, SIAM J. Appl. Math., 15 (1967), pp. 284–293.

## JUSTIFICATION OF THE DESCRIBING FUNCTION METHOD\*

A. R. BERGEN† AND R. L. FRANKS‡

**Abstract.** Explicit conditions are given under which the use of the describing function method to investigate the nature of oscillations in autonomous nonlinear feedback systems is justified. When these conditions are satisfied, bounds are given for the frequency, fundamental magnitude, and higher harmonics of the oscillation based on the describing function approximation.

The feedback systems considered are those which can be decomposed into a linear time-invariant subsystem, not necessarily causal, stable, or finite-dimensional, and a nonlinear frequency independent subsystem, possibly containing hysteresis.

The approach taken is to consider the describing function equation as an approximation to a determining equation for periodic solutions of the autonomous system's operator equation, and to use local degree theory to guarantee the existence of a solution to the determining equation.

**1. Introduction.** The describing function method was introduced in this country by R. J. Kochenburger [6] in 1950. It is used by engineers to investigate the nature of oscillations in nonlinear feedback systems. Frequently it is used as an aid in design to avoid or attenuate such periodic responses. This paper is concerned with justifying the use of the describing function method to investigate periodic responses in autonomous nonlinear feedback systems, and gives error bounds for those responses.

Many authors have discussed the use of the describing function. A good discussion is given in Gelb and Vander Velde [4]. The general approach is as follows:

Consider the autonomous feedback system in Fig. 1.  $\mathcal{N}$  is usually a frequency independent nonlinear subsystem whose characteristic has odd symmetry and  $\mathcal{G}$  is a time-invariant linear subsystem. The input to  $\mathcal{N}$  is taken to be  $x = \text{Re}(re^{j\omega t})$ . The steady state output of  $\mathcal{N}$  and the input to  $\mathcal{G}$  is then

$$y = \text{Re} \left( \sum_{n=1, \text{ odd}}^{\infty} A_n e^{jn\omega t} \right).$$

The system connection requires  $x = -z$ ; i.e.,  $r = -G(j\omega)A_1$  and  $G(jn\omega)A_n = 0$  for  $n$  greater than or equal to 3. The coefficient  $A_1$  ordinarily depends only on  $r$  so it is usually written as  $A_1 = rN(r)$ .  $N(r)$  is called the describing function for  $\mathcal{N}$  and is its steady state first harmonic gain. The first harmonic equation can then be written as

$$(1) \quad 1 + G(j\omega)N(r) = 0$$

and is called the describing function equation. This equation is usually solved graphically by plotting the Nyquist locus  $G(j\omega): \omega \in \mathbb{R}^+$  and the critical locus,

---

\* Received by the editors June 11, 1970, and in revised form January 4, 1971.

† Department of Electrical Engineering and Computer Sciences, Electronic Research Laboratory, University of California, Berkeley, California 94720. This author's work was supported by the National Science Foundation under Grant GK-10656X.

‡ Department of Electrical Engineering and Computer Sciences, Electronics Research Laboratory, University of California, Berkeley, California. Now at Bell Telephone Laboratories, Whippany, New Jersey 07981. This author's work was supported by a NASA Traineeship.

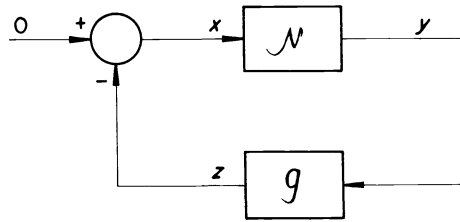


FIG. 1. System (S)

$-1/N(r): r \in \mathbb{R}^+$ . An intersection of these loci gives values of  $\omega$  and  $r$  which satisfy the describing function equation. The basic appeal of the method is the simplicity of this graphical solution.

The condition that  $G(jn\omega)A_n = 0$  for  $n$  greater than or equal to 3 can never be satisfied, but it is ignored. If there is a pair  $(\omega, r)$  which satisfies the describing function equation, then it is felt that the system has a periodic solution near  $x(t) = r \sin \omega t$ ; the heuristic justification is that if  $\mathcal{N}$  is not too nonlinear and  $\mathcal{G}$  is a good low-pass filter, then the higher harmonics  $G(jn\omega)A_n$  are approximately zero and can be ignored. Notice the importance of the assumption of an odd nonlinearity in this connection.

Clearly this approach requires mathematical justification, and error bounds for the describing function approximation would be desirable.<sup>1</sup> Bass [1] considered the problem of justifying the describing function method and gave conditions under which the method was justified. Unfortunately, as he points out, one cannot tell whether or not the conditions are satisfied except in very simple cases.

**2. Outline of method.** The system of interest is shown in Fig. 1, where  $\mathcal{N}$  is a frequency-independent nonlinear subsystem and  $\mathcal{G}$  is a linear time-invariant subsystem. The problem is to show that under certain conditions related to the describing function method, there is a periodic function  $\hat{x}: \mathbb{R} \rightarrow \mathbb{R}$  such that

$$(2) \quad \hat{x} = -\mathcal{G}\mathcal{N}\hat{x}.$$

The approach, motivated by Cesari [3], is first to time scale the problem and then to separate it into two, more tractable problems.

Using time scaling, it is sufficient to consider functions of period  $2\pi$ . If  $\hat{x}$  is a periodic solution of period  $T$ , define the time-scaled function  $x$  by  $x(\omega t) = \hat{x}(t)$ , where  $\omega T = 2\pi$ . Then  $x$  has period  $2\pi$ . The operator  $\mathcal{G}$  is not invariant under time scaling, that is,  $[\mathcal{G}\hat{x}](t) \neq [\mathcal{G}x](\omega t)$ . This means that time scaling requires the introduction of new operators,  $\mathcal{G}_\omega: \omega \in [\omega_*, \omega^*]$ , such that  $[\mathcal{G}\hat{x}](t) = [\mathcal{G}_\omega x](\omega t)$ .

In terms of time-scaled functions and operators, the problem is to find a function  $x$  of period  $2\pi$  and an  $\omega > 0$  such that

$$(3) \quad x = -\mathcal{G}_\omega \mathcal{N}x.$$

Note that if  $\mathcal{N}$  were not frequency independent, an appropriate operator  $\mathcal{N}_\omega$  would be required.

<sup>1</sup> With the paper virtually complete the authors learned of a similar effort by Kudrewicz [7] reported in the 1969 IFAC Conference in Warsaw. The system considered by us includes hysteresis nonlinearities, and the criterion for justification is easier to apply and somewhat less restrictive.

To separate the problem into two, more tractable problems, let  $\mathcal{P}x$  be the projection of a  $2\pi$  periodic function  $x$  onto its first  $M$  harmonics and  $\mathcal{P}^*x$  be the projection of  $x$  onto its remaining harmonics, so that  $x = \mathcal{P}x + \mathcal{P}^*x$ . Then  $x = -\mathcal{G}_\omega \mathcal{N}x$  if and only if both

$$(4) \quad \mathcal{P}^*x = -\mathcal{P}^*\mathcal{G}_\omega \mathcal{N}x,$$

$$(5) \quad \mathcal{P}x = -\mathcal{P}\mathcal{G}_\omega \mathcal{N}x.$$

Under certain conditions, for any  $\omega \in [\omega_*, \omega^*]$  and any set  $\mathbf{r}$  of Fourier coefficients for the first  $M$  harmonics of  $x$ , there is a function  $\tilde{x}(\mathbf{r})$  with those first  $M$  harmonics which satisfies (4). Substituting  $\tilde{x}(\mathbf{r})$  into (5) gives

$$(6) \quad \mathcal{P}\tilde{x}(\mathbf{r}) = -\mathcal{P}\mathcal{G}_\omega \mathcal{N}\tilde{x}(\mathbf{r}).$$

Equation (6) is a determining equation for (3) in the sense that if there are an  $\omega$  and a set  $\mathbf{r}$  of Fourier coefficients such that (6) is satisfied, then  $x = \tilde{x}(\mathbf{r})$  is a solution to (3) for that  $\omega$ . Defining  $\hat{x}(t) = x(\omega t)$  implies  $\hat{x}$  is the required periodic response of the system.

If  $M = 1$ , equation (6) corresponds to an operator  $V: \mathbb{R}^2 \rightarrow \mathbb{R}^2$ . The describing function equation corresponds to an operator  $U: \mathbb{R}^2 \rightarrow \mathbb{R}^2$ .  $U$  is an approximation of  $V$  and has local degree one on a particular set. Under certain conditions  $V$  is homotopically equivalent to  $U$  on that set, so  $V$  has degree one on that set. Therefore (6) has a solution for some  $(\omega, r)$  in that set. This means that (3) has a periodic solution.

This approach is carried out in the next three sections. It is recommended that the reader who is unfamiliar with local degree refer to Appendix II at this time. Theorems II.2 and II.3 in that appendix are fundamental to the motivation of this paper.

**3. Basic result.** We are interested in the solution of (2) in a space of periodic functions (more precisely periodic functions with only odd harmonics). This will be done by solving (3) in the space

$$\mathcal{L}_\pi^2 = \{x: \mathbb{R} \rightarrow \mathbb{R} | x \text{ has period } 2\pi \text{ and only odd harmonics}\},$$

$\mathcal{L}_\pi^2$  is a Hilbert space with the inner product

$$\langle x, y \rangle = \frac{1}{\pi} \int_0^{2\pi} x(\theta)y(\theta) d\theta.$$

The corresponding norm will be referred to as  $\|x\|$ . The elements of  $\mathcal{L}_\pi^2$  are half-wave symmetric periodic functions. We shall refer to them as  $\pi$ -symmetric functions.

In order to avoid proving some later theorems twice we shall introduce new operators  $\mathcal{A}_\omega$  and  $\mathcal{F}$ , and new equations

$$(7) \quad x = -\mathcal{A}_\omega \mathcal{F}x,$$

$$(8) \quad \mathcal{P}^*x = -\mathcal{P}^*\mathcal{A}_\omega \mathcal{F}x,$$

$$(9) \quad \mathcal{P}x = -\mathcal{P}\mathcal{A}_\omega \mathcal{F}x.$$

In § 5,  $\mathcal{A}_\omega$  will be related to  $\mathcal{G}_\omega$ , and  $\mathcal{F}$  will be related to  $\mathcal{N}$ .

The basic result of this section is Theorem 1 which gives sufficient conditions for the existence of a solution to (8) and (9) in the space  $\mathcal{L}_\pi^2$ .

Let

$$\xi_k(\theta) = \begin{cases} \sin k\theta, & k \text{ odd,} \\ \cos (k - 1)\theta, & k \text{ even;} \end{cases}$$

then  $\{\xi_k\}$  is a complete orthonormal basis for  $\mathcal{L}_\pi^2$ . For a fixed positive integer  $\kappa$ , define:

1.  $\mathcal{P} : \mathcal{L}_\pi^2 \rightarrow \mathcal{L}_\pi^2$  such that  $\mathcal{P}x \triangleq \sum_{k=1}^\kappa \langle x, \xi_k \rangle \xi_k$ ,
2.  $\mathcal{P}^* : \mathcal{L}_\pi^2 \rightarrow \mathcal{L}_\pi^2$  such that  $\mathcal{P}^*x = x - \mathcal{P}x$ ,
3.  $y : \mathbb{R}^{\kappa-1} \rightarrow \mathcal{P}\mathcal{L}_\pi^2$  such that  $y(\mathbf{r}) \triangleq \sum_{k=1}^{\kappa-1} r_k \xi_k$ , where  
 $\mathbf{r} = (r_1, r_2, \dots, r_{\kappa-1}) \in \mathbb{R}^{\kappa-1}$ ,
4. the usual Lipschitz norm

$$\|\mathcal{F}\| \triangleq \sup \left\{ \frac{\|\mathcal{F}x - \mathcal{F}y\|}{\|x - y\|} : x, y \in \mathcal{L}_\pi^2, x \neq y \right\}.$$

For the remainder of this section the following hypothesis is required.

**HYPOTHESIS A.**

1. There is an  $\omega_*$  and  $\omega^*$  such that for each  $\omega \in [\omega_*, \omega^*]$ ,  $\mathcal{A}_\omega : \mathcal{L}_\pi^2 \rightarrow \mathcal{L}_\pi^2$  is a bounded linear operator such that  $\|\mathcal{A}_\omega - \mathcal{A}_v\| \rightarrow 0$  as  $\omega \rightarrow v$  for all  $v \in [\omega_*, \omega^*]$ .
2.  $\mathcal{F} : \mathcal{L}_\pi^2 \rightarrow \mathcal{L}_\pi^2$  is a bounded continuous operator with  $\|\mathcal{F}\| < \infty$ .
3. There exists a function  $\rho : [\omega_*, \omega^*] \rightarrow \mathbb{R}$  such that  $\|\mathcal{F}\| \|\mathcal{P}^* \mathcal{A}_\omega\| \leq \rho_\omega$  for all  $\omega \in [\omega_*, \omega^*]$  and  $\sup \{\rho_\omega : \omega \in [\omega_*, \omega^*]\} < 1$ .

Theorem 1 can most conveniently be stated and proved by first proving two lemmas. Lemma 1 gives conditions under which an equation similar to (6) has a solution, and gives important properties of that solution.

**LEMMA 1.** *Given Hypothesis A, for each  $\omega \in [\omega_*, \omega^*]$  and each  $\mathbf{r} \in \mathbb{R}^{\kappa-1}$  the equation  $x = -\mathcal{P}^* \mathcal{A}_\omega \mathcal{F} [y(\mathbf{r}) + x]$  has a unique solution  $x^*(\omega, \mathbf{r})$ . Define  $x(\omega, \mathbf{r}) = y(\mathbf{r}) + x^*(\omega, \mathbf{r})$ . Then:*

- a.  $x : [\omega_*, \omega^*] \times \mathbb{R}^{\kappa-1} \rightarrow \mathcal{L}_\pi^2$  is continuous,
- b.  $\mathcal{P}x(\omega, \mathbf{r}) = y(\mathbf{r})$ ,
- c.  $\mathcal{P}^*x(\omega, \mathbf{r}) = -\mathcal{P}^* \mathcal{A}_\omega \mathcal{F} x(\omega, \mathbf{r})$ ,
- d.  $\|\mathcal{P}^*x(\omega, \mathbf{r})\| \leq \frac{\|\mathcal{P}^* \mathcal{A}_\omega \mathcal{F} y(\mathbf{r})\|}{1 - \rho_\omega}$ .

The proofs of all lemmas and theorems are given in Appendix I.

Before giving conditions under which (6) has a solution it is necessary to define functions  $U : [\omega_*, \omega^*] \times \mathbb{R}^{\kappa-1} \rightarrow \mathbb{R}^\kappa$  and  $V : [\omega_*, \omega^*] \times \mathbb{R}^{\kappa-1} \rightarrow \mathbb{R}^\kappa$  whose

components  $k = 1, 2, \dots, \kappa$  satisfy

$$(10) \quad U_k(\omega, \mathbf{r}) = r_k + \langle \mathcal{A}_\omega \mathcal{F} y(\mathbf{r}), \xi_k \rangle,$$

$$(11) \quad V_k(\omega, \mathbf{r}) = r_k + \langle \mathcal{A}_\omega \mathcal{F} x(\omega, \mathbf{r}), \xi_k \rangle,$$

where  $r_k$  is the  $k$ th component of  $\mathbf{r} \in \mathbb{R}^{\kappa-1}$  for  $k \leq \kappa - 1$  and  $r_\kappa = 0$ .

In Theorem 2, the condition  $U(\omega, \mathbf{r}) = 0$  will be identified with the describing function problem.  $U$  is an approximation for  $V$  in the sense of the following lemma.

LEMMA 2. *Given Hypothesis A,  $U$  and  $V$  are continuous and*

$$\|U(\omega, \mathbf{r}) - V(\omega, \mathbf{r})\|_2 \leq \frac{\|\mathcal{P} \mathcal{A}_\omega\| \cdot \|\mathcal{F}\| \cdot \|\mathcal{P}^* \mathcal{A}_\omega \mathcal{F} y(\mathbf{r})\|}{1 - \rho_\omega}.$$

( $\|\cdot\|_2$  is the Euclidean norm in  $\mathbb{R}^\kappa$ .)

The main result of this section can now be given in terms of  $U, V$ , and  $d[U, 0, \Omega]$ , the local degree of  $U$  with respect to 0 and the set  $\Omega$ . For the reader's convenience, Appendix II contains a definition of local degree and two related theorems.

THEOREM 1. *Given Hypothesis A, if there exists an open bounded set  $\Omega \subset \mathbb{R}^\kappa$  such that:*

1.  $\Omega \subset (\omega_*, \omega^*) \times \mathbb{R}^{\kappa-1}$ ,
2.  $d[U, 0, \Omega] \neq 0$ ,
3.  $|U(\omega, \mathbf{r}) - V(\omega, \mathbf{r})|_2 < |U(\omega, \mathbf{r})|_2$  for each  $(\omega, \mathbf{r}) \in \partial\Omega$ , then there exists a point  $(\omega, \mathbf{r}) \in \Omega$  such that:
  - a.  $\mathcal{P}x(\omega, \mathbf{r}) = -\mathcal{P} \mathcal{A}_\omega \mathcal{F} x(\omega, \mathbf{r})$ ,
  - b.  $\mathcal{P}^*x(\omega, \mathbf{r}) = -\mathcal{P}^* \mathcal{A}_\omega \mathcal{F} x(\omega, \mathbf{r})$ ,
  - c.  $\|\mathcal{P}^*x(\omega, \mathbf{r})\| \leq \frac{\|\mathcal{P}^* \mathcal{A}_\omega \mathcal{F} y(\mathbf{r})\|}{1 - \rho_\omega}$ .

**4. Application to autonomous systems.** In order to use Theorem 1 to investigate the periodic oscillations of the autonomous feedback system in Fig. 1, it is necessary to make some assumptions about  $\mathcal{N}$  and  $\mathcal{G}$ .

HYPOTHESIS B1. The system (S) shown in Fig. 1 satisfies:

1.  $\mathcal{N}: \mathcal{L}_\pi^2 \rightarrow \mathcal{L}_\pi^2$ .
2. If  $x \in \mathcal{L}_\pi^2$  and  $\tilde{x}(t) = x(\omega t)$  for all  $t \in \mathbb{R}$  and some  $\omega > 0$ , then  $(\mathcal{N}x)(\omega t) = (\mathcal{N}\tilde{x})(t)$  for all  $t \in \mathbb{R}$ .
3.  $\mathcal{G}$  is a linear time-invariant system with frequency response function  $G(j\omega)$ .
4. For any odd integer  $k$ ,  $G(j\omega k)$  is continuous at each  $\omega \in [\omega_*, \omega^*]$ .
5.  $|G(j\omega k)| \rightarrow 0$  as  $k \rightarrow \text{odd } \infty$ , uniformly in  $\omega \in [\omega_*, \omega^*]$ .

Hypothesis B1.1 means  $\pi$ -symmetric inputs to  $\mathcal{N}$  result in  $\pi$ -symmetric outputs. Hypothesis B1.2 means  $\mathcal{N}$  is frequency independent for half-wave symmetric periodic inputs. Hypotheses B1.3, B1.4 and B1.5 imply  $\mathcal{G}$  is a low-pass filter, at least for half-wave symmetric periodic inputs which have their fundamental frequencies in a given band.

Rather than consider the operation of  $\mathcal{G}$  on periodic functions of all periods, we shall restrict our attention to an equivalent operation on functions of period  $2\pi$ . The response to  $\mathcal{G}$  is frequency-dependent so it is necessary to consider a continuum of operators,  $\{\mathcal{G}_\omega: \omega \in [\omega_*, \omega^*]\}$ , to obtain this equivalent operation.

DEFINITION (the operator  $\mathcal{G}_\omega$ ). For each  $\omega \in [\omega_*, \omega^*]$ ,  $\mathcal{G}_\omega: \mathcal{L}_\pi^2 \rightarrow \mathcal{L}_\pi^2$  is such that if the Fourier series of  $z$  is  $\sum_{n=-\infty}^\infty c_n e^{jn\theta}$ , then the Fourier series of  $\mathcal{G}_\omega z$  is

$$\sum_{n=-\infty}^\infty c_n G(jn\omega) e^{jn\theta}.$$

The equivalence of the operation of the linear time-invariant system  $\mathcal{G}$  and the operators  $\mathcal{G}_\omega$  is given by the following lemma.

LEMMA 3. If  $x \in \mathcal{L}_\pi^2$  is such that  $x = -\mathcal{G}_\omega \mathcal{N}x$ , then  $\tilde{x} = -\mathcal{G} \mathcal{N} \tilde{x}$ , where  $x(t) \triangleq \tilde{x}(\omega t)$ .

Properties of the operators  $\mathcal{G}_\omega$  are given in the next lemma.

LEMMA 4. Given the system (S) and an even positive integer  $\kappa$ :

- a.  $\|\mathcal{G}_\omega\| = \max \{|G(j\omega k)|: k \text{ an odd integer}\} < \infty$ ,
- b.  $\|\mathcal{P}^* \mathcal{G}_\omega\| = \max \{|G(j\omega k)|: k \text{ an odd integer}, k < \kappa\} < \infty$ ,
- c.  $\|\mathcal{P}^* \mathcal{G}_\omega\| = \max \{|G(j\omega k)|: k \text{ an odd integer}, k > \kappa\} < \infty$ ,
- d.  $\|\mathcal{P}^*(\mathcal{G}_\omega - \mathcal{G}_v)\| \leq \|\mathcal{G}_\omega - \mathcal{G}_v\| \rightarrow 0$  as  $\omega \rightarrow v$  for all  $v \in [\omega_*, \omega^*]$ .

In the most straightforward application of Theorem 1,  $\mathcal{A}_\omega = \mathcal{G}_\omega$  and  $\mathcal{F} = \mathcal{N}$ . In this case the components of the function  $U: [\omega_*, \omega^*] \times \mathbb{R}^{\kappa-1} \rightarrow \mathbb{R}^\kappa$  are given by the following lemma.

LEMMA 5. Given Hypothesis B1, if  $\mathcal{A}_\omega = \mathcal{G}_\omega$ ,  $\mathcal{F} = \mathcal{N}$ , and  $U$  is defined as in (10), then

$$U_k(\omega, \mathbf{r}) = \begin{cases} \text{Re} [r_k + G(j\omega k)\eta_k], & k \text{ odd,} \\ r_k + \text{Im} G[j\omega(k-1)]\eta_{k-1}, & k \text{ even,} \end{cases}$$

where

$$\eta_k \triangleq \frac{j}{\pi} \int_0^{2\pi} \mathcal{N}[y(\mathbf{r})](\theta) e^{-jk\theta} d\theta.$$

If  $\|\mathcal{N}\| < \infty$ , then parts 1 and 2 of Hypothesis A are satisfied with  $\mathcal{A}_\omega \triangleq \mathcal{G}_\omega$  and  $\mathcal{F} \triangleq \mathcal{N}$ . Defining  $\rho_\omega = \|\mathcal{N}\| \max \{|G(j\omega k)|: k \text{ an odd integer}, k > \kappa\}$  for some even positive integer  $\kappa$  gives  $\rho_\omega \geq \|\mathcal{N}\| \|\mathcal{P}^* \mathcal{G}_\omega\|$  for all  $\omega \in [\omega_*, \omega^*]$ . Hypothesis A will be completely satisfied if  $\sup \{\rho_\omega: \omega \in [\omega_*, \omega^*]\} < 1$ , and this can always be accomplished by choosing  $\kappa$  large enough since  $|G(j\omega k)| \rightarrow 0$  as  $k \rightarrow \infty$ , uniformly in  $\omega \in [\omega_*, \omega^*]$ .

Therefore Theorem 1 can be applied to any system (S) which satisfies the relatively weak requirements of Hypothesis B1, if  $\|\mathcal{N}\| < \infty$ , and if a set  $\Omega$  can be found satisfying the theorem's hypothesis.

Searching for an appropriate set  $\Omega$  is greatly simplified if it is possible to take  $\kappa = 2$ . In this case all computations are in  $\mathbb{R}^2$  and can be performed conveniently in the complex plane. In fact, the conditions are closely related to the loci used in the describing function method. The case  $\kappa = 2$  will be considered exclusively henceforth.

### 5. Justification and error bounds for the describing function method.

Hypothesis B2 and Lemma 6 give conditions such that  $d[\omega, 0, \Omega] \neq 0$  for a function  $W$  which will later be identified with the function  $U$  in Theorem 1. The function  $W$  is introduced to allow application of Theorem 1 to systems which do not satisfy  $\|\mathcal{N}\| < \infty$ . With  $\kappa = 2$ ,  $\mathbf{r}$  is a scalar and will henceforth be written as  $r$ .

**HYPOTHESIS B2.**

1.  $N$  has a continuous derivative on  $[r_*, r^*]$ , where  $0 < r_* < r^*$  and

$$N(r) = \frac{j}{\pi r} \int_0^{2\pi} \mathcal{N}(r\xi_1)(\theta)e^{-j\theta} d\theta.$$

2.  $G$  has a continuous derivative on  $[\omega_*, \omega^*]$ , where  $0 < \omega_* < \omega^*$  and  $G$  is abbreviated notation for  $G[j(\cdot)]$ .

3.  $1 + N(r)G(j\omega) = 0$  with  $(\omega, r) \in [\omega_*, \omega^*] \times [r_*, r^*]$  has a unique solution  $(\omega_0, r_0)$  and

- a.  $(\omega_0, r_0) \in (\omega_*, \omega^*) \times (r_*, r^*)$ ,
- b.  $dN/dr|_{r_0} \neq 0$ ,
- c.  $dG/d\omega|_{\omega_0} \neq 0$ ,
- d. the loci of  $-1/N$  and  $G$  are not tangent at  $(\omega_0, r_0)$ .

Usually inspection of  $N, G$  and the loci of  $-1/N$  and  $G$  is enough to check this hypothesis.

LEMMA 6. Given Hypotheses B1 and B2, define  $W: [\omega_*, \omega^*] \times [r_*, r^*] \rightarrow \mathbb{R}^2$  by

$$W(\omega, r) \triangleq (\text{Re}[r + rN(r)G(j\omega)], \text{Im}[r + rN(r)G(j\omega)]);$$

then  $d[W, 0, \Omega] \neq 0$  for any open bounded set  $\Omega$  in  $\mathbb{R}^2$  such that  $(\omega_0, r_0) \in \Omega \subset (\omega_*, \omega^*) \times (r_*, r^*)$ .

Hypotheses B3 and B3' are given below in order to facilitate the statement of Theorem 2, the main result of this paper. Hypothesis B3 covers the case where  $\|\mathcal{N}\| < \infty$ .

**HYPOTHESIS B3.**

- 1. There is a real number  $M$  such that  $\|\mathcal{N}\| \leq M$ .
- 2.  $\rho_\omega = M \max\{|G(j\omega k)| : k = 3, 5, \dots\} < 1$  for all  $\omega \in [\omega_*, \omega^*]$ .

A large class of nonlinear systems which exhibit hysteresis do not satisfy  $\|\mathcal{N}\| < \infty$ . This case is dealt with in Hypothesis B3'. To state it we need another definition.

DEFINITION (the space  $\mathcal{L}_\pi^\infty$ ).

$$\mathcal{L}_\pi^\infty = \{x: \mathbb{R} \rightarrow \mathbb{R} | x \text{ is continuous, } x(\theta) = -x(\theta + \pi)\},$$

where  $\|x\|_\infty = \max_{\theta \in \mathbb{R}} |x(\theta)|$ .

Notice that any function in  $\mathcal{L}_\pi^\infty$  is also in  $\mathcal{L}_\pi^2$ .

**HYPOTHESIS B3'.**

- 1. There is a real number  $M'$  such that  $\|\mathcal{N}x - \mathcal{N}y\| \leq M'\|x - y\|_\infty$  for all  $x, y \in \mathcal{L}_\pi^\infty$ .
- 2.  $k|G(j\omega k)| \rightarrow 0$  as  $k \xrightarrow{\text{odd}} \infty$ , uniformly in  $\omega \in [\omega_*, \omega^*]$ .
- 3.  $\rho'_\omega = (M'\pi/\sqrt{8}) \max\{k|G(j\omega k)| : k = 3, 5, \dots\} < 1$  for all  $\omega \in [\omega_*, \omega^*]$ .

The main result of this paper can now be stated.

THEOREM 2. Given Hypotheses B1, B2, and either B3 or B3', if there exists an open bounded set  $\Omega \subset \mathbb{R}^2$  such that:

- 1.  $(\omega_0, r_0) \in \Omega \subset (\omega_*, \omega^*) \times (r_*, r^*)$ ,
- 2. for all  $(\omega, r) \in \partial\Omega$ ,

$$B(\omega)T(r) < \left| \frac{1}{N(r)} + G(j\omega) \right|,$$



where

$$T(r) \triangleq \left( \frac{\|\mathcal{N}(r\xi_1)\|^2}{|rN(r)|^2} - 1 \right)^{1/2}$$

and

$$B(\omega) = \begin{cases} \frac{|G(j\omega)|\rho_\omega}{1 - \rho_\omega} & \text{if B3 is satisfied,} \\ \frac{|G(j\omega)|\rho'_\omega}{1 - \rho'_\omega} & \text{if B3' is satisfied;} \end{cases}$$

then for some  $(\omega, r) \in \Omega$  there exists an  $x \in \mathcal{L}^2_\pi$  such that :

- a.  $x = -\mathcal{G}_\omega \mathcal{N} x$ ,
- b.  $(\mathcal{P}x)(\theta) = r\xi_1(\theta) = r \sin \theta$ ,
- c.  $\|\mathcal{P}^* x\| \leq |rN(r)|T(r)S(\omega)$ ,

where

$$S(\omega) = \begin{cases} \frac{\rho_\omega}{(1 - \rho_\omega)M} & \text{if B3 is satisfied,} \\ \frac{\rho'_\omega}{(1 - \rho'_\omega)M'} \frac{\sqrt{8}}{\pi} & \text{if B3' is satisfied;} \end{cases}$$

i.e., there exists a periodic function  $\tilde{x}$  with fundamental harmonic  $h(t) = r \sin \omega t$  such that  $\tilde{x} = -\mathcal{G} \mathcal{N} \tilde{x}$  for some initial state of  $\mathcal{G}$  and

$$\frac{\omega}{2\pi} \int_0^{\omega/(2\pi)} |\tilde{x}(t) - h(t)|^2 dt \leq |rN(r)|T(r)S(\omega).$$

Theorem 2 gives not only sufficient conditions for the existence of a periodic response of system (S) but gives error bounds for that response in terms of error bounds on both its fundamental frequency and amplitude, and also on the  $\mathcal{L}^2$ -norm of its higher harmonics.

The heuristic justification of the describing function method stated in § 1 may now be rendered more precise. The statement “ $\mathcal{N}$  is not too nonlinear” is replaced by the requirement of a small  $T(r)$  in Theorem 2; notice that  $T(r) = 0$  if  $\mathcal{N}$  is linear. The statement “ $\mathcal{G}$  is a sufficiently good low-pass filter” is replaced by Hypotheses B1.5, B3.2 (or B3.3’), and the requirement of a small  $B(\omega)$  in Theorem 2. Hypothesis B2 specifies the nature of the intersection of the loci of  $-1/N$  and  $G$  and makes precise the notion that the intersection should be “solid.”

As a further connection with the describing function method we note that the key inequality in Theorem 2 has as its right-hand side a quantity which may be obtained immediately from the graphs of  $G(j\omega)$  and  $-1/N(r)$ .

**6. Examples.** Theorem 2 can be applied to the autonomous nonlinear feedback system shown in Fig. 1 whether or not it contains hysteresis. This section contains an example of the use of Theorem 2 in either case.

The approach taken here is to use a computer to calculate  $-1/N(r)$ ,  $G(j\omega)$ , and  $|1/N(r) + G(j\omega)|/(B(\omega)T(r))$  as functions of  $\omega$  and  $r$ . Plotting the loci of  $-1/N$  and  $G$  gives the values of  $\omega_0$  and  $r_0$  at their intersection. With the printout of

$|1/N(r) + G(j\omega)|/(B(\omega)T(r))$  in a rectangular array, it is a simple matter to find a rectangle in that array with sides  $\omega = \omega_*$ ,  $\omega = \omega^*$ ,  $r = r_*$ , and  $r = r^*$  such that

$$1 < \frac{|1/N(r) + G(j\omega)|}{B(\omega)T(r)}$$

on its sides, and which contains  $(\omega_0, r_0)$ .

After finding such a rectangle it is a simple matter to check the remaining hypotheses of Theorem 2.

**6.1. Without hysteresis.** Consider the system shown in Fig. 1 with  $\mathcal{G}$  having the transfer function

$$G(s) = \frac{1000(s^2 + 3s + 2)}{s^5 + 31s^4 + 259s^3 + 1319s^2 + 1289s + 990}$$

and  $\mathcal{N}$  such that

$$(\mathcal{N}x)(\theta) = \begin{cases} 16 & \text{if } x(\theta) > 1, \\ 16x(\theta) & \text{if } |x(\theta)| \leq 1, \\ -16 & \text{if } x(\theta) < -1. \end{cases}$$

For this  $\mathcal{N}$  the describing function is [4, p. 59]

$$N(r) = \begin{cases} \frac{32}{\pi} \left[ \sin^{-1} \frac{1}{r} + \frac{1}{r} \sqrt{1 - \frac{1}{r^2}} \right], & r > 1, \\ 16, & 0 \leq r \leq 1. \end{cases}$$

Hypothesis B1 is clearly satisfied for this system.

$N$  is continuously differentiable for  $r > 1$  and  $G$  is continuously differentiable on the imaginary axis. The loci  $-1/N(r)$  and  $G(j\omega)$  are shown in Fig. 2. They intersect at  $(\omega_0, r_0) \approx (13.5, 4.6)$  and are clearly not tangent there. A simple computation, or closer examination of the loci, shows that both  $dN/dr$  and  $dG/d\omega$  are not zero at the intersection. Therefore Hypothesis B2 is satisfied.

Hypothesis B3 is satisfied with  $M = 16$ ,  $\omega_* = 10$  and  $\omega^* = 21$ .

Table 1 gives  $|1/N(r) + G(j\omega)|/(B(\omega)T(r))$  for several values of  $\omega$  and  $r$ . This ratio is greater than one for the points tabulated on the boundary of the set

$$\Omega = \{(\omega, r) : 10.5 < \omega < 14.5, 3.5 < r < 7.5\}.$$

It is in fact greater than one for every point on the boundary of this set. Therefore the system has a periodic response with fundamental frequency  $\omega \in (10.5, 14.5)$  and with fundamental magnitude  $r \in (3.5, 7.5)$ .

**6.2. With hysteresis.** Consider the system shown in Fig. 1 with  $\mathcal{G}$  having the transfer function

$$G(s) = 15,000e^{-s/3}/(s + 1)^4$$

and with  $\mathcal{N}$  the backlash nonlinearity shown in Fig. 3. For this  $\mathcal{N}$  the describing

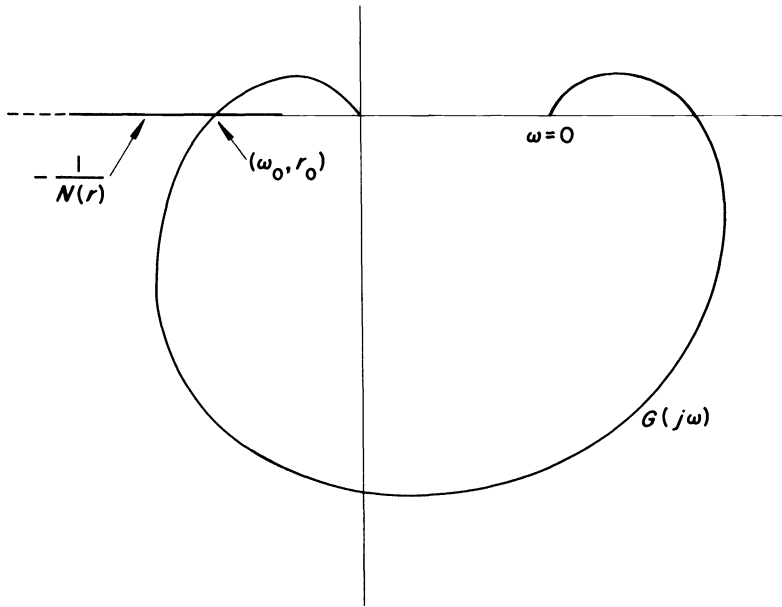


FIG. 2. Loci for the example in § 6.1 (not to scale)

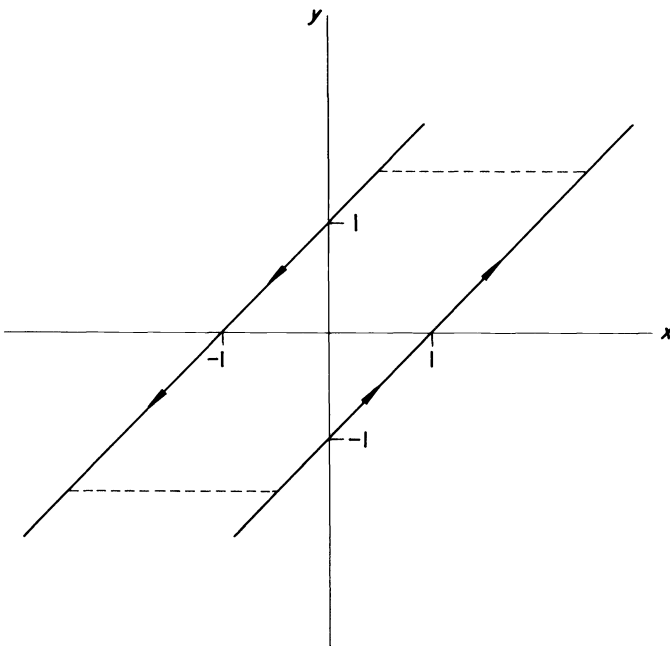


FIG. 3. Backlash nonlinearity for the example in § 6.2

TABLE I  
 $|1/N(r) + G(j\omega)|/B(\omega)T(r)$  versus  $(\omega, r)$  used to find the set  $\Omega$  for the example in § 6.1

$r = 8.0$	.95	1.00	1.08	1.37	1.99	3.01	4.45	6.35	8.77	11.78	15.46	19.91	25.21
7.5	.99	<b>1.01</b>	1.01	1.14	1.59	2.45	3.74	5.47	7.69	<b>10.47</b>	13.88	18.01	22.95
7.0	1.04	1.06	1.00	.98	1.22	1.90	3.02	4.57	6.59	9.14	12.29	16.11	20.69
6.5	1.11	1.14	1.07	.94	.92	1.35	2.28	3.66	5.48	7.80	10.68	14.19	18.41
6.0	1.21	1.27	1.21	1.04	.81	.85	1.54	2.73	4.36	6.45	9.05	12.25	16.12
5.5	1.32	1.42	1.41	1.26	.97	.60	.80	1.78	3.20	5.07	7.41	10.30	13.80
5.0	1.46	1.60	1.65	1.57	1.33	.90	.32	.79	2.03	3.68	5.76	8.34	11.47
4.5	1.61	1.87	1.94	1.94	1.80	1.49	.98	.26	.90	2.30	4.10	6.37	9.14
4.0	1.80	2.08	2.28	2.38	2.36	2.20	1.87	1.36	.85	1.23	2.59	4.47	6.86
3.5	2.03	<b>2.39</b>	2.68	2.90	3.02	3.02	2.88	2.59	2.22	<b>1.79</b>	1.96	3.03	4.84
3.0	2.33	2.79	3.19	3.55	3.83	4.01	4.07	4.01	3.82	3.54	3.27	3.29	3.97
	$\omega = 10.0$	<b>10.5</b>	11.0	11.5	12.0	12.5	13.0	13.5	14.0	<b>14.5</b>	15.0	15.5	16.0

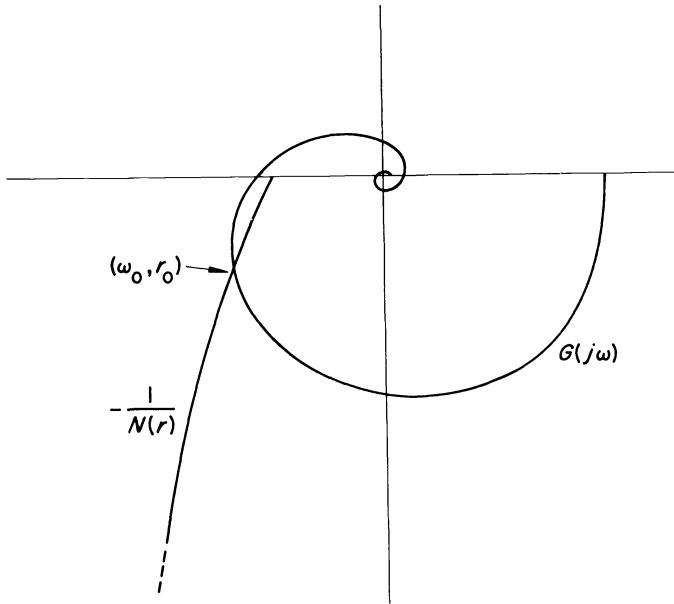


FIG. 4. Loci for the example in § 6.2 (not to scale)

function is [4, p. 69]

$$N(r) = \frac{1}{2} - \frac{1}{\pi} \sin^{-1} \left( 1 - \frac{2}{r} \right) + \frac{1}{\pi} \left( 1 - \frac{2}{r} \right) \sqrt{\frac{4}{r} - \frac{4}{r^2}} - j \frac{4}{\pi} \left( \frac{1}{r} - \frac{1}{r^2} \right), \quad r \geq 1.$$

Hypothesis B1 is clearly satisfied for this system.

$N$  is continuously differentiable for  $r > 1$  and  $G$  is continuously differentiable on the imaginary axis. The loci of  $G(j\omega)$  and  $-1/N(r)$  are shown in Fig. 4. They intersect at  $(\omega_0, r_0) \approx (1.35, 8.2)$  and are clearly not tangent there. Both  $dN/dr$  and  $dG/d\omega$  are not zero at the intersection. Therefore Hypothesis B2 is satisfied.

Hypothesis B3' is satisfied with  $M' = \sqrt{2}$ ,  $\omega_* = 7.5$  and  $\omega^* = 11$ .

Finally the set

$$\Omega = \{(\omega, r) : 8.0 < \omega < 8.4, 1.3 < r < 1.4\}$$

satisfies the remaining hypothesis of Theorem 2. Therefore the system has a periodic response with fundamental frequency  $\omega \in (8.0, 8.4)$  and fundamental magnitude  $r \in (1.3, 1.4)$ .

**7. Conclusion.** Sufficient conditions for the existence of half-wave symmetric oscillations are given and error bounds for the approximate solution found by the use of the describing function method are provided. The heuristic justifications concerning the describing function method are validated and rendered more precise.

**Appendix I.** This appendix consists of the proofs of all theorems and lemmas in this paper. Appendix II contains the definition of local degree and statements of theorems which are used in the proofs here.

*Proof of Lemma 1.* Let

$$F(\omega, \mathbf{r})x \triangleq -\mathcal{P}^* \mathcal{A}_\omega \mathcal{F}[y(\mathbf{r}) + x]$$

for  $(\omega, \mathbf{r}) \in [\omega_*, \omega^*] \times \mathbb{R}^{k-1}$ .

Given  $(\omega, \mathbf{r}) \in [\omega_*, \omega^*] \times \mathbb{R}^{k-1}$  and  $x$  and  $z \in \mathcal{P}^* \mathcal{L}_\pi^2$ , the linearity of  $\mathcal{A}_\omega$  and the bounds on the operator norms assumed in Hypothesis A yield

$$\begin{aligned} \|F(\omega, \mathbf{r})x - F(\omega, \mathbf{r})z\| &= \|\mathcal{P}^* \mathcal{A}_\omega \mathcal{F}[y(\mathbf{r}) + x] - \mathcal{P}^* \mathcal{A}_\omega \mathcal{F}[y(\mathbf{r}) + z]\| \\ &= \|\mathcal{P}^* \mathcal{A}_\omega [\mathcal{F}[y(\mathbf{r}) + x] - \mathcal{F}[y(\mathbf{r}) + z]]\| \\ &\leq \|\mathcal{P}^* \mathcal{A}_\omega\| \cdot \|\mathcal{F}\| \cdot \|x - z\| \\ &\leq \rho_\omega \|x - z\|. \end{aligned}$$

$F(\omega, \mathbf{r}) : \mathcal{P}^* \mathcal{L}_\pi^2 \rightarrow \mathcal{P}^* \mathcal{L}_\pi^2$  is a contraction for each  $(\omega, \mathbf{r}) \in [\omega_*, \omega^*] \times \mathbb{R}^{k-1}$  since  $\rho_\omega < 1$  for each  $\omega \in [\omega_*, \omega^*]$

By the contraction theorem, Theorem II.1 in Appendix II, there is a unique  $x^*(\omega, \mathbf{r}) \in \mathcal{P}^* \mathcal{L}_\pi^2$  such that

$$(I.1) \quad x^*(\omega, \mathbf{r}) = -\mathcal{P}^* \mathcal{A}_\omega \mathcal{F}[y(\mathbf{r}) + x^*(\omega, \mathbf{r})]$$

and

$$(I.2) \quad \|x^*(\omega, \mathbf{r})\| \leq \frac{\|\mathcal{P}^* \mathcal{A}_\omega \mathcal{F} y(\mathbf{r})\|}{1 - \rho_\omega}.$$

$x : [\omega_*, \omega^*] \times \mathbb{R}^{k-1} \rightarrow \mathcal{L}_\pi^2$  is continuous by Theorem II.1, since  $F(\omega, \mathbf{r})$  is continuous.

The remainder of the lemma may be established by noting that  $x(\omega, \mathbf{r}) = y(\mathbf{r}) + x^*(\omega, \mathbf{r})$  implies  $\mathcal{P}x(\omega, \mathbf{r}) = y(\mathbf{r})$ , which is part b.

It also implies

$$\begin{aligned} \mathcal{P}^* x(\omega, \mathbf{r}) &= x^*(\omega, \mathbf{r}) \\ (I.3) \quad &= -\mathcal{P}^* \mathcal{A}_\omega \mathcal{F}[y(\mathbf{r}) + x^*(\omega, \mathbf{r})] \quad (\text{by (I.1)}) \\ &= -\mathcal{P}^* \mathcal{A}_\omega \mathcal{F} x(\omega, \mathbf{r}). \end{aligned}$$

This completes part c.

Equations (I.2) and (I.3) imply

$$\|\mathcal{P}^*x(\omega, \mathbf{r})\| \leq \frac{\|\mathcal{P}^*\mathcal{A}_\omega\mathcal{F}y(\mathbf{r})\|}{1 - \rho_\omega}.$$

*Proof of Lemma 2.* (i) By Lemma 1,  $\mathcal{A}_\omega\mathcal{F}x(\omega, \mathbf{r})$  is continuous in  $(\omega, \mathbf{r})$ . Therefore  $\langle \mathcal{A}_\omega\mathcal{F}x(\omega, \mathbf{r}), \xi_k \rangle$  is continuous in  $(\omega, \mathbf{r})$ , and as a result of (11),  $V$  is continuous. Similarly  $U$  is continuous.

(ii)

$$\begin{aligned} |U_k(\omega, \mathbf{r}) - V_k(\omega, \mathbf{r})|^2 &= |\langle \mathcal{A}_\omega\mathcal{F}y(\mathbf{r}), \xi_k \rangle - \langle \mathcal{A}_\omega\mathcal{F}x(\omega, \mathbf{r}), \xi_k \rangle|^2 \\ &= |\langle \mathcal{A}_\omega[\mathcal{F}y(\mathbf{r}) - \mathcal{F}x(\omega, \mathbf{r})], \xi_k \rangle|^2 \\ &= \langle \mathcal{A}_\omega[\mathcal{F}y(\mathbf{r}) - \mathcal{F}x(\omega, \mathbf{r})], \xi_k \rangle^2. \end{aligned}$$

Therefore,

$$\begin{aligned} |U(\omega, \mathbf{r}) - V(\omega, \mathbf{r})|_2^2 &= \sum_{k=1}^\kappa \langle \mathcal{A}_\omega[\mathcal{F}y(\mathbf{r}) - \mathcal{F}x(\omega, \mathbf{r})], \xi_k \rangle^2 \\ &= \|\mathcal{P}\mathcal{A}_\omega[\mathcal{F}y(\mathbf{r}) - \mathcal{F}x(\omega, \mathbf{r})]\|^2. \end{aligned}$$

Therefore,

$$\begin{aligned} |U(\omega, \mathbf{r}) - V(\omega, \mathbf{r})|_2 &= \|\mathcal{P}\mathcal{A}_\omega[\mathcal{F}y(\mathbf{r}) - \mathcal{F}x(\omega, \mathbf{r})]\| \\ &\leq \|\mathcal{P}\mathcal{A}_\omega\| \cdot \|\mathcal{F}\| \cdot \|y(\mathbf{r}) - x(\omega, \mathbf{r})\| \\ &\leq \|\mathcal{P}\mathcal{A}_\omega\| \cdot \|\mathcal{F}\| \cdot \|x^*(\omega, \mathbf{r})\| \\ &\leq \|\mathcal{P}\mathcal{A}_\omega\| \cdot \|\mathcal{F}\| \cdot \frac{\|\mathcal{P}^*\mathcal{A}_\omega\mathcal{F}y(\mathbf{r})\|}{1 - \rho_\omega}. \end{aligned}$$

*Proof of Theorem 1.* (i) In light of Theorem II.2 in Appendix II, define

$$\Phi(\omega, \mathbf{r}, \mu) = U(\omega, \mathbf{r}) - \mu[U(\omega, \mathbf{r}) - V(\omega, \mathbf{r})] \quad \text{for } \mu \in [0, 1] \quad \text{and } (\omega, \mathbf{r}) \in \bar{\Omega}.$$

Then  $\Phi: \bar{\Omega} \times [0, 1] \rightarrow \mathbb{R}^\kappa$  is continuous since  $U$  and  $V$  are continuous.

If  $(\omega, \mathbf{r}) \in \partial\Omega$ ,

$$\begin{aligned} |\Phi(\omega, \mathbf{r}, \mu)|_2 &= |U(\omega, \mathbf{r}) - \mu[U(\omega, \mathbf{r}) - V(\omega, \mathbf{r})]|_2 \\ &\geq |U(\omega, \mathbf{r})|_2 - \mu|U(\omega, \mathbf{r}) - V(\omega, \mathbf{r})|_2 \\ &\geq |U(\omega, \mathbf{r})| - |U(\omega, \mathbf{r}) - V(\omega, \mathbf{r})|_2 \\ &> 0 \quad (\text{by assumption 3 of Theorem 1}). \end{aligned}$$

Therefore  $(\omega, \mathbf{r}) \in \partial\Omega$  implies  $\Phi(\omega, \mathbf{r}, \mu) \neq 0$  for all  $\mu \in [0, 1]$ .

By Theorem II.2,  $d[\Phi(\cdot, \cdot, 1), 0, \Omega] = d[\Phi(\cdot, \cdot, 0), 0, \Omega]$ . Now  $\Phi(\cdot, \cdot, 1) = V$  and  $\Phi(\cdot, \cdot, 0) = U$  so by assumption 2 of Theorem 1,

$$d[V, 0, \Omega] = d[U, 0, \Omega] \neq 0.$$

Therefore, by Theorem II.3 there is a point  $(\omega, \mathbf{r}) \in \Omega$  such that  $V(\omega, \mathbf{r}) = 0$ ; i.e., using (11), we have

$$-r_k = \langle \mathcal{A}_\omega\mathcal{F}x(\omega, \mathbf{r}), \xi_k \rangle \quad \text{for } k = 1, 2, \dots, \kappa.$$

Therefore,

$$\begin{aligned} \mathcal{P}\mathcal{A}_\omega\mathcal{F}x(\omega, \mathbf{r}) &= -\sum_{k=1}^{\kappa-1} r_k \xi_k \\ &= -y(\mathbf{r}) = -\mathcal{P}x(\omega, \mathbf{r}). \end{aligned}$$

(ii) By Lemma 1,

$$\mathcal{P}^*\mathcal{A}_\omega\mathcal{F}x(\omega, \mathbf{r}) = -\mathcal{P}^*x(\omega, \mathbf{r}).$$

(iii) Also by Lemma 1,

$$\|\mathcal{P}^*x(\omega, \mathbf{r})\| \leq \frac{\|\mathcal{P}^*\mathcal{A}_\omega\mathcal{F}y(\mathbf{r})\|}{1 - \rho_\omega},$$

which proves the theorem.

*Proof of Lemma 3.* Suppose  $x \in \mathcal{L}_\pi^2, x = -\mathcal{G}_\omega \mathcal{N}x, \omega \in [\omega_*, \omega^*], \omega > 0$ , and  $\tilde{x}(t) = x(\omega t)$  for all  $t \in \mathbb{R}$ . Then

$$\begin{aligned} \tilde{x}(t) = x(\omega t) &= -[\mathcal{G}_\omega \mathcal{N}x](\omega t) \\ &= -\left[ \sum_{-\infty}^{\infty} G(j\omega k) \eta_k e^{jk\theta} \right] \Big|_{\theta = \omega t} \\ &= -\sum_{-\infty}^{\infty} G(j\omega k) \eta_k e^{jk\omega t}, \end{aligned}$$

where  $\eta_k$  is the  $k$ th complex Fourier coefficient of  $\mathcal{N}x$ ; i.e.,

$$\begin{aligned} \eta_k &= \frac{1}{\pi} \int_0^{2\pi} [\mathcal{N}x](\theta) e^{-jk\theta} d\theta \\ &= \frac{2}{T} \int_0^T [\mathcal{N}x](\omega t) e^{-jk\omega t} dt. \end{aligned}$$

Therefore  $\eta_k$  is the  $k$ th complex Fourier coefficient of  $\mathcal{N}\tilde{x}$  as well. For the input  $\mathcal{N}x: \mathbb{R} \rightarrow \mathbb{R}$  the output of  $\mathcal{G}$  is

$$[\mathcal{G}\mathcal{N}\tilde{x}](t) = \sum_{k=-\infty}^{\infty} G(j\omega k) \eta_k e^{j\omega k t}, \quad t \in \mathbb{R}.$$

Therefore  $\tilde{x} = -\mathcal{G}\mathcal{N}\tilde{x}$ .

*Proof of Lemma 4.* To prove parts a, b and c, pick any  $x \in \mathcal{L}_\pi^2$ , expand it in a Fourier series, and use the usual manipulations. This same approach gives

$$\|\mathcal{G}_\omega - \mathcal{G}_v\| = \max \{|G(j\omega k) - G(jvk)| : k \text{ odd}\}.$$

By Hypotheses B1.4 and B1.5,  $\|\mathcal{G}_\omega - \mathcal{G}_v\| \rightarrow 0$  as  $\omega \rightarrow v$ .

*Proof of Lemma 5.*

$$U_k(\omega, r) = r_k + \langle \mathcal{G}_\omega \mathcal{N}y(\mathbf{r}), \xi_k \rangle.$$

The complex Fourier coefficients of  $\mathcal{N}y(\mathbf{r})$  for  $l = \pm 1, \pm 3, \pm 5, \dots$  are

$$\frac{1}{2\pi} \int_0^{2\pi} \mathcal{N}[y(\mathbf{r})](\theta) e^{-jl\theta} d\theta = \frac{1}{2j} \eta_l.$$

Therefore  $\mathcal{G}_{\omega, \mathcal{N}} y(\mathbf{r})$  has Fourier series  $(1/(2j)) \sum_{-\infty}^{\infty} G(j\omega l) \eta_l e^{jl\theta}$ . Now for  $n = \pm 1, \pm 3, \dots$ ,

$$\frac{1}{\pi} \int_0^{2\pi} \left[ \frac{1}{2j} \sum_{-\infty}^{\infty} G(j\omega l) \eta_l e^{jl\theta} \right] e^{+jn\theta} d\theta = \frac{1}{j} G(j\omega n) \eta_{-n}.$$

Therefore if  $k$  is odd,  $\xi_k(\theta) = \sin k\theta$  and

$$\begin{aligned} \langle \mathcal{G}_{\omega, \mathcal{N}} y(\mathbf{r}), \xi_k \rangle &= \frac{1}{\pi} \int_0^{2\pi} \left[ \frac{1}{2j} \sum_{-\infty}^{\infty} G(j\omega l) \eta_l e^{jl\theta} \right] \left[ \frac{e^{j\theta k} - e^{-j\theta k}}{2j} \right] d\theta \\ &= -\frac{1}{2} G(-j\omega k) \eta_{-k} + \frac{1}{2} G(j\omega k) \eta_k \\ &= \frac{1}{2} [\bar{G}(j\omega k) \bar{\eta}_k + G(j\omega k) \eta_k] \\ &= \text{Re } G(j\omega k) \eta_k \end{aligned}$$

since

$$\eta_{-k} = \frac{j}{\pi} \int_0^{2\pi} \mathcal{N}(y(\mathbf{r}))(\theta) e^{jk\theta} d\theta = -\bar{\eta}_k,$$

the conjugate of  $-\eta_k$ . Therefore

$$\begin{aligned} U_k(\omega, \mathbf{r}) &= r_k + \langle \mathcal{G}_{\omega, \mathcal{N}} y(\mathbf{r}), \xi_k \rangle \\ &= \text{Re } [r_k + G(j\omega k) \eta_k], \end{aligned} \quad k \text{ odd.}$$

If  $k$  is even,  $\xi_k(\theta) = \cos(k-1)\theta$  and

$$\begin{aligned} \langle \mathcal{G}_{\omega, \mathcal{N}} y(\mathbf{r}), \xi_k \rangle &= \frac{1}{\pi} \int_0^{2\pi} \left[ \frac{1}{2j} \sum_{-\infty}^{\infty} G(j\omega l) \eta_l e^{jl\theta} \right] \left[ \frac{e^{j\theta(k-1)} + e^{-j\theta(k-1)}}{2} \right] \\ &= \frac{1}{2j} G(-j\omega(k-1)) \eta_{-(k-1)} + \frac{1}{2j} G(j\omega(k-1)) \eta_{k-1} \\ &= \frac{1}{2j} [-\bar{G}(j\omega(k-1)) \bar{\eta}_{k-1} + G(j\omega(k-1)) \eta_{k-1}] \\ &= \frac{1}{2j} [2j \text{Im } G(j\omega(k-1)) \eta_{k-1}] \\ &= \text{Im } G(j\omega(k-1)) \eta_{k-1}. \end{aligned}$$

Therefore,

$$\begin{aligned} U_k(\omega, \mathbf{r}) &= r_k + \langle \mathcal{G}_{\omega, \mathcal{N}} y(\mathbf{r}), \xi_k \rangle \\ &= r_k + \text{Im } [G(j\omega(k-1)) \eta_{k-1}], \end{aligned} \quad k \text{ even.}$$

*Proof of Lemma 6.* From assumption 3 of Hypothesis B2, for all  $(\omega, r) \in \Omega \subset [\omega_*, \omega^*] \times [r_*, r^*]$ ,

$$W(\omega, r) \triangleq (\text{Re } [r + rN(r)G(j\omega)], \text{Im } [r + rN(r)G(j\omega)]) = (0, 0)$$

if and only if  $(\omega, r) = (\omega_0, r_0) \in \Omega$ . By assumptions 1 and 2 of Hypothesis B2,  $W(\omega, r)$  is  $C^{(1)}$  on  $\bar{\Omega}$ . Using the analytic definition in Appendix II, we may therefore calculate the local degree of the map  $W$  relative to the point 0 and the set  $\Omega$ ,

$$d[W, 0, \Omega] = \text{sgn det } J_W(\omega_0, r_0),$$



where

$$J_W(\omega_0, r_0) = \begin{bmatrix} \frac{\partial W_1}{\partial \omega} & \frac{\partial W_1}{\partial r} \\ \frac{\partial W_2}{\partial \omega} & \frac{\partial W_2}{\partial r} \end{bmatrix} \Big|_{(\omega_0, r_0)}$$

We next use the remaining assumptions of Hypothesis B2 to show that  $\det J_W(\omega_0, r_0) \neq 0$  and therefore  $d[W, 0, \Omega] \neq 0$ . By direct calculation,

$$J_W(\omega_0, r_0) = r_0 \begin{bmatrix} \operatorname{Re} N(r_0)G'(j\omega_0) & \operatorname{Re} N'(r_0)G(j\omega_0) \\ \operatorname{Im} N(r_0)G'(j\omega_0) & \operatorname{Im} N'(r_0)G(j\omega_0) \end{bmatrix},$$

where the prime indicates the derivative and use has been made of the assumption in B2 that  $1 + N(r_0)G(j\omega_0) = 0$ .

Now if the columns of  $J_W$  are linearly independent,  $\det [J_W(\omega_0, r_0)] \neq 0$ . It is natural to reintroduce the complex plane;  $\det [J_W(\omega_0, r_0)] \neq 0$  if the complex numbers  $z_1 = N(r_0)G'(j\omega_0)$  and  $z_2 = N'(r_0)G(j\omega_0)$  are linearly independent. Now  $N(r_0)G(j\omega_0) = -1$  implies  $N(r_0)$  and  $G(j\omega_0) \neq 0$ . Assumption 3 of B2 specifies  $G'(j\omega_0)$  and  $N'(r_0) \neq 0$ . Therefore neither of the complex numbers  $z_1$  or  $z_2$  is zero.

We now use the relation

$$\frac{d}{dr} \left[ -\frac{1}{N(r)} \right] \Big|_{r_0} = \frac{N'(r_0)}{N^2(r_0)}$$

and find

$$z_2 = N'(r_0)G(j\omega_0) = -N(r_0) \frac{d}{dr} \left[ -\frac{1}{N(r)} \right] \Big|_{r_0}.$$

It is clear that the two (nonzero) complex numbers  $z_1$  and  $z_2$  are linearly independent if  $G'(j\omega_0)$  and  $(d/dr)[-1/N(r)]|_{r_0}$  are linearly independent. Geometrically,  $G'(j\omega_0)$  is a vector tangent to the Nyquist locus at the point  $\omega_0$  while  $(d/dr)[-1/N(r)]|_{r_0}$  is a vector tangent to the critical locus at  $r_0$ . Since by assumption neither vector is zero, the linear independence is assured if the two tangents to the loci at their point of intersection  $(\omega_0, r_0)$  differ. This is assured by part d in assumption 3 of Hypothesis B2.

Therefore  $d[\omega, 0, \Omega] \neq 0$ .

*Proof of Theorem 2.* It is useful first to derive expressions (I.4) and (I.5) below, and then to consider cases B and B' separately. Set  $\kappa = 2$  for defining  $\mathcal{P}$  and  $\mathcal{P}^*$ . Then:

(i)

$$\begin{aligned} rN(r) &= \frac{j}{\pi} \int_0^{2\pi} [\mathcal{N}(r\xi_1)](\theta) e^{-j\theta} d\theta \\ &= \frac{j}{\pi} \int_0^{2\pi} [\mathcal{N}(r\xi_1)](\theta) [\cos \theta - j \sin \theta] d\theta \\ &= \langle \mathcal{N}(r\xi_1), \xi_1 \rangle + j \langle \mathcal{N}(r\xi_1), \xi_2 \rangle, \end{aligned}$$

so

$$[\mathcal{P}\mathcal{N}(r\xi_1)](\theta) \triangleq \langle \mathcal{N}(r\xi_1), \xi_1 \rangle \sin \theta + \langle \mathcal{N}(r\xi_1), \xi_2 \rangle \cos \theta,$$

$$(I.4) \quad [\mathcal{P}\mathcal{N}(r\xi_1)](\theta) = |rN(r)| \sin [\theta + \angle N(r)] \quad \text{for } r > 0.$$

(ii)

$$\begin{aligned} \|\mathcal{P}^*\mathcal{N}(r\xi_1)\|^2 &= \|\mathcal{N}(r\xi_1)\|^2 - \|\mathcal{P}\mathcal{N}(r\xi_1)\|^2 \\ &= \|\mathcal{N}(r\xi_1)\|^2 - |rN(r)|^2 \end{aligned}$$

so

$$\begin{aligned} T(r) &= \sqrt{\frac{\|\mathcal{N}(r\xi_1)\|^2}{|rN(r)|^2} - 1} = \sqrt{\frac{\|\mathcal{N}(r\xi_1)\|^2 - |rN(r)|^2}{|rN(r)|^2}} \\ &= \frac{\|\mathcal{P}^*\mathcal{N}(r\xi_1)\|}{|rN(r)|}. \end{aligned}$$

Therefore  $B(\omega)T(r) < |1/N(r) + G(j\omega)|$  implies

$$(I.5) \quad B(\omega)\|\mathcal{P}^*\mathcal{N}(r\xi_1)\| < |r + rN(r)G(j\omega)| = |W(\omega, r)|_2.$$

Case B3. (i) Hypothesis A is satisfied with  $\mathcal{A}_\omega \triangleq \mathcal{G}_\omega, \mathcal{F} \triangleq \mathcal{N}, \kappa \triangleq 2$  and  $y(r) = r\xi_1$ .

*Proof.* Clearly part 1 is satisfied by Hypothesis B1 and Lemma 4. Parts 2 and 3 are satisfied by Hypothesis B3 and Lemma 4.

(ii)

$$\begin{aligned} \langle \mathcal{G}_\omega \mathcal{N}(r\xi_1), \xi_1 \rangle &= \langle \mathcal{G}_\omega \mathcal{P}\mathcal{N}(r\xi_1), \xi_1 \rangle \\ &= \frac{1}{\pi} \int_0^{2\pi} |G(j\omega)| |rN(r)| \sin(\theta + \angle N + \angle G) \sin \theta \, d\theta \\ &= |G(j\omega)| \cdot |rN(r)| \cos(\angle N + \angle G) \\ &= \text{Re} [rN(r)G(j\omega)], \end{aligned} \quad r \geq r_*.$$

Similarly,

$$\langle \mathcal{G}_\omega \mathcal{N}(r\xi_1), \xi_2 \rangle = \text{Im} [rN(r)G(j\omega)].$$

Therefore by (7),

$$\begin{aligned} U(\omega, r) &= (r + \text{Re} [rN(r)G(j\omega)], \text{Im} [rN(r)G(j\omega)]) \\ &= (\text{Re} [r + rN(r)G(j\omega)], \text{Im} [r + rN(r)G(j\omega)]) \\ &= W(\omega, r). \end{aligned}$$

So  $d[U, 0, \Omega] \neq 0$  by Lemma 4.

(iii)

$$\begin{aligned} |U(\omega, r) - V(\omega, r)|_2 &\leq \frac{\|\mathcal{P}\mathcal{G}_\omega\| \cdot \|\mathcal{N}\| \cdot \|\mathcal{P}^*\mathcal{G}_\omega\mathcal{N}(r\xi_1)\|}{1 - \rho_\omega} \\ &\leq \frac{\|\mathcal{P}\mathcal{G}_\omega\| \cdot \|\mathcal{N}\| \cdot \|\mathcal{P}^*\mathcal{G}_\omega\| \cdot \|\mathcal{P}^*\mathcal{N}(r\xi_1)\|}{1 - \rho_\omega} \\ &= \frac{|G(j\omega)|\rho_\omega}{1 - \rho_\omega} \|\mathcal{P}^*\mathcal{N}(r\xi_1)\| \\ &= B(\omega)\|\mathcal{P}^*\mathcal{N}(r\xi_1)\| \\ &< |W(\omega, r)|_2 = |U(\omega, r)|_2, \end{aligned}$$

by (I.5) for all  $(\omega, r) \in \partial\Omega$ . Therefore,

$$|U(\omega, r) - V(\omega, r)|_2 < |U(\omega, r)| \quad \text{for all } (\omega, r) \in \partial\Omega.$$

By Theorem 1, there is an  $(\omega, r) \in \Omega$  such that  $x = -\mathcal{G}_\omega\mathcal{N}x$ , where  $x = x(\omega, r)$ . Also for  $x \in \mathcal{L}_\pi^2$ ,

$$\begin{aligned} \mathcal{P}x = r\xi_1 \quad \text{and} \quad \|\mathcal{P}^*x\| = \|x^*(\omega, r)\| &\leq \frac{\|\mathcal{P}^*\mathcal{G}_\omega\mathcal{N}(r\xi_1)\|}{1 - \rho_\omega} \\ &\leq \frac{\|\mathcal{P}^*\mathcal{G}_\omega\| \cdot \|\mathcal{P}^*\mathcal{N}(r\xi_1)\|}{1 - \rho_\omega} \\ &= \left[ \frac{\rho_\omega}{1 - \rho_\omega} \|\mathcal{P}^*\mathcal{N}(r\xi_1)\| \right] \frac{1}{M} \\ &= \frac{\rho_\omega}{1 - \rho_\omega} \frac{|rN(r)|}{M} T(r). \end{aligned}$$

*Case B3'*: The proof in this case is very similar to the preceding one. The essential difference is that hypothesis B3.1 is not satisfied. This problem is avoided by defining an operator  $\mathcal{C} : \mathcal{L}_\pi^2 \rightarrow \mathcal{L}_\pi^2$  such that  $y$  is the derivative of  $\mathcal{C}y$ , and by defining an operator  $\mathcal{D} : \mathcal{G}_\omega\mathcal{L}_\pi^2 \rightarrow \mathcal{L}_\pi^2$  such that  $\mathcal{D}y$  is the derivative of  $y$ , i.e.,  $\mathcal{D}\mathcal{C}y = y$ .

The proof now proceeds in essentially the same way as before but with  $\mathcal{A}_\omega = \mathcal{D}\mathcal{G}_\omega$  and  $\mathcal{F} = \mathcal{N}\mathcal{C}$ .

(i) Define  $\mathcal{C} : \mathcal{L}_\pi^2 \rightarrow \mathcal{L}_\pi^2$  such that

$$(\mathcal{C}x)(\theta) = \begin{cases} \frac{1}{2} \int_0^\pi \text{sgn}(\theta - \tau)x(\tau) d\tau, & 0 \leq \theta < \pi, \\ -(\mathcal{C}x)(\theta - \pi), & \pi \leq \theta < 2\pi. \end{cases}$$

Note 1.

$$(\mathcal{C}x)(\theta) = \int_0^\theta x(\tau) d\tau - \frac{1}{2} \int_0^\pi x(\tau) d\tau, \quad 0 \leq \theta < \pi,$$

so  $(\mathcal{C}x)(\theta) = -(\mathcal{C}x)(\pi) = (\mathcal{C}x)(2\pi)$  and  $\mathcal{C}x$  is continuous for all  $x \in \mathcal{L}_\pi^2$ .

Note 2.

$$\mathcal{C}[e^{jn(\cdot)}] = \frac{1}{jn} e^{jn(\cdot)}, \quad n \text{ odd.}$$

Note 3.

$$\begin{aligned} |(\mathcal{C}x)(\theta)| &= \frac{1}{2} \left| \int_0^\pi \operatorname{sgn}(\theta - \tau)x(\tau) d\tau \right| \quad (0 \leq \theta < 2\pi) \\ &\leq \frac{1}{2} \int_0^\pi |x(\tau)| d\tau \\ &\leq \frac{1}{2} \left( \int_0^\pi 1^2 d\tau \right)^{1/2} \left( \int_0^\pi |x(\tau)|^2 d\tau \right)^{1/2} \\ &= \frac{\pi}{2\sqrt{2}} \|x\|. \end{aligned}$$

Therefore,

$$\|\mathcal{C}x\|_\infty \leq \frac{\pi}{2\sqrt{2}} \|x\| \quad \text{for all } x \in \mathcal{L}_\pi^2.$$

Now

$$\begin{aligned} \|\mathcal{N}\mathcal{C}x - \mathcal{N}\mathcal{C}y\| &\leq M' \|\mathcal{C}x - \mathcal{C}y\|_\infty \\ &\leq M' \frac{\pi}{2\sqrt{2}} \|x - y\| \quad \text{for all } x, y \in \mathcal{L}_\pi^2, \end{aligned}$$

so  $\|\mathcal{N}\mathcal{C}\| \leq \pi M'/(2\sqrt{2})$ .

(ii) Define  $\mathcal{D}$  such that  $\mathcal{D}[e^{jn(\cdot)}] = jne^{jn(\cdot)}$  and

$$\mathcal{D} \sum_{k=1}^\infty a_k \xi_k = \sum_{k=1}^\infty a_k \mathcal{D}\xi_k.$$

Now  $\mathcal{D}$  does not map  $\mathcal{L}_\pi^2$  into  $\mathcal{L}_\pi^2$ , but  $\mathcal{D}\mathcal{G}_\omega$  does.  $[\mathcal{D}\mathcal{G}_\omega][e^{jn(\cdot)}] = jnG(jn\omega) \cdot e^{jn(\cdot)}$ , and  $[jnG(jn\omega)]$  satisfies Hypothesis B1, by Hypothesis B3', so by Lemma 4,

$$\begin{aligned} \mathcal{D}\mathcal{G}_\omega &: \mathcal{L}_\pi^2 \rightarrow \mathcal{L}_\pi^2, \\ \|\mathcal{D}\mathcal{G}_\omega\| &= \max \{n|G(jn\omega)| : n \text{ odd}\}, \\ \|\mathcal{P}^*\mathcal{D}\mathcal{G}_\omega\| &= \max \{k|G(jk\omega)| : k = 3, 5, \dots\}, \\ \|\mathcal{P}\mathcal{D}\mathcal{G}_\omega\| &= |G(j\omega)|, \\ \|\mathcal{P}^*(\mathcal{D}\mathcal{G}_\omega = \mathcal{D}\mathcal{G}_\nu)\| &\leq \|\mathcal{D}\mathcal{G}_\omega - \mathcal{D}\mathcal{G}_\nu\|_{\nu \rightarrow \omega} \rightarrow 0 \quad \text{for all } \omega \in [\omega_*, \omega^*]. \end{aligned}$$

(iii)

$$\begin{aligned} \|\mathcal{P}^*\mathcal{D}\mathcal{G}_\omega\| \cdot \|\mathcal{N}\mathcal{C}\| &\leq \frac{\pi M'}{2\sqrt{2}} \max \{k|G(jk\omega)| : k = 3, 5, \dots\} \\ &\triangleq \rho'_\omega < 1 \quad \text{for all } \omega \in [\omega_*, \omega^*] \end{aligned}$$

by Hypothesis B3'. Therefore Hypothesis A is satisfied with  $\mathcal{A}_\omega \mathcal{D}\mathcal{G}_\omega, \mathcal{F} = \mathcal{N}\mathcal{C}$ , and  $\kappa = 2$ .

(iv) Hypothesis A is independent of the basis chosen for  $\mathcal{L}_\pi^2$  so define a new orthonormal basis for  $\mathcal{L}_\pi^2, \{z_k\}_1^\infty$ , such that

$$z_1 = \xi_2, \quad z_2 = -\xi_1, \quad z_k = \xi_k \quad \text{for } k > 2.$$

Notice that  $\mathcal{D}\xi_1 = z_1, \mathcal{D}\xi_2 = z_2, \mathcal{C}z_1 = \xi_1$  and  $\mathcal{C}z_2 = \xi_2$ . So,

$$\begin{aligned} \langle \mathcal{D}y, z_i \rangle &= \langle \mathcal{D}[\langle y, \xi_1 \rangle \xi_1 + \langle y, \xi_2 \rangle \xi_2], z_i \rangle \\ &= \langle [\langle y, \xi_1 \rangle z_1 + \langle y, \xi_2 \rangle z_2], z_i \rangle \\ &= \langle y, \xi_i \rangle \quad (i = 1 \text{ and } 2) \quad \text{for all } y \in \mathcal{L}_\pi^2, \end{aligned}$$

and

$$\begin{aligned} \langle \mathcal{D}\mathcal{G}_\omega \mathcal{N}\mathcal{C}y(r), z_i \rangle &= \langle \mathcal{G}_\omega \mathcal{N}\mathcal{C}y(r), \xi_i \rangle \\ &= \langle \mathcal{G}_\omega \mathcal{N}\mathcal{C}(rz_1), \xi_i \rangle \\ &= \langle \mathcal{G}_\omega \mathcal{N}\mathcal{C}(r\xi_1), \xi_i \rangle \quad (i = 1 \text{ and } 2). \end{aligned}$$

Therefore,

$$\begin{aligned} U(\omega, r) &\triangleq (r + \langle \mathcal{D}\mathcal{G}_\omega \mathcal{N}\mathcal{C}y(r), z_1 \rangle, \langle \mathcal{D}\mathcal{G}_\omega \mathcal{N}\mathcal{C}y(r), z_2 \rangle) \\ &= (r + \langle \mathcal{G}_\omega \mathcal{N}\mathcal{C}(r\xi_1), \xi_1 \rangle, \langle \mathcal{G}_\omega \mathcal{N}\mathcal{C}(r\xi_1), \xi_2 \rangle) \\ &= W(\omega, r) \end{aligned}$$

as in Case B3. Therefore  $d[u, 0, \Omega] \neq 0$  by Lemma 4.

(v)

$$\begin{aligned} |U(\omega, r) - V(\omega, r)|_2 &\leq \frac{\|\mathcal{P}\mathcal{D}\mathcal{G}_\omega\| \cdot \|\mathcal{N}\mathcal{C}\| \cdot \|\mathcal{P}^*\mathcal{D}\mathcal{G}_\omega \mathcal{N}\mathcal{C}(r\xi_1)\|}{1 - \rho_\omega} \quad (\text{by Lemma 2}). \\ &\leq \frac{\|\mathcal{P}\mathcal{D}\mathcal{G}_\omega\| \cdot \|\mathcal{N}\mathcal{C}\| \cdot \|\mathcal{P}^*\mathcal{D}\mathcal{G}_\omega\| \cdot \|\mathcal{P}^*\mathcal{N}\mathcal{C}(r\xi_1)\|}{1 - \rho_\omega} \\ &\leq \frac{|G(j\omega)|\rho_\omega}{1 - \rho_\omega} \|\mathcal{P}^*\mathcal{N}\mathcal{C}(r\xi_1)\| \\ &= B(\omega)\|\mathcal{P}^*\mathcal{N}\mathcal{C}(r\xi_1)\|, \end{aligned}$$

so

$$|U(\omega, r) - V(\omega, r)|_2 < B(\omega)\|\mathcal{P}^*\mathcal{N}\mathcal{C}(r\xi_1)\| \quad \text{for all } (\omega, r) \in \bar{\Omega}.$$

By assumption 2 of Theorem 2,  $(\omega, r) \in \partial\Omega$  implies  $B(\omega)T(r) < |1/N(r) + G(j\omega)|$ , and that implies

$$B(\omega)\|\mathcal{P}^*\mathcal{N}\mathcal{C}(r\xi_1)\| < |W(\omega, r)|_2 = |U(\omega, r)|_2$$

by the comments at the beginning of the proof. Therefore

$$|U(\omega, r) - V(\omega, r)|_2 < |U(\omega, r)|_2 \quad \text{for all } (\omega, r) \in \partial\Omega.$$

Therefore the hypotheses of Theorem 1 are satisfied, so there exists  $(\omega, r) \in \Omega$  such that

$$x(\omega, r) = -\mathcal{D}\mathcal{G}_\omega \mathcal{N} \mathcal{C}x(\omega, r), \quad \mathcal{C}x(\omega, r) = -\mathcal{G}_\omega \mathcal{N} \mathcal{C}x(\omega, r).$$

Let  $z \triangleq \mathcal{C}x(\omega, r)$ . Then  $z = -\mathcal{G}_\omega \mathcal{N} z$ , which proves part a.

Also

$$\mathcal{P}z = \mathcal{P} \mathcal{C}x(\omega, r) = \mathcal{C}\mathcal{P}x(\omega, r) = \mathcal{C}y(r) = \mathcal{C}rz_1 = r\xi_1,$$

which proves part b.

$$\begin{aligned} \|\mathcal{P}^*z\| &= \|\mathcal{P}^* \mathcal{C}x(\omega, r)\| = \|\mathcal{C}\mathcal{P}^*x(\omega, r)\| \leq \|\mathcal{P}^*x(\omega, r)\| \\ &\leq \frac{\|\mathcal{P}^* \mathcal{D}\mathcal{G}_\omega \mathcal{N} \mathcal{C}y(r)\|}{1 - \rho'_\omega} \quad (\text{by Theorem 1}) \\ &\leq \frac{\|\mathcal{P}^* \mathcal{D}\mathcal{G}_\omega\| \cdot \|\mathcal{P}^* \mathcal{N} \mathcal{C}y(r)\|}{1 - \rho'_\omega} = \frac{\rho'_\omega}{1 - \rho'_\omega} \frac{\sqrt{8}}{(\pi M')} \|\mathcal{P}^* \mathcal{N}(r\xi_1)\| \\ &= \frac{\rho'_\omega \sqrt{8T(r)}|rN(r)|}{(1 - \rho'_\omega)\pi M'}. \end{aligned}$$

**Appendix II.** This appendix is included for the reader's convenience. It contains a statement of the contraction theorem, a definition of local degree, and two theorems related to local degree which are used in the proof of Theorem 1. The material related to contraction is taken from [5] and material related to local degree is taken from [2].

**THEOREM II.1** (fixed points by contraction). *Let  $\mathcal{L}$  be a Banach space, let  $S$  be a set  $\mathbb{R}^n$ , and for each  $s \in S$ , let  $T_s: \mathcal{L} \rightarrow \mathcal{L}$  be continuous and such that  $\|T_s x - T_s y\| \leq \rho_s \|x - y\|$  for all  $x, y \in \mathcal{L}$ , where  $\rho_s < 1$ . Then:*

- a. *There exists a unique  $x_s^* \in \mathcal{L}$  such that  $T_s x_s^* = x_s^*$ .*
- b.  $\|x_s^* - x\| \leq \|T_s x - x\| / (1 - \rho_s)$  for all  $x \in \mathcal{L}$ .  
*In particular  $\|x_s^*\| \leq \|T\mathbf{0}\| / (1 - \rho_s)$ , where  $\mathbf{0}$  is the zero of  $\mathcal{L}$ .*
- c. *If in addition,  $s \rightarrow s_0$  implies  $\|T_s - T_{s_0}\| \rightarrow 0$ , then  $s \rightarrow s_0$  implies  $\|x_s - x_{s_0}\| \rightarrow 0$ .*

**DEFINITION** (local degree for a special case). Let:

- 1.  $\Omega$  be an open bounded set in  $\mathbb{R}^n$ ,
- 2.  $f: \bar{\Omega} \rightarrow \mathbb{R}^n$  have continuous first partial derivatives,
- 3.  $f(x) \neq p$  for all  $x \in \partial\Omega$ , the boundary of  $\Omega$ ,
- 4.  $\det J_f(x) \triangleq \det(\partial f_i / \partial x_j) \neq 0$  for all  $x \in f^{-1}(p)$ .

Then the *degree of  $f$  relative to  $p$  and  $\Omega$*  is

$$d[f, p, \Omega] \triangleq \sum_{x \in f^{-1}(p) \cap \Omega} \text{sgn} [\det J_f(x)].$$

Berger and Berger continue the definition to the case where  $f$  is continuous but not differentiable. This part of the definition has been omitted as it is not needed here.

**THEOREM II.2** (homotopic invariance of  $d$ ). *If  $\Phi: \bar{\Omega} \times [0, 1] \rightarrow \mathbb{R}^n$  is continuous and  $\Phi(x, \mu) \neq p$  for all  $x \in \partial\Omega$  and all  $\mu \in [0, 1]$ , then  $d[\Phi(\cdot, \mu), p, \Omega]$  is the same for all  $\mu \in [0, 1]$ .*

THEOREM II.3 (existence of preimages). *If  $d[f, p, \Omega] \neq 0$ , then there is an  $x \in \Omega$  such that  $f(x) = p$ .*

## REFERENCES

- [1] R. W. BASS, *Mathematical legitimacy of equivalent linearization by describing functions*, Proc. 1960 International Federation of Automatic Control Congress, Butterworths, London, pp. 2074–2084.
- [2] M. BERGER AND M. BERGER, *Perspectives in Nonlinearity*, W. A. Benjamin, New York, 1968.
- [3] L. CESARI, *Functional Analysis and Periodic Solutions of Nonlinear Differential Equations*, Contributions to Differential Equations, vol. 1, Interscience, New York, 1962.
- [4] A. GELB AND W. E. VANDER VELDE, *Multiple-Input Describing Functions and Nonlinear System Design*, McGraw-Hill, New York, 1968.
- [5] L. V. KANTOROVICH AND G. P. AKILOV, *Functional Analysis in Normal Spaces*, Pergamon Press, Oxford, 1964.
- [6] R. J. KOCHENBURGER, *A frequency response method for analyzing and synthesizing contractor servomechanisms, Part 1*, Trans. AIEE, 69 (1950), pp. 270–284.
- [7] J. KUDREWICZ, *Theorems on the existence of periodic vibrations based upon the describing function method*, Proc. International Federation on Automatic Control, Warsaw, Poland, 1969.

## AN EXISTENCE THEOREM FOR LAGRANGE PROBLEMS WITH UNBOUNDED CONTROLS AND A SLENDER SET OF EXCEPTIONAL POINTS\*

L. CESARI,† J. R. LA PALM‡ AND D. A. SANCHEZ§

**Abstract.** In a number of existence theorems for Lagrange problems with unbounded controls certain growth conditions are required, which generally can be omitted for bounded controls. In the present paper the authors show that even for unbounded controls the growth conditions can be dispensed with in an arbitrary subset in the  $tx$ -space  $E_{1+n}$ , provided such a subset is "slender" according to suitable definitions. Any set contained in finitely many or countably many smooth curves  $x = \phi(t)$ ,  $a \leq t \leq b$ , is certainly slender. Many examples are given.

1. We deal here with Lagrange problems of control concerning the minimum of an integral of the form

$$(1) \quad I[x, u] = \int_{t_1}^{t_2} f_0(t, x(t), u(t)) dt$$

in classes  $\Omega$  of admissible pairs  $x(t) = (x^1, \dots, x^n)$ ,  $u(t) = (u^1, \dots, u^m)$ ,  $t_1 \leq t \leq t_2$ ,  $x$  absolutely continuous (AC),  $u$  measurable, satisfying a system of ordinary differential equations

$$(2) \quad dx/dt = f(t, x(t), u(t)), \quad t \in [t_1, t_2] \quad \text{a.e.},$$

$f = (f_1, \dots, f_n)$ , constraints of the forms

$$(3) \quad (t, x(t)) \in A \subset E_{n+1}, \quad t \in [t_1, t_2],$$

$$(4) \quad u(t) \in U(t, x(t)) \subset E_m, \quad t \in [t_1, t_2] \quad \text{a.e.},$$

and boundary conditions usually written in the McShane form

$$(5) \quad (t_1, x(t_1), t_2, x(t_2)) \in B \subset E_{2n+1},$$

where  $A, B$  are given fixed sets and  $U(t, x)$  is a given set which may depend on  $t$  and  $x$ .

We do not exclude that  $A$  may coincide with the whole space  $E_{n+1}$  and that  $U$  may coincide with the space  $E_m$ . For  $m = n$ ,  $f = u$ ,  $U = E_m$  these problems reduce to free problems of the calculus of variations, that is, problems concerning the minimum of a functional

$$(6) \quad I[x] = \int_{t_1}^{t_2} f_0(t, x(t), x'(t)) dt$$

\* Received by the editors February 11, 1970, and in final revised form February 4, 1971.

† Department of Mathematics, University of Michigan, Ann Arbor, Michigan 48104. This author's research was supported in part by AFOSR Research Grant 69-1662, and in part by NONR 233 (76).

‡ Department of Mathematics, Wayne State University, Detroit, Michigan 48202. This author's research was completed in the frame of AFOSR Research Project 942-65.

§ Department of Mathematics, University of California, Los Angeles, California 90024. This author's research was supported in part by the National Science Foundation under Research Grant GP-9658, and in part by NONR 233 (76).



in classes  $\Omega$  of AC vector-valued functions  $x(t) = (x^1, \dots, x^n)$ ,  $t_1 \leqq t \leqq t_2$ , satisfying constraint and boundary conditions

$$(7) \quad (t, x(t)) \in A \subset E_{n+1}, \quad t \in [t_1, t_2],$$

$$(8) \quad (t_1, x(t_1), t_2, x(t_2)) \in B \subset E_{2n+1}.$$

In previous papers Cesari [lab] has given existence theorems for optimal solutions with unbounded controls, based on convexity properties of the relevant sets, Helly's selection theorem, recent Filippov's type implicit function theorems, and ensuing lower closure theorems. In such an analysis Ascoli's selection theorem is still used on the trajectories. When the sets  $U(t, x)$  are unbounded, or are not uniformly bounded, in order to guarantee the compactness of a minimizing sequence of trajectories, growth conditions concerning  $f_0$  and  $f$  are usually requested. Cesari [1b] and Cesari, La Palm and Nishiura [2] proposed the following growth condition ( $\varepsilon$ ): Given  $\varepsilon > 0$  there is some locally  $L$ -integrable scalar function  $\psi_\varepsilon(t) \geqq 0$  such that  $|f(t, x, u)| \leqq \psi_\varepsilon(t) + \varepsilon f_0(t, x, u)$  for all  $t, x, u$ . As shown by Cesari [1b] and La Palm [3b] this condition is quite general, and contains growth conditions previously proposed, in particular, those of Tonelli [8] and McShane [4ab] for free problems.

Tonelli [8ab], McShane [4abd], and Turner [9] proved that for free problems growth conditions can be dispensed with in an arbitrary set  $G$  of points  $(t, x)$  in  $A$ , for which the only requirement is to be "slender" in the terminology of Turner [9]. For instance, Tonelli [8ab] and McShane [4ab] required  $G$  to be either finite, or countable, or contained in a countable set of straight lines parallel to the  $t$ -axis, or of curves each of the type  $\Gamma : x = x(t), t' \leqq t \leqq t'', x$  being AC on  $[t', t'']$ , and other analogous cases. McShane [4d] and Turner [9] proposed general definitions of "slender" sets, which contain all these cases.

The purpose of this paper is to show that also for general Lagrange problems growth conditions can be dispensed with in arbitrary slender subsets of  $A$ . We shall use the same definition of "slender" sets that Turner used for free problems. In § 2 we state an existence theorem for Lagrange problems (1)–(5) with a slender set  $G \subset A$  of exceptional sets. In § 3 we prove the theorem, and in § 4 we state and prove a slight variant of the same existence theorem. In § 5 we summarize a number of points concerning the detailed conditions which we are using (upper semi-continuity of sets, property (Q), growth conditions). In § 6 we derive a number of existence theorems for free problems (6)–(8) as corollaries of Cesari's previous existence theorem, the present existence theorem, and the previous analysis concerning growth conditions. In particular, we prove that Tonelli's existence theorems I–V of [8a] can be obtained as corollaries of analogous results for Lagrange problems. In § 7 we list a few examples of Lagrange and free problems satisfying the main hypotheses of the present existence theorem.

**2. Statement of the theorem.** In the formulation of the Lagrange problem (1)–(5) above we assume that  $A$  and  $B$  are given subsets of the  $tx$ -space  $E_{n+1}$  and of the  $t_1x_1t_2x_2$ -space  $E_{2n+1}$  respectively, and that for every  $(t, x) \in A$  a given subset  $U(t, x)$  of  $E_m$  is given (control space). We denote by  $M$  the set of all  $(t, x, u) \in E_{n+m+1}$  with  $(t, x) \in A, u \in U(t, x)$ , and we assume that  $f_0$  and  $f = (f_1, \dots, f_n)$  are given real functions on  $M$ . A pair  $x(t) = (x^1, \dots, x^n), u(t) = (u^1, \dots, u^m)$ ,

$t_1 \leqq t \leqq t_2$ ,  $x$  AC,  $u$  measurable, is said to be admissible provided (2), (3), (4), (5) hold, and, in addition,  $f_0(t, x(t), u(t))$  is  $L$ -integrable in  $[t_1, t_2]$ .

Let  $G$  be a fixed subset of  $A$  in the  $tx$ -space  $E_{n+1}$ ,  $G \subset A \subset E_{n+1}$ . For every set  $S$  of the  $t$ -axis and  $i = 1, \dots, n$ , we shall denote by  $G^i(S)$  the set of all real  $\xi$  such that for some  $(t, x)$  we have  $(t, x) \in G$ ,  $t \in S$ ,  $x^i = \xi$ . We shall say that  $G^i(S)$  is the "image" of  $S$  on the  $x^i$ -axis by means of the set  $G$ . Note that, when  $G$  is the graph of a curve  $x = g(t)$ ,  $t_1 \leqq t \leqq t_2$ , in the  $tx$ -space, or  $x^i = g_i(t)$ ,  $t_1 \leqq t \leqq t_2$ ,  $i = 1, \dots, n$ , then  $G^i(S) = g_i(S)$  is exactly the image of  $S$  on the  $x^i$ -axis by means of the component  $g_i$  of  $g$ . A set of  $G \subset A$  is said to be "slender" if the following property holds:

- (s) For every set  $S$  of measure zero on the  $t$ -axis, the sets  $G^i(S)$  also have measure zero on the  $x^i$ -axis,  $i = 1, \dots, n$ .

In other words,  $G$  is slender provided  $|S| = 0$  implies  $|G^i(S)| = 0$ ,  $i = 1, \dots, n$ .

As usual we denote by  $M$  the subset of all  $(t, x, u) \in E_{n+m+1}$  with  $(t, x) \in A$ ,  $u \in U(t, x)$  (graph of  $U(t, x)$ ). For every  $(\bar{t}, \bar{x}) \in A$  and  $\delta > 0$  we denote by  $N_\delta(\bar{t}, \bar{x})$ , or  $\delta$ -neighborhood of  $(\bar{t}, \bar{x})$  in  $A$ , the set of all  $(t, x) \in A$  at a distance  $\leqq \delta$  from  $(\bar{t}, \bar{x})$ . Also, we denote by  $U(\bar{t}, \bar{x}; \delta)$  the union of all  $U(t, x)$  with  $(t, x) \in N_\delta(\bar{t}, \bar{x})$ , and then the usual Kuratowski upper semicontinuity property of the sets  $U(t, x)$  at  $(\bar{t}, \bar{x})$  is expressed by the following requirement, or property  $(U)$  at  $(\bar{t}, \bar{x})$ :

$$U(\bar{t}, \bar{x}) = \bigcap_\delta \text{cl } U(\bar{t}, \bar{x}; \delta).$$

It is well known that, if  $A$  is closed, then  $U(t, x)$  has property  $(U)$  at every point  $(t, x)$  of  $A$  if and only if  $M$  is closed.

As usual, we denote by  $\tilde{Q}(t, x)$  the set of all  $\bar{z} = (z^0, z) \in E_{n+1}$  with  $z^0 \geqq f_0(t, x, u)$ ,  $z = f(t, x, u)$ ,  $u \in U(t, x)$ . We then say that the sets  $\tilde{Q}(t, x)$  satisfy property  $(Q)$  at a point  $(\bar{t}, \bar{x})$  of  $A$  provided:

$$\tilde{Q}(\bar{t}, \bar{x}) = \bigcap_\delta \text{cl co } \tilde{Q}(\bar{t}, \bar{x}; \delta).$$

Cesari showed in [1c] that this condition can be thought of as an extension for Lagrange problems of Tonelli's and McShane's property of seminormality for free problems, and we shall summarize the main concepts in § 4 below. Property  $(Q)$  has been used extensively by Cesari, Olech, Lasota, La Palm, Baum, Angell, and others. (For the comparison with a different interpretation of this property in Olech [7ab] see Cesari [1d].)

**EXISTENCE THEOREM 2.1** (For Lagrange problems (1)–(5) with unbounded controls and an exceptional slender set). *Let  $A$  be compact, and  $B$  and  $M$  closed. Let  $f_0(t, x, u)$ ,  $f(t, x, u) = (f_1, \dots, f_n)$  be continuous on  $M$ , and let us assume that the sets  $\tilde{Q}(t, x)$  are all convex, closed, and satisfy property  $(Q)$  at all points  $(t, x) \in A$  with exception perhaps of a set of points whose  $t$ -coordinate lies on a set of measure zero on the  $t$ -axis.*

*Let us assume that: (A) For every point  $(\bar{t}, \bar{x}) \in A$  there are a neighborhood  $N_\delta(\bar{t}, \bar{x})$ , real constants  $r, b = (b_1, \dots, b_n)$ , and  $v > 0$  such that  $(t, x) \in N_\delta(\bar{t}, \bar{x})$ ,  $u \in U(t, x)$  implies  $f_0(t, x, u) \geqq r + b \cdot f(t, x, u) + v|f(t, x, u)|$ . Let us assume that a closed set  $G \subset A$  is given satisfying property (s):  $|G^i(S)| = 0$ ,  $i = 1, \dots, n$ , for every set  $S$  of measure zero on the  $t$ -axis. Let us assume that the following growth condition holds: (ε) For every point  $(\bar{t}, \bar{x}) \in A - G$  there is a neighborhood  $N_\delta(t, x)$  such that given  $\varepsilon > 0$  we can determine an  $L$ -integrable function  $\psi_\varepsilon(t) \geqq 0$ ,  $\bar{t} - \delta \leqq t \leqq \bar{t} + \delta$ ,*

such that  $(t, x) \in N_\delta(\bar{t}, \bar{x})$ ,  $u \in U(t, x)$  implies  $|f(t, x, u)| \leq \psi_\varepsilon(t) + \varepsilon f_0(t, x, u)$ . If  $\Omega$  is any closed nonempty family of admissible pairs, then functional (1) has an absolute minimum in  $\Omega$ .

In § 6 we shall show that growth condition  $(\varepsilon)$  is a generalization of the growth condition used by Tonelli and McShane for free problems. In view of the terminology in common use for free problems, condition  $(\Lambda)$  can be called a “condition of normality” for the pair  $f_0, f$ . Details concerning these concepts will be given in § 5. Finally, in § 6, we shall show that Tonelli’s existence theorems I–V of [8a] can be derived from the existence theorem above for Lagrange problems, from Cesari’s previous existence theorems [1abcd], and results in [3ab] and [2]. A growth condition similar to condition  $(\varepsilon)$  above has been recently proposed also by C. Olech [7a].

**3. Proof of Existence Theorem 2.1.** We shall denote by  $|x|$  the usual absolute value for real numbers  $x$ , and  $|x| = \max[|x^i|, i = 1, \dots, n]$  for vectors  $x = (x^1, \dots, x^n)$ . If  $X$  denotes the set of all continuous vector functions  $x(t) = (x^1, \dots, x^n)$ , defined in arbitrary intervals  $a \leq t \leq b$ , we shall make  $X$  a metric space by means of the  $\rho$ -metric defined as follows. For any two elements  $x(t) = (x^1, \dots, x^n)$ ,  $a \leq t \leq b$ , and  $y(t) = (y^1, \dots, y^n)$ ,  $c \leq t \leq d$ , let  $\rho(x, y)$  be the distance function defined by

$$\rho(x, y) = |a - c| + |b - d| + \max |x(t) - y(t)|,$$

where  $\max$  is taken for  $-\infty < t < \infty$  and  $x$  and  $y$  are extended by continuity and constancy outside the original intervals  $(a, b)$  and  $(c, d)$ .

From hypothesis  $(\Lambda)$ , to each point  $(\bar{t}, \bar{x}) \in A$  we may associate an expression  $\xi = r + b \cdot f + v|f|$  and a constant  $\bar{\delta} > 0$  such that  $f_0(t, x, u) \geq \xi$  for all  $(t, x, u) \in M$  with  $|t - \bar{t}| \leq 5\bar{\delta}$ ,  $|x - \bar{x}| \leq 5\bar{\delta}$ . A finite number of the cubes  $\{(t, \bar{x}) \mid |t - \bar{t}| \leq \bar{\delta}, |x - \bar{x}| \leq \bar{\delta}\}$  cover  $A$ . Let  $\delta$  be the minimum of the numbers  $\bar{\delta}$  associated with the members of such a finite covering. Divide  $E_{n+1}$  into cubes  $S$  whose sides have length  $\delta$  by means of the hyperplanes  $t = h\delta$ ,  $x^i = l_i\delta$  for  $h, l_i = 0, \pm 1, \pm 2, \dots, i = 1, 2, \dots, n$ . Each cube  $S$  which intersects  $A$  intersects one of the cubes of the covering. Let us associate with this cube  $S$  the expression  $\xi = r + b \cdot f + v|f|$  which is associated with the cube of the finite covering. In this manner a set of cubes  $S_{hl}$  of the form  $\{(t, x) \mid (h - 1)\delta \leq t \leq h\delta, (l_i - 1)\delta \leq x^i \leq l_i\delta, i = 1, \dots, n\}$  are obtained which cover  $A$ , and with each cube  $S_{hl}$  we have associated an expression

$$\xi_{hl} = r_{hl} + b_{hl} \cdot f(t, x, u) + v_{hl}|f(t, x, u)|, \quad l = (l_1, \dots, l_n),$$

with  $f_0(t, x, u) \geq \xi_{hl}$  for all  $(t, x) \in A \cap S_{hl}$ ,  $u \in U(t, x)$ , and this is true even if we replace  $S_{hl}$  by any one of the  $3^{n+1} - 1$  adjacent cubes. The  $t$ -coordinates of the vertices of these cubes define a partition of a section of the  $t$ -axis; thus, there are two integers  $p, q$  such that each vertex of the cubes has  $t$ -coordinate of the form  $s\delta$  with  $p \leq s \leq q$ .

For every  $t_0 \in [p\delta, q\delta]$ , the set  $G^i(\{t_0\})$  has measure zero. Given  $\eta$ ,  $0 < \eta < \delta$ , for every  $i = 1, \dots, n$ , the set  $G^i(\{t_0\})$  can be covered by an open set of measure  $< \eta$ . The Cartesian product  $F_0$  of the open sets is open in the hyperplane  $H(t_0): t = t_0$ . Now the set  $(H(t_0) - F_0) \cap A$  is a compact subset of  $A$  free of points of the exceptional closed set  $G$ . By hypothesis  $(\varepsilon)$ , given  $N > 0$  integer and any point

$(t_0, x_0) \in (H(t_0) - F_0) \cap A$ , there is a number  $\rho_0 > 0$  and an  $L$ -integrable function  $\psi_{0N}(t)$ ,  $t_0 - \rho_0 \leq t \leq t_0 + \rho_0$ , such that  $|f(t, x, u)| \leq \psi_{0N}(t) + N^{-1} \cdot f_0(t, x, u)$  when  $(t, x) \in A, |t - t_0| \leq \rho_0, |x - x_0| \leq \rho_0, u \in U(t, x)$ .

The compact set  $(H(t_0) - F) \cap A$  may be covered by a finite number of these hypercubes. Let  $\bar{\rho} > 0$  be less than  $\delta$  and less than the minimum of the numbers  $\rho_0$  of the given finite covering. Thus, there is an interval  $(t_0 - \bar{\rho}, t_0 + \bar{\rho})$  such that  $|f(t, x, u)| \leq \psi_N(t) + N^{-1}f_0(t, x, u)$  when  $(t, x) \in A, (t, x)$  is in the union  $K$  of these open hypercubes, and  $|t - t_0| < \bar{\rho}$ . Here  $\psi_N(t) = \max\{\psi_{0N}(t)\}$ , this maximum being taken on the finitely many functions  $\psi_{0N}(t)$  associated with the chosen finite covering. The complement  $H$  of  $K$  in the section  $(t_0 - \bar{\rho}, t_0 + \bar{\rho})$  has the property that each of its projections on the coordinate axes  $x^i$  has measure  $< \eta$ .

In this manner we have associated an open interval of the form  $(t - \rho, t + \rho)$  with each  $t \in [p\delta, q\delta]$ , thereby obtaining an open covering of  $[p\delta, q\delta]$ . By a suitable contraction these intervals may be used to define a partition  $P: p\delta = t_0 < t_1 < \dots < t_R = q\delta$  of  $[p\delta, q\delta]$ , and it may be assumed without loss of generality that the points  $s\delta$  for  $p \leq s \leq q$  are used in this partition. Now let us refine the previous partition of  $A$  into parts  $S_{ml}$  by means of the hyperplanes  $t = t_j, j = 1, \dots, R$ . The new parts are intervals, say  $Q_{jl}$ ; we shall still call them cubes, for the sake of simplicity. These parts  $Q_{jl}$  are of the form

$$Q_{jl} = \{(t, x) | t_{j-1} \leq t \leq t_j, (l_i - 1)\delta \leq x^i \leq l_i\delta, i = 1, \dots, n\},$$

$$j = 1, \dots, R, \quad t_j - t_{j-1} \leq \delta, \quad l = (l_1, \dots, l_n).$$

Let  $\xi_{jl} = r_{jl} + b_{jl} \cdot f + v_{jl}|f|$  denote the function associated with the cube of the former partition which contains  $Q_{jl}$ .

Summarizing, the following type of partition may be obtained: Given  $\eta > 0$  and  $N > 0$ , there are expressions  $\xi_{jl}$  as above, and a partition of  $A$  into cubes  $Q_{jl}$  as above, whose edges in the  $x^i$  direction have length  $\delta > 0$  independent of  $\eta$  and  $N$ , such that  $f_0(t, x, u) \geq \xi_{jl}$  for all  $(t, x) \in Q_{jl} \cap A$ , or  $(t, x) \in A$  in any of the  $3^{n+1} - 1$  cubes adjacent to  $Q_{jl}$ , and  $|f(t, x, u)| \leq \psi_N(t) + N^{-1}f_0(t, x, u)$  for all  $(t, x)$  with  $t_{j-1} \leq t \leq t_j$ , except for a set  $H_j$  (made up of cubes  $Q_{jl}$ ) whose projections on the  $x^i$ -axes have measures  $< \eta, i = 1, \dots, n$ . Moreover, the constants  $r_{jl}$  and  $b_{jl}$  are independent of  $\eta$  and  $N$ .

Let  $r = \max|r_{jl}|, b = \max|b_{jl}|, v = \min v_{jl}$ , and take  $\eta < \delta/2, N > 2b \cdot (1 + 4\sqrt{n+1})$ .

Let  $\Omega^*$  denote the class of all "admissible pairs"  $x(t), u(t), a \leq t \leq b$ , defined on an arbitrary interval  $[a, b]$ , in the sense that we require  $x$  to be AC in  $[a, b], u$  to be measurable in  $[a, b]$ , we require  $(t, x(t)) \in A$  for all  $t \in [a, b], u(t) \in U(t, x(t))$  and  $x'(t) = f(t, x(t), u(t))$  a.e. in  $[a, b]$ , and we require  $f_0(t, x(t), u(t))$  to be  $L$ -integrable in  $[a, b]$ . Thus, we do not require any boundary condition for inclusion in  $\Omega^*$ , and thus  $\Omega^* \supset \Omega$ . Let  $C_j: x = x(t), t_{j-1} \leq t \leq t_j$ , be the part of the trajectory  $x$  (if any) defined on  $[t_{j-1}, t_j]$ . Divide  $C_j$  into more subarcs  $C_{j1}, \dots, C_{jT_j}$  as follows: The first endpoint of  $C_{1j}$  is  $x(t_{j-1})$  (or  $x(a)$  if  $t_{j-1} < a < t_j$ ); the second endpoint is either the first point where  $C_j$  leaves one of the  $3^n - 1$  cubes in the section  $\{(t, x) | (t, x) \in A, t_{j-1} \leq t \leq t_j\}$  adjacent to the cube containing  $x(t_{j-1})$ , or  $x(t_j)$  if  $C_j$  does not leave these  $3^n - 1$  cubes (or  $x(b)$  if  $t_{j-1} < b < t_j$ ). The first point of  $C_{j2}$  is the endpoint of  $C_{j1}$ , and the endpoint of  $C_{j2}$  is either the first point of  $C_j$

which leaves the  $3^n - 1$  cubes adjacent to the cube containing the first endpoint of  $C_{j_2}$  or  $x(t_j)$  if  $C_j$  does not leave these cubes (or  $x(b)$  if  $t_{j-1} < b < t_j$ ). Continuing in this manner,  $C_j$  is broken up into arcs  $C_{j\sigma}$ ,  $\sigma = 1, \dots, T_j$ . This process must terminate after a finite number of steps, since each arc  $C_{j\sigma}$  except the last has length  $> \delta$ .

Let  $\Lambda_{j\sigma}$  be the set of all  $t$  in the domain of  $C_{j\sigma}$ , where  $x(t) \in H_j$ ; let  $\Lambda'_{j\sigma}$  be the complement of  $\Lambda_{j\sigma}$  in this domain. Let  $\lambda_{j\sigma} = \int_{\Lambda_{j\sigma}} |x'(t)| dt$ , and  $\lambda'_{j\sigma} = \int_{\Lambda'_{j\sigma}} |x'(t)| dt$ . Let the initial point of  $C_{j\sigma}$  be in  $Q_{jlk}$ . Then

$$\begin{aligned} I[x, u] &\geq \sum_{j=1}^R \sum_{\sigma=1}^{T_j} \left\{ \int_{\Lambda_{j\sigma}} [r_{jl\sigma} + b_{jl\sigma} \cdot x'(t) + v_{j\sigma}|x'(t)|] dt \right. \\ &\quad \left. + \int_{\Lambda'_{j\sigma}} f_0(t, x(t), u(t)) dt \right\} \\ &\geq \sum_{j=1}^R \sum_{\sigma=1}^{T_j} \left\{ -r|\Lambda_{j\sigma}| - b \left| \int_{\Lambda_{j\sigma}} x'(t) dt \right| + v_{j\sigma}\lambda_{j\sigma} \right. \\ &\quad \left. + \int_{\Lambda'_{j\sigma}} f_0(t, x(t), u(t)) dt \right\}. \end{aligned}$$

For  $t \in \Lambda'_{j\sigma}$  we have  $|f(t, x, u)| \leq \psi_N(t) + N^{-1}f_0(t, x, u)$  for all  $u \in U(t, x(t))$ . Then

$$\begin{aligned} I[x, u] &\geq -r(b - a) + \sum_{j=1}^R \sum_{\sigma=1}^{T_j} \left\{ -b \left| \int_{\Lambda_{j\sigma}} x'(t) dt \right| + v\lambda_{j\sigma} \right. \\ &\quad \left. + \int_{\Lambda_{j\sigma}} [f_0(t, x(t), u(t)) + N\psi_N(t)] dt - \int_{\Lambda_{j\sigma}} N\psi_N(t) dt \right\} \\ &\geq -r(b - a) - \int_a^b N\psi_N(t) dt \\ &\quad + \sum_{j=1}^R \sum_{\sigma=1}^{T_j} \left\{ -b \left| \int_{\Lambda_{j\sigma}} x'(t) dt \right| + v\lambda_{j\sigma} + N\lambda'_{j\sigma} \right\}. \end{aligned}$$

Now  $|\left(\int_{\Lambda_{j\sigma}} + \int_{\Lambda'_{j\sigma}}\right)x'(t) dt| \leq 2\delta\sqrt{n + 1}$ , and thus

$$\begin{aligned} \left| \int_{\Lambda_{j\sigma}} x'(t) dt \right| - \left| \int_{\Lambda'_{j\sigma}} x'(t) dt \right| &\leq 2\delta\sqrt{n + 1}, \\ \left| \int_{\Lambda_{j\sigma}} x'(t) dt \right| &\leq 2\delta\sqrt{n + 1} + \lambda'_{j\sigma} \quad \text{for all } j \text{ and } \sigma. \end{aligned}$$

Also,  $\lambda'_{j\sigma} > \delta - \eta > \delta - \delta/2 = \delta/2$  for  $\sigma = 1, \dots, T_j - 1$ . Let

$$\lambda = \sum_{j=1}^R \sum_{\sigma=1}^{T_j} \lambda_{j\sigma}, \quad \lambda' = \sum_{j=1}^R \sum_{\sigma=1}^{T_j} \lambda'_{j\sigma},$$

and let  $D$  be the diameter of  $A$ . Then

$$\begin{aligned}
 I[x, u] &> -r(b - a) - \int_a^b N\psi_N(t) dt + v\lambda + N\lambda' + \sum_{j=1}^R \sum_{\sigma=1}^{T_j} \{-b(\lambda'_{j\sigma} + 2\delta\sqrt{n+1})\} \\
 &\geq -rD - \int_{p\delta}^{q\delta} N\psi_N(t) dt + v\lambda + N\lambda' - b \sum_{j=1}^R \sum_{\sigma=1}^{T_j-1} (\lambda'_{j\sigma} + 4\sqrt{n+1}\lambda'_{j\sigma}) \\
 &\quad - b \sum_{j=1}^R (\lambda'_{jT_j} + 2\delta\sqrt{n+1}) \\
 &\geq -rD - \int_{p\delta}^{q\delta} N\psi_N(t) dt - 2\delta bR\sqrt{n+1} + v\lambda + [N - b(1 + 4\sqrt{n+1})]\lambda' \\
 &\geq -rD - \int_{p\delta}^{q\delta} N\psi_N(t) dt - 2\delta bR\sqrt{n+1} + v\lambda + 2^{-1}N\lambda'.
 \end{aligned}$$

Let  $V = \min(v, N/2) > 0$ . Then

$$I[x; u] > -rD - \int_{p\delta}^{q\delta} N\psi_N(t) dt - 2\delta bR\sqrt{n+1} + V \int_a^b |x'(t)| dt.$$

Thus, given any constant  $M_0 > 0$ , for any admissible pair  $x, u$  with  $I[x, u] \leq M_0$ , the trajectory  $x$  has uniformly bounded total variation, and then uniformly bounded length. Moreover, this inequality shows that

$$Z = \inf I[x, u] \geq -rD - \int_{p\delta}^{q\delta} N\psi_N(t) dt - 2\delta bR\sqrt{n+1},$$

where inf is taken in the class  $\Omega^* \supset \Omega$ .

It will now be proved that for the same admissible pairs  $x(t), u(t), a \leq t \leq b$ , with  $I[x, u] \leq M_0$ , the trajectories  $x$  are equicontinuous.

If they are not, then there is an  $\varepsilon > 0$  such that for every positive integer  $k$  there is an admissible pair  $x_k(t), u_k(t), a_k \leq t \leq b_k$ , and two points  $t_{k1}, t_{k2} \in [a_k, b_k]$ , such that  $0 < t_{k2} - t_{k1} < k^{-1}, I[x_k, u_k] \leq M_0$ , and  $|x_k(t_{k2}) - x_k(t_{k1})| > \varepsilon$ . Suppose without loss of generality that  $t_{k1} \rightarrow t_0, t_{k2} \rightarrow t_0, x_k(t_{k1}) \rightarrow x_1, x_k(t_{k2}) \rightarrow x_2$  as  $k \rightarrow \infty$ . Then  $|x_2 - x_1| \geq \varepsilon$ . Let  $L$  be a bound for the lengths of the trajectories  $x_k$ .

The sets  $G^i(\{t_0\})$  have measure zero. Hence, they may be covered by open sets  $F_i$  of measure  $< \eta, 0 < \eta \leq \varepsilon/4$ . Let  $F$  denote the set  $F = \{(t_0, x) | x^i \in F_i, i = 1, \dots, n\}$ . Then  $F$  is open in the hyperplane  $H(t_0): t = t_0$ . Let  $N = (4/\varepsilon)[M_0 + |2Z - 2r\delta - bL - 1|]$ , where  $r, \delta$  and  $b$  are the constants defined above. The set  $(H(t_0) - F) \cap A$  is compact, and for every point  $(t_0, \bar{x}) \in (H(t_0) - F) \cap A$  there is a  $\bar{\rho} > 0$  and an  $L$ -integrable function  $\psi_N(t)$  such that  $|f(t, x, u)| \leq \psi_N(t) + N^{-1}f_0(t, x, u)$  when  $|t - t_0| < \bar{\rho}$ . A finite number of these intervals  $\{(t, x) | |t - t_0| < \bar{\rho}, |x - \bar{x}| < \bar{\rho}\}$  covers  $(H(t_0) - F) \cap A$ . Let  $\rho$  be the minimum  $\bar{\rho}$  for such a covering.

Divide the curve  $C_k: x = x_k(t), a_k \leq t \leq b_k$ , into three parts:  $C_{k1}, C_{k2}, C_{k3}$  according as  $a_k \leq t \leq t_{k1}, t_{k1} \leq t \leq t_{k2}, t_{k2} \leq t \leq b_k$ , respectively. Divide the interval  $[t_{k1}, t_{k2}]$  into two subsets, say  $E_2 = [t, x(t) \in F], E_1 = [t_{k1}, t_{k2}] - E_2$ .

Then, for some  $k_0$  and all  $k > k_0$  we have  $|t_{k1} - t_0| < \rho$ ,  $|t_{k2} - t_0| < \rho$ , and  $|x_k(t_{k1}) - x_k(t_{k2})| > \varepsilon/2$ . Then for  $k > k_0$  we have also

$$\begin{aligned} I[x_k, u_k] &= I_{k1} + I_{k2} + I_{k3} \geq 2Z + I_{k2} \\ &\geq 2Z + \left( \int_{E_1} + \int_{E_2} \right) f_0(t, x_k(t), u_k(t)) dt \\ &\geq 2Z + \int_{E_1} (-N\psi_N(t) + N|f(t, x_k(t), u_k(t))|) dt \\ &\quad + \int_{E_2} [-r - b|x'_k(t)|] dt \\ &\geq 2Z - 2r\delta - bL + N \int_{E_1} |x'_k(t)| dt - \int_{t_{1k}}^{t_{2k}} N\psi_N(t) dt. \end{aligned}$$

Now

$$\begin{aligned} \int_{E_2} |x'_k(t)| dt &< \eta, \\ \int_{E_1} |x'_k(t)| dt &= \int_{t_{1k}}^{t_{2k}} |x'_k(t)| dt - \int_{E_2} |x'_k(t)| dt \\ &\geq (\varepsilon/2) - \eta \geq \varepsilon/4. \end{aligned}$$

On the other hand, for  $k$  sufficiently large, the integral of  $N\psi_N(t)$  over  $[t_{1k}, t_{2k}]$  is certainly  $< 1$  since  $t_{2k} - t_{1k} \rightarrow 0$ . Finally,

$$\begin{aligned} I[x_k, u_k] &> 2Z - 2r\delta - bL - 1 + N\varepsilon/4 \\ &= 2Z - 2r\delta - bL - 1 + [M_0 + |2Z - 2r\delta - bL - 1|] \geq M_0, \end{aligned}$$

a contradiction. Thus, for admissible pairs  $x, u$  with  $I[x, u] \leq M$ , the trajectories  $x$  are equicontinuous.

Let  $i$  be the infimum of  $I[x, u]$  in the original family  $\Omega$  of admissible pairs  $x(t), u(t)$ ,  $t_1 \leq t \leq t_2$ . Thus, for these pairs, the trajectories satisfy the required boundary conditions  $(t_1, x(t_1), t_2, x(t_2)) \in B$ , together with  $(t, x(t)) \in A$  for all  $t \in [t_1, t_2]$ . Since  $\Omega \subset \Omega^*$ , we have  $i \geq Z > -\infty$ ; thus,  $i$  is finite. Let  $x_k(t), u_k(t)$ ,  $t_{1k} \leq t \leq t_{2k}$ ,  $k = 1, 2, \dots$ , be a minimizing sequence of admissible pairs, all in  $\Omega$ , with  $I[x_k, u_k] < i + k^{-1}$ ,  $k = 1, 2, \dots$ . Then the trajectories  $x_k$  are equibounded and equicontinuous. By a suitable extraction we obtain a subsequence, which we still call  $[k]$ , such that  $x_k$  converges in the  $\rho$ -metric toward a continuous vector function  $x(t)$ ,  $t_1 \leq t \leq t_2$ , and  $(t, x(t)) \in A$  for all  $t_1 \leq t \leq t_2$ , and  $(t_1, x(t_1), t_2, x(t_2)) \in B$ . Also, the curve  $C: x = x(t)$ ,  $t_1 \leq t \leq t_2$ , has finite length  $l(C)$ , since  $l(C) \leq \liminf l(C_k) \leq L$  as  $k \rightarrow \infty$ .

Let us prove that  $x$  is AC. Suppose that this is not the case. Let  $s$ ,  $0 \leq s \leq l = l(C)$ , denote the usual arc length parameter along  $C$ , thought of as a path curve in the  $tx$ -space  $E_{n+1}$ . Note that given any measurable set  $E \subset [a, b]$ , the usual Lebesgue measure  $|E|$  of  $E$  is the infimum of the numbers  $\sum_i (\beta_i - \alpha_i)$  for any countable covering  $(\alpha_i, \beta_i)$ ,  $i = 1, 2, \dots$ , of  $E$ . Analogously, we can define another measure, or length measure  $l(E)$  of  $E$ , by taking the infimum of the numbers  $\sum_i (s(\beta_i) - s(\alpha_i))$  for all the same open coverings of  $E$ . Obviously,  $|E| \leq l(E)$ . If  $x$

is not AC then there is some set  $E$  of Lebesgue measure zero on  $[a, b]$  which has positive length measure, or  $|E| = 0, l(E) = \lambda > 0$ . Now the  $n$  sets  $G^i(E), i = 1, \dots, n$ , have all zero Lebesgue measure. Thus, if  $P = \{(t, x)|t \in E, x = x(t)\}$ , then  $P \cap G$  has projection of zero Lebesgue measure on each coordinate axis. As a consequence, there is some closed set  $E' \subset E$  such that  $|E'| = 0, l(E') > \lambda/2$ , and  $(t, x(t)) \notin G$  for  $t \in E'$ .

Let  $N = (2/\lambda)(|i| + 2 + rD + bL)$ , where  $L$  is a bound on the lengths of the trajectories  $x$  in  $\Omega$ . Then, there is some  $\rho > 0$  and an  $L$ -integrable function  $\psi_N(t) \geq 0$  such that  $|f(t, x, u)| \leq \psi_N(t) + N^{-1}f_0(t, x, u)$  when  $(t, x)$  has a distance  $< \rho$  from the compact set  $P' = \{(t, x)|t \in E', x = x(t)\}$ . Since  $E'$  is compact,  $|E'| = 0$ , it may be covered by a finite set of open intervals  $(\alpha_j, \beta_j), j = 1, \dots, R$ , such that if  $F = \cup_{j=1}^R (\alpha_j, \beta_j)$ , we have  $\int_F N \psi_N(t) dt < 1$ , and  $x$  maps  $F$  into the  $\rho$ -neighborhood of  $P'$ . Let  $k_0$  be such that  $\int_F |x'_k(t)| dt > \lambda/2$  for all  $k \geq k_0$ . Finally,

$$\begin{aligned} I[x_k, u_k] &= \int_{t_{1k}}^{t_{2k}} f_0(t, x_k(t), u_k(t)) dt \\ &\geq -rD - bL + \int_F f_0(t, x_k(t), u_k(t)) dt \\ &\geq -rD - bL + \int_F (-N\psi_N(t) + N|x'_k(t)|) dt \\ &\geq -rD - bL - 1 + N\lambda/2 \\ &= -rD - bL - 1 + (|i| + 2 + rD + bL) = |i| + 1 \geq i + 1, \end{aligned}$$

a contradiction. We have proved that  $x$  is AC. By the lower closure theorem (see, for instance, Cesari [1d]) we know that there is a measurable function  $u(t), t_1 \leq t \leq t_2$ , such that the pair  $x(t), u(t), t_1 \leq t \leq t_2$ , is admissible, and

$$I[x, u] \leq \liminf I[x_k, u_k] = i.$$

Here  $(x, u)$  belongs to  $\Omega$  since  $\Omega$  is closed; thus  $I[x, u] \geq i$ , and hence  $I[x, u] = i$ . This proves Existence Theorem 2.1.

*Remark 1.* In the proof of Existence Theorem 2.1, we actually prove that there is a positive constant  $m_0$  and a real number  $m_1$  such that  $I[x, u] \geq m_0 BV(x) + m_1$ , where  $BV(x)$  is the total variation of  $x$ . Hence, the trajectories of a minimizing sequence for  $I$  in  $\Omega$  have uniformly bounded variation. The proof of this inequality is similar to the one given by Turner [9] for free problems in  $E_n$ . Also, McShane [4b] noticed an analogous inequality for free problems for  $n = 1$ .

*Remark 2.* Condition (A) in Existence Theorem 2.1 can be replaced by the following one: ( $\bar{A}$ ) For every point  $(\bar{t}, \bar{x}) \in A$  there are a neighborhood  $N_\delta(\bar{t}, \bar{x})$ , real constants  $r, b = (b_1, \dots, b_n)$ , and  $v > 0$ , and a point  $\bar{z} \in E_n$  such that  $(t, x) \in N_\delta(\bar{t}, \bar{x}), u \in U(t, x)$  implies  $f_0(t, x, u) \geq r + b \cdot f(t, x, u) + v|f(t, x, u) - \bar{z}|$ . The proof is a simple modification of the proof given above. In particular, we can take  $\bar{z} = f(\bar{t}, \bar{x}, \bar{u})$  for some  $\bar{u} \in U(\bar{t}, \bar{x})$ .



**4. A variant of Existence Theorem 2.1.**

EXISTENCE THEOREM 4.1. *This is the same as Theorem 2.1 with condition (A) satisfied only at the points  $(\bar{i}, \bar{x})$  of the closed set  $G$ .*

*Proof.* Let us perform the initial subdivision in cubes  $Q_{hl}$  as at the beginning of § 3 by taking  $\delta > 0$  so small that the following holds:

(a) If  $Q_{hl} \cap G \neq \emptyset$ , then  $Q_{hl}$  and the  $3^{n+1} - 1$  adjacent cubes  $Q$  form a set  $F$  with  $F \cap A \subset N_\delta$ , where  $N_\delta$  denotes one of the neighborhoods of property (A). (b) If  $Q_{hl}$  is any one of the cubes  $Q$  with  $Q_{hl} \cap A \neq \emptyset$  and not covered by any of the sets  $F$  above, then  $Q_{hl}$  and the  $3^{n+1} - 1$  adjacent cubes  $Q$  form a set  $F$  with  $F \cap A \subset N_\delta$ , where  $N_\delta$  denotes one of the neighborhoods of property ( $\epsilon$ ).

We may denote by  $Q_1, \dots, Q_\nu$ , the cubes  $Q_{hl}$  of the first kind, and by  $Q'_1, \dots, Q'_\mu$ , the cubes  $Q_{hl}$  of the second kind. Finally, for every curve  $C: x = x(t)$ ,  $a \leqq t \leqq b$ , with  $(t, x(t)) \in A$  there is a subdivision into finitely many arcs  $C_1, \dots, C_\nu$ ,  $C'_1, \dots, C'_\mu$ , each  $C_1, \dots, C_\nu$  contained in a cube of the first kind, and each  $C'_1, \dots, C'_\mu$  in a cube of the second kind. Each of these arcs has diameter  $> \delta$ . We shall now apply to the arcs  $C_1, \dots, C_\nu$  the arguments of § 3, and to the arcs  $C'_1, \dots, C'_\mu$ , the usual arguments of [1a, pp. 391–395], and [1b, Remark 10, p. 539]. These arguments show that for  $I[x, u] \leqq M_0$  the curves  $C: x = x(t)$ ,  $t_1 \leqq t \leqq t_2$ , have total length  $l(C) \leqq L$  for some fixed  $L$ , and this yields bounds for the numbers  $\nu$  and  $\mu$ . The concluding argument of the proof is then the same as in § 3.

*Remark.* In Existence Theorem 4.1, condition (A) at the points of  $G$  can be replaced by a condition more geometric in character which we shall state in § 5 after Theorem 5.5.

**5. Remarks.** (a) Given any two functions  $g(t, x, u)$ ,  $f_0(t, x, u)$ ,  $f_0$  scalar, defined and continuous on the same set  $M$  as in § 2, we say that  $g(t, x, u)$  is “of slower growth than  $f_0(t, x, u)$  with respect to  $u$ ” in a neighborhood  $N_\delta(\bar{i}, \bar{x})$  of  $(\bar{i}, \bar{x}) \in A$ , provided, given  $\epsilon > 0$ , there is some  $\bar{u} = \bar{u}(\delta, \epsilon) \geqq 0$  such that  $(t, x, u) \in M$ ,  $(t, x) \in N_\delta(t, x)$ ,  $|u| \geqq \bar{u}$ , implies  $|g(t, x, u)| \leqq \epsilon f_0(t, x, u)$ . We proved in [1b, (2.1)(i), (2.2) (ii)] and in [1c, (2.i)] the following criterion for property (Q).

CRITERION 5.1. *If  $g$  and  $f(t, x, u)$  are of slower growth than  $f_0(t, x, u)$  with respect to  $u$  in a neighborhood  $N_\delta(\bar{i}, \bar{x})$  of  $(\bar{i}, \bar{x})$  in  $A$ , then the sets  $\tilde{Q}(t, x)$ , if convex, satisfy property (Q) at  $(\bar{i}, \bar{x})$ .*

(b) Let  $f(t, x, u) = (f_1, \dots, f_n)$  and  $f_0(t, x, u)$  be continuous and defined on  $M$  as in § 2. As usual, we denote by  $Q(t, x) \subset E_n$  the sets defined by  $Q(t, x) = f(t, x, U(t, x)) = \{z | z = f(t, x, u), u \in U(t, x)\} \subset E_n$ . As usual we denote by  $R$  the linear manifold of  $E_n$  of minimum dimension  $r$ ,  $0 \leqq r \leqq n$ , containing  $Q(t, x)$ , and we denote by  $\text{Rint } Q(t, x)$  the set of all points  $z \in Q(t, x)$  which are interior to  $Q(t, x)$  with respect to  $R$ . We say that  $f(t, x, u) = (f_1, \dots, f_n)$  and  $f_n(t, x, u)$  satisfy property ( $\alpha$ ) at a point  $(\bar{i}, \bar{x}) \in A$  provided:

$$(\alpha) \text{ If } (z^0, z) \in \bigcap_\delta \text{cl co } \tilde{Q}(\bar{i}, \bar{x}; \delta), \text{ then } z \in Q(\bar{i}, \bar{x}).$$

As we noticed in [1c] this condition is necessary for the sets  $\tilde{Q}(t, x)$  to satisfy property (Q) at  $(\bar{i}, \bar{x})$ .

For free problems (that is, when  $m = n$ ,  $f = u$ ,  $U = E_n$ ; see § 5 below), we have  $Q(t, x) = E_n$  for all  $(t, x) \in A$ , and then condition ( $\alpha$ ) is trivially satisfied. We

consider again continuous functions  $f(t, x, u) = (f_1, \dots, f_n)$  and  $f_0(t, x, u)$  scalar as in § 2. We say that  $f$  and  $f_0$  satisfy property (X) at the point  $(\bar{t}, \bar{x}) \in A$  provided, for every  $\bar{z} \in Q(\bar{t}, \bar{x})$  there is at least one point  $\bar{u} \in U(\bar{t}, \bar{x})$  such that (i)  $\bar{z} = f(\bar{t}, \bar{x}, \bar{u})$ , and (ii) given  $\varepsilon > 0$  there are numbers  $\delta > 0$  and  $r, b = (b_1, \dots, b_n)$  such that

$$(X') f_0(t, x, u) \geq r + b \cdot f(t, x, u) \quad \text{for all } (t, x) \in N_\delta(\bar{t}, \bar{x}) \quad \text{and } u \in U(t, x);$$

$$(X'') f_0(\bar{t}, \bar{x}, \bar{u}) \leq r + b \cdot f(\bar{t}, \bar{x}, \bar{u}) + \varepsilon.$$

We say that  $f$  and  $f_0$  satisfy property (X\*) at the point  $(\bar{t}, \bar{x}) \in A$  provided, for every  $\bar{z} \in Q(\bar{t}, \bar{x})$ , there is at least one point  $\bar{u} \in U(\bar{t}, \bar{x})$  such that (i)  $\bar{z} = f(\bar{t}, \bar{x}, \bar{u})$ , and (ii) given  $\varepsilon > 0$  there are numbers  $\delta > 0, \nu > 0$ , and  $r, b = (b_1, \dots, b_n)$  real such that

$$(X'*) f_0(t, x, u) \geq r + b \cdot f(t, x, u) + \nu |f(t, x, u) - f(\bar{t}, \bar{x}, \bar{u})|$$

$$\text{for all } (t, x) \in N_\delta(\bar{t}, \bar{x}) \quad \text{and } u \in U(t, x);$$

$$(X'') f_0(\bar{t}, \bar{x}, \bar{u}) \leq r + b \cdot \bar{f}(\bar{t}, \bar{x}, \bar{u}) + \varepsilon.$$

Here, by  $b \cdot f$  we mean as usual  $b_1 f_1 + \dots + b_n f_n$ .

For free problems, that is,  $m = n, f = u, U = E_n$ , condition (X) reduces to the following condition which concerns only the function  $f_0$ , or condition (X<sub>f</sub>): For every  $\bar{u} = (\bar{u}^1, \dots, \bar{u}^n) \in E_n$  and  $\varepsilon > 0$  there are numbers  $\delta > 0$  and  $r, b = (b_1, \dots, b_n)$  real such that

$$(X'_f) f_0(t, x, u) \geq r + b \cdot u \quad \text{for all } (t, x) \in N_\delta(\bar{t}, \bar{x}) \quad \text{and all } u \in E_n,$$

$$(X''_f) f_0(\bar{t}, \bar{x}, \bar{u}) \leq r + b \cdot \bar{u} + \varepsilon.$$

When this condition is satisfied, we say that  $f_0(t, x, u), (t, x) \in A \times E_n$ , is *semi-normal* at  $(\bar{t}, \bar{x}) \in A$ . This is an equivalent formulation of definitions due to Tonelli [8ab] and McShane [4]. Again, for free problems, condition (X\*) reduces to the following condition concerning  $f_0$ , or condition (SN): For every  $\bar{u} = (\bar{u}^1, \dots, \bar{u}^n) \in E_n$  and  $\varepsilon > 0$ , there are numbers  $\delta > 0, \nu > 0$ , and  $r, b = (b_1, \dots, b_n)$  real such that

$$(SN') f_0(t, x, u) \geq r + b \cdot u + \nu |u - \bar{u}| \quad \text{for all } (t, x) \in N_\delta(\bar{t}, \bar{x}) \quad \text{and } u \in E_n;$$

$$(SN'') f_0(\bar{t}, \bar{x}, \bar{u}) \leq r + b \cdot \bar{u} + \varepsilon.$$

When this condition is satisfied, we say that  $f_0(t, x, u), (t, x) \in A \times E_n$ , is *normal* at  $(\bar{t}, \bar{x}) \in A$ . This is an equivalent formulation of definitions due to Tonelli [8ab] and McShane [4b].

We see, therefore, that conditions (X) and (X\*) above appear as natural extensions to Lagrange problems of the usual conditions of seminormality and normality respectively for free problems.

We refer in [1c, (4.i)] the following second criterion for property (Q) of the sets  $\tilde{Q}(t, x)$ .

CRITERION 5.2. *If conditions (α) and (X) hold at the point  $(\bar{t}, \bar{x}) \in A$ , then the set  $\tilde{Q}(\bar{t}, \bar{x})$  is closed and convex, and the sets  $\tilde{Q}(t, x)$  satisfy property (Q) at  $(\bar{t}, \bar{x})$ .*

(c) The case in which  $f$  is linear in  $u$ , or  $f(t, x, u) = B(t, x)u + C(t, x)$  ( $B, C$  matrices of dimensions  $n \times m, n \times 1$  with continuous entries), is of interest. The

free problems are of this kind, with  $m = n, f = u, B$  the identity matrix of order  $n$ , and  $C = 0$ . We proved in [1c, (7.i)] the following criterion for property (Q).

CRITERION 5.3. *If  $A$  is closed,  $U = E_m, M = A \times E_m, f_0$  is continuous on  $M$ , convex in  $u$ , and seminormal in  $u$ , at a point  $(\bar{t}, \bar{x}) \in A$ , and if  $f(t, x, u) = B(t, x)u + C(t, x)$ , then the sets  $\tilde{Q}(t, x)$  satisfy property (Q) at  $(\bar{t}, \bar{x})$ .*

(d) Given  $f(t, x, u) = (f_1, \dots, f_n)$  and  $f_0(t, x, u)$  continuous on  $M$  as usual, and for any point  $z \in Q(t, x) = f(t, x, U(t, x)) \subset E_n$ , let us define  $T(z; t, x)$  by taking

$$\begin{aligned} T(z; t, x) &= \inf [z^0 | (z^0, z) \in \tilde{Q}(t, x)] \\ &= \inf [z^0 | z^0 \geq f_0(t, x, u), z = f(t, x, u) \text{ for some } u \in U(t, x)]. \end{aligned}$$

Thus, for every  $(\bar{t}, \bar{x}) \in A$ , the scalar function  $T(z; \bar{t}, \bar{x}), z \in Q(\bar{t}, \bar{x})$ , is defined, with  $-\infty \leq T(z; \bar{t}, \bar{x}) < +\infty$ . It was proved in [1c, (8.i)] that, if  $\tilde{Q}(t, x)$  is convex, then either  $T(z; \bar{t}, \bar{x}) = -\infty$  for all  $z \in \text{Rint } Q(\bar{t}, \bar{x})$ , or  $T(z; \bar{t}, \bar{x}) > -\infty$  for all  $z \in Q(\bar{t}, \bar{x})$ . In the latter case,  $T(z; \bar{t}, \bar{x})$  is finite everywhere and a convex function of  $z$  in  $Q(\bar{t}, \bar{x})$ ,  $T(z; \bar{t}, \bar{x})$  is bounded below on every bounded subset of  $Q(\bar{t}, \bar{x})$ , and  $T(z; \bar{t}, \bar{x})$  is continuous in the convex set  $\text{Rint } Q(\bar{t}, \bar{x})$ . Finally, if  $\tilde{Q}(\bar{t}, \bar{x})$  is convex and closed, and  $T(z; \bar{t}, \bar{x}) > -\infty$  in  $Q(\bar{t}, \bar{x})$ , then  $T(z; \bar{t}, \bar{x})$  is a lower semicontinuous function of  $z$  at every point  $z \in Q(\bar{t}, \bar{x}) - \text{Rint } Q(\bar{t}, \bar{x})$ . We proved in [1c, (9.i)] and [1e, (4.vii)] the following statements.

THEOREM 5.4. *If  $T(z; \bar{t}, \bar{x}) > -\infty$  in  $Q(\bar{t}, \bar{x})$ , then the sets  $\tilde{Q}(t, x)$  have property (Q) at  $(\bar{t}, \bar{x})$  if and only if properties (α) and (X) hold at the point  $(\bar{t}, \bar{x})$ .*

THEOREM 5.5. *If  $T(z; \bar{t}, \bar{x}) > -\infty$  in  $Q(\bar{t}, \bar{x})$  and if the sets  $\tilde{Q}(t, x)$  contain no straight line, then the sets  $\tilde{Q}(t, x)$  have property (Q) at  $(\bar{t}, \bar{x})$  if and only if properties (α) and (X\*) hold.*

As a consequence of Theorem 5.5 we can now replace condition (Λ) in Existence Theorem 4.1 by the following condition more geometric in character: (Λ') At every point  $(\bar{t}, \bar{x}) \in G$  the sets  $\tilde{Q}(\bar{t}, \bar{x})$  satisfy property (Q),  $T(z; \bar{t}, \bar{x}) > -\infty$ , and  $\tilde{Q}(\bar{t}, \bar{x})$  contains no straight line.

Indeed, condition (Λ') implies that both conditions (α) and (X\*) hold at the same points, and the latter implies (Λ) for  $\bar{z} = 0$ .

Theorem 5.5 reduces for free problems ( $m = n, f = u, U = E_n$ ) to the following well-known statement:

THEOREM 5.6. *If  $f_0(t, x, u)$  is continuous in  $A \times E_n$ , then  $f_0$  is normal in  $u$  at  $(\bar{t}, \bar{x})$  if and only if  $f_0(\bar{t}, \bar{x}, u)$  is convex in  $u$ , and for no  $\bar{u}, u_1 \in E_n, u_1 \neq 0$ , it occurs that  $f_0(\bar{t}, \bar{x}, \bar{u}) = 2^{-1}[f_0(\bar{t}, \bar{x}, \bar{u} + \lambda u_1) + f_0(\bar{t}, \bar{x}, \bar{u} - \lambda u_1)]$  for all  $\lambda \geq 0$ .*

This last statement was proved by Tonelli [8a] under smoothness hypotheses, and then extended by Turner [9] under the sole continuity hypotheses above.

**6. Free problems.** We consider now free problems, that is, Lagrange problems with  $m = n, f = u, U = E_n$ , as stated in § 1. Thus, we are interested in problems concerning the minimum of an integral

$$(6) \quad I[x] = \int_{t_1}^{t_2} f_0(t, x(t), x'(t)) dt$$

in classes  $\Omega$  of admissible AC vector-valued functions  $x(t) = (x^1, \dots, x^n)$ ,

$t_1 \leqq t \leqq t_2$ , satisfying constraint and boundary conditions

$$(7) \quad (t, x(t)) \in A \subset E_{n+1}, \quad t \in [t_1, t_2],$$

$$(8) \quad (t_1, x(t_1), t_2, x(t_2)) \in B \subset E_{2n+1}.$$

The only requirement for admissibility (§ 3), which has not yet been restated, is now that  $f_0(t, x(t), x'(t))$  be  $L$ -integrable in  $[t_1, t_2]$ . We shall show in this section that existence theorems I to V of the memoir [8a] of Tonelli for free problems can be derived as particular cases from the existence theorems above for Lagrange problems, from the previous Cesari existence theorems [1abcd], and results in [3ab] and [2].

**THEOREM 6.1.** *Let  $A$  be compact,  $B$  closed,  $f_0(t, x, x')$  continuous in  $A \times E_n$  and convex in  $u$  for every  $(t, x) \in A$ , and let us assume that: ( $\phi$ ) There exists a continuous real-valued function  $\phi(\xi)$ ,  $0 \leqq \xi < +\infty$ , with  $f_0(t, x, x') \geqq \phi(|x'|)$ ,  $\phi(\xi)/\xi \rightarrow +\infty$  as  $\xi \rightarrow +\infty$ . Then, integral (6) has an absolute minimum in every closed nonempty class of admissible AC curves  $x = x(t)$ ,  $t_1 \leqq t \leqq t_2$ , satisfying (7) and (8).*

Indeed, here  $f = u$ ,  $m = n$ ,  $U = E_n$ , and by force of ( $\phi$ ) both 1 and  $f$  are of slower growth than  $f_0$  in  $A$ , so that the convex sets  $\tilde{Q}(t, x) = [(z^0, u)|z^0 \geqq f_0(t, x, u), u \in E_n]$  satisfy property (Q) at every  $(\bar{t}, \bar{x}) \in A$  by force of Criterion 5.1. Theorem 6.1 is now a corollary of Cesari's Existence Theorem 1 for Lagrange problems in [1a, p. 390]. Of course, we could use Theorem 2.1 instead. Indeed, by force of ( $\phi$ ), first  $\phi(\xi) \geqq c + d\xi$  for suitable constants  $c$  real and  $d > 0$ , and then condition ( $\Lambda$ ) is trivially satisfied with  $r = c$ ,  $b = 0$ ,  $v = d$ . On the other hand, condition ( $\phi$ ) certainly implies condition ( $\epsilon$ ). Thus, Existence Theorem 2.1 applies with  $G = \phi$ . Note that Theorem 6.1 is essentially Tonelli's Theorem I of [8a, p. 208] for any  $n \geqq 1$  instead of  $n = 1$  only.

**THEOREM 6.2.** *This is the same as Theorem 6.1 with ( $\phi$ ) replaced by*

$$(l) \quad \lim_{|u| \rightarrow \infty} f_0(t, x, u)/|u| = +\infty \quad \text{for every } (t, x) \in A.$$

Indeed, for  $f_0$  convex, conditions ( $\phi$ ) and (l) are equivalent (see [8a]). Note that Criterion 5.2 is essentially Tonelli's Theorem II of [8a, p. 211] for any  $n \geqq 1$  instead of  $n = 1$  only.

**THEOREM 6.3.** *This is the same as Theorem 6.1 with ( $\phi$ ) replaced by: ( $\epsilon_0$ ) Given  $\epsilon > 0$  there is a locally integrable function  $\psi_\epsilon(t)$  which may depend on  $\epsilon$ , such that  $|u| \leqq \psi_\epsilon(t) + \epsilon f_0(t, x, u)$  for all  $(t, x, u) \in A \times E_n$ .*

Indeed, condition ( $\epsilon_0$ ) implies that, for all  $(\bar{t}, \bar{x}) \in A$ , with exception perhaps of a set of points whose  $t$ -coordinate lies in a set of measure zero on the  $t$ -axis, the figurative  $z = f_0(\bar{t}, \bar{x}, u)$ ,  $u \in E_n$ , does not contain any straight line, and hence, because of Theorem 5.6,  $f_0$  is normal in  $u$  at  $(\bar{t}, \bar{x})$ . This in turn implies, because of Criterion 5.3, that the sets  $\tilde{Q}(t, x) = [(z^0, u)|z^0 \geqq f_0(t, x, u), u \in E_n]$  satisfy condition (Q) at the same points  $(\bar{t}, \bar{x}) \in A$ . We are now in a position to apply Cesari's existence theorem for Lagrange problems [1b, Cor. 2, pp. 545–546, and Remark 10, p. 539]. Theorem 6.3 is thereby proved.

Concerning Tonelli's Theorem III of [8a, p. 213], let us assume that there is a closed subset  $E$  of  $A$  such that  $f_0$  satisfies condition (l) at every point of  $A - E$ , and the condition ( $\beta$ ) stated by Tonelli in [8a, p. 213] and restated in [3b, p. 380]

under the name of condition (T). As proved by J. R. La Palm in [3b, pp. 380–382],  $f_0$  then satisfies condition  $(\varepsilon_0)$  in A. and Theorem 6.3 applies.

**THEOREM 6.4.** *Let A be compact, B closed,  $f_0$  continuous in  $A \times E_n$ , convex in u for every  $(t, x) \in A$ , and normal. Let G be a closed subset of A contained in a countable family of AC curves  $\Gamma: x = x(t)$ ,  $t' \leqq t \leqq t''$ , or in a collection of straight lines  $x = \bar{x} = (\bar{x}^1, \dots, \bar{x}^n)$  such that each set  $\{\bar{x}^i\}$ ,  $i = 1, \dots, n$ , of real numbers is of measure zero on the real line. Let us assume that at every point  $(t, x) \in A - G$  we have: (l)  $f_0(t, x, u)/|u| \rightarrow +\infty$  as  $|u| \rightarrow +\infty$ . Then, the integral (6) has an absolute minimum in every closed nonempty class  $\Omega$  of AC curves  $x = x(t)$ ,  $t_1 \leqq t \leqq t_2$ , satisfying (7) and (8).*

Since  $f_0$  is normal,  $f_0$  satisfies conditions (SN'), (SN'') of § 5, part (b). In particular, for  $\bar{u} = 0$  we obtain  $f_0(t, x, u) \geqq r + b \cdot u + v|u|$  for all  $(t, x) \in N_\delta(\bar{t}, \bar{x})$ , and all  $u \in E_n$ . In other words,  $f_0$  satisfies property (Λ) of Theorem 2.1 with  $f = u$ ,  $m = n$ ,  $U = E_n$ . Furthermore, again by force of the normality, and because of Criterion 5.3, the sets

$$\tilde{Q}(t, x) = [(z^0, u)|z^0 \geqq f_0(t, x, u), u \in E_n]$$

satisfy property (Q) for every  $(t, x) \in A$ . Finally, for every  $(\bar{t}, \bar{x}) \in A - G$ , there is a neighborhood  $N_\delta(\bar{t}, \bar{x})$  at all points of which condition (l) holds; hence, by Tonelli's argument,  $f_0$  satisfies a condition (ϕ) in  $N_\delta(\bar{t}, \bar{x})$ , and, a fortiori, a condition (ε) as in Theorem 2.1 with  $f = u$ ,  $m = n$ ,  $U = E_n$ . Theorem 2.1 applies and Theorem 6.4 is proved. Note that Theorem 6.4 is essentially Tonelli's Theorem IV with  $n \geqq 1$  instead of  $n = 1$ .

**THEOREM 6.5.** *Let A be compact, B closed,  $f_0$  continuous in  $A \times E_n$ , and convex in u for every  $(t, x) \in A$ . Let G be any subset of A as described in Theorem 6.4. Let us assume that (i) at every point  $(\bar{t}, \bar{x}) \in G$  the figurative  $Q(\bar{t}, \bar{x}) = [(z^0, u), z^0 = f_0(\bar{t}, \bar{x}, u), u \in E_n]$  contains no straight line; (ii) for every  $(\bar{t}, \bar{x}) \in A - G$  there is a neighborhood  $N_\delta(\bar{t}, \bar{x})$  of  $(\bar{t}, \bar{x})$  in A and, for every  $\varepsilon > 0$ , a locally integrable nonnegative function  $\psi_\varepsilon(t)$  such that  $|u| \leqq \psi_\varepsilon(t) + \varepsilon f_0(t, x, u)$  for all  $(t, x) \in N_\delta(\bar{t}, \bar{x})$  and  $u \in E_n$ . Then, the same conclusion holds as in Theorem 6.4.*

For any point  $(\bar{t}, \bar{x}) \in G$  the figurative  $Q(\bar{t}, \bar{x})$  contains no straight line, and then the same holds at every point  $(t, x)$  of some neighborhood  $N_\delta(\bar{t}, \bar{x})$  of  $(\bar{t}, \bar{x})$  in A. Thus,  $f_0$  is normal in  $N_\delta(\bar{t}, \bar{x})$ , and consequently satisfies property (Λ) in  $N_\delta(\bar{t}, \bar{x})$ . Also, the sets  $\tilde{Q}(t, x) = [(z^0, u)|z^0 \geqq f_0(t, x, u), u \in E_n]$  satisfy property (Q) in all of  $N_\delta(\bar{t}, \bar{x})$  by force of Criterion 5.3. Note that for every  $(\bar{t}, \bar{x}) \in A - G$  the growth condition described in hypothesis (ii) holds in a neighborhood  $N_\delta(\bar{t}, \bar{x})$  of  $(\bar{t}, \bar{x})$  in A. As a consequence, property (l) holds at all  $(t, x) \in N_\delta(\bar{t}, \bar{x})$  with exception perhaps of a set of points whose t-coordinate lies in a set of measure zero on the t-axis. Hence, the figurative  $Q(t, x)$  contains no straight line and  $f_0$  is normal at all  $(t, x) \in N_\delta(\bar{t}, \bar{x})$  with perhaps the same exception as above.

By force of Criterion 5.3, property (Q) holds at all points of A-G with the same possible exception. Thus, property (Q) holds in all of A, with the same possible exception, as required in Theorem 2.1. Conditions (s) and (ε) of Theorem 2.1 hold, and condition (Λ) of Theorem 4.1 also holds. Thus, Theorem 4.1 applies and Theorem 6.5 is proved.

Concerning Tonelli's Theorem V of [8a, p. 225], let us assume that G is given as in Theorem 6.5 and that hypothesis (i) is satisfied. Also, let us assume that there

is a closed set  $E$  of  $A$  such that  $f_0$  satisfies condition (I) at every point of  $A-G-E$ , and that condition ( $\beta$ ) stated by Tonelli in [8a, p. 213] holds. This condition was restated in [3b, p. 380] under the name of Tonelli's condition (T). As proved by J. R. La Palm in [3b, pp. 380-382],  $f_0$  then satisfies hypothesis (ii) in A-G, and Theorem 6.5 applies. The reader may have noticed that the terminology used in the present paper is slightly different from the one used in [1c]. This change was made in order to be closer to Tonelli's and McShane's terminology.

**7. Examples.** (a) For free problems the functions

$$1. f_0 = |x|u^2 + (1 + u^2)^{1/2}, \quad m = n = 1, \quad U = E_1, \quad f = u;$$

$$2. f_0 = |x||u|^{1+\alpha} + (1 + u^2)^{1/2}, \quad \alpha > 0, \quad m = n = 1, \quad U = E_1, \quad f = u;$$

$$3. f_0 = (x^2 + y^2)(u^2 + v^2) + (1 + u^2 + v^2)^{1/2}, \quad m = n = 2, \quad U = E_2, \\ f_1 = u, \quad f_2 = v;$$

$$4. f_0 = (x - t)^2 u^2 + (1 + u^2)^{1/2}, \quad n = m = 1, \quad U = E_1, \quad f = u;$$

$$5. f_0 = (x \sin 1/x)^2 u^2 + (1 + u^2)^{1/2}, \quad m = n = 1, \quad U = E_1, \quad f = u;$$

$$6. f_0 = ((x^2 - t^2)^2 + y^2)(u^2 + v^2) + (1 + u^2 + v^2)^{1/2}, \quad m = n = 2, \\ U = E_2, \quad f_1 = u, \quad f_2 = v;$$

$$7. f_0 = |x|u^2 + (1 + u^2)^{1/2} - 2u, \quad m = n = 1, \quad U = E_1, \quad f = u;$$

$$8. f_0 = (x^2 + y^2)(u^2 + v^2) + (1 + u^2 + v^2)^{1/2} - 3u - 5v - 1, \\ m = n = 2, \quad U = E_2, \quad f_1 = u, \quad f_2 = v,$$

all satisfy the condition of the Existence Theorem 2.1. The exceptional set is respectively, the single straight line  $x = 0$  in the first, second, and seventh examples, the single straight line  $x = y = 0$  in the third and eighth, the straight line  $x = t$  in the fourth, the countable family of straight lines  $x = 0$  and  $x = \pm(k\pi)^{-1}$ ,  $k = 1, 2, \dots$ , in the fifth, and the two straight lines,  $x = \pm t$ ,  $y = 0$  in the sixth. In the seventh and eighth examples condition ( $\Lambda$ ) is satisfied by taking respectively  $b = 2$ , and  $b = (3, 5)$  in the whole of  $A$ .

(b) For general Lagrange problems the following functions:

$$1. f_0 = |x|u^4 + (1 + u^4)^{1/2}, \quad f = u^2, \quad n = 1, \quad m = 1, \quad U = E_1;$$

$$2. f_0 = |x|u^4 + (1 + u^4)^{1/2} - 3u^2, \quad f = u^2, \quad n = 1, \quad m = 1, \quad U = E_1;$$

$$3. f_0 = x^3 u^2, \quad f = x^2 u, \quad n = 1, \quad m = 1, \quad x \geq 0, \quad U = E_1;$$

$$4. f_0 = x^3 u^2 - x^2 u, \quad f = x^2 u, \quad n = 1, \quad m = 1, \quad x \geq 0, \quad U = E_1,$$

satisfy the conditions of the Existence Theorem 2.1.

#### REFERENCES

- [1a] L. CESARI, *Existence theorems for weak and usual optimal solutions in Lagrange problems with unilateral constraints. I, II*, Trans. Amer. Math. Soc., 124 (1966), pp. 369-412, 413-429.

- [1b] ———, *Existence theorems for optimal controls of the Mayer type*, this Journal, 6 (1968), pp. 517–552.
- [1c] ———, *Seminormality and upper semicontinuity in optimal control*, J. Optimization Theory Appl., 6 (1970), pp. 114–137.
- [1d] ———, *Closure, lower closure, and semicontinuity theorems in problems of optimal control*, this Journal, 9 (1971), pp. 287–301.
- [1e] ———, *Lagrange problems of optimal control and convex sets not containing any straight line*, J. Optimization Theory Appl., to appear.
- [2] L. CESARI, J. R. LA PALM AND T. NISHIURA, *Remarks on some existence theorems for optimal control*, Ibid., 3 (1969), pp. 296–305.
- [3a] J. R. LA PALM, *Remarks on certain growth conditions in problems of optimal control*, Ibid., 4 (1969), pp. 321–329.
- [3b] ———, *A recent significant growth condition in optimal control theory*, Ibid., 4 (1969), pp. 378–385.
- [4a] E. J. MCSHANE, *On the semicontinuity of integrals in the calculus of variations*, Ann. of Math., 33 (1932), pp. 460–484.
- [4b] ———, *Existence theorems for ordinary problems of the calculus of variations*, Ann. Scuola Norm. Sup. Pisa (2), 3 (1934), pp. 181–211, 287–315.
- [4c] ———, *Semicontinuity of integrals in the calculus of variations*, Duke Math. J., 2 (1936), pp. 597–616.
- [4d] ———, *Some existence theorems for problems in the calculus of variations*, Ibid., 4 (1938), pp. 132–156.
- [4e] ———, *A navigation problem in the calculus of variations*, Amer. Math. J., 59 (1937), pp. 327–334.
- [4f] ———, *Relaxed controls and variational problems*, this Journal, 5 (1967), pp. 438–485.
- [5] E. J. MCSHANE AND R. B. WARFIELD, *On Filippov's implicit functions lemma*, Proc. Amer. Math. Soc., 18 (1969), pp. 41–47.
- [6] T. NISHIURA, *On an existence theorem for optimal control*, this Journal, 5 (1967), pp. 532–544.
- [7a] C. OLECH, *Existence theorems for optimal problems with vector valued cost functions*, Trans. Amer. Math. Soc., 136 (1967), pp. 157–180.
- [7b] ———, *Existence theorems for optimal control problems involving multiple integrals*, J. Differential Equations, 6 (1969), pp. 512–526.
- [8a] L. TONELLI, *Sugli integrali del calcolo delle variazioni in forma ordinaria*, Ann. Scuola Norm. Sup. Pisa (2), 3 (1934), pp. 401–450 = Opere Scelte, Edizioni Cremonese, Roma, 3 (1962), pp. 192–254.
- [8b] ———, *Fondamenti di Calcolo delle Variazioni*, 2 vols., Zanichelli, Bologna, 1921–23.
- [9] L. TURNER, *The direct method in the calculus of variations*, Doctoral thesis, Purdue University, Lafayette, Ind., 1957.

## NONLINEAR PROGRAMMING IN COMPLEX SPACE: NECESSARY CONDITIONS\*

ROBERT A. ABRAMS† AND ADI BEN-ISRAEL‡

**Abstract.** Necessary conditions of the Kuhn–Tucker type are given for two classes of nonlinear programming problems over polyhedral cones in finite-dimensional complex space. The first class consists of problems of the form:

$$\text{Minimize } \operatorname{Re} f(z) \quad \text{subject to } g(z) \in S,$$

where  $S$  is a polyhedral cone in  $C^m$  and  $f: C^n \rightarrow C, g: C^n \rightarrow C^m$  are analytic functions. A necessary condition for a feasible point  $z^0$  to be optimal is that there exist a vector  $u \in S^*$  such that  $\overline{\nabla f}(z^0) = [D_z^H g(z^0)]u$  and  $\operatorname{Re}(g(z^0), u) = 0$ . The second class consists of problems of the form:

$$\text{Minimize } \operatorname{Re} f(z, \bar{z}) \quad \text{subject to } g(z, \bar{z}) \in S,$$

where  $f: C^{2n} \rightarrow C, g: C^{2n} \rightarrow C^m$  are analytic. A necessary condition for a feasible point  $z^0$  to be optimal is that there exist a  $u \in S^*$  such that  $\overline{\nabla_z f}(z^0, \bar{z}^0) + \nabla_{\bar{z}} f(z^0, \bar{z}^0) = [D_z^H g(z^0, \bar{z}^0)]u + [D_{\bar{z}}^T g(z^0, \bar{z}^0)]\bar{u}$  and  $\operatorname{Re}(g(z^0, \bar{z}^0), u) = 0$ . The derivation of necessary conditions for problems of the first class is analogous to that used in the real case. For problems of the second class, necessary conditions are obtained by considering an equivalent problem in the form of the first class.

**1. Introduction.** In this paper first order necessary conditions of the Kuhn–Tucker type are given for nonlinear programming problems over polyhedral cones in finite-dimensional complex space. Two classes of problems are considered. In the first class, both the objective and constraint functions are assumed to be analytic functions of  $n$  complex variables, i.e., the problems are of the form:

$$\text{Minimize } \operatorname{Re} f(z) \quad \text{subject to } g(z) \in S,$$

where  $S$  is a polyhedral cone in  $C^m$  and the functions  $f: C^n \rightarrow C$  and  $g: C^n \rightarrow C^m$  are analytic, at least in a neighborhood of a feasible point  $z^0$  being tested for optimality. A necessary condition for optimality of  $z^0$ , given in Theorem 1, is that there exist a  $u$  in a subcone of  $S^*$  such that

$$\overline{\nabla_z f}(z^0) = [D_z^H g(z^0)]u \quad \text{and} \quad \operatorname{Re}(g(z^0), u) = 0.$$

In the second, more general, class of problems, both objective and constraint functions may involve  $z$  and  $\bar{z}$  and are assumed to be analytic functions of  $2n$  complex variables in a neighborhood of  $(z^0, \bar{z}^0)$ , where  $z^0$  is a feasible point being tested for optimality. Problems in this class are of the form:

$$\text{Minimize } \operatorname{Re} f(z, \bar{z}) \quad \text{subject to } g(z, \bar{z}) \in S,$$

where  $f: C^{2n} \rightarrow C$  and  $g: C^{2n} \rightarrow C^m$  are analytic. A necessary condition for the optimality of  $z^0$  is, by Theorem 3, that there exist a vector  $u$  in a certain subcone

\* Received by the editors February 23, 1970, and in revised form October 8, 1970.

† Department of Industrial Engineering and Management Sciences, Northwestern University, Evanston, Illinois 60201.

‡ Departments of Engineering Science and of Industrial Engineering and Management Sciences, Northwestern University, Evanston, Illinois 60201, and Department of Applied Mathematics, Technion-Israel Institute of Technology, Haifa, Israel. The work of this author was supported in part by the National Science Foundation under Project GP 13546.



of  $S^*$  such that

$$\overline{\nabla_z f(z^0, \bar{z}^0)} + \nabla_{\bar{z}} f(z^0, \bar{z}^0) = [D_z^H g(z^0, \bar{z}^0)]u + [D_{\bar{z}}^T g(z^0, \bar{z}^0)]\bar{u}$$

and  $\text{Re}(g(z^0, \bar{z}^0), u) = 0$ .

Sufficient conditions for optimality, and duality theorems based on the complex analogues of convexity, are given in [1], [2] and [11].

Complex linear programming was studied in [13] (where complex mathematical programming was introduced), [5] and [3] (where the natural extension to quadratic programming is made). Other references are given in the survey article [4], which also contains some applications of complex programming, and in which some of the results of the present paper and of [2] are announced without proof.

This paper contains five sections and an Appendix. The preliminaries of § 2 are used to derive the necessary conditions for problems of the first and second classes in § 3 and § 4 respectively. An example is given in § 5. Some results on derivatives of complex functions which are used here and in [2] are derived in the Appendix.

**2. Preliminaries.**

NOTATIONS 2.1.

$C^n [R^n]$  denotes the  $n$ -dimensional complex [real] vector space;

$C^{m \times n} [R^{m \times n}]$ , the  $m \times n$  complex [real] matrices;

$R_+^n \equiv \{x \in R^n : x_i \geq 0, i = 1, \dots, n\}$ , the nonnegative orthant of  $R^n$ ;

$x \geq y$  denotes  $x - y \in R_+^n$  for  $x, y \in R^n$ .

For  $A = (a_{ij}) \in C^{m \times n}$ ,

$\bar{A} \equiv (\bar{a}_{ij})$  denotes the conjugate,

$A^T \equiv (a_{ji})$  denotes the transpose,

$A^H \equiv \bar{A}^T$  denotes the conjugate transpose.

For  $x = (x_i) \in C^n, y \in C^n$ ,

$(x, y) \equiv y^H x$  denotes the inner product of  $x$  and  $y$ ,

$\bar{x} \equiv (\bar{x}_i)$  denotes the conjugate,

$\text{Re } x \equiv (\text{Re } x_i) \in R^n$  denotes the real part,

$\text{Im } x \equiv (\text{Im } x_i) \in R^n$  denotes the imaginary part,

$\arg x \equiv (\arg x_i)$  denotes the argument of  $x$ .

For a subspace  $L \subset C^n, L^\perp \equiv \{y \in C^n : l \in L \Rightarrow (y, l) = 0\}$  denotes the orthogonal complement of  $L$ .

For a nonempty set  $S \subset C^n, S^* \equiv \{y \in C^n : x \in S \Rightarrow \text{Re}(y, x) \geq 0\}$  denotes the dual (also polar) of  $S$ .

For a nonempty set  $S \subset R^n$ ,

$$S^* = \{y \in R^n : x \in S \Rightarrow (x, y) \geq 0\}.$$

For an analytic function  $f : C^n \rightarrow C$  and a point  $z^0 \in C^n, \nabla_z f(z^0) \equiv ((\partial f / \partial z_i)(z^0)), i = 1, \dots, n$ , denotes the gradient of  $f$  at  $z^0$ .

For a complex function  $f(w^1, w^2)$  analytic in the  $2n$  variables  $(w^1, w^2)$  at the point  $(z^0, \bar{z}^0) \in C^n \times C^n$ ,

$$\nabla_z f(z^0, \bar{z}^0) \equiv \left( \frac{\partial f}{\partial w_i^1}(z^0, \bar{z}^0) \right), \quad i = 1, \dots, n,$$

and

$$\nabla_z f(z^0, \bar{z}^0) \equiv \left( \frac{\partial f}{\partial w_i^2}(z^0, \bar{z}^0) \right), \quad i = 1, \dots, n.$$

For an analytic function  $g: C^n \rightarrow C^m$ ,

$$D_z g(z^0) \equiv \left( \frac{\partial g_i}{\partial z_j}(z^0) \right), \quad i = 1, \dots, m, \quad j = 1, \dots, n.$$

Similarly, for a function  $g: C^n \times C^n \rightarrow C^m$  analytic in the  $2n$  variables  $(w^1, w^2)$  at  $(z^0, \bar{z}^0) \in C^n \times C^n$ ,

$$D_z g(z^0, \bar{z}^0) = \left( \frac{\partial g_i}{\partial w_j^1}(z^0, \bar{z}^0) \right), \quad i = 1, \dots, m, \quad j = 1, \dots, n,$$

and

$$D_z g(z^0, \bar{z}^0) = \left( \frac{\partial g_i}{\partial w_j^2}(z^0, \bar{z}^0) \right), \quad i = 1, \dots, m, \quad j = 1, \dots, n.$$

Also

$$\begin{aligned} D_z^T g(z^0, \bar{z}^0) &= (D_z g(z^0, \bar{z}^0))^T, \\ D_z^H g(z^0, \bar{z}^0) &= (D_z g(z^0, \bar{z}^0))^H. \end{aligned}$$

**DEFINITIONS 2.2.** A nonempty set  $S \subset C^n$  is:

- (a) *convex* if  $0 \leq \lambda \leq 1 \Rightarrow \lambda S + (1 - \lambda)S \subset S$ ,
- (b) a *cone* if  $0 \leq \lambda \Rightarrow \lambda S \subset S$ ,
- (c) a *polyhedral cone* if for some positive integer  $k$  and  $A \in C^{n \times k}$ ,

$$S = AR_+^k = \{Ax : x \in R_+^k\};$$

i.e.,  $S$  is generated by finitely many vectors (the columns of  $A$ ).

The following results are needed in the sequel.

*Result 2.3.* A polyhedral cone in  $C^n$  is a closed convex cone.

*Result 2.4.* A nonempty set  $S \subset C^n$  is a closed convex cone if and only if  $S = S^{**}$  (for proof see, e.g., [5, Theorem 1.5]).

*Result 2.5.* If  $S, T$  are polyhedral cones, then  $S \times T$  is a polyhedral cone.

*Result 2.6.* For any nonempty sets  $S, T: (S \times T)^* = S^* \times T^*$ .

*Result 2.7.* Let  $A \in C^{m \times n}, b \in C^m$  and  $S \subset C^n$  be a polyhedral cone. Then the following are equivalent:

- (a)  $Ax = b, x \in S$  is consistent.
- (b)  $A^H y \in S^* \Rightarrow \text{Re}(b, y) \geq 0$ . (See [5, Theorem 3.5].)

*Result 2.8.* The nonnegative orthant  $R_+^n$  is a self-dual set in  $R^n: (R_+^n)^* = R_+^n$ .

*Result 2.9.* Let  $S$  be a polyhedral cone in  $C^n$ . Then  $S$  is the intersection of finitely many closed half-spaces, each including the origin in its boundary:

$$S = \bigcap_{k=1}^p H_{u_k},$$

where  $H_{u_k} = \{z \in C^n : \text{Re}(z, u_k) \geq 0\}$ . (It is proved similarly to the real case, e.g., [14].)

*Result 2.10.* Let  $S = \bigcap_{k=1}^p H_{u_k}$  be a polyhedral cone in  $C^n$  or  $R^n$  and let  $z^0 \in S$ . Then  $S(z^0)$  is defined to be the intersection of those closed half-spaces  $H_{u_k}$  which include  $z^0$  in their boundaries, i.e.,

$$S(z^0) = \bigcap_{k \in B(z^0)} H_{u_k},$$

where  $B(z^0) \equiv \{k: \text{Re}(z^0, u_k) = 0\}$ . If  $z^0$  is in the interior of  $S$ , then  $S(z^0) = C^n$ .

*Result 2.11.* Let  $\emptyset \neq S \subset T \subset C^n$ . Then  $T^* \subset S^*$ .

*Result 2.12.* Let  $\{S_i: i = 1, \dots, p\}$  be closed convex cones in  $C^n$ . Then

$$\left( \bigcap_{i=1}^p S_i \right)^* = \text{cl} \sum_{i=1}^p S_i^*$$

(which follows from [5, Corollary 1.7]).

**3. Necessary conditions for problems of the first class.** Necessary conditions for optimal points of analytic nonlinear programming problems will now be derived. In the real case the necessary conditions are obtained by using the Farkas theorem to show that the gradient of the objective function lies in a cone generated by the gradients of the constraint functions, e.g., [10, § 2.1]. The generalized Farkas theorem, Result 2.7, is used here in a similar manner.

Let  $f: C^n \rightarrow C$  and  $g: C^n \rightarrow C^m$  both be analytic in a neighborhood of a point  $z^0$ . ( $g: C^n \rightarrow C^m$  is analytic if each of its components  $g_i: C^n \rightarrow C, i = 1, \dots, m$ , is analytic.) Let  $S$  be a polyhedral cone in  $C^m$ . The complex nonlinear programming problem of the first class is:

(1)  $\text{Minimize } \text{Re } f(z) \text{ subject to } g(z) \in S.$

An analogue of the Kuhn–Tucker constraint qualification provides a regularity condition which the constraint is assumed to satisfy at a point being checked for optimality. Let  $z^0$  be a feasible point of (1). Then the (Kuhn–Tucker) *constraint qualification* holds at  $z^0$ , or the point  $z^0$  is *qualified*, if every  $z \in C^n$  such that  $[D_z g(z^0)]z \in S(g(z^0))$  is tangent to a once differentiable arc  $a(\theta)$  beginning at  $z^0$  and leading into the feasible region, i.e.,  $a(0) = z^0, g(a(\theta)) \in S$  for  $0 \leq \theta < \varepsilon$  and  $a'(0) = kz$  for some  $k > 0, \varepsilon > 0$ .

**THEOREM 1.** *Let  $S$  be a polyhedral cone in  $C^m$ . Let  $f: C^n \rightarrow C$  and  $g: C^n \rightarrow C^m$  both be analytic in a neighborhood of a qualified point  $z^0$ . Then a necessary condition for  $z^0$  to be a local minimum of the problem:*

$$\text{Minimize } \text{Re } f(z) \text{ subject to } g(z) \in S,$$

*is that there exist a vector  $u \in [S(g(z^0))]^* \subset S^*$  such that*

(2) 
$$\overline{\nabla_z f(z^0)} = [D_z^H g(z^0)]u$$

*and*

(3) 
$$\text{Re}(g(z^0), u) = 0.$$

*Proof.* Define the set

(4) 
$$Z = \{z \in C^n: [D_z g(z^0)]z \in S(g(z^0)), \text{Re}(\overline{\nabla_z f(z^0)}, z) < 0\}.$$

If  $Z$  is empty, then

(5) 
$$[D_z g(z^0)]z \in S(g(z^0)) \Rightarrow \text{Re}(\overline{\nabla_z f(z^0)}, z) \geq 0.$$

The cone  $S(g(z^0))$  is polyhedral and hence closed and convex. Therefore by (2.4),  $S(g(z^0))$  may be replaced by  $([S(g(z^0))]^*)^*$  in (5). By Result 2.7, the complex Farkas theorem, (5) is equivalent to the existence of a  $u \in [S(g(z^0))]^*$  such that

$$\overline{\nabla_z f(z^0)} = [D_z^H g(z^0)]u.$$

Thus in order to prove (2), it is sufficient to show that  $Z$  is empty.

Consider any  $z$  such that  $[D_z g(z^0)]z \in S(g(z^0))$ . Since  $z^0$  is a qualified point, there is a feasible arc  $a(\theta)$  with  $a(0) = z^0$  and  $a'(0) = kz$  for some  $k > 0$ . By assumption,  $z^0$  is a local minimum of  $\text{Re } f(z)$  along  $a(\theta)$ . Therefore

$$(6) \quad \frac{d}{d\theta} \text{Re } [f(a(\theta))]_{\theta=0} \geq 0$$

which is equivalent to

$$(7) \quad \text{Re } \frac{d}{d\theta} f(a(\theta))_{\theta=0} \geq 0.$$

Using the chain rule, we find that (7) becomes  $\text{Re } [\nabla_z^T f(z^0)a'(0)] \geq 0$ , or with  $a'(0) = kz$ ,  $\text{Re } (\overline{\nabla f(z)}, z) \geq 0$ , which proves that  $Z$  is empty.

To prove (3), note that  $u \in [S(g(z^0))]^*$ . Then

$$\begin{aligned} [S(g(z^0))]^* &= \left( \bigcap_{k \in B(g(z^0))} H_{u_k} \right)^* = \text{cl } \sum_{k \in B(g(z^0))} (H_{u_k})^* \\ &= \text{cl } \sum_{k \in B(g(z^0))} \{ \alpha_k u_k : \alpha_k \geq 0 \} = \sum_{k \in B(g(z^0))} \{ \alpha_k u_k : \alpha_k \geq 0 \}, \end{aligned}$$

where the second equality follows from (2.12) and the last follows from the fact that  $\sum_{k \in B(g(z^0))} \{ \alpha_k u_k : \alpha_k \geq 0 \}$  is a polyhedral cone and hence closed (Result 2.3). Therefore  $u \in [S(g(z^0))]^*$  implies

$$(8) \quad u = \sum_{i \in B(g(z^0))} \beta_i u_i \quad \text{with } \beta_i \geq 0.$$

By the definition of  $B(g(z^0))$ ,  $\text{Re } (u_i, g(z^0)) = 0$  for  $i \in B(g(z^0))$ , and hence (3) follows from (8). This completes the proof.

For comparison with the real case, a real version of Theorem 1, i.e., the Kuhn–Tucker theorem for constraints over polyhedral cones, is stated now. The proof follows from the proof of Theorem 1 by replacing complex spaces with real spaces and replacing analytic functions with differentiable functions. The definition of a qualified point is analogous to that given above for the complex case.

**THEOREM 2.** *Let  $S$  be a polyhedral cone in  $R^m$ . Let  $f: R^n \rightarrow R$  and  $g: R^n \rightarrow R^m$  both be (Fréchet) differentiable in a neighborhood of a qualified point  $x^0$ . Then a necessary condition for  $x^0$  to be a local minimum of the problem:*

$$\text{Minimize } f(x) \quad \text{subject to } g(x) \in S,$$

*is that there exist a  $u \in [S(g(x^0))]^*$  such that*

$$\nabla_x f(x^0) = [D_x^T g(x^0)]u$$

*and*

$$(g(x^0), u) = 0.$$

The usual form of the Kuhn–Tucker theorem for  $l$  inequality constraints and  $p$  equality constraints is obtained from the above theorem by taking

$$S = R_+^l \times \mathbf{0}^{m-l},$$

where  $\mathbf{0}^{m-l}$  is the zero element of the  $(m - l)$ -dimensional real vector space. The representation of  $S$  as an intersection of half-spaces is

$$S = \bigcap_{k=1}^m H_{e_k} \bigcap_{k=l+1}^m H_{-e_k},$$

where  $e_k$  is the  $k$ th unit vector in  $R^m$ . Suppose that at the point  $x^0$  the first  $r$  inequality constraints are binding and that  $l - r$  are not, i.e.,

$$(g(x^0), e_k) = 0, \quad k = 1, \dots, r, \quad \text{and} \quad (g(x^0), e_k) > 0, \quad k = r + 1, \dots, l.$$

Then

$$S(g(x^0)) = \bigcap_{k=1}^r H_{e_k} \bigcap_{k=l+1}^m H_{e_k} \bigcap_{k=l+1}^m H_{-e_k}$$

and

$$[S(g(x^0))]^* = \sum_{k=1}^r \{\alpha_k e_k : \alpha_k \geq 0\} + \sum_{k=l+1}^m \{(\alpha_k - \alpha'_k)e_k : \alpha_k \geq 0, \alpha'_k \geq 0\}.$$

Thus  $u \in [S(g(x^0))]^*$  implies

$$u = \sum_{k=1}^r \beta_k e_k + \sum_{k=l+1}^m \gamma_k e_k,$$

where  $\beta_k \geq 0, k = 1, \dots, r$ , and  $\gamma_k$  is unrestricted,  $k = l + 1, \dots, m$ . The  $\beta_k$  and the  $\gamma_k$  are multipliers of the inequality constraints and the equality constraints respectively.

**4. Necessary conditions for problems of the second class.** The class of problems under consideration is now extended to those in which  $f$  and  $g$  are analytic functions of the  $2n$  variables  $(w^1, w^2) \in C^{2n}$  and in which the additional constraint  $w^2 = \overline{w^1}$  is added. Thus let  $f : C^{2n} \rightarrow C$  and  $g : C^{2n} \rightarrow C^m$  both be analytic in a neighborhood of a point  $(z^0, \overline{z^0}) \in C^{2n}$  and let  $S$  be a polyhedral cone in  $C^m$ . Consider the problem:

$$(9) \quad \text{Minimize} \quad \text{Re } f(w^1, w^2) \quad \text{subject to } g(w^1, w^2) \in S \quad \text{and} \quad w^2 = \overline{w^1},$$

which is rewritten in the more suggestive form:

$$(10) \quad \text{Minimize} \quad \text{Re } f(z, \bar{z}) \quad \text{subject to } g(z, \bar{z}) \in S.$$

From (A.20) it follows that a function explicitly involving the variable  $\bar{z}$ , such as  $f(z, \bar{z})$  and  $g(z, \bar{z})$  in (10), cannot be an analytic function of the variable  $z$ . Therefore Theorem 1 is not directly applicable to problems of the form (10). However, it is possible to recast (10) in the form (1), and thus optimality conditions for (10) can be obtained from Theorem 1.

Define the set  $Q$  by

$$(11) \quad Q \equiv \left\{ \begin{pmatrix} w^1 \\ w^2 \end{pmatrix} \in C^{2n} : w^2 = \overline{w^1} \right\}.$$

$Q$  is a polyhedral cone; it is generated by the set of vectors

$$(12) \quad \left\{ \bigcup_{j=1}^n \begin{pmatrix} e_j \\ e_j \end{pmatrix} \bigcup_{j=1}^n \begin{pmatrix} -e_j \\ -e_j \end{pmatrix} \bigcup_{j=1}^n \begin{pmatrix} ie_j \\ -ie_j \end{pmatrix} \bigcup_{j=1}^n \begin{pmatrix} -ie_j \\ ie_j \end{pmatrix} \right\},$$

where  $e_j$  is the  $j$ th unit vector in  $R^n$ . The dual of  $Q$  is easily seen to be the polyhedral cone

$$(13) \quad Q^* = \left\{ \begin{pmatrix} w^1 \\ w^2 \end{pmatrix} \in C^{2n} : w^2 = -\overline{w^1} \right\}.$$

Since  $S$  and  $Q$  are polyhedral cones,  $S \times Q$  is also (Result 2.5), and therefore (10) may be written in the form (1) as:

$$(14) \quad \begin{array}{l} \text{Minimize} \quad \text{Re } f(w^1, w^2) \\ \text{subject to} \quad \begin{pmatrix} g(w^1, w^2) \\ w^1 \\ w^2 \end{pmatrix} \in S \times Q. \end{array}$$

A qualified point for (10) is defined using (14) and the definition of a qualified point for (1). Thus  $(v^1, v^2)$  is a qualified point for (14) if, for every  $(w^1, w^2) \in C^{2n}$  such that

$$(15) \quad D_{\begin{pmatrix} w^1 \\ w^2 \end{pmatrix}} \begin{pmatrix} g(v^1, v^2) \\ v^1 \\ v^2 \end{pmatrix} \begin{pmatrix} w^1 \\ w^2 \end{pmatrix} \in [S \times Q] \begin{pmatrix} g(v^1, v^2) \\ v^1 \\ v^2 \end{pmatrix},$$

there exists an arc  $a(\theta) \subset C^{2n}$  such that  $a(0) = (v^1, v^2)$ ,  $a'(\theta) = k \begin{pmatrix} w^1 \\ w^2 \end{pmatrix}$ ,  $k > 0$ , and

$$\begin{pmatrix} g(a(\theta)) \\ a(\theta) \end{pmatrix} \in S \times Q, \quad 0 \leq \theta < \varepsilon.$$

Condition (15) is equivalent to

$$(16) \quad \begin{pmatrix} D_{\begin{pmatrix} w^1 \\ w^2 \end{pmatrix}} g(v^1, v^2) \\ I \end{pmatrix} \begin{pmatrix} w^1 \\ w^2 \end{pmatrix} \in S(g(v^1, v^2)) \times Q(v^1, v^2).$$

For any point  $(v^1, v^2)$  in the cone  $Q$ ,  $Q(v^1, v^2) = Q$ , and therefore (16) and (15) become

$$D_{\begin{pmatrix} w^1 \\ w^2 \end{pmatrix}} g(v^1, v^2) \begin{pmatrix} w^1 \\ w^2 \end{pmatrix} \in S(g(v^1, v^2)) \quad \text{and} \quad \begin{pmatrix} w^1 \\ w^2 \end{pmatrix} \in Q.$$

Thus we may say that a point  $(z^0, \overline{z^0})$  is qualified for (10) if for all  $(z, \overline{z})$  such that

$$[D_{\begin{pmatrix} z \\ \overline{z} \end{pmatrix}} g(z^0, \overline{z^0})] \begin{pmatrix} z \\ \overline{z} \end{pmatrix} \in S(g(z^0, \overline{z^0})),$$

there exists an arc  $(\beta(\theta), \overline{\beta(\theta)}) \subset C^{2n}$  such that  $\beta(0) = z^0, \beta'(0) = kz$  for some  $k > 0$  and  $g(\beta(\theta), \overline{\beta(\theta)}) \in S$  for  $0 \leq \theta < \epsilon$ .

**THEOREM 3.** *Let  $S$  be a polyhedral cone in  $C^m$ . Let  $f: C^{2n} \rightarrow C$  and  $g: C^{2n} \rightarrow C^m$  both be analytic in a neighborhood of a qualified point  $(z^0, \overline{z^0})$ . A necessary condition for  $z^0, \overline{z^0}$  to be a local minimum of the problem:*

$$\text{Minimize } \operatorname{Re} f(z, \overline{z}) \text{ subject to } g(z, \overline{z}) \in S,$$

is that there exist a  $u \in [S(g(z^0, \overline{z^0}))]^* \subset S^*$  such that

$$(17) \quad \overline{\nabla_z f(z^0, \overline{z^0})} + \nabla_{\overline{z}} f(z^0, \overline{z^0}) = [D_z^H g(z^0, \overline{z^0})]u + [D_{\overline{z}}^T g(z^0, \overline{z^0})]\overline{u}$$

and

$$(18) \quad \operatorname{Re}(g(z^0, \overline{z^0}), u) = 0.$$

*Proof.* The problem under consideration has been shown to be equivalent to (14) which is of the form (1). We apply Theorem 1 to (14). Thus a necessary condition for qualified point  $(w^1, w^2)$  to be an optimal point of (14) is that there exist a

$$\begin{pmatrix} u \\ \lambda \end{pmatrix} \in \left\{ [S \times Q] \begin{pmatrix} g(w^1, w^2) \\ w^1 \\ w^2 \end{pmatrix} \right\}^* = [S(g(w^1, w^2))]^* \times Q^*$$

such that

$$(19) \quad \overline{\nabla_{\begin{pmatrix} w^1 \\ w^2 \end{pmatrix}} f(w^1, w^2)} = D_{\begin{pmatrix} w^1 \\ w^2 \end{pmatrix}}^H \begin{pmatrix} g(w^1, w^2) \\ w^1 \\ w^2 \end{pmatrix} \begin{pmatrix} u \\ \lambda \end{pmatrix}$$

and

$$(20) \quad \operatorname{Re} \left( \begin{pmatrix} g(w^1, w^2) \\ w^1 \\ w^2 \end{pmatrix}, \begin{pmatrix} u \\ \lambda \end{pmatrix} \right) = 0.$$

From (13),  $\lambda \in Q^*$  implies  $\lambda = \begin{pmatrix} \pi \\ -\overline{\pi} \end{pmatrix}$ . Therefore (19) becomes

$$(21) \quad \begin{aligned} \overline{\nabla_{w^1} f(w^1, w^2)} &= D_{w^1}^H g(w^1, w^2)u + \pi, \\ \overline{\nabla_{w^2} f(w^1, w^2)} &= D_{w^2}^H g(w^1, w^2)u - \overline{\pi}. \end{aligned}$$

Conjugating the second equation of (21) and adding to the first yields

$$\overline{\nabla_{w^1} f(w^1, w^2)} + \nabla_{w^2} f(w^1, w^2) = D_{w^1}^H g(w^1, w^2)u + D_{w^2}^T g(w^1, w^2)\overline{u}$$

which is equivalent to (17) since for any feasible point of (14),  $w^2 = \overline{w^1}$ .

Substituting  $\lambda = \begin{pmatrix} \pi \\ -\overline{\pi} \end{pmatrix}$ , and  $w^2 = \overline{w^1}$  in (20) yields

$$(22) \quad \operatorname{Re}(g(w^1, \overline{w^1}), u) + \operatorname{Re}[w^{1H}\pi - w^{1T}\overline{\pi}] = 0.$$

The second term of (22) is zero, and therefore (23) follows. This completes the proof.

**COROLLARY.** *In addition to the hypothesis of Theorem 3, assume that  $g(z, \bar{z})$  is real-valued for all  $z \in C^n$ . Then  $S$  may be assumed to be real, and a necessary condition for a qualified point  $(z^0, \bar{z}^0)$  to be a local minimum of (10) is that there exist a  $v \in [S(g(z^0, \bar{z}^0))]^*$ , where  $S(g(z^0, \bar{z}^0))$  and the dual cone are calculated in  $R^m$ , such that*

$$(23) \quad \overline{\nabla_z f(z^0, \bar{z}^0)} + \nabla_{\bar{z}} f(z^0, \bar{z}^0) = D_z^H g(z^0, \bar{z}^0)v$$

and

$$(24) \quad (g(z^0, \bar{z}^0), v) = 0.$$

*Proof.* If  $S$  contains any complex vectors, it may be replaced by the polyhedral cone  $\hat{S} \equiv \{z : z \in S, \text{Im } z = 0\}$  with no loss in generality. Since  $g(z, \bar{z})$  is real-valued, (A.33) implies  $D_{\bar{z}} g(z^0, \bar{z}^0) = D_z g(z^0, \bar{z}^0)$ . Therefore (17) yields the existence of a vector  $u$  in the complex dual of  $[\hat{S}(g(z^0, \bar{z}^0))]^*$  such that

$$(25) \quad \overline{\nabla_z f(z^0, \bar{z}^0)} + \nabla_{\bar{z}} f(z^0, \bar{z}^0) = D_z^H g(z^0, \bar{z}^0)(u + \bar{u}).$$

The cone  $\hat{S} \subset C^m$  may be represented as the intersection of  $S$  with the cone

$$(26) \quad K = \bigcap_{j=1}^m H_{ie_j} \bigcap_{j=1}^m H_{-ie_j},$$

where  $e_j$  is the  $j$ th unit vector in  $R^m$ . Therefore

$$\hat{S}(g(z^0, \bar{z}^0)) = (S \cap K)(g(z^0, \bar{z}^0)) = S(g(z^0, \bar{z}^0)) \cap K(g(z^0, \bar{z}^0)) = S(g(z^0, \bar{z}^0)) \cap K$$

is a real polyhedral cone. The cones  $\hat{S} \subset C^m$  and  $\hat{S}(g(z^0, \bar{z}^0)) \subset C^m$  which contain only real vectors will be identified with their isomorphic images in  $R^m$ . The same notation,  $\hat{S}$  and  $\hat{S}(g(z^0, \bar{z}^0))$ , will be used in either case, but the meaning will be clear from the context.

If  $\hat{S} \subset C^m$  is written as an intersection of half-spaces; i.e.,  $\hat{S} = \bigcap_{k=1}^p H_{u_k}$ , then in  $R_m$  we have  $\hat{S} = \bigcap_{k=1}^p H_{\text{Re } u_k}$ . For any point  $g(z^0, \bar{z}^0) \in \hat{S}$ ,  $\text{Re}(u_k, g(z^0, \bar{z}^0)) = 0$  if and only if  $(\text{Re } u_k, g(z^0, \bar{z}^0)) = 0$ . Thus  $\hat{S}(g(z^0, \bar{z}^0)) = \bigcap_{k \in B(g(z^0, \bar{z}^0))} H_{u_k}$  is the same as  $\bigcap_{k \in B(g(z^0, \bar{z}^0))} H_{\text{Re } u_k}$ , where the  $H_{\text{Re } u_k}$  are calculated in  $R^m$ .

It follows easily that if  $u \in [\hat{S}(g(z^0, \bar{z}^0))]^*$  calculated in  $C^m$ , then  $v = 2 \text{Re } u$  is an element of  $[\hat{S}(g(z^0, \bar{z}^0))]^*$  calculated in  $R^m$ . Thus (23) holds. Since both  $v$  and  $g(z^0, \bar{z}^0)$  are real, (24) follows from (18). This completes the proof.

If the objective function  $f(z, \bar{z})$  is real-valued for  $z \in C^n$ , then using (A.33) we may replace the sum  $\overline{\nabla_z f(z^0, \bar{z}^0)} + \nabla_{\bar{z}} f(z^0, \bar{z}^0)$  in (17) or (23) by either of its terms. If both  $f(z, \bar{z})$  and  $g(z, \bar{z})$  were obtained from analytic functions  $f(z)$  and  $g(z)$  by substituting  $z = x + iy$  and then  $x = (z + \bar{z})/2$  and  $y = (z - \bar{z})/(2i)$ , then (A.20) implies that  $\nabla_{\bar{z}} f(z, \bar{z}) = D_z g(z, \bar{z}) = 0$ . Using (A.22) and (A.34) we obtain  $\nabla_z f(z, \bar{z}) = \nabla_z f(z)$  and  $D_z g(z, \bar{z}) = D_z g(z)$ , and thus, in this case, (17) reduces to (2).

**5. Example.** A simple example, with  $n = m = 1$ , is used to illustrate the above results:

$$(27) \quad \text{Minimize } \text{Re } z^2 \quad \text{subject to } |z| \leq 1.$$



To put (27) in the form (1), note that  $|z| \leq 1$  may be rewritten :

$$(28) \quad \log z \equiv \log |z| + i \arg z \in S \equiv \text{left half-plane.}$$

However, in (28) the point  $z = 0$  has been mapped to the point at infinity and must therefore be considered separately. By Theorem 1, a necessary condition for the problem :

$$\text{Minimize } \operatorname{Re} z^2 \quad \text{subject to } \log z \in S,$$

to have a minimum at a point  $z$  is that there exist a  $u \in S^*$  such that

$$(29) \quad \overline{2z} = \frac{1}{z}u \quad \text{and} \quad \operatorname{Re}(\log z, u) = 0.$$

In this case  $S^*$  is the negative real axis and thus (29) becomes

$$(30) \quad \bar{z}^2 = \frac{u}{2} \leq 0 \quad \text{and} \quad (\log |z|, u) = 0.$$

The second condition implies that either  $|z| = 1$  or  $u = 0$ . If  $|z| = 1$ , the first condition implies that  $\bar{z}^2$  is negative and therefore  $\operatorname{Re} z = 0$  and  $\operatorname{Im} z = +1$ , which are clearly optimal points of the problem. If  $u = 0$ , then the first condition implies  $z = 0$ , which is a saddle point of the problem.

By writing  $|z| \leq 1$  as  $1 - z\bar{z} \geq 0$ , (27) may alternatively be put in the form of (10) Using the corollary to Theorem 3 and the following paragraph, the necessary conditions are that there exist a  $u \in [S(g(z, \bar{z}))]^* \subset S^* = R_+$  such that

$$(31) \quad \bar{z} = -zu \quad \text{and} \quad (1 - z\bar{z}, u) = 0.$$

If  $u = 0$ , then  $z = 0$ . If  $u > 0$ , then the second constraint gives  $|z| = 1$  and the first implies

$$\bar{z}^2 = -z\bar{z}u = -u < 0,$$

and hence (31) is equivalent to (30).

**Appendix. Identities for differentiation of complex functions.** The definition of analyticity of a function of several variables is usually given in terms of a power series or by means of the Cauchy–Riemann equation in each variable, e.g., [6], [7] and [8]. The power series definition will be used in this paper.

**DEFINITION.** Let  $f: C^n \rightarrow C$  and let  $D$  be a domain in  $C^n$ . Then  $f(z)$  will be called *analytic* in  $D$  if in some neighborhood of every point of  $D$  it may be represented by an absolutely convergent power series about that point in the  $n$  complex variables.

If  $f(z)$  is analytic in  $n$  variables, then it is analytic in each variable separately when the others are held fixed, and hence the Cauchy–Riemann equations hold in each variable, i.e.,

$$(A.1) \quad \nabla_x F^R(x, y) = \nabla_y F^I(x, y), \quad \nabla_x F^I(x, y) = -\nabla_y F^R(x, y),$$

where  $F^R(x, y)$  and  $F^I(x, y)$  are the real and imaginary parts of

$$(A.2) \quad F(x, y) \equiv f(z(x, y)).$$

When  $f(z)$  is analytic it also follows that all partials of all orders of  $F^R(x, y)$  and  $F^I(x, y)$  exist and are continuous. Therefore both  $F^R(x, y)$  and  $F^I(x, y)$  have Fréchet derivatives at points  $(x, y)$ , where  $z = x + iy$  is a point at which  $f(z)$  is analytic.

Applying the chain rule to (A.2) gives

$$(A.3a) \quad \frac{\partial F(x, y)}{\partial x_k} = \frac{\partial f(z)}{\partial z_k} \frac{\partial z_k}{\partial x_k} = \frac{\partial f(z)}{\partial z_k}, \quad k = 1, \dots, n,$$

or

$$(A.3b) \quad \nabla_x F(x, y) = \nabla_z f(z)$$

and

$$(A.4a) \quad \frac{\partial F(x, y)}{\partial y_k} = \frac{\partial f(z)}{\partial z_k} \frac{\partial z_k}{\partial y_k} = i \frac{\partial f(z)}{\partial z_k}, \quad k = 1, \dots, n,$$

$$(A.4b) \quad \nabla_y F(x, y) = i \nabla_z f(z).$$

If in the functions  $F^R(x, y)$  and  $F^I(x, y)$ ,  $(z + \bar{z})/2$  and  $(z - \bar{z})/(2i)$  are formally substituted for  $x$  and  $y$  respectively, “functions” of the “variables”  $z$  and  $\bar{z}$  result. It will be convenient at times to formally differentiate these “functions” with respect to  $z$  or  $\bar{z}$ . The following lemma guarantees that these functions are analytic.

LEMMA A.1. *Let  $f: C^n \rightarrow C$  be analytic in  $D$ . Then there exist functions  $f^R: C^{2n} \rightarrow C$  and  $f^I: C^{2n} \rightarrow C$ , both analytic in a neighborhood of  $(z, \bar{z})$  for any  $z \in D$ , such that on the manifold  $\{(w^1, w^2) \in C^{2n}: w^2 = \overline{w^1}\}$  both functions are real-valued, and*

$$(A.5) \quad f(z) = f^R(z, \bar{z}) + if^I(z, \bar{z}).$$

*Proof.* Let  $z^0 \in D$ . Then  $f(z)$  may be represented as a power series about  $z^0$  which is absolutely convergent in some polycylinder  $P = \{z \in C^n: |z_k - z_k^0| < r_k, k = 1, \dots, n\}$ . Define  $h: C^{2n} \rightarrow C^n$  by  $h(u_1, \dots, u_n, v_1, \dots, v_n) = (u_1 + iv_1, \dots, u_n + iv_n)$ . The function  $h(u, v)$  is clearly analytic and therefore the composition of  $f(z)$  and  $h(u, v)$  is analytic in  $h^{-1}(P)$ , i.e.,

$$(A.6) \quad E(u, v) \equiv f(h(u, v)), \quad \text{where } (u, v) \in C^{2n},$$

is analytic in  $h^{-1}(P)$ . The point  $(u^0, v^0) \in C^{2n}$ , where  $u^0 = \text{Re } z^0$  and  $v^0 = \text{Im } z^0$ , is an element of  $h^{-1}(P)$ , and hence  $E(u, v)$  may be expanded in an absolutely convergent power series about this point; i.e.,

$$(A.7) \quad E(u, v) = \sum_j a_j (u_1 - u_1^0)^{j_1} \dots (u_n - u_n^0)^{j_n} (v_1 - v_1^0)^{j_{n+1}} \dots (v_n - v_n^0)^{j_{2n}},$$

where  $j = (j_1, \dots, j_{2n})$  and each  $j_k$  is a positive integer, converges in a neighborhood of the (real) point  $(u^0, v^0)$  contained in  $h^{-1}(P)$ .

The analytic functions  $E^R(u, v)$  and  $E^I(u, v)$  are defined by replacing  $a_j$  by  $\text{Re } a_j$  and  $\text{Im } a_j$  respectively in (A.7), i.e.,

$$(A.8) \quad E^R(u, v) = \sum_j \text{Re}(a_j) (u_1 - u_1^0)^{j_1} \dots (v_n - v_n^0)^{j_{2n}}$$

and

$$(A.9) \quad E^I(u, v) = \sum_j \text{Im}(a_j) (u_1 - u_1^0)^{j_1} \dots (v_n - v_n^0)^{j_{2n}}.$$

The series (A.8) and (A.9) converge absolutely at any point at which (A.7) converges absolutely, and hence they represent analytic functions at these points. The function  $p: C^{2n} \rightarrow C^{2n}$  defined by

$$(A.10) \quad p(z^1, z^2) = \left( \frac{z^1 + z^2}{2}, \frac{z^1 - z^2}{2i} \right)$$

is now composed with (A.8) and (A.9) to give

$$(A.11) \quad f^R(z^1, z^2) \equiv E^R(p(z^1, z^2)),$$

$$(A.12) \quad f^I(z^1, z^2) \equiv E^I(p(z^1, z^2)).$$

Since  $p(z^1, z^2)$  is analytic in each component on the whole plane,  $f^R(z^1, z^2)$  and  $f^I(z^1, z^2)$  are analytic in the inverse image under  $p$  of the domain of analyticity of  $E^R(u, v)$  and  $E^I(u, v)$  which as noted above contains the domain of analyticity of  $E(u, v)$ . Thus if  $f(z)$  is analytic in a neighborhood of  $z^0$ ,  $E^R(u, v)$  and  $E^I(u, v)$  are analytic in a neighborhood of  $(\text{Re } z^0, \text{Im } z^0) \in C^{2n}$ . The inverse image of this point under the map  $p$  is the point  $(z^0, \bar{z}^0) \in C^{2n}$ , and therefore  $f^R(z^1, z^2)$  and  $f^I(z^1, z^2)$  are analytic in a neighborhood of this point. Since  $p(z, \bar{z})$  is real-valued,  $f^R(z, \bar{z})$  and  $f^I(z, \bar{z})$  are both real-valued. Also  $h(p(z, \bar{z})) = z$ . Therefore

$$(A.13) \quad \begin{aligned} f^R(z, \bar{z}) + if^I(z, \bar{z}) &= E^R(p(z, \bar{z})) + iE^I(p(z, \bar{z})) \\ &= E(p(z, \bar{z})) = f[h(p(z, \bar{z}))] = f(z). \end{aligned}$$

This completes the proof.

A number of identities will now be derived for the analytic function  $f(z)$  and the related functions defined above. With  $f^R(z^1, z^2)$  and  $f^I(z^1, z^2)$  defined by (A.11) and (A.12), let

$$(A.14) \quad f(z^1, z^2) \equiv f^R(z^1, z^2) + if^I(z^1, z^2).$$

Derivatives of left- and right-hand terms of (A.14) will usually be evaluated on the manifold  $\{(z^1, z^2): z^2 = \bar{z}^1\}$ . Thus as in Notations 2.1,  $\partial f(z, \bar{z})/\partial z_k^1$  and  $\partial f(z, \bar{z})/\partial z_k^2$  will be denoted by  $\partial f(z, \bar{z})/\partial z_k$  and  $\partial f(z, \bar{z})/\partial \bar{z}_k$  respectively; i.e., derivatives will be written as though  $z$  and  $\bar{z}$  were independent variables. Derivatives of  $f^R(z^1, z^2)$  and  $f^I(z^1, z^2)$  are written similarly. From (A.11) and (A.12) we have

$$(A.15) \quad f^R(z, \bar{z}) = E^R(p(z, \bar{z})), \quad f^I(z, \bar{z}) = E^I(p(z, \bar{z})), \quad f(z, \bar{z}) = E(p(z, \bar{z})).$$

Denote the components of  $p(z^1, z^2)$  (defined in (A.10)) by  $(p_1^1(z^1, z^2) \cdots p_n^1(z^1, z^2), p_1^2(z^1, z^2) \cdots p_n^2(z^1, z^2))$ . Differentiating the last of equations (A.15) using the chain rule (which may be applied to the composition of analytic functions) gives

$$(A.16) \quad \frac{\partial f(z, \bar{z})}{\partial z_k} = \frac{\partial E}{\partial u_k} \frac{\partial p_k^1}{\partial z_k} + \frac{\partial E}{\partial v_k^2} \frac{\partial p_k^2}{\partial z_k}, \quad k = 1, \dots, n,$$

or

$$(A.17) \quad \nabla_z f(z, \bar{z}) = \frac{1}{2} \nabla_u E(x, y) - \frac{i}{2} \nabla_v E(x, y).$$

The functions  $E(u, v)$ ,  $E^R(u, v)$  and  $E^I(u, v)$  are analytic continuations of  $F(x, y)$ ,  $F^R(x, y)$  and  $F^I(x, y)$  respectively. Therefore  $\nabla_u E(x, y) = \nabla_x F(x, y)$  and  $\nabla_v E(x, y) = \nabla_y F(x, y)$  with similar equations holding for  $E^R(u, v)$  and  $E^I(u, v)$ . Thus (A.17) becomes

$$(A.18) \quad \nabla_z f(z, \bar{z}) = \frac{1}{2} \nabla_x F(x, y) - \frac{i}{2} \nabla_y F(x, y).$$

Taking partials with respect to  $z^2$  instead of  $z^1$  as was done above yields

$$(A.19) \quad \nabla_{\bar{z}} f(z, \bar{z}) = \frac{1}{2} \nabla_x F(x, y) + \frac{i}{2} \nabla_y F(x, y).$$

Expanding the right-hand side of (A.19) gives, as an alternate form of the Cauchy–Riemann equations,

$$(A.20) \quad \begin{aligned} \nabla_{\bar{z}} f(z, \bar{z}) &= \frac{1}{2} [\nabla_x F^R(x, y) + i \nabla_x F^I(x, y) + i \nabla_y F^R(x, y) - \nabla_y F^I(x, y)] \\ &= 0 \quad (\text{by (A.1)}). \end{aligned}$$

A similar application of the chain rule to  $f^R(z, \bar{z})$  instead of  $f(z, \bar{z})$  gives

$$(A.21) \quad 2 \nabla_z f^R(z, \bar{z}) = \nabla_x F^R(x, y) - i \nabla_y F(x, y),$$

which, with (A.1) and (A.3), yields the interesting formula

$$(A.22) \quad 2 \nabla_z f^R(z, \bar{z}) = \nabla_z f(z).$$

For the applications of this paper the domain of interest of a function of  $2n$  complex variables will be the manifold  $\{(w^1, w^2) : w^2 = \overline{w^1}\}$ . The following lemma for functions of  $2n$  complex variables is analogous to Lemma A.1.

LEMMA A.2. *Let  $f: C^{2n} \rightarrow C$  be analytic in a domain  $D$  containing a point  $(z^0, \bar{z}^0) \in C^{2n}$ . Then there exist functions  $f^R: C^{2n} \rightarrow C$  and  $f^I: C^{2n} \rightarrow C$  analytic in a neighborhood of  $(z^0, \bar{z}^0)$  such that on the manifold  $\{(w^1, w^2) : w^2 = \overline{w^1}\}$ ,  $f^R(w^1, w^2)$  and  $f^I(w^1, w^2)$  are real-valued and*

$$(A.23) \quad f(w^1, w^2) = f^R(w^1, w^2) + if^I(w^1, w^2).$$

*Outline of proof.* Expand  $f(w^1, w^2)$  in an absolutely convergent series about  $(z^0, \bar{z}^0)$ . Define  $\bar{f}(w^1, w^2)$  by

$$(A.24) \quad \bar{f}(w^1, w^2) = \overline{f(w^2, w^1)}.$$

Then the series  $\bar{f}(w^1, w^2)$  also converges, and hence  $\bar{f}(w^1, w^2)$  is analytic in a neighborhood of  $(z^0, \bar{z}^0)$ . The required functions are then defined by

$$(A.25) \quad f^R(w^1, w^2) = \frac{f(w^1, w^2) + \bar{f}(w^1, w^2)}{2}$$

and

$$(A.26) \quad f^I(w^1, w^2) = \frac{f(w^1, w^2) - \bar{f}(w^1, w^2)}{2i}.$$

They are analytic and, since  $\overline{f(z, \bar{z})} = \bar{f}(z, \bar{z})$ , real-valued on  $\{(w^1, w^2) : w^2 = \overline{w^1}\}$ . This completes the proof.

Identities for analytic functions of  $2n$  variables are now obtained. Let  $f: C^{2n} \rightarrow C$  be analytic in a neighborhood of  $(z, \bar{z})$ . Define  $h: R^{2n} \rightarrow C^n$  by  $h(u, v) = u + iv$  and  $H: R^{4n} \rightarrow C$  by  $H(x^1, y^1, x^2, y^2) = f(h(x^1, y^1), h(x^2, y^2))$ . In addition, define  $g_1(x) = x$  and  $g_2(x) = -x$  for  $x \in R^n$ . Then, what is written informally as  $F(x, y) = f(z(x, y), \bar{z}(x, y))$  may be defined more precisely as

$$(A.27) \quad F(x, y) \equiv H(x, y, g_1(x), g_2(y)) = f(h(x, y), h(g_1(x), g_2(y))).$$

$F(x, y)$  is the composition of a complex function of complex variables and a complex function of real variables both of which are differentiable. Hence the chain rule may be applied to (A.27) to give

$$(A.28) \quad \begin{aligned} \frac{\partial F(x, y)}{\partial x_k} &= \frac{\partial f(z, \bar{z})}{\partial w_k^1} \frac{\partial h}{\partial u_k} + \frac{\partial f(z, \bar{z})}{\partial w_k^2} \frac{\partial h}{\partial u_k} \frac{\partial g}{\partial x_k} \\ &= \frac{\partial f(z, \bar{z})}{\partial z} + \frac{\partial f(z, \bar{z})}{\partial \bar{z}}, \end{aligned} \quad k = 1, \dots, n,$$

where  $\partial f(z, \bar{z})/\partial z_k$  and  $\partial f(z, \bar{z})/\partial \bar{z}_k$  have been used to denote  $\partial f(z, \bar{z})/\partial z_k^1$  and  $\partial f(z, \bar{z})/\partial z_k^2$  respectively.

In a similar manner we obtain

$$(A.29) \quad \frac{\partial F(x, y)}{\partial y_k} = i \frac{\partial f(z, \bar{z})}{\partial z_k} - i \frac{\partial f(z, \bar{z})}{\partial \bar{z}_k}, \quad k = 1, \dots, n,$$

or rewriting (A.28) and (A.29) in gradient notation,

$$(A.30a) \quad \nabla_x F(x, y) = \nabla_z f(z, \bar{z}) + \nabla_{\bar{z}} f(z, \bar{z}),$$

$$(A.30b) \quad \nabla_y F(x, y) = i \nabla_z f(z, \bar{z}) - i \nabla_{\bar{z}} f(z, \bar{z}).$$

Multiplying (A.30a) by  $i$  and adding to (A.30b) gives

$$(A.31) \quad \nabla_z f(z, \bar{z}) = \frac{1}{2} \nabla_x F(x, y) - \frac{i}{2} \nabla_y F(x, y),$$

while multiplying (A.30b) by  $i$  and adding to (A.30a) gives

$$(A.32) \quad \nabla_{\bar{z}} f(z, \bar{z}) = \frac{1}{2} \nabla_x F(x, y) + \frac{i}{2} \nabla_y F(x, y).$$

Since (A.30)–(A.32) hold for any analytic function of  $2n$  variables, they hold for the  $f^R(w^1, w^2)$  and  $f^I(w^1, w^2)$  defined in Lemma A.1, and hence (A.31) and (A.32) agree with (A.18) and (A.19). If the function is real-valued, e.g.,  $f^R(z, \bar{z})$  or  $f^I(z, \bar{z})$  as defined in Lemmas A.1 or A.2, then  $F(x, y)$  is real-valued; and therefore, (A.31) and (A.32) imply

$$(A.33) \quad \overline{\nabla_z f(z, \bar{z})} = \nabla_{\bar{z}} f(z, \bar{z}).$$

A final identity is now obtained. First differentiate (A.23) with respect to  $w^1$  and then with respect to  $w_2$  to give

$$(a) \quad \nabla_z f(z, \bar{z}) = \nabla_z f^R(z, \bar{z}) + i \nabla_z f^I(z, \bar{z})$$

and

$$(b) \quad \nabla_z f(z, \bar{z}) = \nabla_z f^R(z, \bar{z}) + i \nabla_z f^I(z, \bar{z}).$$

Conjugating (a), adding (b) and using (A.33) yields

$$(A.34) \quad \overline{\nabla_z f(z, \bar{z})} + \nabla_z f(z, \bar{z}) = 2 \nabla_z f^R(z, \bar{z}).$$

#### REFERENCES

- [1] R. A. ABRAMS, *Nonlinear programming in complex space*, Doctoral thesis in Applied Mathematics, Northwestern Univ., Evanston, Ill., 1969.
- [2] ———, *Nonlinear programming in complex space: Sufficient conditions and duality*, J. Math. Anal. Appl., to appear.
- [3] R. A. ABRAMS AND A. BEN-ISRAEL, *A duality theorem for complex quadratic programming*, J. Optimization Theory Appl., 4 (1969), pp. 244–252.
- [4] ———, *Complex mathematical programming*, Developments in Operations Research, B. Avi-Itzhak, ed., Gordon and Breach, New York, 1971.
- [5] A. BEN-ISRAEL, *Linear equations and inequalities on finite dimensional, real or complex, vector spaces: A unified theory*, J. Math. Anal. Appl., 27 (1969), pp. 367–389.
- [6] S. BOCHNER AND W. T. MARTIN, *Several Complex Variables*, Princeton University Press, Princeton, N.J., 1948.
- [7] H. J. BREMERMAN, *Several complex variables*, Studies in Real and Complex Analysis, I. Hirschman, ed., The Mathematical Association of America, Washington, D.C., 1965.
- [8] J. DIEUDONNÉ, *Foundations of Modern Analysis*, Academic Press, New York, 1960.
- [9] J. FARKAS, *Über die Theorie der einfachen Ungleichungen*, J. Reine Angew. Math., 124 (1902), pp. 1–24.
- [10] A. V. FIACCO AND G. P. MCCORMICK, *Nonlinear Programming: Sequential Unconstrained Minimization Techniques*, John Wiley, New York, 1968.
- [11] M. A. HANSON AND B. MOND, *Duality for nonlinear programming in complex space*, J. Math. Anal. Appl., 28 (1969), pp. 52–58.
- [12] H. W. KUHN AND A. W. TUCKER, *Non-linear programming*, Proc. Second Berkeley Symposium on Mathematical Statistics and Probability, J. Neyman, ed., University of California Press, Berkeley, Calif., 1951, pp. 481–493.
- [13] N. LEVINSON, *Linear programming in complex space*, J. Math. Anal. Appl., 14 (1966), pp. 44–62.
- [14] H. WEYL, *The elementary theory of convex polyhedra*, Contributions to the Theory of Games, vol. I, Annals of Mathematics Studies No. 24, H. W. Kuhn and A. W. Tucker, eds., Princeton University Press, Princeton, N.J., 1950, pp. 3–18.

## REPRESENTATION OF MARTINGALES, QUADRATIC VARIATION AND APPLICATIONS\*

E. WONG†

**Abstract.** In this paper, we present two related results. First, we shall obtain a sufficient condition under which a second order sample-continuous martingale can be represented as a stochastic integral in terms of a Brownian motion. Secondly, we shall show that if  $X$  and  $Y$  are sample-continuous local martingales (not necessarily with respect to the same family of  $\sigma$ -algebras) and if either  $X + Y$  or  $X - Y$  is almost surely of bounded variation, then the quadratic variations of the two martingales are equal. This rather simple result has some surprising consequences.

**1. Introduction.** Let  $\{X_t, t \geq 0\}$  be a sample-continuous second order martingale. Then  $\{X_t^2, t \geq 0\}$  is a sample-continuous first order submartingale and the conditions for Meyer's decomposition [1] are always satisfied; thus we can write

$$(1) \quad X_t^2 = M_t + A_t, \quad t \geq 0,$$

where  $M$  is a martingale,  $A$  is an increasing process, and both are sample-continuous. The decomposition is unique if we set  $M_0 = X_0^2$ . If  $X$  is a local martingale [2], (1) remains valid except now  $M$  is a local martingale. Following Kunita and Watanabe [2], we shall adopt the suggestive notation  $\langle X \rangle_t$  for the increasing process  $A_t$ .

In this paper, we present two related results. First, we shall obtain a sufficient condition under which a second order sample-continuous martingale can be represented as a stochastic integral in terms of a Brownian motion. Second, we shall show that if  $X$  and  $Y$  are sample-continuous local martingales (not necessarily with respect to the same family of  $\sigma$ -algebras) and if either  $X + Y$  or  $X - Y$  is almost surely of bounded variation, then  $\langle X \rangle_t = \langle Y \rangle_t$ . This rather simple result has some surprising consequences.

**2. Martingales and stochastic integrals.** Let  $(\Omega, \mathcal{A}, \mathcal{P})$  be a probability space, and let  $\{\mathcal{A}_t, t \geq 0\}$  be an increasing family of sub- $\sigma$ -algebras. A process  $\{X_t, t \geq 0\}$  is said to be *adapted* to  $\{\mathcal{A}_t\}$  if, for each  $t$ ,  $X_t$  is  $\mathcal{A}_t$ -measurable. We say that  $\{X_t, \mathcal{A}_t\}$  is a *martingale* if  $X$  is adapted to  $\{\mathcal{A}_t\}$  and for every  $t > s$ ,

$$(2) \quad E^{\mathcal{A}_s} X_t = X_s$$

almost surely. If  $\{X_t, \mathcal{A}_t\}$  is a sample-continuous second order martingale, then the increasing process  $\langle X \rangle_t$  introduced earlier is well-defined and  $E\langle X \rangle_t < \infty$ .

A nonnegative random variable  $\tau$  is said to be a *stopping time* of  $\{\mathcal{A}_t\}$  if  $\{\omega : \tau(\omega) \leq t\} \in \mathcal{A}_t$  for every  $t$ . A process  $\{X_t, \mathcal{A}_t\}$  is said to be a *local martingale* if there exists an increasing sequence of stopping times  $\{\tau_n\}$  such that  $\tau_n \uparrow \infty$

---

\* Received by the editors September 15, 1970, and in revised form January 20, 1971.

† Department of Electrical Engineering and Computer Sciences and the Electronics Research Laboratory, University of California, Berkeley, California 94720. This work was supported by the United States Army Research Office, Durham, under Contract DAHC04-67-0046 and the National Science Foundation under Grant GK-1065X.

almost surely and  $\{X_{\min(t, \tau_n)}, \mathcal{A}_t\}$  is a second order martingale. Let  $t \wedge s$  denote  $\min(t, s)$  and set

$$X_{n,t} = X_{t \wedge \tau_n}.$$

Kunita and Watanabe have proved [2] that if  $\{X_t, \mathcal{A}_t\}$  is a sample-continuous local martingale, then there exists a sample-continuous increasing process  $\langle X \rangle_t$  such that

$$\langle X \rangle_{t \wedge \tau_n} = \langle X_n \rangle_t.$$

Of course, now we need not have  $E\langle X \rangle_t < \infty$ .

If  $\{W_t, \mathcal{A}_t\}$  is a sample-continuous second order martingale such that for  $t > s$ ,

$$(3) \quad E^{\mathcal{A}_s}(W_t - W_s)^2 = t - s,$$

then  $W$  is necessarily a Brownian motion and for each  $s$ ,  $\{W_t - W_s, t \geq s\}$  is independent of  $\mathcal{A}_s$  [3, p. 384]. We describe this situation by saying that  $\{W_t, \mathcal{A}_t\}$  is a Brownian motion. Let  $\{W_t, \mathcal{A}_t\}$  be a Brownian motion and let  $\{\phi_t, t \geq 0\}$  be a measurable process adapted to  $\{\mathcal{A}_t\}_t$  such that

$$(4) \quad \int_0^t E\phi_s^2 ds < \infty$$

for each  $t$ . The stochastic integral  $\int_0^t \phi_s dW_s$  is well-defined as the quadratic limit of a sequence of sums  $\sum_v \phi_{t_v^{(n)}}[W_{t_{v+1}^{(n)}} - W_{t_v^{(n)}}]$ , where  $\{t_v^{(n)}\}$  is a sequence of partitions of  $[0, t]$  such that

$$\max_v (t_{v+1}^{(n)} - t_v^{(n)}) \xrightarrow{n \rightarrow \infty} 0.$$

If we define

$$(5) \quad X_t = \int_0^t \phi_s dW_s$$

and choose a separable version for  $X$ , then  $\{X_t, \mathcal{A}_t\}$  is a second order sample-continuous martingale, with

$$(6) \quad E^{\mathcal{A}_s}(X_t - X_s)^2 = \int_s^t E^{\mathcal{A}_s}\phi_\tau^2 d\tau, \quad 0 \leq s \leq t.$$

If, instead of (4),  $\phi$  merely satisfies

$$(7) \quad \int_0^t \phi_s^2 ds < \infty, \quad \text{a.s.},$$

then  $\int_0^t \phi_s dW_s$  can be defined as follows: Let  $\tau_n(\omega)$  be defined by

$$(8) \quad \tau_n(\omega) = \begin{cases} \inf \left\{ t : \int_0^t \phi_s^2(\omega) ds \geq n \right\} \\ \infty \quad \text{if } \int_0^t \phi_s^2(\omega) ds < n \text{ for all } t, \end{cases}$$



and set

$$(9) \quad \phi_{n,s}(\omega) = \begin{cases} \phi_s(\omega), & s \leq \tau_n(\omega), \\ 0, & s > \tau_n(\omega). \end{cases}$$

For each  $n$ ,  $\int_0^t \phi_{n,s} dW_s$  is well-defined. It can be shown that  $\int_0^t \phi_{n,s} dW_s$  converges in probability as  $n \rightarrow \infty$ , and we define  $\int_0^t \phi_s dW_s$  as this limit. Now, the process  $X_t = \int_0^t \phi_s dW_s$  need no longer be second order or a martingale, but it is still sample-continuous if a separable version is chosen. Moreover, for each  $n$ ,  $\{X_{t \wedge \tau_n}, \mathcal{A}_t\}$  is a second order martingale. By definition,  $X$  is a local martingale.

If  $X$  is a stochastic integral of the form

$$(10) \quad X_t = X_0 + \int_0^t \phi_s dW_s$$

and  $f$  is a twice continuously differentiable function of a real variable, then Ito's differentiation formula [4] yields

$$(11) \quad f(X_t) = f(X_0) + \int_0^t f'(X_s)\phi_s dW_s + \frac{1}{2} \int_0^t f''(X_s)\phi_s^2 ds.$$

In particular,

$$(12) \quad X_t^2 - X_0^2 = 2 \int_0^t X_s \phi_s dW_s + \int_0^t \phi_s^2 ds.$$

We can now identify  $X_0^2 + 2 \int_0^t X_s \phi_s dW_s$  as the local martingale term in the decomposition (1) of  $X_t^2$ , and thus the increasing process is given by

$$(13) \quad \langle X \rangle_t = \int_0^t \phi_s^2 ds.$$

If  $X$  is of the form (10), then we can define the stochastic integral  $\int_0^t \psi_s dX_s$  by

$$\int_0^t \psi_s dX_s = \int_0^t \psi_s \phi_s dW_s,$$

provided that  $\int_0^t \psi_s^2 \phi_s^2 ds < \infty$  almost surely. If  $X$  is a local martingale, not necessarily of the form (10), the stochastic integral  $\int_0^t \phi_s dX_s$  can still be defined, provided that  $\int_0^t \phi_s^2 d\langle X \rangle_s < \infty$  almost surely [2]. Finally, if  $X = Y + Z$  where  $Y$  is a local martingale and  $\{Z_t, t \geq 0\}$  is a process with sample functions almost surely of bounded variation, then we can define  $\int_0^t \psi_s dX_s$  by

$$\int_0^t \psi_s dX_s = \int_0^t \psi_s dY_s + \int_0^t \psi_s dZ_s,$$

provided that the first integral exists as a stochastic integral and the second as a Stieltjes integral. For the case where  $Y$  and  $Z$  are sample continuous and  $f$  is a

twice continuously differentiable function, Kunita and Watanabe has extended Ito's differentiation rule to read

$$f(X_t) = f(X_0) + \int_0^t f'(X_s) dX_s + \frac{1}{2} \int_0^t f''(X_s) d\langle Y \rangle_s.$$

In particular, we note that

$$X_t^2 - X_0^2 - 2 \int_0^t X_s dX_s = \langle Y \rangle_t$$

is independent of  $Z$  and serves to define  $\langle X \rangle_t$ .

**3. Representation of martingales.** Not every sample-continuous second order martingale can be represented as a stochastic integral in the form of (10). It is clear from (13) that for such a representation to be possible the increasing process  $\langle X \rangle(\omega, t)$  must be an absolutely continuous function of  $t$  (w.r.t. the Lebesgue measure) for almost all  $\omega$ . As Fisk has observed [5], this condition is also sufficient by virtue of a theorem of Doob [3, p. 449], but it may be necessary to enlarge the underlying probability space by the adjunction of a Brownian motion. Specifically, Doob proved the following theorem.

**THEOREM 3.1 (Doob).** *Let  $\{X_t, \mathcal{A}_t, 0 \leq t \leq T\}$  be a sample-continuous second order martingale. Suppose that there exists a nonnegative measurable process  $\{\psi_t, 0 \leq t \leq T\}$  adapted to  $\{\mathcal{A}_t\}$  such that for  $t > s$ ,*

$$(14) \quad E^{\mathcal{A}_s}(X_t - X_s)^2 = \int_s^t E^{\mathcal{A}_s} \psi_\tau d\tau.$$

*If the set  $\{(\omega, t) : \psi(\omega, t) = 0\}$  has zero  $d\mathcal{P} dt$  measure, then there exists a Brownian motion  $\{W_t, \mathcal{A}_t, 0 \leq t \leq T\}$  such that*

$$(15) \quad X_t = X_0 + \int_0^t \psi_s^{1/2} dW_s$$

*with probability 1. Without the hypothesis that  $\psi$  vanishes almost nowhere, representation (15) is still valid with the adjunction of a Brownian motion to the probability space.*

The condition that  $\langle X \rangle$  be almost surely continuous with respect to the Lebesgue measure is both somewhat stringent and difficult to verify. Perhaps, it is more natural to consider representations of a more general form

$$(16) \quad X_t(\omega) = X_0(\omega) + \int_0^t \phi_s(\omega) dW_{F(s)}(\omega),$$

where  $W$  is a Brownian motion and  $F$  is a continuous increasing function.

**THEOREM 3.2.** *Let  $\{X_t, \mathcal{A}_t, 0 \leq t \leq T\}$  be a sample-continuous local martingale. We assume that  $\{\mathcal{A}_t\}$  is right-continuous (i.e.  $\bigcap_{s>t} \mathcal{A}_s = \mathcal{A}_t$ ) and each  $\mathcal{A}_t$  is completed. A representation of the form (16) exists if and only if  $\langle X \rangle$  is absolutely continuous with respect to  $F$  with probability 1.*

*Proof.* We only need to prove the theorem for the case  $F(t) = t$ , because  $\{X_t, \mathcal{A}_t\}$  can be transformed into a sample-continuous local martingale  $\{\tilde{X}_t, \tilde{\mathcal{A}}_t\}$  with  $E\tilde{X}_t^2 = t$  by defining

$$\begin{aligned} F^{-1}(t) &= \inf \{s : F(s) = t\}, \\ \tilde{X}_t &= X_{F^{-1}(t)}, \\ \tilde{\mathcal{A}}_t &= \mathcal{A}_{F^{-1}(t)}. \end{aligned}$$

Even though  $F^{-1}$  may be discontinuous,  $\{\tilde{\mathcal{A}}_t\}$  is right-continuous and  $\tilde{X}$  is still sample-continuous, because  $F(t) = F(s)$  implies  $X_t = X_s$  almost surely. Since  $X_t = \tilde{X}_{F(t)}$  with probability 1, a representation

$$\tilde{X}_t = \tilde{X}_0 + \int_0^t \tilde{\phi}_s dW_s,$$

where  $\{W_t, \tilde{\mathcal{A}}_t\}$  is a Brownian motion, implies a representation

$$\begin{aligned} X_t &= \tilde{X}_{F(t)} = X_0 + \int_0^{F(t)} \tilde{\phi}_s dW_s \\ &= X_0 + \int_0^t \phi_s dW_{F(s)}, \end{aligned}$$

which is just (16).

To prove Theorem 3.2 for the case  $F(t) = t$ , we first note that necessity follows from (13). To prove sufficiency, we assume that  $\langle X \rangle$  is absolutely continuous with probability 1 and write

$$(17) \quad \langle X \rangle_t(\omega) = \int_0^t \psi_s(\omega) ds,$$

where  $\psi$  can always be chosen to be a measurable process because  $\langle X \rangle$  is a measurable process. For each  $t$ ,  $\psi_t$  is measurable with respect to  $\cap_{s>t} \mathcal{A}_s$ , which is equal to  $\mathcal{A}_t$  by assumption.

We now follow Doob [3, p. 449] and define

$$(18) \quad g_s(\omega) = \begin{cases} \psi_s^{-1/2}(\omega) & \text{if } \psi_s(\omega) > 0, \\ 0 & \text{if } \psi_s(\omega) = 0. \end{cases}$$

Since  $\int_0^t g_s^2 d\langle X \rangle_s \leq t < \infty$ , the integral  $\int_0^t g_s dX_s = W_t$  is well-defined as a local martingale. If  $\psi_s(\omega) > 0$  for almost all  $(\omega, s)$ , then  $\langle W \rangle_t = t$  and it follows from Theorem 2.3 of [2] that  $W$  is a Brownian motion. If not, we adjoin an independent Brownian motion  $B$  to the underlying probability space and define

$$W_t = \int_0^t g_s dX_s + \int_0^t [1 - g_s \sqrt{\psi_s}] dB_s.$$

Then,  $W$  is a Brownian motion. In either case, we have

$$X_t = X_0 + \int_0^t \sqrt{\psi_s} dW_s$$

and the proof is complete.

Theorem 3.2 is basically the same as Theorem 2.1 of [5], except for the introduction of the increasing function  $F$  and the generalization to local martingales. We now come to the main result of the section, namely, a sufficient condition for the representation (16) that can be verified in terms of two-dimensional distributions of a martingale  $X$ .

**THEOREM 3.3.** *Let  $\{X_t, \mathcal{A}_t, 0 \leq t \leq T\}$  be a sample-continuous second order martingale and let  $F(t) = E(X_t - X_0)^2 = E\langle X \rangle_t$ . Suppose that there exist finite positive constants  $\alpha$  and  $\beta$  such that*

$$(19) \quad \sup_{0 < F(t) - F(s) \leq \beta} \frac{E|X_t - X_s|^{2+2\alpha}}{[F(t) - F(s)]^{1+\alpha}} < \infty.$$

Then  $\langle X \rangle$  is almost surely absolutely continuous with respect to the  $F$ -measure, and  $X$  has a representation of the form of (16).

*Proof.* By virtue of the Lebesgue decomposition, we can always write

$$(20) \quad \langle X \rangle_t(\omega) = \int_0^t \psi_s(\omega) dF(s) + \mu_t(\omega),$$

where  $\mu$  is almost surely singular with respect to  $F$ . Now,  $E\psi_s = 1$  implies

$$E\mu_t = - \int_0^t dF(t) + E\langle X \rangle_t = 0,$$

which in turn implies that  $\mu \equiv 0$  almost surely since  $\mu$  is nonnegative and sample-continuous. Therefore, we only need to prove that (19) implies  $E\psi_s = 1, 0 \leq s \leq T$ .

Let  $T_n = \{t_v^{(n)}, v = 0, 1, \dots, n\}$  be a sequence of nested (i.e.,  $T_{n+1} \supset T_n$ ) partitions of the interval  $[0, T]$  such that

$$\max_v [F(t_{v+1}^{(n)}) - F(t_v^{(n)})] \xrightarrow{n \rightarrow \infty} 0.$$

Define  $\psi_{n,t}, 0 \leq t \leq T$ , as follows:

$$(21) \quad \psi_{n,t} = \frac{\langle X \rangle_{t_{v+1}^{(n)}} - \langle X \rangle_{t_v^{(n)}}}{F(t_{v+1}^{(n)}) - F(t_v^{(n)})}, \quad t_v^{(n)} \leq t < t_{v+1}^{(n)}.$$

It is well known (see, e.g., [3, pp. 346–347]) that for each  $\omega, \psi_{n,t}$  converges for almost all  $t$  ( $F$ -measure) to the Radon-Nikodym derivative of the absolutely continuous component of  $\langle X \rangle$  with respect to  $F$ . That is,  $\psi_{n,t} \rightarrow \psi_t$  for almost all  $(\omega, t)$ . Since it is obvious that  $E\psi_{n,t} = 1$ , the desired result  $E\psi_t = 1$  will follow if for each  $t, \{\psi_{n,t}\}$  is a uniformly integrable family of random variables.

Now, it is known [6] that for any  $p > 1/2$  there exists a constant  $\kappa_p$  such that

$$E|\langle X \rangle_t - \langle X \rangle_s|^p \leq \kappa_p E|X_t - X_s|^{2p}.$$

Therefore, if we let  $N$  be the smallest  $n$  such that

$$\max_v [F(t_{v+1}^{(n)}) - F(t_v^{(n)})] \leq \beta,$$

then

$$\sup_{n \geq N} E\psi_{n,t}^{1+\alpha} \leq \kappa_{1+\alpha} \sup_{0 < F(t) - F(s) \leq \beta} \left\{ \frac{E|X_t - X_s|^{2+2\alpha}}{[F(t) - F(s)]^{1+\alpha}} \right\} < \infty,$$

so that  $\{\psi_{n,t}\}$  is a uniformly integrable family of random variables. This, together with Theorem 3.2, completes the proof.

Theorem 3.2 is reminiscent of Kolmogorov's condition for sample continuity and has similar advantages, the primary one being that it can be verified in terms of the two-dimensional distributions of  $X$ .

**4. Quadratic variation.** Let  $T_n = \{t_v^{(n)}\}$  be a nested sequence of partitions of  $[0, T]$  such that  $\max_v (t_v^{(n)} - t_{v+1}^{(n)}) \xrightarrow{n \rightarrow \infty} 0$ . Let  $X$  be a sample-continuous second order martingale, and let  $t \wedge s$  denote  $\min(t, s)$ . Fisk [5] has shown that the sequence of sums

$$(22) \quad Q_n(X, t) = \sum_v [X_{t \wedge t_{v+1}^{(n)}} - X_{t \wedge t_v^{(n)}}]^2$$

converges to  $\langle X \rangle_t$  in  $L^1$ -mean, i.e.,

$$(23) \quad E|Q_n(X, t) - \langle X \rangle_t| \xrightarrow{n \rightarrow \infty} 0.$$

For this reason,  $\langle X \rangle_t$  is said to be the quadratic variation of  $X$  on  $[0, t]$ . Now, suppose that  $\{Z_t, 0 \leq t \leq T\}$  is a sample-continuous process, the sample functions of which are almost surely of bounded variation. Then there exists an almost surely finite random variable  $A$  such that

$$\sup_n \sum_v |Z_{t_{v+1}^{(n)}} - Z_{t_v^{(n)}}| \leq A.$$

Therefore,

$$Q_n(Z, T) = \sum_v |Z_{t_{v+1}^{(n)}} - Z_{t_v^{(n)}}|^2 \leq A \max_v |Z_{t_{v+1}^{(n)}} - Z_{t_v^{(n)}}| \xrightarrow[n \rightarrow \infty]{a.s.} 0$$

so that  $Z$  has zero quadratic variation on  $[0, T]$ .

**THEOREM 4.1.** *Let  $\{\mathcal{A}_t\}$  and  $\{\tilde{\mathcal{A}}_t\}$  be two increasing families of  $\sigma$ -algebras and let  $\{X_t, \mathcal{A}_t\}$  and  $\{\tilde{X}_t, \tilde{\mathcal{A}}_t\}$  be sample-continuous local martingales. If  $X + \tilde{X}$  or  $X - \tilde{X}$  is of bounded variation, then  $\langle X \rangle_t = \langle \tilde{X} \rangle_t$  almost surely for every  $t$ .*

*Proof.* First, suppose that  $X$  and  $\tilde{X}$  are second order martingales. Let  $B$  denote the process of bounded variation given by  $X - \tilde{X}$  or  $X + \tilde{X}$ . From the inequality

$$\left( \sqrt{\sum_k a_k^2} - \sqrt{\sum_k b_k^2} \right)^2 \leq \sum_k (a_k \pm b_k)^2 \leq \left( \sqrt{\sum_k a_k^2} + \sqrt{\sum_k b_k^2} \right)^2,$$

we have

$$\left( \sqrt{Q_n(\tilde{X}, t)} - \sqrt{Q_n(B, t)} \right)^2 \leq Q_n(X, t) \leq \left( \sqrt{Q_n(\tilde{X}, t)} + \sqrt{Q_n(B, t)} \right)^2.$$

Since  $Q_n(\tilde{X}, t)$ ,  $Q_n(X, t)$  and  $Q_n(B, t)$  converge respectively to  $\langle \tilde{X} \rangle_t$ ,  $\langle X \rangle_t$  and 0, we have  $\langle X \rangle_t = \langle \tilde{X} \rangle_t$  almost surely for all  $t \in [0, T]$ .

If  $X$  and  $\tilde{X}$  are local martingales, then there exist sequences of stopping times  $\{\tau_m\}$  and  $\{\tilde{\tau}_m\}$  increasing to  $\infty$ , so that  $X_{m,t} = X_{\tau_m \wedge t}$  and  $\tilde{X}_{m,t} = \tilde{X}_{\tilde{\tau}_m \wedge t}$  are second order martingales for each  $m$ . Let

$$\Omega_{m,t} = \{\omega : t < \tau_m \wedge \tilde{\tau}_m\}.$$

Since  $\tau_m \uparrow \infty$  and  $\tilde{B}_m \uparrow \infty$ , we have  $\Omega_{m,t} \uparrow \Omega$ . For  $\omega \in \Omega_{m,t}$ , we have

$$\begin{aligned} X_{m,s}(\omega) &= X_s(\omega), & 0 \leq s \leq t, \\ \tilde{X}_{m,s}(\omega) &= \tilde{X}_s(\omega), & 0 \leq s \leq t. \end{aligned}$$

Therefore, for  $\omega \in \Omega_{m,t}$ , either  $X_m(\omega) + \tilde{X}_m(\omega)$  or  $X_m(\omega) - \tilde{X}_m(\omega)$  is of bounded variation on  $[0, t]$ , and it follows from the same inequality as before that

$$\langle X_m \rangle_t(\omega) = \langle \tilde{X}_m \rangle_t(\omega), \quad \omega \in \Omega_{m,t},$$

or

$$\langle X \rangle_t(\omega) = \langle \tilde{X} \rangle_t(\omega), \quad \omega \in \Omega_{m,t}.$$

The proof is completed by letting  $m \uparrow \infty$ .

We should note that in Theorem 4.1 we do not assume that  $X$  and  $\tilde{X}$  are local martingales with respect to the same family of  $\sigma$ -algebras. If they are, then  $X + \tilde{X}$  and  $X - \tilde{X}$  are both local martingales. If one of them is also of bounded variation, say  $X + \tilde{X}$ , then  $X_t + \tilde{X}_t = X_0 + \tilde{X}_0$  with probability 1 for all  $t \in [0, T]$ , in which case the result of Theorem 4.1 trivially follows. The more interesting cases arise when neither  $X + \tilde{X}$  nor  $X - \tilde{X}$  is a local martingale.

An interesting application of Theorem 4.1 is in connection with quasi-martingales [7]. A process  $\{X_t, 0 \leq t \leq T\}$  is said to be a quasi-martingale with respect to  $\{\mathcal{A}_t\}$  if there exist  $\{B_t, 0 \leq t \leq T\}$  and  $\{M_t, 0 \leq t \leq T\}$  both adapted to  $\{\mathcal{A}_t\}$  such that  $X = M + B$ ,  $B$  is of bounded variation, and  $\{M_t, \mathcal{A}_t\}$  is a martingale. We shall be interested only in those cases where both  $M$  and  $B$  are sample-continuous and where the total variation of  $B$  has a finite expectation. Under these assumptions, if  $\{X_t, \mathcal{A}_t\}$  is a quasi-martingale, then  $\{X_t, \mathcal{A}_{xt}\}$  is always a quasi-martingale, where  $\mathcal{A}_{xt}$  denotes the  $\sigma$ -algebra generated by  $\{X_s, 0 \leq s \leq t\}$ . More generally, let  $\{\mathcal{A}_t\}$  be any increasing family of  $\sigma$ -algebras such that  $\mathcal{A}_{xt} \subseteq \tilde{\mathcal{A}}_t \subseteq \mathcal{A}_t$  for every  $t$ . Then,  $\{X_t, \tilde{\mathcal{A}}_t\}$  is a quasi-martingale.<sup>1</sup> That is, there exist  $\tilde{B}$  and  $\tilde{M}$  adapted to  $\{\tilde{\mathcal{A}}_t\}$  such that  $B$  is of bounded variation.  $\{\tilde{M}_t, \tilde{\mathcal{A}}_t\}$  is a martingale, both  $\tilde{M}$  and  $\tilde{B}$  are sample-continuous, and the total variation of  $\tilde{B}$  has finite expectation. It follows from Theorem 4.1 that  $\langle M \rangle_t = \langle \tilde{M} \rangle_t$  for every  $t$ .

An important class of quasi-martingales is made up of Ito processes, which are processes having the representation

$$(24) \quad X_t = X_0 + \int_0^t \psi_s ds + \int_0^t \phi_s dW_s, \quad 0 \leq t \leq T,$$

where  $\psi$  and  $\phi$  are measurable processes adapted to an increasing family of  $\sigma$ -algebras  $\{\mathcal{A}_t\}$ ,  $\{W_t, \mathcal{A}_t\}$  is a Brownian motion, and  $X_0$  is  $\mathcal{A}_0$ -measurable. In addition, we assume

$$(25) \quad \int_0^T E|\psi_s| ds < \infty,$$

$$(26) \quad \int_0^T E\phi_s^2 ds < \infty.$$

---

<sup>1</sup> This fact is easily proved by verifying the two conditions of Theorem 3.3 in [7].

It is clear that  $\{X_t, \mathcal{A}_t\}$  is a quasi-martingale. Thus,  $\{X_t, \mathcal{A}_{xt}\}$  is also a quasi-martingale and the representation of  $X$  as a quasi-martingale with respect to  $\{\mathcal{A}_{xt}\}$  is given in the following theorem.

THEOREM 4.2. Let  $X_t, 0 \leq t \leq T$ , be an Ito process satisfying (24), (25) and

$$(27) \quad \int_0^T \phi_s^2 ds < \infty \quad \text{almost surely.}$$

Then there exists a representation of  $X$  in the form

$$(28) \quad X_t = X_0 + \int_0^t E^{\tilde{\mathcal{A}}_s} \psi_s ds + \int_0^t |\psi_s| d\tilde{W}_s,$$

where  $\{\tilde{\mathcal{A}}_t\}$  is any increasing family of  $\sigma$ -algebras satisfying

$$(29) \quad \mathcal{A}_{xt} \subseteq \tilde{\mathcal{A}}_t \subseteq \mathcal{A}_t \quad \text{for each } t$$

and  $\tilde{W}$  is a Brownian motion.

Remark. If we take  $\tilde{\mathcal{A}}_t = \mathcal{A}_{xt}$ , then we have

$$X_t = X_0 + \int_0^t E^{\mathcal{A}_{xs}} \psi_s ds + \int_0^t |\phi_s| d\tilde{W}_s,$$

where both integrands  $E^{\mathcal{A}_{xs}} \psi_s$  and  $|\phi_s|$  are  $\mathcal{A}_{xs}$ -measurable functions for almost all  $s$ . The latter stems from the fact that

$$(30) \quad |\phi_t| = \left[ \frac{d}{dt} \langle X \rangle_t \right]^{1/2},$$

where  $d/dt$  stands for the Radon-Nikodym derivative and  $\langle X \rangle_t$  is defined by

$$(31) \quad \langle X \rangle_t = X_t^2 - X_0^2 - 2 \int_0^t X_s dX_s.$$

Proof. First, suppose that, in place of (27), the stronger condition (26) is satisfied. Then,  $\{X_t, \tilde{\mathcal{A}}_t\}$  is a quasi-martingale and we can write

$$X_t = X_0 + \tilde{B}_t + \tilde{\mathcal{M}}_t, \quad 0 \leq t \leq T,$$

where  $\tilde{B}$  is of bounded variation and  $\{\tilde{M}_t, \tilde{\mathcal{A}}_t\}$  is a martingale. Since

$$\tilde{M}_t - \int_0^t \phi_s dW_s = \int_0^t \psi_s ds - \tilde{B}_t$$

and  $M_t = \int_0^t \phi_s dW_s$  is a martingale under condition (26), we have from Theorem 4.1:

$$\langle \tilde{M} \rangle_t = \langle M \rangle_t = \int_0^t \phi_s^2 ds, \quad 0 \leq t \leq T;$$

Theorem 3.1 yields the representation

$$\langle \tilde{M} \rangle_t = \int_0^t |\phi_s| d\tilde{W}_s,$$

where  $\tilde{W}$  is a Brownian motion.

To get an integral representation for  $\tilde{B}$ , let  $T_n = \{t_v^{(n)}\}$ ,  $n = 1, 2, \dots$ , be a nested sequence of partitions of  $[0, T]$  such that  $\max_v(t_{v+1}^{(n)} - t_v^{(n)}) \xrightarrow{n \rightarrow \infty} 0$ . Define

$$\tilde{B}_{n,t} = \sum_v E[X_{t \wedge t_{v+1}^{(n)}} - X_{t \wedge t_v^{(n)}} | \tilde{\mathcal{A}}_{t_v^{(n)}}].$$

Fisk [7] has shown that

$$E|\tilde{B}_{n,t} - \tilde{B}_t| \xrightarrow{n \rightarrow \infty} 0.$$

From (24) we get

$$E[X_{t_{v+1}^{(n)}} - X_{t_v^{(n)}} | \tilde{\mathcal{A}}_{t_v^{(n)}}] = \int_{t_v^{(n)}}^{t_{v+1}^{(n)}} E(\psi_s | \tilde{\mathcal{A}}_{t_v^{(n)}}) ds.$$

If we denote  $f_s = E^{\tilde{\mathcal{A}}_s} \psi_s$ , then

$$\begin{aligned} & \sup_{0 \leq t \leq T} E \left| \int_0^t E^{\tilde{\mathcal{A}}_s} \psi_s ds - \tilde{B}_t \right| \\ & \leq \lim_{n \rightarrow \infty} \left\{ E \left| \sum_v \int_{t_v^{(n)}}^{t_{v+1}^{(n)}} [f_s - f_{t_v^{(n)}}] ds \right| \right. \\ & \quad \left. + E \left| \sum_v \int_{T^{(n)}}^{t_{v+1}^{(n)}} E[\psi_s - \psi_{t_v^{(n)}} | \tilde{\mathcal{A}}_{t_v^{(n)}}] ds \right| \right\} \\ & \leq \lim_{n \rightarrow \infty} E \sum_v \int_{t_v^{(n)}}^{t_{v+1}^{(n)}} [|f_s - f_{t_v^{(n)}}| + |\psi_s - \psi_{t_v^{(n)}}|] ds. \end{aligned}$$

For a suitable sequence of partitions  $\{t_v^{(n)}\}$ , the last limit can always be made zero (Doob [3, pp. 63–65]). Since the first expression is independent of  $\{t_v^{(n)}\}$ , we must have

$$\sup_{0 \leq t \leq T} E \left| \int_0^t E^{\tilde{\mathcal{A}}_s} \psi_s ds - \tilde{B}_t \right| = 0.$$

If  $\phi$  satisfies (27) and not (26), we define  $\{\tau_n\}$  by (8). Since  $\langle X \rangle_t = \int_0^t \psi_s^2 ds$ ,  $\{\tau_n\}$  is a sequence of stopping times not only for  $\{\mathcal{A}_t\}$  but also for  $\{\tilde{\mathcal{A}}_t\}$ . Therefore, for each  $n$ ,  $\{X_{t \wedge \tau_n}, \tilde{\mathcal{A}}_t\}$  is a quasi-martingale with a representation

$$X_{t \wedge \tau_n} = X_0 + \int_0^{t \wedge \tau_n} E^{\tilde{\mathcal{A}}_s} \psi_s ds + \tilde{M}_{t \wedge \tau_n}.$$

Hence,  $\tilde{M}_t = X_t - X_0 - \int_0^t E^{\tilde{\mathcal{A}}_s} \psi_s ds$  is a local martingale. It follows from Theorem 4.1 and (24) that

$$\langle \tilde{M} \rangle_t = \int_0^t \phi_s^2 ds,$$

and the proof is completed by using Theorem 3.2.

Theorem 4.2 can be generalized to a vector Ito process practically without change. Let  $X$  be an  $n$ -vector-valued process satisfying

$$(32) \quad X_t = X_0 + \int_0^t \psi_s ds + \int_0^t \phi_s dW_s,$$



where  $W$  is an  $m$ -vector process, the components of which are independent Brownian motions and  $\phi$  is an  $n \times m$  matrix. Instead of (25) and (27), we now assume

$$(33) \quad \int_0^T E \|\psi_s\| ds < \infty$$

and

$$(34) \quad \int_0^T \|\phi_s\|^2 ds < \infty \quad \text{almost surely,}$$

where  $\|\cdot\|$  denotes the Euclidean norm.

Let  $\alpha$  be an  $n$ -vector and let prime denote transpose. For every  $\alpha \in R^n$ ,  $\{\alpha' X_t\}$  is a scalar-valued Ito process of the form (24), and we can write

$$\alpha' X_t = \int_0^t \alpha' E^{\tilde{\mathcal{A}}_s} \psi_s ds + M_{\alpha,t},$$

where  $\{M_{\alpha,t}, \tilde{\mathcal{A}}_t\}$  is a local martingale for any  $\{\tilde{\mathcal{A}}_t\}$  satisfying (29) and

$$\langle M_\alpha \rangle_t = \int_0^t \alpha' \phi_s \phi_s' \alpha ds.$$

It follows that  $\int_0^t \phi_s \phi_s' ds$  is  $\mathcal{A}_{xt}$ -measurable for every  $t$ , and the positive semi-definite matrix  $\phi_s \phi_s'$  is  $\mathcal{A}_{xs}$ -measurable for almost all  $s$ . Let  $\phi_s \phi_s'$  be diagonalized so that

$$(35) \quad \phi_s \phi_s' = A_s \Lambda_s A_s',$$

where  $\Lambda_s$  is diagonal and the orthogonal matrix  $A_s$  is  $\mathcal{A}_{xt}$ -measurable for almost all  $s$ . Now, if we define

$$Y_t = \int_0^t A_s' dX_s$$

and apply Theorem 4.2 to the components of  $Y$ , we get the following theorem.

**THEOREM 4.3.** *Let  $X$  be an  $n$ -vector Ito process of the form (32) satisfying (33) and (34). Let  $A_s$  and  $\Lambda_s$  be matrices defined by the diagonalization (35). Then, for any increasing family of  $\sigma$ -algebras satisfying (29), there exists an  $n$ -vector Brownian motion  $\tilde{W}$  so that*

$$(36) \quad X_t = X_0 + \int_0^t E^{\tilde{\mathcal{A}}_s} \psi_s ds + \int_0^t A_s \Lambda_s^{1/2} d\tilde{W}_s.$$

**5. Applications.** Equation (32) is widely used to model a dynamical system disturbed by Gaussian white noise. Theorem 4.3 has some interesting and surprising consequences for such models. For example, in filtering problems we often interpret  $X$  in (32) as the observed process,  $\psi$  as the process to be estimated, and the stochastic integral term as the noise. If we identify  $\hat{\psi}_t = E^{\tilde{\mathcal{A}}_t} \psi_t$  as the estimator, then Theorem 4.3 implies that

$$\int_0^t (\psi_s - \hat{\psi}_s) ds = - \int_0^t \phi_s dW_s + \int_0^t A_s \Lambda_s^{1/2} d\tilde{W}_s.$$

The stochastic integral  $-\int_0^t \phi_s dW_s$  can always be re-represented as  $\int_0^t A_s \Lambda_s^{1/2} dV_s$ , where  $V$  is now an  $n$ -dimensional Brownian motion, so that

$$\int_0^t (\psi_s - \hat{\psi}_s) ds = \int_0^t A_s \Lambda_s^{1/2} (dV_s + d\tilde{W}_s).$$

Formally, the estimation error can be expressed as

$$\psi_t - \hat{\psi}_t = A_t \Lambda_t^{1/2} (\dot{V}_t + \dot{\tilde{W}}_t).$$

One may be tempted to say that the estimation error is "white," but that would be misleading. For the special case  $\phi(\omega, t) = K$ , a constant, we have

$$\psi_t - \hat{\psi}_t = A \Lambda^{1/2} (\dot{V}_t + \dot{\tilde{W}}_t),$$

which is the sum of Gaussian white noise (albeit dependent) processes. This was originally observed by Wonham [8] in a more limited context and generalized by Kailath [9]. We note that the same observation is true if  $\phi\phi'$  is a constant even if  $\phi$  is not.

Another observation worthy of note is that for almost all  $t$ ,  $\phi_t\phi_t'$  is  $\mathcal{A}_{xt}$ -measurable. Thus, for example, if  $\psi_t$  is a function only of  $t$  and  $\{\phi_s\phi_s', 0 \leq s \leq t\}$ , then the estimation error is necessarily zero for almost all  $t$ . Such a result would be difficult to prove outside of the context of quadratic variations. A similar observation can be made with respect to singular detection [10].

Finally, as a simple example of possible applications to control problems, consider a scalar equation

$$dX_t = f(X_0^t, U_0^t, t) dk + U_t dW_t,$$

where  $X$  represents the state and  $U$  the control. Suppose that  $U_t \geq 0$  for all  $t$ . Then Theorem 4.2 implies that for any optimization problem whatever, the optimizing control can always be implemented in state feedback form. This means that even if we observe the past of the noise process  $W$  and use it in constructing the control, performance cannot be improved. Surely, this is an unexpected result.

**Acknowledgment.** I am grateful to Professor P. P. Varaiya for many valuable suggestions. In particular, it was at his persuasion and with his help that I extended the result of § 4 to local martingales and to the multidimensional case.

#### REFERENCES

- [1] P. A. MEYER, *A decomposition theorem for supermartingales*, Illinois J. Math., 6 (1962), pp. 193–205.
- [2] H. KUNITA AND S. WATANABE, *On square integrable martingales*, Nagoya Math. J., 30 (1967), pp. 209–245.
- [3] J. L. DOOB, *Stochastic Processes*, John Wiley, New York, 1953.
- [4] K. ITO, *On a formula concerning stochastic differentials*, Nagoya Math. J., 3 (1951), pp. 55–65.
- [5] D. L. FISK, *Sample quadratic variation of sample continuous second order martingales*, Z. Wahrscheinlichkeitstheorie und Verw. Gebiete, 6 (1966), pp. 273–278.
- [6] P. A. MEYER, *Une majoration du processus croissant naturel associé à une surmartingale*, Seminaire de probabilités, II, Springer-Verlag, Berlin, 1967.
- [7] D. L. FISK, *Quasi-martingales*, Trans. Amer. Math. Soc., 120 (1965), pp. 369–389.
- [8] W. M. WONHAM, *Random differential equations in control theory*, Probabilistic Methods in Applied Mathematics, vol. 2, A. T. Bharucha-Reid, ed., Academic Press, New York, 1970.

- [9] T. KAILATH, *An innovation approach to least-squares estimation. Part I: Linear filtering in additive white noise*, IEEE Trans. Automatic Control, AC-13 (1968), pp. 646–655.
- [10] E. WONG AND M. ZAKAI, *The oscillation of stochastic integrals*, Z. Wahrscheinlichkeitstheorie und Verw. Gebiete, 4 (1965), pp. 103–112.