

# Stepsize Restrictions for Stability of One-Step Methods in the Numerical Solution of Initial Value Problems

By M. N. Spijker

**Abstract.** This paper deals with the analysis of general one-step methods for the numerical solution of initial (-boundary) value problems for stiff ordinary and partial differential equations. Restrictions on the stepsize are derived that are necessary and sufficient for the rate of error growth in these methods to be of moderate size. These restrictions are related to disks contained in the stability region of the method, and the errors are measured with arbitrary norms (not necessarily generated by an inner product).

The theory is illustrated in the numerical solution of a diffusion-convection problem where the error growth is measured with the maximum norm.

## 1. Introduction.

1.1. *The Relevance of  $A[D]$ -Stability.* In 1963 Dahlquist [4] introduced the concept of  $A$ -stability in the analysis of numerical methods for solving initial value problems for ordinary differential equations. This concept is based on the error propagation in a numerical method when it is applied to the simple scalar test equation  $dU(t)/dt = \lambda U(t)$  with  $\lambda \in \mathbb{C}$ . Between 1976-1979 the criterion of  $A$ -stability proved to be of great relevance in a rigorous analysis of error propagation in methods for solving problems that are essentially more general than the scalar test problem (cf. [2], [8], [1], [13], [5]). A priori estimates were obtained of error propagation in general one-step methods applied to arbitrarily stiff linear systems of ordinary differential equations and linear partial differential equations (cf. [2], [8], [1]).

After 1963, the concept of a stability region  $S$  in the complex plane (cf. Subsection 1.2) was studied, which led to several weaker versions of  $A$ -stability (cf., e.g., [14]). One of these versions,  $A(\alpha)$ -stability, can also be applied successfully in a rigorous analysis analogous to the above (see [2], [20]).

In the present paper, we focus on some other weaker version of  $A$ -stability namely  $A[D]$ -stability (cf. [14], [9], [10] and Subsection 1.2). Here  $D$  denotes a disk in the complex plane (bounded by a circle passing through the origin) and it is required that  $D \subset S$ . In contrast to the requirements of  $A$ -stability and  $A(\alpha)$ -stability, the requirement  $D \subset S$  can be fulfilled by explicit methods.

Under the assumption of  $A[D]$ -stability, we give in this paper an analysis of error propagation in a framework that has similarity to those mentioned above (cf. [2], [8],

---

Received March 26, 1984; revised February 18, 1985.  
1980 *Mathematics Subject Classification.* Primary 65L20, 65M10.

[1], [20]). We shall arrive at conditions on the stepsize that are (necessary and) sufficient for the rate of error growth, in a given method, to be of moderate size. These conditions will be illustrated in the solution of a diffusion-convection problem.

Further, we shall relate the concept of  $A[D]$ -stability to some other known stability concepts (spectral stability condition, von Neumann stability condition, weak and strong stability, contractivity).

We note that, under the assumption of  $A[D]$ -stability, already rigorous results on error propagation were obtained for linear  $k$ -step methods by Nevanlinna [13] and Dahlquist [5]. These results are more general than ours in that the differential equations considered were allowed to be nonlinear, and  $k$  was allowed to be  $> 1$ . On the other hand, the norms we shall deal with need not be generated by an inner product, and the number of stages in the numerical schemes that we shall consider is not restricted to 1—as was the case in [13], [5].

1.2. *Notations and Definitions.* With  $\phi$  we denote a given rational function,  $\phi(\zeta) = P_1(\zeta)/P_0(\zeta)$ , where  $P_1, P_0$  are polynomials with real coefficients, no common zero and  $P_1(0) = P_0(0) = 1$ . For any  $s \times s$  matrix  $T$  we say that  $\phi(T)$  exists, and we write  $\phi(T) = P_1(T)P_0(T)^{-1}$  whenever the matrix  $P_0(T)$  is regular.

We shall be concerned with the rate of growth of vectors  $u_n \in \mathbf{R}^s$  that are computed from the recurrence relation

$$(1.1) \quad u_n = \phi(hA)u_{n-1} \quad (n = 1, 2, 3, \dots).$$

Here  $h > 0$  is the *stepsize*, and  $A$  denotes a given real  $s \times s$  matrix.

In many applications  $u_n$  stands for an approximation to  $U(nh)$ . Here  $U(t) \in \mathbf{R}^s$  denotes the solution to a given initial value problem for a system of ordinary differential equations

$$(1.2) \quad \frac{d}{dt}U(t) = AU(t) \quad (t \geq 0), \quad U(0) = u_0.$$

Many known numerical methods for solving ordinary differential equations, such as *Runge-Kutta* methods and *Rosenbrock* methods, result, when applied to (1.2), in a procedure of type (1.1).

Further, many numerical schemes for solving *initial-boundary value problems* in *partial differential equations* can be written in the form (1.1) (see, e.g., Section 7, and [17], [8], [19], [23], [26]). Here  $s$  will often stand for the (large) number of gridpoints involved at any fixed time level  $t = nh$ .

In this paper, the rate of growth of the  $u_n$  will be related to the size of  $h$  and of quantities  $r$  and  $R$  that we are now going to define.

The *stability region*  $S$  of  $\phi$  is defined by

$$S = \{ \zeta \mid \zeta \in \mathbf{C} \text{ and } \phi \text{ is regular at } \zeta \text{ with } |\phi(\zeta)| \leq 1 \}.$$

The disk  $D = \mathcal{D}(\xi, \rho)$  in the complex plain is defined by

$$D = \{ \zeta \mid \zeta \in \mathbf{C} \text{ and } |\zeta - \xi| \leq \rho \}$$

for  $\xi \in \mathbf{C}$ ,  $0 \leq \rho < \infty$ . Further, the procedure (1.1) is called  $A[D]$ -stable if  $D \subset S$ .

We define the *stability radius*  $r \in [0, \infty]$  of  $\phi$  by

$$r = \sup \{ \rho \mid 0 \leq \rho < \infty \text{ and } \mathcal{D}(-\rho, \rho) \subset S \}.$$

Finally,  $R \in [0, \infty]$  is defined by

$$R = \sup\{\rho \mid \rho = 0, \text{ or } 0 < \rho < \infty \text{ and } \phi \text{ absolutely monotonic on } [-\rho, 0]\}$$

(a function is called absolutely monotonic on an interval if the values of the function and of all its derivatives are finite and  $\geq 0$  on that interval). For reasons becoming evident in the subsequent we call  $R$  the *contractivity radius* of  $\phi$ .

In this paper, we deal mainly with matrices  $A$  such that (1.2) is *stable* in the sense that the solution to (1.2) remains bounded as  $t \rightarrow \infty$  (for each starting vector  $u_0 \in \mathbf{R}^s$ ). We aim at transparent conditions on the stepsize  $h > 0$  that guarantee an analogous stability-behavior for the vectors  $u_n$  satisfying (1.1). Similarly to the frameworks in [2], [8], [1], [20], and motivated by the applications to partial differential equations mentioned above, we shall focus on stability results for (1.1) which hold uniformly with respect to the dimension  $s$ .

In most applications, it is the growth of the difference between two solutions, say  $u_n$  and  $\tilde{u}_n$ , to (1.1) which is significant. For instance,  $\tilde{u}_n$  may stand for the numerical approximation obtained in the presence of a rounding error  $v_0 = \tilde{u}_0 - u_0$ . Since the resulting error  $v_n = \tilde{u}_n - u_n$  then also satisfies (1.1), our stability results on the growth of  $u_n$  will also be relevant to the *growth of errors*  $v_n$  in the application of method (1.1).

1.3. *Organization of the Paper.* In Section 2 we introduce the classes of matrices  $A$  for which we shall analyze the stability of process (1.1). We prove a theorem on these classes that was already applied in [22].

In Section 3 we deal with a recurrence relation for  $u_n$  that can be viewed as originating from an application of method (1.1) to a simple (but nonscalar) testproblem of type (1.2). This recurrence relation enables us to relate in a natural way  $A[D]$ -stability to a number of well-known other stability concepts.

In Theorem 4.1 of Section 4, we present a condition (in terms of the radius  $r$ ) on the stepsize  $h > 0$  which is necessary and sufficient for (a weak version of) stability of the general process (1.1). At the end of Section 4, we also present stability results for the case where  $h$  in (1.1) is replaced by a variable stepsize  $h_n > 0$  ( $n = 1, 2, 3, \dots$ ).

Theorem 5.1 in Section 5 is concerned with an analogous condition (in terms of the radius  $R$ ) on the stepsize  $h > 0$  which is necessary and sufficient for the stronger version of stability called contractivity.

In Section 6, we discuss modifications of Theorems 4.1, 5.1.

Section 7 contains an illustration of the material of Sections 2, 4, 5 in the numerical solution of a diffusion-convection problem.

We note that the tools used in the proofs in Sections 4, 5 mainly consist in power series expansions for matrix-valued functions and Parseval's formula for complex functions. Our arguments are therefore more elementary than those used in the proofs of the important paper by Brenner and Thomée [1] referred to in Subsection 1.1.

**2. Two Classes of Matrices  $A$ .** Let  $\mathbf{R}^s$  denote the  $s$ -dimensional real vector space equipped with an arbitrary norm  $|\cdot|$ . Let  $\omega, \tau, \alpha \in \mathbf{R}$  with  $\tau > 0$ .

Following [23], we denote by  $\mathcal{L}(\mathbf{R}^s, \omega, \tau)$ , the collection of all real  $s \times s$  matrices  $A$  such that

$$(2.1) \quad \|A + \tau^{-1}\| \leq \tau^{-1} + \omega.$$

Here  $\|\cdot\|$  denotes the (lub-) matrix norm induced by  $|\cdot|$ , i.e.,

$$\|T\| = \sup\{|Tx|: x \in \mathbf{R}^s \text{ with } |x| = 1\}$$

for any  $s \times s$  matrix  $T$ .

Further, we denote by  $\mathcal{B}(\mathbf{R}^s, \omega, \alpha)$ , the collection of all real  $s \times s$  matrices  $A$  such that

$$(2.2) \quad \|A\| \leq \alpha, \quad \mu[A] \leq \omega.$$

Here  $\mu[\cdot]$  denotes the logarithmic norm induced by  $|\cdot|$ , i.e.,

$$\mu[A] = \lim_{t \rightarrow 0^+} t^{-1}(\|I + tA\| - 1)$$

(cf. [3];  $I$  denotes the identity matrix).

The inequality (2.1) implies that the spectrum of the matrix  $A$  is contained in the disk  $\mathcal{D}(-\tau^{-1}, \tau^{-1} + \omega)$ . For normal  $A$  and the Euclidean norm this spectral property is even equivalent to (2.1). In [10], [23] property (2.1) was called a *circle condition*. For further interpretations of the class  $\mathcal{L}(\mathbf{R}^s, \omega, \tau)$ , see [10], [23].

The inequalities in (2.2) only involve the (logarithmic) norm of  $A$  itself so that the intuitive meaning of the class  $\mathcal{B}(\mathbf{R}^s, \omega, \alpha)$  is even more easily grasped than that of  $\mathcal{L}(\mathbf{R}^s, \omega, \tau)$ . Therefore, the following Theorem 2.1, relating  $\mathcal{L}(\mathbf{R}^s, \omega, \tau)$  to  $\mathcal{B}(\mathbf{R}^s, \omega, \alpha)$  is of importance.

We refer to Section 7 for an example involving the classes  $\mathcal{L}(\mathbf{R}^s, \omega, \tau)$ ,  $\mathcal{B}(\mathbf{R}^s, \omega, \alpha)$ .

With  $\mathbf{R}_p^s$ , we shall denote  $\mathbf{R}^s$  when equipped with the  $p$ th Hölder norm  $|x|_p = (\sum_j |\xi_j|^p)^{1/p}$  (when  $1 \leq p < \infty$ ),  $|x|_p = \max_j |\xi_j|$  (when  $p = \infty$ ) for  $x = (\xi_1, \xi_2, \dots, \xi_s)^T \in \mathbf{R}^s$ . The corresponding matrix norm is denoted by  $\|\cdot\|_p$ .

**THEOREM 2.1.** (Relations between  $\mathcal{L}(\mathbf{R}^s, \omega, \tau)$  and  $\mathcal{B}(\mathbf{R}^s, \omega, \alpha)$ .) *Let  $\omega, \alpha \in \mathbf{R}$  with  $\alpha > \omega$ . Then*

(i)  $\mathcal{B}(\mathbf{R}^s, \omega, \alpha) \supset \mathcal{L}(\mathbf{R}^s, \omega, 2(\alpha - \omega)^{-1})$  for each norm  $|\cdot|$  in  $\mathbf{R}^s$ ;

(ii)  $\mathcal{B}(\mathbf{R}^s, \omega, \alpha) = \mathcal{L}(\mathbf{R}^s, \omega, 2(\alpha - \omega)^{-1})$  when the norm in  $\mathbf{R}^s$  is given by  $|x| = |Qx|_p$  where  $Q$  is a regular  $s \times s$  matrix and  $p = 1$  or  $p = \infty$ .

*Proof.* 1. Defining  $\tau = 2(\alpha - \omega)^{-1}$ , we have for any  $A \in \mathcal{L}(\mathbf{R}^s, \omega, \tau)$ , the inequalities

$$\|A\| - \tau^{-1} \leq \|A + \tau^{-1}\| \leq \omega + \tau^{-1},$$

and therefore,

$$\|A\| \leq \omega + 2\tau^{-1} = \alpha.$$

This proves (i) since  $\mu[A] \leq \tau^{-1}(\|I + \tau A\| - 1) \leq \tau^{-1}(\tau(\omega + \tau^{-1}) - 1) = \omega$  (cf. [3]).

2. In order to prove (ii) we first assume  $p = \infty$ ,  $Q = I$  (the identity). Let  $A \in \mathcal{B}(\mathbf{R}_\infty^s, \omega, \alpha)$ ,  $A = (\alpha_{ij})$ .

For  $x = (\xi_1, \dots, \xi_s)^T$ ,  $y = (\eta_1, \dots, \eta_s)^T \in \mathbf{R}_\infty^s$  with  $|x|_\infty = 1$ ,  $y = (I + \tau A)x$  we have

$$|\eta_i| \leq |1 + \tau\alpha_{ii}| \cdot |\xi_i| + \tau \sum_{j \neq i} |\alpha_{ij}| \cdot |\xi_j| \quad (1 \leq i \leq s).$$

Hence,

$$|\eta_i| \leq \theta_i(1 + \tau\alpha_{ii}) + (1 - \theta_i)(-1 - \tau\alpha_{ii}) + \tau \sum_{j \neq i} |\alpha_{ij}| \quad (1 \leq i \leq s)$$

with  $\theta_i = 0$  (when  $1 + \tau\alpha_{ii} < 0$ ) or  $\theta_i = 1$  (when  $1 + \tau\alpha_{ii} \geq 0$ ). Consequently (cf. [3]),

$$\begin{aligned} |\eta_i| &\leq \theta_i \left\{ 1 + \tau\alpha_{ii} + \tau \sum_{j \neq i} |\alpha_{ij}| \right\} + (1 - \theta_i) \left\{ -1 - \tau\alpha_{ii} + \tau \sum_{j \neq i} |\alpha_{ij}| \right\} \\ &\leq \theta_i(1 + \tau\omega) + (1 - \theta_i)(-1 + \tau\|A\|) \\ &\leq \theta_i(1 + \tau\omega) + (1 - \theta_i)(-1 + \tau \cdot (\omega + 2\tau^{-1})) = 1 + \tau\omega. \end{aligned}$$

Therefore,  $\|(I + \tau A)x\|_\infty = \|y\|_\infty \leq 1 + \tau\omega$ , which proves  $A \in \mathcal{L}(\mathbf{R}_\infty^s, \omega, \tau)$ . Hence,  $\mathcal{B}(\mathbf{R}_\infty^s, \omega, \alpha) = \mathcal{L}(\mathbf{R}_\infty^s, \omega, \tau)$ .

3. We note that for any  $s \times s$  matrix  $A = (\alpha_{ij})$  we have  $\|A\|_1 = \|A^T\|_\infty$ ,  $\mu_1[A] = \mu_\infty[A^T]$ . From  $\mathcal{B}(\mathbf{R}_\infty^s, \omega, \alpha) = \mathcal{L}(\mathbf{R}_\infty^s, \omega, \tau)$ , we thus conclude that also  $\mathcal{B}(\mathbf{R}_1^s, \omega, \alpha) = \mathcal{L}(\mathbf{R}_1^s, \omega, \tau)$ .

4. Let  $|\cdot|$  be any norm in  $\mathbf{R}^s$  for which we have  $\mathcal{B}(\mathbf{R}^s, \omega, \alpha) = \mathcal{L}(\mathbf{R}^s, \omega, \tau)$ . Let  $Q$  be a regular  $s \times s$  matrix and put  $|x|^* = |Qx|$  (for  $x \in \mathbf{R}^s$ ). The matrix norm and logarithmic norm induced by  $|\cdot|^*$  are denoted by  $\|\cdot\|^*$ ,  $\mu^*[\cdot]$ .

For any  $s \times s$  matrix  $A$  we have

$$\mu^*[A] = \mu[QAQ^{-1}], \quad \|A\|^* = \|QAQ^{-1}\| \quad \text{and} \quad \|A + \tau^{-1}\|^* = \|QAQ^{-1} + \tau^{-1}\|.$$

From these relations it easily follows that also  $\mathcal{B}(\mathbf{R}^s, \omega, \alpha) = \mathcal{L}(\mathbf{R}^s, \omega, \tau)$  with respect to the norm  $|\cdot|^*$ .

In view of the parts 2, 3 the proof is now complete.  $\square$

*Remarks.* 1. By part (i) of the above theorem, we have for each  $A \in \mathcal{L}(\mathbf{R}^s, \omega, \tau)$ , the inequality  $\mu[A] \leq \omega$ , and therefore (cf. [3]),

$$|U(t)| \leq \exp(\omega t)|u_0| \quad (\text{for } t \geq 0)$$

for the solution to (1.2). If  $\omega \leq 0$ , the solution to (1.2) is thus bounded uniformly for  $0 \leq t < \infty$ .

2. One easily verifies that statement (ii) in the above theorem is not valid in case  $Q = I$  and  $p = 2$ .

3. Statement (ii) with  $Q = I$ ,  $p = \infty$ , is a basic means for proving the expression for the stepsize threshold that was presented, without full proof, in [22].

**3. Relating  $A[D]$ -Stability to Other Stability Concepts.** In order to obtain insight into the possible rate of growth of vectors  $u_n$  computed from (1.1) when  $A$  belongs to one of the classes defined in Section 2, it is useful to consider the case

$$(3.1) \quad u_n = \phi(A_s)u_{n-1} \quad (n = 1, 2, 3, \dots).$$

Here  $A_s$  denotes a bi-diagonal  $s \times s$  matrix of the form

$$(3.2) \quad A_s = \begin{pmatrix} \xi & \rho & 0 & \cdots & 0 \\ 0 & \xi & \rho & \cdot & \cdot \\ \vdots & \cdot & \cdot & \cdot & \cdot \\ \vdots & \cdot & \cdot & \cdot & \cdot \\ 0 & \cdots & 0 & \rho & \xi \end{pmatrix}$$

with  $\xi, \rho \in \mathbf{R}, \rho \geq 0$ . Moreover, the subsequent stability analysis will relate  $A[D]$ -stability in a natural fashion to a number of other stability concepts.

**THEOREM 3.1.** (Necessary condition for weak stability of (3.1).) *Let  $\phi(\zeta)$  be regular at  $\zeta = \xi$ , and let  $1 \leq p \leq \infty$ . Assume there exist  $\gamma, q < \infty$ , such that*

$$(3.3) \quad |u_n|_p \leq \gamma n^q \cdot |u_0|_p \quad (\text{for all } n \geq 1, s \geq 1, u_0 \in \mathbf{R}^s)$$

whenever  $u_n$  satisfies (3.1). Then,

$$(3.4) \quad \mathcal{D}(\xi, \rho) \subset S.$$

*Proof.* Apply Lemma 3.2, to be given at the end of this section.  $\square$

The above theorem, constituting the main result of this section, shows that there is *weak stability* (in the sense of (3.3)) only if there is  $A[D]$ -stability with disk  $D = \mathcal{D}(\xi, \rho)$ . This contrasts with the familiar *spectral condition for stability*, which for (1.1) requires that the spectrum of the matrix  $hA$  is contained in (the interior of)  $S$ . The latter condition reduces for the case of (3.1) to the requirement  $|\phi(\xi)| < 1$ . This is, generally, a much weaker requirement than (3.4). Similar shortcomings of the spectral condition for stability (uniformly with respect to  $s$ ) were stated, e.g. in [16], [17, p. 152], [11, pp. 258–261], [12], [19], [26].

We note that Theorem 3.1 can also be proved, slightly differently, by using the *Godunov-Ryabenkii criterion for stability* (cf. [17, p. 153]). We have preferred the proof via Lemma 3.2 since it is more direct and shorter.

It is interesting, and not surprising in view of the Godunov-Ryabenkii criterion, that condition (3.4) can be arrived at in a heuristic fashion by putting  $s = \infty$  in (3.1). The components  $u_{n,j}$  of the vector  $u_n$  would then satisfy

$$u_{n,j} = \sum_{k=0}^{\infty} \frac{\phi^{(k)}(\xi)}{k!} \rho^k u_{n-1,j+k} \quad (j \in \mathbf{Z}, n \geq 1),$$

and the *von Neumann condition for stability* (cf. [17]) applied to this difference scheme, reads

$$(3.5) \quad \left| \sum_{k=0}^{\infty} \frac{\phi^{(k)}(\xi)}{k!} \rho^k e^{ikt} \right| \leq 1 \quad (\text{for } -\infty < t < \infty).$$

Clearly (3.5) is equivalent to (3.4).

We conclude this section with

**LEMMA 3.2.** *Let the rational function  $\phi(\zeta)$  be regular at  $\zeta = \xi$ , and let  $1 \leq p \leq \infty$ . Then*

(i)  $\lim_{s \rightarrow \infty} \|\phi(A_s)\|_p^n = \infty$  (for each  $n \geq 1$ ) when  $\phi(\zeta)$  has a singularity in  $\mathcal{D}(\xi, \rho)$ ;

(ii)  $2 \cdot \|\phi(A_s)\|_p^n \geq M^n - g(s, n)$  (for each  $n \geq 1$ ) with  $\lim_{s \rightarrow \infty} g(s, n) = 0$ ,  $M = \max\{|\phi(\zeta)| : \zeta \in \mathcal{D}(\xi, \rho)\}$  when  $\phi(\zeta)$  is regular on  $\mathcal{D}(\xi, \rho)$ .

*Proof.* With  $E$  we denote the  $s \times s$  matrix all of whose entries  $\epsilon_{ij}$  vanish with the exception of  $\epsilon_{i,i+1} = 1$  ( $1 \leq i \leq s - 1$ ), and we define  $\gamma_j = (j!)^{-1} \phi^{(j)}(\xi) \cdot \rho^j$  (for  $j = 0, 1, 2, \dots$ ).

For  $n \geq 1$ , we have  $\phi(A_s)^n = \sum \gamma_{j_1} \gamma_{j_2} \cdot \dots \cdot \gamma_{j_n} \cdot E^{j_1+j_2+\dots+j_n}$  where the summation is for all integers  $j_1 \geq 0, j_2 \geq 0, \dots, j_n \geq 0$  with  $j_1 + j_2 + \dots + j_n \leq s - 1$  (cf. [6]).

Denoting also for complex vectors  $y \in \mathbf{C}^s$  the  $p$ th Hölder norm by  $|y|_p$ , we have for  $1 \leq p < \infty$  and for any  $y = (\eta_1, \eta_2, \dots, \eta_s)^T \in \mathbf{C}^s$  with  $|y|_p \leq 1$ , the inequality

$$2 \cdot \|\phi(A_s)^n\|_p \geq |\phi(A_s)^n y|_p = \left\{ \sum_{j=1}^s \left| \sum_{k=0}^{s-j} \beta_k \eta_{j+k} \right|^p \right\}^{1/p},$$

where

$$\beta_k = \sum_{j_1 + \dots + j_n = k} \gamma_{j_1} \cdot \gamma_{j_2} \cdot \dots \cdot \gamma_{j_n} \quad (k \geq 0).$$

Let  $\theta \in (0, 1]$  be such that  $\phi(\zeta)$  is regular on  $\mathcal{D}(\xi, \theta\rho)$ . Choosing  $y = (\eta_1, \eta_2, \dots, \eta_s)^T \in \mathbf{C}^s$  with  $\eta_k = \eta_0 \theta^k \exp(ikt)$  (for  $1 \leq k \leq s$ ),  $\eta_0 = s^{-1/p}$  (when  $\theta = 1$ ),  $\eta_0 = \theta^{-1}(1 - \theta^p)^{1/p}$  (when  $0 < \theta < 1$ ),  $t \in \mathbf{R}$ , we thus obtain

$$2 \cdot \|\phi(A_s)^n\|_p \geq \eta_0 \left\{ \sum_{m=0}^{s-1} \theta^{p(s-m)} \left| \sum_{k=0}^m \beta_k \theta^k \exp(ikt) \right|^p \right\}^{1/p}.$$

Writing  $\zeta = \xi + \theta\rho \exp(it)$ , we have

$$\lim_{m \rightarrow \infty} \sum_{k=0}^m \beta_k \theta^k \exp(ikt) = \phi(\zeta)^n.$$

From this relation, we can conclude that, for  $1 \leq p < \infty$ ,

$$2 \cdot \|\phi(A_s)^n\|_p \geq |\phi(\zeta)^n| - g(s, n, \theta)$$

with  $\lim_{s \rightarrow \infty} g(s, n, \theta) = 0$ . Since  $\|\phi(A_s)^n\|_\infty = \|\phi(A_s)^n\|_1$ , this result is also valid for  $p = \infty$ .

Statement (i) of the lemma follows easily by varying  $t$  and  $\theta$  appropriately, and statement (ii) follows by choosing  $\theta = 1$  and  $t$  such that  $|\phi(\xi + \rho \cdot \exp(it))| = M$ .  $\square$

*Remark.* Although it is not essential in the stability considerations concerning process (3.1), it is worth mentioning that the factor 2 occurring in statement (ii) (Lemma 3.2) can be omitted. This follows from a proof communicated to us recently by M. Crouzeix.

#### 4. On the Relevance of the Radius $r$ to the Stability of (1.1).

4.1. *A Necessary and Sufficient Condition for Stability of (1.1).* In this section, we give a stability analysis of (1.1) based on  $A[D]$ -stability with  $D = \mathcal{D}(-r, r)$ , where  $r$  is the stability radius defined in Subsection 1.2.

In the present subsection, we state the main result of this analysis in Theorem 4.1. The proof of the theorem will be based on Theorem 3.1 and on two corollaries to the technical Lemma 4.2. This lemma will be presented in Subsection 4.2 while its corollaries are given in Subsection 4.3. The corollaries in the latter subsection also contain some conclusions that have not been incorporated into the main Theorem 4.1, among other things a conclusion on the variable stepsize version (4.1) of (1.1).

The subsequent Theorem 4.1 provides a restriction on the stepsize  $h > 0$  (see (s1)) which is sufficient for weak stability of the process (1.1) (uniformly with respect to  $s$ ; see (s3)). Further, the theorem implies that restriction (s1) is also a necessary condition, already for a weaker version (namely (s2)) of the weak stability property (s3). Finally, the theorem shows that, under some additional condition on  $A$ , the

stepsize restriction (s1) is even sufficient for (strong) stability of (1.1) (uniformly with respect to  $s$ ; see (s4)).

An illustration of Theorem 4.1 is presented in Section 7.

**THEOREM 4.1.** *Let  $h, \tau, p$  be given with  $0 < h < \infty, 0 < \tau < \infty, 1 \leq p \leq \infty$ . Then the following four statements (s1)–(s4) are equivalent to each other.*

(s1)  $0 < h \leq r\tau$ ;

(s2)  $\phi(hA)$  exists and (1.1) implies  $|u_n|_p \leq \gamma n^q |u_0|_p$  (for each  $n \geq 1, s \geq 1, u_0 \in \mathbf{R}^s, A = A_s$  (see (3.2)) with  $-\xi = \rho = \tau^{-1}$ ). Here  $\gamma, q$  are constants independent of  $n, s, u_0$ ;

(s3)  $\phi(hA)$  exists and (1.1) implies  $|u_n| \leq \gamma n^{1/2} |u_0|$  (for each  $n \geq 1, s \geq 1, u_0 \in \mathbf{R}^s, A \in \mathcal{L}(\mathbf{R}^s, 0, \tau)$  and each norm  $|\cdot|$  in  $\mathbf{R}^s$ ). Here  $\gamma$  is a constant independent of  $n, s, u_0, A, |\cdot|$ ;

(s4)  $\phi(hA)$  exists and (1.1) implies  $|u_n| \leq [(-\omega\tau)(2 + \omega\tau)]^{-1/2} |u_0|$  (for each  $n \geq 1, s \geq 1, u_0 \in \mathbf{R}^s, \omega < 0, A \in \mathcal{L}(\mathbf{R}^s, \omega, \tau)$  and each norm  $|\cdot|$  in  $\mathbf{R}^s$ ).

*Proof.* 1. Assume (s1). Then (s3) follows easily from Corollary 4.3 stated in the next Subsection 4.3. We only have to apply this corollary with  $h_n \equiv h, \omega = 0$  and  $\rho = r$  (when  $r < \infty$ ) or  $\rho \in [h\tau^{-1}, \infty)$  (when  $r = \infty$ ).

Further, (s4) follows immediately by applying Corollary 4.4.

2. Since  $A = A_s$  with  $-\xi = \rho = \tau^{-1}$  (see (3.2)) belongs to the class  $\mathcal{L}(\mathbf{R}_p^s, 0, \tau)$ , one easily sees that (s3) implies (s2) with  $q = \frac{1}{2}$ .

Applying Theorem 3.1 (with  $\xi, \rho$  replaced by  $h\xi, h\rho$ ), we see that statement (s2) implies

$$\mathcal{D}(-h\tau^{-1}, h\tau^{-1}) = \mathcal{D}(h\xi, h\rho) \subset S.$$

Hence, (s2) implies (s1).

3. It remains to show that (s4) implies (s1).

Let  $\omega < 0, 1 + \omega\tau \geq 0$  and  $A = A_s$  with  $\xi = -\tau^{-1}, \rho = \tau^{-1} + \omega$ . Since  $A \in \mathcal{L}(\mathbf{R}_p^s, \omega, \tau)$ , statement (s4) implies that  $\phi(\zeta)$  is regular at  $\zeta = -h\tau^{-1}$ , and that

$$\|[\phi(hA_s)]^n\|_p \leq \gamma \quad (\text{for all } s \geq 1, n \geq 1)$$

with  $\gamma = [(-\omega\tau)(2 + \omega\tau)]^{-1/2}$ . An application of Theorem 3.1 (with  $q = 0, \xi = -h\tau^{-1}, \rho = h(\tau^{-1} + \omega)$ ) shows that

$$\mathcal{D}(-h\tau^{-1}, h\tau^{-1} + h\omega) \subset S.$$

Since  $S$  is closed, it follows, by letting  $\omega \rightarrow 0^-$ , that

$$\mathcal{D}(-h\tau^{-1}, h\tau^{-1}) \subset S.$$

Hence,  $h\tau^{-1} \leq r$  which implies (s1), and completes the proof of the theorem.  $\square$

*Remarks.* 1. In the majorization (s4) the factor  $n^{1/2}$ , which is present in (s3), has disappeared. This is compensated by the factor  $\tau^{-1/2}$  in (s4), which satisfies

$$\tau^{-1/2} \leq r^{1/2} h^{-1/2} = (r/t)^{1/2} \cdot n^{1/2} \quad \text{with } t = nh.$$

2. The constant  $\gamma$  in (s3) can be chosen, to some degree, independently of the parameters  $h \in (0, \infty)$  and  $\tau \in (0, \infty)$ .

Assume (s1) and  $r < \infty$ . Then (s3) holds with  $\gamma$  only depending on  $\phi$ .



Assume (s1) and  $r = \infty$ . Let  $h\tau^{-1} \leq \rho < \infty$ . Then (s3) holds with  $\gamma$  only depending on  $\phi$  and on  $\rho$ .

These two conclusions easily follow from Corollary 4.3.

4.2. *Formulation and Proof of Lemma 4.2.* In this subsection we deal with a slightly generalized version of procedure (1.1). We consider the recurrence relation

$$(4.1) \quad u_n = \phi(h_n A) u_{n-1} \quad (n = 1, 2, 3, \dots)$$

with arbitrary stepsizes

$$h_n > 0 \quad (n = 1, 2, \dots).$$

The subsequent lemma is a convenient means for obtaining upper bounds for  $|u_n|$  under various conditions on  $A$  and  $h_n$ .

We introduce some notations needed in the formulation of the lemma.

For  $\xi \in \mathbf{C}$ ,  $0 \leq \rho < \infty$ , we define  $\mathcal{M}(\xi, \rho) = \infty$  when  $\phi$  has a singularity in  $\mathcal{D}(\xi, \rho)$ , and  $\mathcal{M}(\xi, \rho) = \max\{|\phi(\zeta)| : \zeta \in \mathcal{D}(\xi, \rho)\}$  when  $\phi$  is regular on  $\mathcal{D}(\xi, \rho)$ .

We assume  $\tau, \omega, \theta \in \mathbf{R}$  with

$$(4.2.a) \quad \tau > 0, \quad 1 + \omega\tau \geq 0, \quad \theta > 0,$$

and we define  $\lambda_n, \mu_n$  (for  $n = 1, 2, 3, \dots$ ) by

$$(4.2.b) \quad \lambda_n = \mathcal{M}(-h_n\tau^{-1}, (1 + \omega\tau)h_n\tau^{-1}),$$

$$(4.2.c) \quad \mu_n = \mathcal{M}(-h_n\tau^{-1}, \theta^{-1}h_n\tau^{-1}).$$

LEMMA 4.2. *Assume (4.2) and  $\lambda_n < \infty, \mu_n < \infty$  ( $n = 1, 2, 3, \dots$ ). Let  $m$  be an arbitrary integer  $\geq 0$ . Then, for each  $n \geq 1, s \geq 1, A \in \mathcal{L}(\mathbf{R}^s, \omega, \tau)$  and each norm  $|\cdot|$  in  $\mathbf{R}^s$ , the following statements (4.3), (4.4) are valid.*

$$(4.3) \quad \phi(h_n A) \text{ exists,}$$

$$(4.4) \quad \begin{aligned} & \|\phi(h_n A) \cdot \dots \cdot \phi(h_2 A) \cdot \phi(h_1 A)\| \\ & \leq \left( \prod_{j=1}^n \lambda_j \right) \left\{ \sum_{k=0}^{m-1} (1 + \omega\tau)^{-2k} \|(I + \tau A)^k\|^2 \right\}^{1/2} \\ & \quad + \left( \prod_{j=1}^n \mu_j \right) \left\{ \sum_{k=m}^{\infty} \theta^{2k} \|(I + \tau A)^k\|^2 \right\}^{1/2}. \end{aligned}$$

The first sum in the right-hand member of (4.4) stands for zero when  $m = 0$ , and it stands for 1 when  $m > 0, (1 + \omega\tau) = 0$ .

*Proof.* Defining  $\xi_n = -h_n\tau^{-1}$  we have, in view of (2.1), for any  $A \in \mathcal{L}(\mathbf{R}^s, \omega, \tau)$ , the representation  $h_n A = \xi_n - \xi_n B$  where the  $s \times s$  matrix  $B$  satisfies

$$B = I + \tau A, \quad \|B\| \leq 1 + \omega\tau.$$

From (4.2.b),  $\lambda_n < \infty$ , we see that  $\phi$  is regular on  $\mathcal{D}(\xi_n, |\xi_n(1 + \omega\tau)|)$ . Therefore, using the above representation for  $h_n A$ , it follows (cf. [6]) that (4.3) holds.

From (4.2.b), (4.2.c),  $\lambda_j < \infty, \mu_j < \infty$  (for  $1 \leq j \leq n$ ), we see that the rational function  $f$  defined by

$$f(\zeta) = \phi(\xi_n(1 - \zeta)) \cdot \dots \cdot \phi(\xi_2(1 - \zeta)) \cdot \phi(\xi_1(1 - \zeta))$$

is regular for all  $\zeta \in \mathbf{C}$  with  $|\zeta| \leq \max[(1 + \omega\tau), \theta^{-1}]$ . Denoting the coefficients of the Maclaurin expansion of  $f(\zeta)$  by  $\gamma_k$ , we thus have

$$(4.5) \quad f(\zeta) = \gamma_0 + \gamma_1\zeta + \gamma_2\zeta^2 + \dots \quad (\text{for } |\zeta| \leq \max[(1 + \omega\tau), \theta^{-1}]).$$

Consequently (cf. [6]),

$$(4.6) \quad \phi(h_n A) \cdot \dots \cdot \phi(h_2 A) \cdot \phi(h_1 A) = f(B) = \gamma_0 + \gamma_1 B + \gamma_2 B^2 + \dots.$$

From (4.6) we obtain, by two applications of Schwarz' inequality, the upper bound

$$\begin{aligned} \|f(B)\| \leq & \left\{ \sum_{k=0}^{m-1} (\gamma_k(1 + \omega\tau)^k)^2 \right\}^{1/2} \left\{ \sum_{k=0}^{m-1} (1 + \omega\tau)^{-2k} \|B^k\|^2 \right\}^{1/2} \\ & + \left\{ \sum_{k=m}^{\infty} (\gamma_k\theta^{-k})^2 \right\}^{1/2} \left\{ \sum_{k=m}^{\infty} \theta^{2k} \|B^k\|^2 \right\}^{1/2}. \end{aligned}$$

For any  $\sigma$  with  $0 < \sigma \leq \max[1 + \omega\tau, \theta^{-1}]$ , we obtain from (4.5), by applying Parseval's formula, the relation

$$\sum_{k=0}^{\infty} (\gamma_k\sigma^k)^2 = (2\pi)^{-1} \int_0^{2\pi} |f(\sigma e^{it})|^2 dt.$$

In view of the definition of  $f$ , it thus follows that

$$\sum_{k=0}^{\infty} (\gamma_k\sigma^k)^2 \leq \prod_{j=1}^n [\mathcal{M}(\xi_j, |\xi_j\sigma|)]^2.$$

Applying the last inequality successively with  $\sigma = 1 + \omega\tau$  and with  $\sigma = \theta^{-1}$ , we easily obtain from the above upper bound for  $\|f(B)\|$  the inequality (4.4).  $\square$

4.3. *Upper bounds for  $|u_n|$ .* We now state three interesting corollaries to the above lemma, the first two corollaries of which were essential in our proof of Theorem 4.1.

**COROLLARY 4.3.** *Let  $0 < \tau < \infty$ ,  $0 < \rho < \infty$ ,  $\rho \leq r$  and  $0 < h_n \leq \rho\tau$ . Assume  $\omega \geq 0$  and  $\phi$  regular on  $\mathcal{D}(-\rho, (1 + \omega\tau)\rho)$ . Then  $\phi(h_n A)$  exists and (4.1) implies  $|u_n| \leq \gamma n^{1/2} L^n |u_0|$  (for each  $n \geq 1$ ,  $s \geq 1$ ,  $u_0 \in \mathbf{R}^s$ ,  $A \in \mathcal{L}(\mathbf{R}^s, \omega, \tau)$  and each norm  $|\cdot|$  in  $\mathbf{R}^s$ ). Here  $\gamma$  only depends on  $\phi$ ,  $\rho$  and on the product  $\omega\tau$ . Further,  $L = \mathcal{M}(-\rho, (1 + \omega\tau)\rho)$ , and  $L = 1$  when  $\omega = 0$ .*

*Proof.* We shall apply Lemma 4.2 with  $\theta > 0$  such that

$$1 + \omega\tau < \theta^{-1}, \quad \mathcal{M}(-\rho, \theta^{-1}\rho) < \infty.$$

Clearly, the relations (4.2) are fulfilled with

$$\lambda_n \leq L = \mathcal{M}(-\rho, (1 + \omega\tau)\rho) < \infty, \quad \mu_n \leq M = \mathcal{M}(-\rho, \theta^{-1}\rho) < \infty.$$

Note that  $L = 1$  when  $\omega = 0$ , since  $\phi(0) = 1$ ,  $\rho \leq r$ .

Let  $n \geq 1$ ,  $s \geq 1$ ,  $u_0 \in \mathbf{R}^s$ ,  $u_0 \neq 0$ ,  $A \in \mathcal{L}(\mathbf{R}^s, \omega, \tau)$  and  $|\cdot|$  be given. By Lemma 4.2 we then have (4.3), and (4.1) then implies

$$|u_n|/|u_0| \leq L^n m^{1/2} + M^n \left\{ \sum_{k=m}^{\infty} \theta^{2k} (1 + \omega\tau)^{2k} \right\}^{1/2}.$$

Choosing  $m = nj$ , we thus have

$$|u_n|/|u_0| \leq L^n (jn)^{1/2} + \{M[\theta(1 + \omega\tau)]^j\}^n \cdot \{1 - \theta^2(1 + \omega\tau)^2\}^{-1/2}.$$

Taking the integer  $j$  so large that  $M[\theta(1 + \omega\tau)]^j \leq L\{1 - \theta^2(1 + \omega\tau)^2\}^{1/2}$  there follows  $|u_n|/|u_0| \leq (j^{1/2} + 1)L^n n^{1/2}$ , which proves the corollary since  $\gamma = (j^{1/2} + 1)$  only depends on  $\phi, \rho, \omega\tau$ .  $\square$

**COROLLARY 4.4.** *Let  $0 < \tau < \infty, 0 < h < \infty$  and  $0 < h \leq r\tau$ . Assume  $\omega < 0$ . Then  $\phi(hA)$  exists and (1.1) implies*

$$|u_n| \leq [(-\omega\tau)(2 + \omega\tau)]^{-1/2} |u_0|$$

(for each  $n \geq 1, s \geq 1, u_0 \in \mathbf{R}^s, A \in \mathcal{L}(\mathbf{R}^s, \omega, \tau)$  and each norm  $|\cdot|$  in  $\mathbf{R}^s$ ).

*Proof.* When  $A \in \mathcal{L}(\mathbf{R}^s, \omega, \tau)$ , we can apply Lemma 4.2 with  $h_n \equiv h, \theta = 1, m = 0$ . By (4.3) the matrix  $\phi(hA)$  thus exists, and by (4.4) we easily arrive at

$$\|\phi(hA)^n\| \leq \left\{ \sum_{k=0}^{\infty} \|(I + \tau A)^k\|^2 \right\}^{1/2}.$$

In view of (2.1), we have  $\|I + \tau A\| \leq 1 + \omega\tau$ , and, consequently,

$$\|\phi(hA)^n\|^2 \leq \sum_{k=0}^{\infty} (1 + \omega\tau)^{2k} = [(-\omega\tau)(2 + \omega\tau)]^{-1}. \quad \square$$

The next corollary is obtained by applying Lemma 4.2 similarly as in the proof of Corollary 4.4 with  $h_n \equiv h, m = 0$ , but with  $\theta = (1 + \omega\tau/2)^{-1}$ . Since the proof is analogous to the above, we omit it.

**COROLLARY 4.5.** *Let  $0 < \tau < \infty, 0 < h < \infty$  and  $0 < h \leq r\tau$ . Assume  $\omega < 0$ . Then  $\phi(hA)$  exists and (1.1) implies  $|u_n| \leq \gamma M^n |u_0|$  (for each  $n \geq 1, s \geq 1, u_0 \in \mathbf{R}^s, A \in \mathcal{L}(\mathbf{R}^s, \omega, \tau)$  and each norm  $|\cdot|$  in  $\mathbf{R}^s$ ). Here  $\gamma = (2 + \omega\tau)[- \omega\tau(4 + 3\omega\tau)]^{-1/2}$ , and  $M = \mathcal{M}(-h\tau^{-1}, (1 + \omega\tau/2)h\tau^{-1})$  satisfies  $M < 1$  (provided the rational function  $\phi$  is no constant).*

This corollary implies that when  $\omega < 0, A \in \mathcal{L}(\mathbf{R}^s, \omega, \tau)$  and  $h$  is restricted as in Theorem 4.1, any  $u_n$  satisfying (1.1) damps out exponentially when  $n \rightarrow \infty$ . This asymptotic behavior of the approximations  $u_n \approx U(nh)$  is desirable since an analogous behavior is shown by the true solution  $U(t)$  to (1.2) (cf. Remark 1 of Section 2).

**5. On the Relevance of the Radius  $R$  to the Stability of (1.1).** In this section the parameter  $R$ , defined in Subsection 1.2, will be compared with the radius  $r$  and related to the stability of the process (1.1).

The definition of  $R$  implies that, for any  $\rho \in (0, R)$ , the Taylor coefficients  $\gamma_k = (k!)^{-1}\phi^{(k)}(-\rho)$  satisfy  $\gamma_k \geq 0$  and

$$\gamma_0 + \gamma_1\rho + \gamma_2\rho^2 + \dots = \phi(0) = 1.$$

Consequently, procedure (1.1) is  $A[D]$ -stable for any  $D = \mathcal{D}(-\rho, \rho)$  with  $\rho \in (0, R)$ , so that

$$(5.1) \quad R \leq r.$$

The subsequent Theorem 5.1 has a structure similar to the one of Theorem 4.1 in the above section. Comparing both theorems we see that, in view of (5.1), the stepsize restriction (S1) occurring in Theorem 5.1 is, generally, more severe than the analogous restriction (s1) in Theorem 4.1. This is in agreement with the fact that the corresponding stability result (S3) in the subsequent theorem is much stronger than

the analogous result (s3) in Theorem 4.1. The property  $|u_n| \leq |u_0|$  occurring in (S3) is often called *contractivity* (cf., e.g., [23], [26]) and is related to what sometimes is called *practical stability* (cf., e.g., [12]).

An illustration of Theorem 5.1 is presented in Section 7.

**THEOREM 5.1.** *Let  $h, \tau, p$  be given with  $0 < h < \infty, 0 < \tau < \infty$ , and  $p = 1$  or  $p = \infty$ . Then, the following three statements are equivalent to each other.*

(S1)  $0 < h \leq R\tau$ ;

(S2)  $\phi(hA)$  exists and (1.1) implies  $|u_n|_p \leq |u_0|_p$  (for each  $n \geq 1, s \geq 1, u_0 \in \mathbf{R}^s, A = A_s$  (see (3.2)) with  $-\xi = \rho = \tau^{-1}$ );

(S3)  $\phi(hA)$  exists and (1.1) implies  $|u_n| \leq |u_0|$  (for each  $n \geq 1, s \geq 1, u_0 \in \mathbf{R}^s, A \in \mathcal{L}(\mathbf{R}^s, 0, \tau)$  and each norm  $|\cdot|$  in  $\mathbf{R}^s$ ).

*Proof.* 1. (S3) implies (S2) since  $A = A_s$  with  $-\xi = \rho = \tau^{-1}$  belongs to  $\mathcal{L}(\mathbf{R}^s, 0, \tau)$ .

2. Assuming (S2) we prove (S1). With  $A$  as in (S2) we have  $\phi(hA) = \gamma_0 + \gamma_1(\beta E) + \dots + \gamma_{s-1}(\beta E)^{s-1}$ , where  $\beta = h\rho, \gamma_k = (k!)^{-1}\phi^{(k)}(-\beta)$  and  $E$  is the matrix defined in the proof of Lemma 3.2. Consequently,

$$1 \geq \|\phi(hA)\|_p = |\gamma_0| + |\gamma_1|\beta + \dots + |\gamma_{s-1}|\beta^{s-1}$$

for each  $s \geq 1$ . We thus obtain

$$1 \geq \sum_{k=0}^{\infty} |\gamma_k|\beta^k \geq \sum_{k=0}^{\infty} \gamma_k\beta^k = \phi(0) = 1.$$

It follows that  $\phi^{(k)}(-\beta) \geq 0$  and  $\phi(\zeta)$  is regular on  $\mathcal{D}(-\beta, \beta)$ .

For any  $t \in [-\beta, 0]$  we thus have

$$\phi^{(k)}(t) = \phi^{(k)}(-\beta) + (t + \beta)\phi^{(k+1)}(-\beta) + (2!)^{-1}(t + \beta)^2\phi^{(k+2)}(-\beta) + \dots \geq 0$$

(for all  $k \geq 0$ ). Hence  $\beta \leq R$ , which proves (S1).

3. Assume (S1) and let  $A \in \mathcal{L}(\mathbf{R}^s, 0, \tau)$ . In view of (5.1) we have, similarly as in the beginning of the proof of Lemma 4.2 (with  $n = 1$ ),

$$\phi(hA) = \gamma_0 + \gamma_1 B + \gamma_2 B^2 + \dots$$

with  $\gamma_k = (k!)^{-1}(h\tau^{-1})^k\phi^{(k)}(-h\tau^{-1}), \|B\| \leq 1$ . Hence,

$$\|\phi(hA)\| \leq |\gamma_0| + |\gamma_1| + |\gamma_2| + \dots = \gamma_0 + \gamma_1 + \gamma_2 + \dots = \phi(0) = 1,$$

which completes the proof of (S3).  $\square$

*Remark.* The implication (S1)  $\Rightarrow$  (S3) follows immediately from [23, Theorem 3.3]. We have included the above proof of this implication since it is very short and keeps the paper self-contained.

**6. Modifications of Theorems 4.1, 5.1.** Suppose the norm in  $\mathbf{R}^s$  is generated by an inner product, i.e.,  $|x| = \langle x, x \rangle^{1/2}$  (for all  $x \in \mathbf{R}^s$ ). Then, an interesting modification of Theorems 4.1, 5.1 is possible.

One arrives at this modification by a straightforward application of Theorem 3.1 and a theorem of J. von Neumann (see [18, p. 432], [8], [2]). The latter theorem, applied to the situation at hand, says that for any  $s \times s$  matrix  $T$  with lub-norm  $\|T\| \leq 1$  and rational  $f$  mapping  $\mathcal{D}(0, 1)$  into  $\mathcal{D}(0, 1)$ , one has the bound  $\|f(T)\| \leq 1$ .

By choosing  $f(\zeta) = \phi(-h\tau^{-1}(1 - \zeta))$ ,  $T = I + \tau A$ , one thus obtains

**THEOREM 6.1.** *Let  $h > 0$ ,  $\tau > 0$  be given. Then, the following three statements are equivalent to each other.*

- ( $\sigma 1$ )  $0 < h \leq r\tau$ ;
- ( $\sigma 2$ )  $\phi(hA)$  exists and (1.1) implies  $|u_n|_2 \leq |u_0|_2$  (for each  $n \geq 1$ ,  $s \geq 1$ ,  $u_0 \in \mathbf{R}^s$ ,  $A = A_s$  (see (3.2)) with  $-\xi = \rho = \tau^{-1}$ );
- ( $\sigma 3$ )  $\phi(hA)$  exists and (1.1) implies  $|u_n| \leq |u_0|$  (for each  $n \geq 1$ ,  $s \geq 1$ ,  $u_0 \in \mathbf{R}^s$ ,  $A \in \mathcal{L}(\mathbf{R}^s, 0, \tau)$  and each norm  $|\cdot|$  in  $\mathbf{R}^s$  generated by an inner product).

By this theorem we thus have contractivity (see ( $\sigma 3$ )) for stepsizes  $h$  subject to ( $\sigma 1$ ). This stepsize restriction is, generally, weaker than restriction (S1) in Theorem 5.1 (cf. (5.1)). On the other hand, ( $\sigma 1$ ) equals restriction (s1) of Theorem 4.1, but ( $\sigma 3$ ) cannot be deduced from the stability statement (s3) of that theorem.

From Theorem 6.1 we conclude that in Theorem 5.1 one cannot allow  $p = 2$ . Only for  $\phi$  with  $R = r$ , Theorem 5.1 remains valid with  $p = 2$ , but this is exceptional (cf. (5.1)).

We conclude this section by briefly discussing the possibility of replacing the factor  $(\gamma n^{1/2})$  in statement (s3) of Theorem 4.1 by a factor, say  $\gamma_0$ , only depending on  $\phi$ . For  $A \in \mathcal{L}(\mathbf{R}^s, 0, \tau)$  the true solution  $U(t)$  to (1.2) is, by Remark 1 of Section 2, bounded (uniformly for  $0 \leq t < \infty$ ). Therefore, one might hope for an analogous boundedness behavior in (s3) for the approximations  $u_n \approx U(nh)$  (uniformly for  $n \geq 0$ ).

Unfortunately, without additional conditions on  $\phi$ , Theorem 4.1 does not allow such a modification. In fact,  $\phi(\zeta) = (1 - \zeta/2)^{-1}(1 + \zeta/2)$  provides a counterexample with  $r = \infty$ . Using arguments taken from [25, pp. 280–287], it can be proved that for  $A = E - I \in \mathcal{L}(\mathbf{R}^\infty, 0, 1)$  (with  $E$  as in the proof of Lemma 3.2) one has  $\sup\|\phi(hA)^n\|_\infty = \infty$ , the supremum being for  $s \geq 1$ ,  $n \geq 1$ ,  $h > 0$ .

**7. Illustration in a Diffusion-Convection Problem.** We turn to the application of the above to the stability analysis of difference methods for solving partial differential equations. Due to the framework we have been using, applications seem possible with arbitrary norms and any number of space variables, any boundary conditions, variable coefficients and variable space discretizations. On the other hand, due to the generality of the above theorems, one would not expect that in any actual application sharp, or refined, stability results can be obtained.

As an illustration, we consider the problem

$$(7.1.a) \quad \frac{\partial}{\partial t} U(x, t) = \frac{\partial^2}{\partial x^2} U(x, t) + b(x) \frac{\partial}{\partial x} U(x, t) + c(x)U(x, t),$$

$$(7.1.b) \quad U(0, t) = U(1, t) = 0,$$

$$(7.1.c) \quad U(x, 0) = U_0(x),$$

where  $0 < x < 1$ ,  $t > 0$  and  $U_0, b, c$  are given bounded real functions (which need not be smooth) with

$$c(x) \leq \omega \quad (0 < x < 1).$$

The following finite-difference scheme has been constructed according to well-known principles (cf. [12], [7], [26], [24], [15]).

$$(7.2.a) \quad h^{-1}(u_j^n - u_j^{n-1}) = \delta^{-2}(1 - \beta_j + \varepsilon_j|\beta_j|)(\theta u_{j-1}^n + (1 - \theta)u_{j-1}^{n-1}) \\ + \delta^{-2}(-2 - 2\varepsilon_j|\beta_j| + c_j\delta^2)(\theta u_j^n + (1 - \theta)u_j^{n-1}) \\ + \delta^{-2}(1 + \beta_j + \varepsilon_j|\beta_j|)(\theta u_{j+1}^n + (1 - \theta)u_{j+1}^{n-1}),$$

$$(7.2.b) \quad u_0^{n-1} = u_{s+1}^{n-1} = 0,$$

$$(7.2.c) \quad u_j^0 = U_0(j\delta),$$

where  $j = 1, 2, \dots, s$  and  $n = 1, 2, 3, \dots$ . In (7.2) we use the notations  $h = \Delta t > 0$ ,  $\delta = \Delta x > 0$ ,  $(s + 1)\delta = 1$ ,  $\beta_j = \delta \cdot b(j\delta)/2$ ,  $c_j = c(j\delta)$ ,  $u_j^n = U(j\delta, nh)$ , and  $\varepsilon_j \in [0, 1]$ ,  $\theta \in [0, 1]$  are parameters specifying the method. The choices  $\varepsilon_j \equiv 0$  and  $\varepsilon_j \equiv 1$  yield central finite-difference and fully upwinded finite-difference approximations to  $\partial U(x, t)/\partial x$ , respectively. The choices  $\theta = 0, \frac{1}{2}, 1$  correspond to the explicit, the Crank-Nicolson and the fully implicit schemes, respectively.

With the definitions  $u_n = (u_1^n, u_2^n, \dots, u_s^n)^T$  and

$$\phi(\xi) = (1 + (1 - \theta)\xi)(1 - \theta\xi)^{-1},$$

the relations (7.2.a), (7.2.b) become equivalent to the recurrence relation (1.1) provided the matrix  $A = (\alpha_{ij})$  is given by  $\alpha_{ij} = 0$  ( $|i - j| > 1$ ),

$$\alpha_{j,j-1} = \delta^{-2}(1 - \beta_j + \varepsilon_j|\beta_j|) \quad (2 \leq j \leq s),$$

$$\alpha_{jj} = \delta^{-2}(-2 - 2\varepsilon_j|\beta_j| + c_j\delta^2) \quad (1 \leq j \leq s),$$

$$\alpha_{j,j+1} = \delta^{-2}(1 + \beta_j + \varepsilon_j|\beta_j|) \quad (1 \leq j \leq s - 1).$$

Using the maximum norm in  $\mathbf{R}^s$ , we have (cf. [3], [26])

$$\|A\|_\infty = \max_i \sum_j |\alpha_{ij}|, \quad \mu_\infty[A] = \max_i \left( \alpha_{ii} + \sum_{j \neq i} |\alpha_{ij}| \right).$$

Consequently, choosing  $\varepsilon_j$  such that

$$(7.3) \quad \varepsilon_j \geq 1 - \frac{2}{\delta|b(j\delta)|} \quad (1 \leq j \leq s),$$

we have  $\mu_\infty[A] \leq \omega$ ,  $\|A\|_\infty \leq \alpha = \delta^{-2}(4 + 2\delta\lambda + \delta^2|c|)$ , where

$$\lambda = \max_j |\varepsilon_j b(j\delta)|, \quad |c| = \sup_x |c(x)|.$$

In view of (2.2), we have  $A \in \mathcal{B}(\mathbf{R}_\infty^s, \omega, \alpha)$ , and from Theorem 2.1, we thus obtain

$$(7.4) \quad A \in \mathcal{L}(\mathbf{R}_\infty^s, \omega, \tau), \quad \tau = \delta^2 \cdot (2 + \delta\lambda + \delta^2(|c| - \omega)/2)^{-1}.$$

One easily verifies that, with  $\phi$  as defined above,

$$(7.5.a) \quad r = (1 - 2\theta)^{-1} \quad (0 \leq \theta < \frac{1}{2}), \quad r = \infty \quad (\frac{1}{2} \leq \theta \leq 1),$$

$$(7.5.b) \quad R = (1 - \theta)^{-1} \quad (0 \leq \theta < 1), \quad R = \infty \quad (\theta = 1).$$

From Theorem 4.1 it thus follows that the solution to (7.2) satisfies the stability estimate

$$(7.6.a) \quad \max_j |u_j^n| \leq \delta^{-1} \cdot |\omega|^{-1/2} (2 + \delta\lambda + \delta^2(|c| - \omega)/2)^{1/2} \cdot \max_j |u_j^0|$$

whenever  $\omega < 0$ , (7.3) and

$$(7.7.a) \quad \begin{cases} h\delta^{-2} \leq (1 - 2\theta)^{-1}(2 + \delta\lambda + \delta^2(|c| - \omega)/2)^{-1} & (\text{for } 0 \leq \theta < \frac{1}{2}), \\ h\delta^{-2} < \infty & (\text{for } \frac{1}{2} \leq \theta \leq 1). \end{cases}$$

Similarly, Theorem 5.1 yields the estimate

$$(7.6.b) \quad \max_j |u_j^n| \leq \max_j |u_j^0|$$

whenever  $\omega \leq 0$ , (7.3) and

$$(7.7.b) \quad \begin{cases} h\delta^{-2} \leq (1 - \theta)^{-1}(2 + \delta\lambda + \delta^2(|c| - \omega)/2)^{-1} & (\text{for } 0 \leq \theta < 1), \\ h\delta^{-2} < \infty & (\text{for } \theta = 1). \end{cases}$$

Further, Corollary 4.3 yields estimates when  $\omega \geq 0$  and the stepsizes  $h_n > 0$  vary under restrictions similar to (7.7.a).

Of course, (7.6) can be used in a standard way (cf. [17], [11]) to obtain maximum-norm bounds for the global error  $u_j^n - U(j\delta, nh)$ . Here the factor  $\delta^{-1}$  in the right-hand member of (7.6.a) (as well as the factor  $n^{1/2}$  appearing in Corollary 4.3) need not stand in the way to obtain bounds for the global error that are of the same order as the local discretization errors (cf., e.g., [17, pp. 124–130], [21]).

We conclude with comparing the above estimates to some stability estimates obtainable from the literature.

For  $h\delta^{-2} < \infty$ ,  $\frac{1}{2} < \theta \leq 1$ , an estimate that is essentially sharper than (7.6.a) follows from the general theory in [1]. But, for  $0 \leq \theta \leq \frac{1}{2}$ , the estimate (7.6.a) is no direct consequence of that theory. Also, the variable stepsize result mentioned above does not follow from [1].

For the case  $\theta = 0$  and the pure diffusion-convection equation (7.1.a), with  $c(x) \equiv 0$ ,  $b(x) \equiv \text{constant}$ , Griffiths et al. [7] presented restrictions on  $\epsilon_j \equiv \epsilon$  and on  $h\delta^{-2}$  that are necessary and sufficient for (7.6.b) to be valid. It is interesting to note that for this case the conditions (7.3), (7.7.b) neatly reduce to the restrictions presented in [7].

**Acknowledgment.** I wish to thank J. M. Sanz-Serna for a stimulating discussion on the topic of this paper.

Institute of Applied Mathematics and Computer Science  
 University of Leiden  
 P. O. Box 9512  
 2300 RA Leiden, The Netherlands

1. P. BRENNER & V. THOMÉE, "On rational approximations of semigroups," *SIAM J. Numer. Anal.*, v. 16, 1979, pp. 683–694.
2. M. CROUZEIX & P. A. RAVIART, "Approximation d'équations d'évolution linéaires par des méthodes multipas," in *Etude Numérique des Grands Systèmes*, Rencontres IRIA-Novosibirsk, 1976, Dunod, Paris.
3. G. DAHLQUIST, "Stability and error bounds in the numerical integration of ordinary differential equations," *Trans. Roy. Inst. Techn.*, No. 130, Stockholm, 1959.
4. G. DAHLQUIST, "A special stability problem for linear multistep methods," *BIT*, v. 3, 1963, pp. 27–43.
5. G. DAHLQUIST, "G-stability is equivalent to A-stability," *BIT*, v. 18, 1978, 384–401.
6. N. DUNFORD & J. T. SCHWARTZ, *Linear Operators, Part I*, Interscience, New York, 1958.

7. D. F. GRIFFITHS, I. CHRISTIE & A. R. MITCHELL, "Analysis of error growth for explicit difference schemes in conduction-convection problems," *Internat. J. Numer. Methods Engrg.*, v. 15, 1980, pp. 1075–1081.
8. R. HERSH & T. KATO, "High-accuracy stable difference schemes for well-posed initial-value problems," *SIAM J. Numer. Anal.*, v. 16, 1979, pp. 670–682.
9. R. JELTSCH & O. NEVANLINNA, "Largest disk of stability of explicit Runge-Kutta methods," *BIT*, v. 18, 1978, pp. 500–502.
10. R. JELTSCH & O. NEVANLINNA, "Stability of explicit time discretizations for solving initial value problems," *Numer. Math.*, v. 17, 1977, pp. 61–91.
11. A. R. MITCHELL & D. F. GRIFFITHS, *The Finite Difference Method in Partial Differential Equations*, Wiley, Chichester, 1980.
12. K. W. MORTON, "Stability of finite difference approximations to a diffusion-convection equation," *Internat. J. Numer. Methods Engrg.*, v. 15, 1980, pp. 677–683.
13. O. NEVANLINNA, "On the numerical integration of nonlinear initial value problems by linear multistep methods," *BIT*, v. 17, 1977, pp. 58–71.
14. F. ODEH & W. LINIGER, "A note on unconditional fixed  $-h$  stability of linear multistep formulae," *Computing*, v. 7, 1971, pp. 240–253.
15. S. PAOLUCCI & D. R. CHENOWETH, "Stability of the explicit finite differenced transport equation," *J. Comput. Phys.*, v. 47, 1982, pp. 489–496.
16. S. V. PARTER, "Stability, convergence, and pseudo-stability of finite-difference equations for an over-determined problem," *Numer. Math.*, v. 4, 1962, pp. 277–292.
17. R. D. RICHTMYER & K. W. MORTON, *Difference Methods for Initial-Value Problems*, 2nd ed., Wiley, New York, 1967.
18. F. RIESZ & B. SZ-NAGY, *Leçons d'Analyse Fonctionnelle*, 2nd ed., Akademiai Kiado, Budapest, 1953.
19. J. M. SANZ-SERNA, *Convergent Approximations to Partial Differential Equations and Stability Concepts of Methods for Stiff Systems of Ordinary Differential Equations*, Report, University of Valladolid, 1983.
20. B. SCHMITT, "Norm bounds for rational matrix functions," *Numer. Math.*, v. 42, 1983, pp. 379–389.
21. M. N. SPIJKER, *Equivalence Theorems for Nonlinear Finite-Difference Methods*, Lecture Notes in Math., Vol. 267, Springer-Verlag, Berlin and New York, 1972, pp. 233–266.
22. M. N. SPIJKER, "Numerical contractivity in the solution of initial value problems," in *Proc. Zweites Seminar über Numerische Behandlung von Differentialgleichungen*, Halle (DDR), 1983.
23. M. N. SPIJKER, "Contractivity in the numerical solution of initial value problems," *Numer. Math.*, v. 42, 1983, pp. 271–290.
24. G. STOYAN, "Monotone difference schemes for diffusion-convection problems," *Z. Angew. Math. Mech.*, v. 59, 1979, pp. 361–372.
25. V. THOMÉE, "Stability of difference schemes in the maximum-norm," *J. Differential Equations*, v. 1, 1965, pp. 273–292.
26. J. G. VERWER & K. DEKKER, *Step-by-Step Stability in the Numerical Solution of Partial Differential Equations*, Report NW 161/83, Mathematical Centre, Amsterdam, 1983.