

# A Predictor-Corrector Method for a Certain Class of Stiff Differential Equations

By Karl G. Guderley and Chen-Chi Hsu\*

**Abstract.** In stiff systems of linear ordinary differential equations, certain elements of the matrix describing the system are very large. Sometimes, e.g., in treating partial differential equations, the problem can be formulated in such a manner that large elements appear only in the main diagonal. Then the elements causing stiffness can be taken into account analytically. This is the basis of the predictor-corrector method presented here. The truncation error can be estimated in terms of the difference between predicted and corrected values in nearly the same manner as for the customary predictor-corrector method. The question of stability, which is crucial for stiff equations, is first studied for a single equation; as expected, the method is much more stable than the usual predictor-corrector method. For systems of equations, sufficient conditions for stability are derived which require less work than a detailed stability analysis. The main tool is a matrix norm which is consistent with a weighted infinity vector norm. The choice of the weights is critical. Their determination leads to the question whether a certain matrix has a positive inverse.

**1. Introduction.** The present paper studies a predictor-corrector method for stiff systems of differential equations which have the following form:

$$(1) \quad \frac{dy}{dx} + \Lambda y = Ay + \Gamma(x) \equiv f(x, y),$$

where  $\Lambda$  is a diagonal matrix which may have some large elements, and the right-hand side is considered as nonstiff. The operator on the left is inverted and the right-hand side is approximated by Lagrangian interpolation polynomials at grid points. The point at which the unknown vector is to be computed is excluded in the predictor phase, but it is included in the corrector phase. The integration of the exponential functions is done analytically.

In general, it would be too costly to bring a stiff system into the form (1). However, systems of this kind arise naturally if one reduces certain partial differential equations into a system of ordinary differential equations by a Galerkin procedure which uses approximating functions closely related to the partial differential equation, see for instance [1]. The ideas of usual predictor-corrector methods can be applied almost immediately to (1). In particular, it is possible to express the truncation error in terms of the difference between predicted and corrected values and to use this relation for step control.

For stiff systems, the question of stability is crucial. We discuss it, first for a single equation, and then for a system, by means of a simplified analysis. Actually, this

Received July 21, 1971.

AMS 1969 subject classifications. Primary 6561, 6580.

Key words and phrases. Predictor-corrector method, stiff differential equations, interval control, stability, weighted infinity vector norm.

\* This work was carried out while one of the authors (C. C. Hsu) was pursuing an NRC-OAR postdoctoral resident research associateship.

latter discussion may be more interesting for the insight which it gives than for the practical work. In practice, the step control will probably suffice. It is true that a step control is only based on accuracy requirements, but the step selected by the method usually lies close to the stability limit, if a step of this size is compatible with the desired accuracy. The reason lies in the fact that instabilities lead to a deterioration of accuracy which leads to a reduction of the step size. The method would be  $A$ -stable if the matrix  $A$  vanishes. Therefore, one expects to find rather wide stability limits if  $A$  is sufficiently small.

Whenever one wants to use a large integration step for the solution of a stiff system, one is confronted with matrix inversions. A method can be made more effective if, within a certain interval which contains a number of integration steps, the operator which determines the differential equation is decomposed into a constant stiff part and a remainder. Then the matrix inversion is needed only once in each interval. The iteration steps which would be needed to invert the complete matrix (in our case  $A - \Lambda$ ) are combined with measures for increasing the order of the integration method; they do not appear explicitly even in the derivation of the formulae. In Gear's work [2], this idea is expressed clearly; in other methods it would probably appear if they are implemented in an efficient manner. In the present procedure the matrix inversion is, of course, trivial. Inversions of (well-conditioned) matrices will appear at values  $x$  where the functions used in Galerkin's method are changed.

The inversion of the operator on the left of (1) leads to the equation

$$(2) \quad \mathbf{y}(x) = e^{-\Lambda(x-x_0)}\mathbf{y}(x_0) + e^{-\Lambda x} \int_{x_0}^x e^{\Lambda\tau}\mathbf{f}(\tau, \mathbf{y}(\tau)) d\tau.$$

This equation is the point of departure for some papers [3], [10]. If one approximates  $\mathbf{f}$  by a polynomial, then, after integration by parts, one is led to matrices  $\exp(-\Lambda h)$ , where  $h$  is the integration step. For a nondiagonal  $\Lambda$ , an approximation must be used at this stage. For stability reasons, it is advisable to use a rational function for the approximation of the exponential function [3], [4]. Here, a matrix inversion is required too. This additional work does not occur in the present approach, for  $\Lambda$  is a diagonal matrix. In Gear's work, an approximation of this kind is not needed for he makes the stronger assumption that  $\mathbf{y}$  and  $\mathbf{f}$ , rather than  $\mathbf{f}$  only, can be approximated by a polynomial. The difference has only minor importance; it would be felt in regions where those contributions of the first term on the right of (2) which are related to large elements of  $\Lambda$  are important and if  $A$  is small. But this happens only rarely and can be handled by taking a smaller integration step.

In summary, we can say that the present paper deals with stiff systems of a form which allows us to use a modified version of the predictor-corrector method for nonstiff systems. One obtains a fairly simple computational procedure and a convenient characterization of the truncation error.

**2. Integration Formulae.** Let the dimension of the vector  $\mathbf{y}$  be  $N$ , and assume that  $\mathbf{f}$  is approximated by a polynomial of degree  $k$ . Assume that the integration has progressed to a station  $x_n = x_0 + nh$  and that we want to compute the value of  $\mathbf{y}$  at  $x_{n+1}$ . The values of  $\mathbf{y}$  and  $\mathbf{f}$  at equidistant stations,  $\mathbf{y}_{n-k+j}$  and  $\mathbf{f}_{n-k+j}$  for  $j = 0, 1, \dots, k$ , have been retained in the memory. Then the prediction polynomial for

$\mathbf{f}$  is obtained by Lagrangian interpolation. One finds

$$(3) \quad \mathbf{f}(x, \mathbf{y}) \simeq F^P B^P \mathbf{g}(\xi), \quad \xi = (x - x_n)/h,$$

where  $F^P$  is an  $N$  by  $(k + 1)$  matrix whose  $j$ th column is given by  $\mathbf{f}_{n-k+i}$ , for short  $F^P = (\mathbf{f}_{n-k}, \mathbf{f}_{n-k+1}, \dots, \mathbf{f}_n)$ ,  $B^P$  is a  $(k + 1)$  by  $(k + 1)$  constant matrix, and  $\mathbf{g}(\xi)$  is a  $(k + 1)$ -dimensional vector whose  $j$ th component is given by  $\xi^{j-1}$ . The value of  $\mathbf{y}$  at  $x_{n+1}$  is predicted by substituting (3) into (2) and carrying out a number of integrations by parts. One obtains

$$(4) \quad \mathbf{y}_{n+1}^P = e^{-\Lambda h} \mathbf{y}_n + h[\mathbf{q}_{n+1} - e^{-\Lambda h} \mathbf{q}_n],$$

in which the vector  $\mathbf{q}_{n+\xi}$  is given by

$$(5) \quad \mathbf{q}_{n+\xi} = (\Lambda h)^{-1} F^P B^P \mathbf{g}(\xi) - (\Lambda h)^{-2} F^P B^P \frac{d\mathbf{g}}{d\xi} + \dots + (-1)^k (\Lambda h)^{-k-1} F^P B^P \frac{d^k \mathbf{g}}{d\xi^k}.$$

For small  $\Lambda h$ , including  $\Lambda h = 0$ , an alternate expression for  $\mathbf{q}_{n+\xi}$  is found by expanding  $\exp(\Lambda h \xi)$  in (2) before the integration is carried out. One finds

$$(6) \quad \mathbf{q}_{n+\xi} = e^{-\Lambda h \xi} \sum_j (\Lambda h)^j F^P B^P \mathbf{s}_j(\xi), \quad j = 0, 1, \dots,$$

where the column vector  $\mathbf{s}_j(\xi)$  is

$$(7) \quad \mathbf{s}_j(\xi) = \frac{1}{j!} \int_0^\xi \tau^j \mathbf{g}(\tau) d\tau.$$

The prediction formula (4) can be written in the form

$$(8) \quad \mathbf{y}_{n+1}^P = e^{-\Lambda h} \mathbf{y}_n + h[V_0 \mathbf{f}_n + V_1 \mathbf{f}_{n-1} + \dots + V_k \mathbf{f}_{n-k}],$$

where the diagonal matrices  $V_i$  depend only upon  $\Lambda h$ .

From the predicted value  $\mathbf{y}_{n+1}^P$  so obtained, one computes the predicted value of  $\mathbf{f}_{n+1}$ . The correction polynomial of degree  $k$  for  $\mathbf{f}$  is obtained in a similar manner:

$$(9) \quad \mathbf{f}(x, \mathbf{y}) \simeq F^C B^C \mathbf{g}(\xi), \quad \xi = (x - x_n)/h.$$

Here,  $F^C$  is an  $N$  by  $(k + 1)$  matrix whose  $j$ th column is given by  $\mathbf{f}_{n-k+j+1}$ ,  $F^C = (\mathbf{f}_{n-k+1}, \mathbf{f}_{n-k+2}, \dots, \mathbf{f}_{n+1})$ ,  $B^C$  is a  $(k + 1)$  by  $(k + 1)$  matrix. Proceeding in complete analogy, one finds

$$(10) \quad \mathbf{y}_{n+1}^C = e^{-\Lambda h} \mathbf{y}_n + h[W_0 \mathbf{f}_{n+1} + W_1 \mathbf{f}_n + \dots + W_k \mathbf{f}_{n-k+1}],$$

where  $W_i$  are also diagonal matrices which depend only upon  $\Lambda h$ . Specific expressions for  $V_i$ ,  $W_i$ ,  $B^P$  and  $B^C$  for the case  $k = 4$  are given in the Appendix.

To start or restart the integration procedure, Picard's iteration method is employed. Assume that the initial value  $\mathbf{y}_i$  at the station  $x_i$  is given, then the values  $\mathbf{y}_{i+j}$ , for  $j = 1, 2, \dots, k$ , are computed simultaneously by iteration. Assume that the  $m$ th approximation for these values of  $\mathbf{y}$  has been found, one then computes the corresponding values of  $\mathbf{f}$  and obtains an interpolation polynomial for the  $m$ th approximation. The recurrence relation for generating these values of  $\mathbf{y}$  are obtained from (2) and (3) with  $n$  in the definition of  $\xi$  and  $F^P$  replaced by  $i + k$ . One finds

$$(11) \quad \mathbf{y}_{i+j}^{(m+1)} = e^{-\Lambda h} \mathbf{y}_{i+j-1}^{(m+1)} + h[\mathbf{q}_{i+j}^{(m)} - e^{-\Lambda h} \mathbf{q}_{i+j-1}^{(m)}], \quad j = 1, \dots, k,$$

where the superscript denotes the order of approximation and the function  $\mathbf{q}$  is defined by (5). The first approximation for  $\mathbf{y}_{i+j}$  is obtained from (11) by neglecting the last term of the equation.

**3. Truncation Error.** In problems of the kind considered, initial perturbations die out with increasing  $x$ . The same is true for the propagation of errors unless the numerical procedure is not stable. Accordingly, most of the truncation error is generated locally, and the local truncation error should be a useful measure for the entire truncation error. In the present method, the local truncation error is caused by the fact that  $\mathbf{f}$  is approximated by a polynomial of degree  $k$ . According to [5], the error caused by this approximation in the corrector phase is given by

$$(12) \quad \epsilon(\xi) = \frac{\mathbf{f}^{k+1}(\bar{\xi})}{(k+1)!} (\xi - \xi_{1-k})(\xi - \xi_{2-k}) \cdots (\xi - \xi_0)(\xi - \xi_1), \quad \xi = \frac{x - x_n}{h},$$

where  $\mathbf{f}^{k+1}(\bar{\xi})$  is the  $(k+1)$ th derivative evaluated at some station  $\bar{\xi}$  in the interval  $\xi_{1-k} \leq \xi \leq \xi_1$ . The value  $\bar{\xi}$  depends on the value of  $\xi$  for which  $\epsilon(\xi)$  is evaluated. The local truncation error for the corrector is obtained from (2) and (12). One finds

$$(13) \quad \mathbf{y}(x_{n+1}) - \mathbf{y}_{n+1}^C = \mathbf{t} = h e^{-\Lambda h} \int_0^1 e^{\Lambda h \xi} \epsilon(\xi) d\xi.$$

To compute  $\mathbf{t}$ , an estimate for  $\mathbf{f}^{k+1}(\bar{\xi})$  is required. In general,  $\mathbf{f}^{k+1}$  is not available, therefore, the additional assumption is made that  $\mathbf{f}$  is exactly given by a polynomial of degree  $(k+1)$ . Then the derivative  $\mathbf{f}^{k+1}$  is a constant vector; it can be expressed in terms of the values  $\mathbf{f}_{n-k}, \mathbf{f}_{n-k+1}, \cdots, \mathbf{f}_n$  and  $\mathbf{f}_{n+1}$ , and (13) can be evaluated.

If  $\mathbf{f}^{k+1}$  is constant, the difference between the predicted and the corrected values is simply related to the local truncation error. This relation is used for step control. One has

$$(14) \quad \mathbf{y}_{n+1}^C - \mathbf{y}_{n+1}^P = G\mathbf{t},$$

where  $G$  is a diagonal matrix which depends on the degree of the polynomial  $k$  and  $\Lambda h$ . This matrix  $G$  can be determined from any convenient example, since it does not depend on  $\mathbf{f}$ . To derive (14), we first observe that  $\mathbf{f}$  as a polynomial of degree  $(k+1)$  depends linearly upon  $(k+2)$  parameters. We choose for these parameters the values  $\mathbf{f}_{n-k+1}, \cdots, \mathbf{f}_n, \mathbf{f}_{n+1}$  and the derivative  $\mathbf{f}^{k+1}$ . In the corrector formula,  $\mathbf{f}$  is computed from  $\mathbf{f}_{n-k+1}, \cdots, \mathbf{f}_n, \mathbf{f}_{n+1}$ , the correction polynomial can therefore be obtained from the exact polynomial by setting  $\mathbf{f}^{k+1} = 0$ . One finds by an integration that

$$(15) \quad \mathbf{t} = hG_1\mathbf{f}^{k+1},$$

where  $G_1$  is a diagonal matrix which depends on  $k$  and  $\Lambda h$ . This result is, of course, already implied by (13). In the predictor phase, the approximation polynomial is determined by  $\mathbf{f}_{n-k}, \mathbf{f}_{n-k+1}, \cdots, \mathbf{f}_n$ . But in the "exact" polynomial the constant vector  $\mathbf{f}^{k+1}$  is a linear function of  $\mathbf{f}_{n-k}, \mathbf{f}_{n-k+1}, \cdots, \mathbf{f}_n$  and  $\mathbf{f}_{n+1}$ . Thus, the prediction polynomial can be considered as a linear function of the parameters  $\mathbf{f}_{n-k+1}, \cdots, \mathbf{f}_n, \mathbf{f}_{n+1}$  and  $\mathbf{f}^{k+1}$ . Now, for  $\mathbf{f}^{k+1} = 0$ , the prediction polynomial is the same as the correction polynomial, therefore

$$(16) \quad \mathbf{y}_{n+1}^C - \mathbf{y}_{n+1}^P = hG_2 \mathbf{f}^{k+1}.$$

Again,  $G_2$  is a diagonal matrix depending on  $k$  and  $\Delta h$ . Equations (15) and (16) then lead to (14).

To find the matrix  $G$ , we choose for convenience  $\mathbf{f}_{n-k+j} = 0$  for  $j = 0, 1, \dots, k$  and  $\mathbf{f}_{n+1} =$  unit vector. One then obtains

$$(17) \quad G = (k + 1) \frac{\int_0^1 e^{\Delta h \xi} \xi(\xi + 1) \cdots (\xi + k - 1) d\xi}{\int_0^1 e^{\Delta h \xi} (\xi - 1)\xi(\xi + 1) \cdots (\xi + k - 1) d\xi}.$$

For large values of  $\Delta h$  one has

$$(18) \quad G \sim -(k + 1)\Delta h.$$

The values of  $G$  for  $k = 1, 2, 3$  and  $4$  are given in Fig. 1. In practice, a simple bound for  $G$  will be sufficient. For  $k = 4$ , the matrix

$$(19) \quad \tilde{G} = -0.95(10 + 5\Delta h)(2 + \Delta h)/(1 + \Delta h)$$

serves this purpose, it satisfies

$$(20) \quad 0.9 |G| \leq |\tilde{G}| \leq |G|.$$

**4. Stability.** We shall use the usual stability definition, namely, that the solution should remain bounded as  $x$  tends to infinity. If the matrix  $A$  on the right-hand side of (1) is zero, the present integration method is perfect, except for errors in the

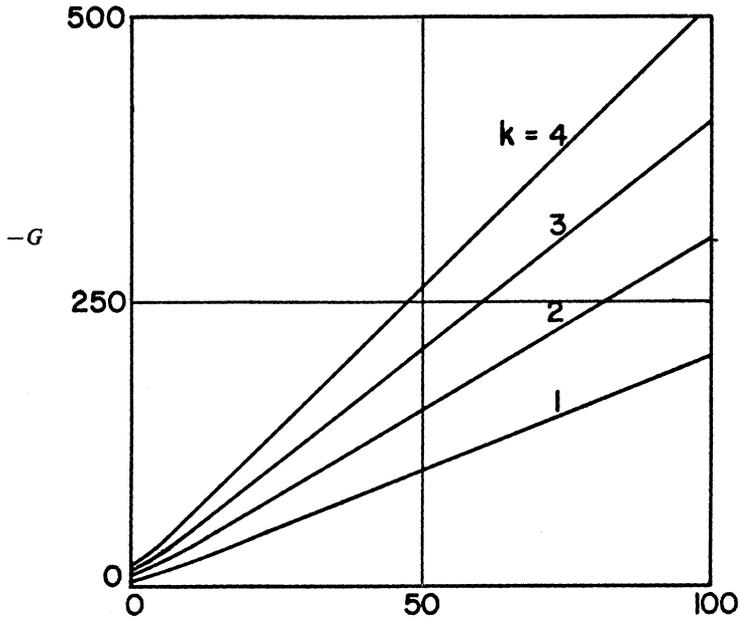


FIGURE 1. Ratio Between the Difference of Corrected and Predicted  $y$  and the Truncation Error

evaluation of integrals. Therefore, one would expect the method to be stable for rather large values of  $h$  if  $A$  is sufficiently small in comparison to  $\Lambda$ . For integration methods in which the matrix governing the system is treated as a whole, the stability analysis is greatly simplified by applying a similarity transformation which brings the matrix into its diagonal form; the results found from integrating the original system and the transformed system are exactly the same, except for round-off errors. Therefore, only a single equation needs to be considered in the stability discussion. However, if several matrices are considered simultaneously, as in the present method and also in Gear's approach [2], this simplification does not materialize. Nevertheless, we discuss the stability of a single equation in this section, for the results so obtained can be considered as an indication of the usefulness of the method. (This is the same as assuming that  $A$  is a diagonal matrix, too.) In the next section, we shall derive sufficient stability conditions for a system by a simplified analysis; even these discussions are still too complicated for practical use. In practice, one will depend upon the step control, via the control of truncation error, as a means of controlling the stability.

Assume that  $A$  in (1) is a diagonal matrix, then we examine the stability of a single equation of the form

$$(21) \quad y' + \lambda y = \gamma y.$$

In the integration formulae, (8) and (10),  $f_{n-k+j}$  is now replaced by  $\gamma y_{n-k+j}$ . Substituting the predicted value of  $y$  into the corrector formula, one obtains

$$(22) \quad y_{n+1}^c = \alpha_0 y_n + \alpha_1 y_{n-1} + \cdots + \alpha_k y_{n-k},$$

where the coefficients  $\alpha_j$  depend only upon  $\gamma h$  and  $\lambda h$ ; for  $k = 4$ , one has specifically

$$(23) \quad \begin{aligned} \alpha_0 &= e^{-\lambda h} + \gamma h(w_0 e^{-\lambda h} + w_1 + \gamma h w_0 v_0), \\ \alpha_1 &= \gamma h(w_2 + \gamma h w_0 v_1), & \alpha_2 &= \gamma h(w_3 + \gamma h w_0 v_2), \\ \alpha_3 &= \gamma h(w_4 + \gamma h w_0 v_3), & \alpha_4 &= (\gamma h)^2 w_0 v_4. \end{aligned}$$

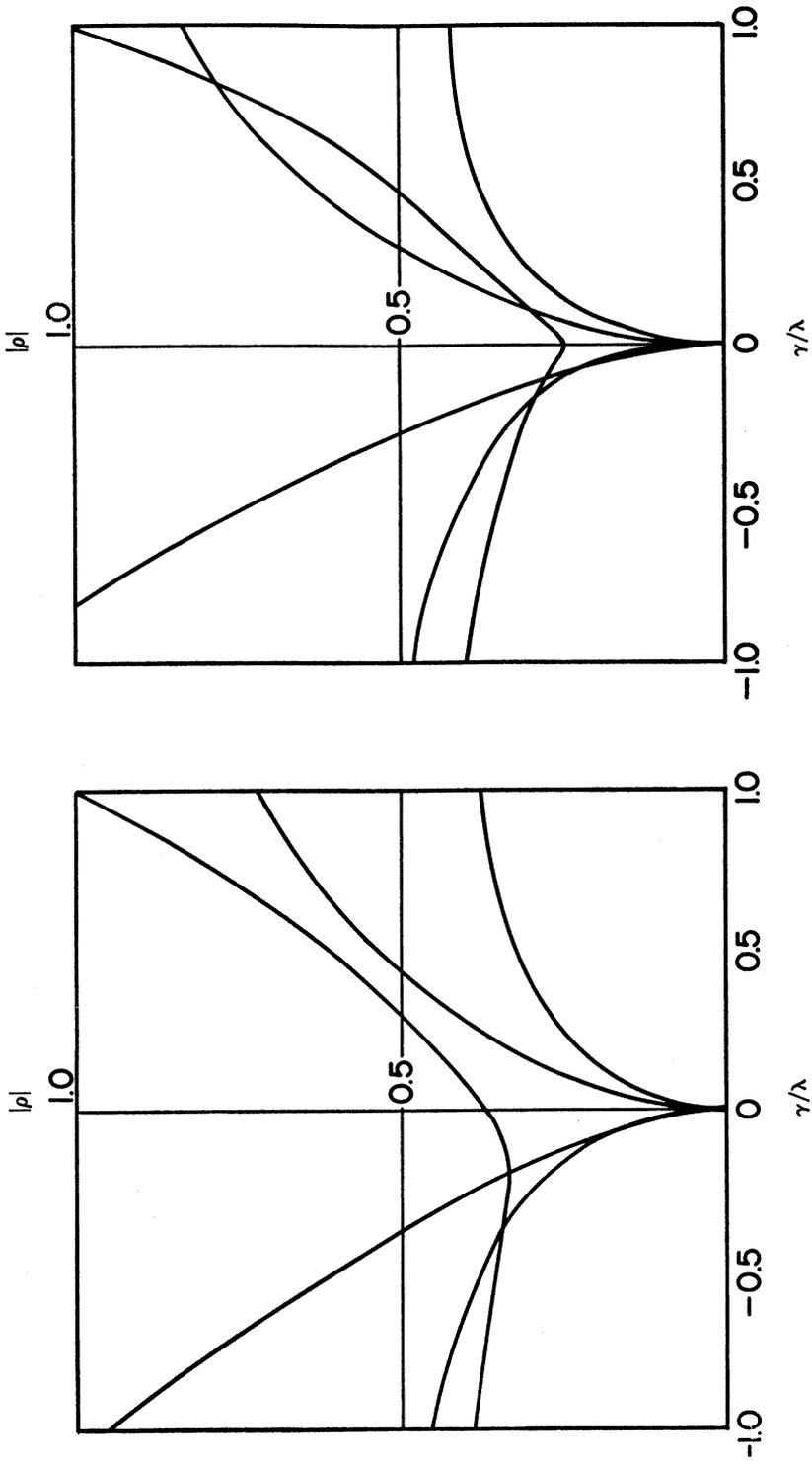
Here,  $w_j$  and  $v_j$  are elements of the matrices given in the Appendix. Now, consider the following characteristic polynomial:

$$(24) \quad \rho^{k+1} - \alpha_0 \rho^k - \alpha_1 \rho^{k-1} - \cdots - \alpha_{k-1} \rho - \alpha_k = 0.$$

The method is stable if the maximum of the absolute values of the roots of (24),  $|\rho|_{\max}$ , is less than one.

The function  $|\rho|_{\max}$  is continuous in  $\lambda h$  and  $\gamma h$ , but its first derivative need not be continuous since  $|\rho|_{\max}$  is not always attained by the same root. This can be seen from Fig. 2 which shows the relative magnitude of different roots in dependence upon  $\lambda h$  and  $\gamma h$ .

Fig. 3 presents for  $k = 4$  the relation between  $|\rho|_{\max}$  and the ratio  $\gamma/\lambda$  for various values of  $\lambda h$ . For  $\lambda h < 1.0$  and positive  $\gamma h$ ,  $|\rho|_{\max}$  is almost exactly approximated by the principal root,  $\exp[-(\lambda - \gamma)h]$ . But for negative  $\gamma h$ ,  $|\rho|_{\max}$  is sometimes attained by spurious roots which are not related to the ideal solution, even if  $\lambda h$  is small. Fig. 4 gives the same curves as Fig. 3, but for larger values of  $\lambda h$ . Here, no resemblance with the curve  $|\rho| = \exp[-(\lambda - \gamma)h]$  exists. For  $\gamma = 0$ , the method is, of course,



(a)  $\lambda h = 1.0$  and  $k = 4$

(b)  $\lambda h = 1.4$  and  $k = 4$

FIGURE 2. Effect of  $\lambda h$  and  $\gamma h$  on the Magnitude of the Roots

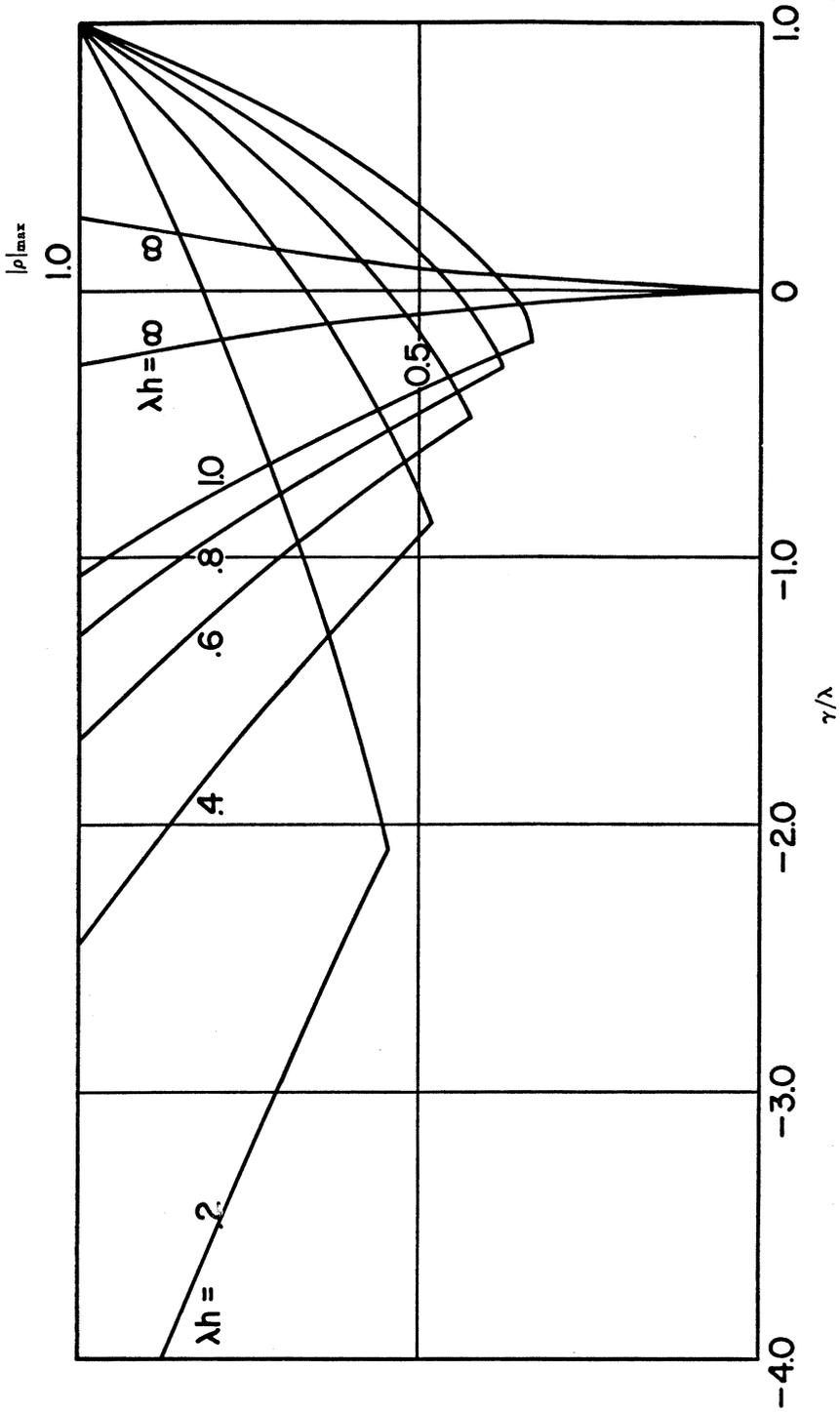


FIGURE 3.  $|\rho|_{\max}$  versus  $\gamma/\lambda$  for  $k = 4$  and Various  $\lambda h$

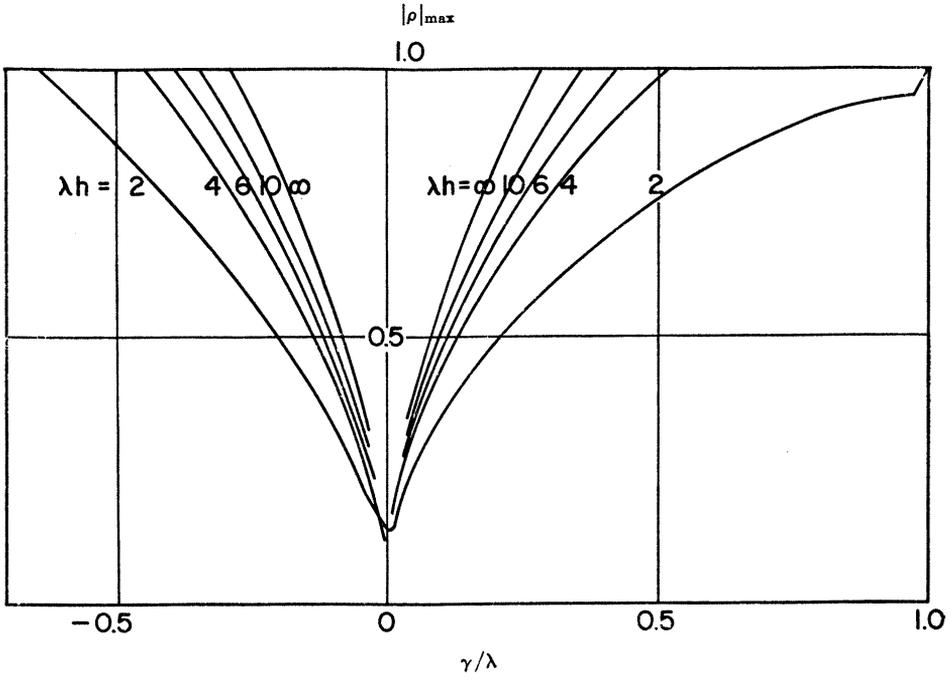


FIGURE 4.  $|\rho|_{\max}$  versus  $\gamma/\lambda$  for  $k = 4$  and  $\lambda h$  Larger than One

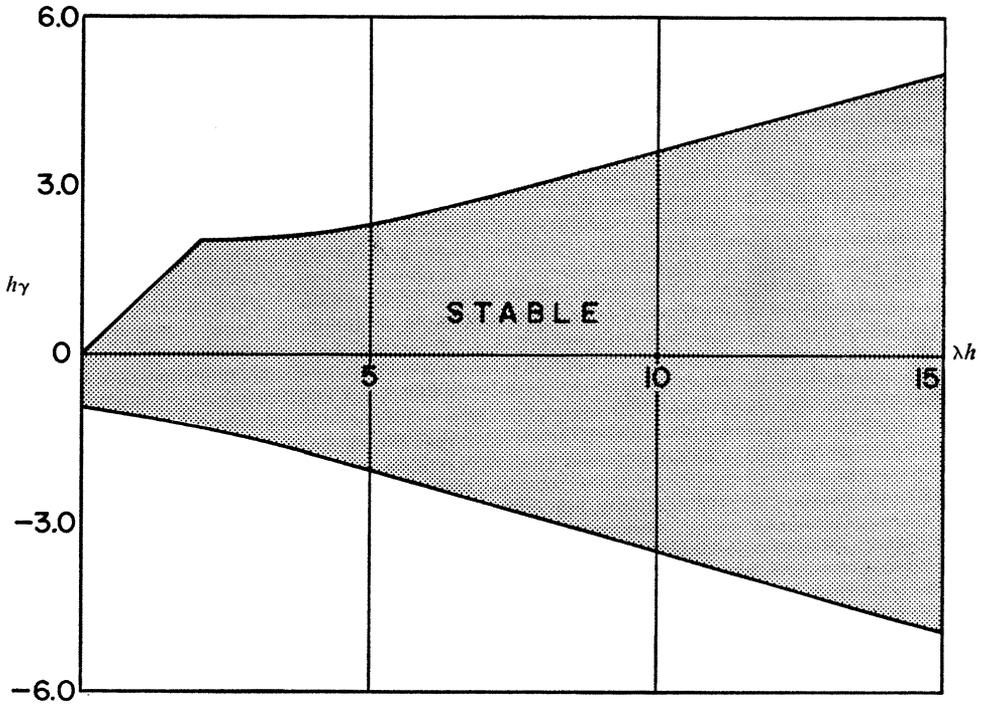


FIGURE 5. Region of Stability for  $k = 4$

exact and  $|\rho|_{\max}$  is given by the principal root,  $\exp(-\lambda h)$ . But, even in a close vicinity of  $\gamma = 0$ , very marked deviations from the correct principal root are encountered. For the exact solution, the right-hand side of (21) is  $\gamma \exp[-(\lambda - \gamma)h]$ ; in the integration scheme, this expression is replaced by a polynomial of degree  $k$ . The deviation between the correct principal root and the actual value of  $|\rho|_{\max}$  must be attributed to this approximation. For  $k = 4$ , the present method is stable for  $|\gamma| \leq 0.28\lambda$ , even if  $\lambda$  approaches infinity.

Fig. 5 shows the region of stability in  $\lambda$ - $\gamma$  plane for  $k = 4$ . Assume that one treats a number of equations simultaneously and that  $\gamma$  has about the same magnitude in all equations. Then, according to Fig. 5, the step size which is admissible from the point of view of stability is determined by those equations with small values of  $\lambda$ ; under these circumstances, the stiffness of the system, which expresses itself by the presence of large values of  $\lambda$ , is no longer critical.

Figs. 6-9 show the effects of the degree  $k$  of the polynomial on the stability. The admissible step size is smaller for higher values of  $k$ . For  $\lambda = 0$ , the method reduces to the usual predictor-corrector method and the results shown in Fig. 6 are the stability criteria.

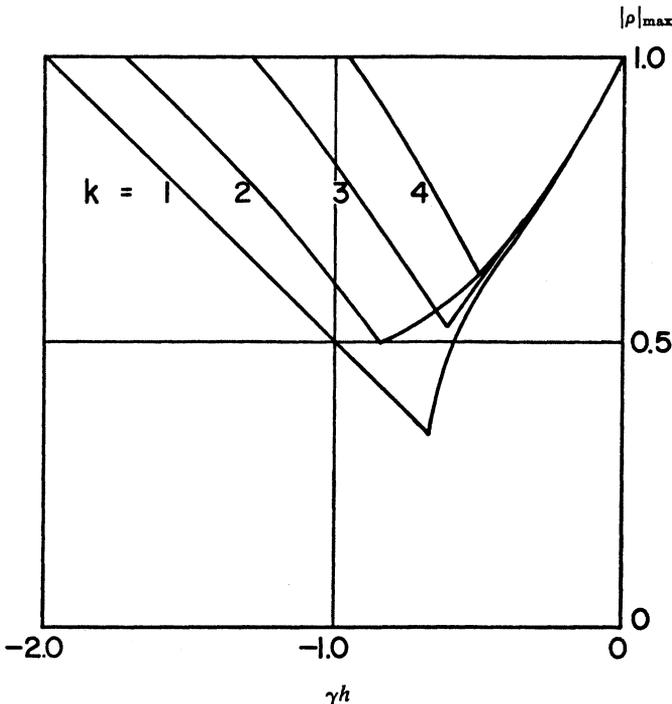


FIGURE 6. Effect of  $k$  on  $|\rho|_{\max}$  for  $\lambda h = 0$

5. Stability Criteria for a System. For a more general stability discussion of the system (1) we consider  $A$  as a constant matrix. The inhomogeneous term  $r(x)$  has no influence in the stability analysis, therefore  $f_i$  in (8) and (10) is replaced by  $Ay_i$ . In the following analysis, the degree of the predictor and corrector polynomials is

assumed to be  $k = 4$ . Substituting the predictor formula (8) into the corrector formula (10), one finds

$$(25) \quad \mathbf{y}_{n+1}^c = Q_0 \mathbf{y}_n + Q_1 \mathbf{y}_{n-1} + Q_2 \mathbf{y}_{n-2} + Q_3 \mathbf{y}_{n-3} + Q_4 \mathbf{y}_{n-4},$$

where the matrices  $Q_i$  are given by

$$(26) \quad \begin{aligned} Q_0 &= e^{-\Lambda h} + h W_0 A (e^{-\Lambda h} + h V_0 A) + h W_1 A, \\ Q_1 &= h W_2 A + h^2 W_0 A V_1 A, & Q_2 &= h W_3 A + h^2 W_0 A V_2 A, \\ Q_3 &= h W_4 A + h^2 W_0 A V_3 A, & Q_4 &= h^2 W_0 A V_4 A. \end{aligned}$$

Now one sets

$$(27) \quad \mathbf{y}_j = \rho^j \mathbf{z}, \quad \mathbf{y}_{n+1}^c = \mathbf{y}_{n+1},$$

where  $\mathbf{z}$  is a vector which does not depend on  $j$ , and  $\rho$  is a complex number. Substitution of (27) into (25) then yields a linear, homogeneous system

$$(28) \quad Q \mathbf{z} = 0,$$

with the  $N$  by  $N$  matrix  $Q$  given by

$$(29) \quad Q = \rho^5 I - \rho^4 Q_0 - \rho^3 Q_1 - \rho^2 Q_2 - \rho Q_3 - Q_4.$$

The system (28) has a nontrivial solution for  $\mathbf{z}$  if and only if the determinant  $|Q|$  vanishes. This requirement leads to a characteristic polynomial of degree  $5N$  in  $\rho$ . The integration method is stable if all zeros of this polynomial lie within the unit circle.

For an assumed value of the step size  $h$ , the elements of the matrices  $Q_0, \dots, Q_4$  can be computed without too much effort. However, if  $N$  is not small, the evaluation of the determinant requires many multiplications of polynomials. A test which avoids the explicit evaluation of the characteristic polynomial would, therefore, be preferable. We make certain simplifications which allow us to derive a sufficient condition for the stability of the system which is somewhat easier to apply. Because of the simplification made, the test may sometimes fail, although the system is stable. Assume that the diagonal part of the matrix  $Q$  is nonsingular and let

$$(30) \quad Q = Q_d + Q_{nd},$$

where  $Q_d$  and  $Q_{nd}$ , respectively, are the diagonal part and the nondiagonal part of  $Q$ . Then the system (28) can be written as

$$(31) \quad (I + Q_d^{-1} Q_{nd}) \mathbf{z} = 0.$$

Stability is guaranteed if we can show that for  $|\rho| \geq 1$  only the trivial solution for  $\mathbf{z}$  exists. Because of the triangular inequality, this is the case if, for  $|\rho| \geq 1$ , the following condition is satisfied

$$(32) \quad \|Q_d^{-1} Q_{nd}\| < 1.$$

Of course, the matrix norm used here must be consistent with the definition used for the vector norm.

For the vector norm, we propose to use a weighted infinity norm where the weights will be determined such that (32) is as permissive as possible. For such a norm defi-

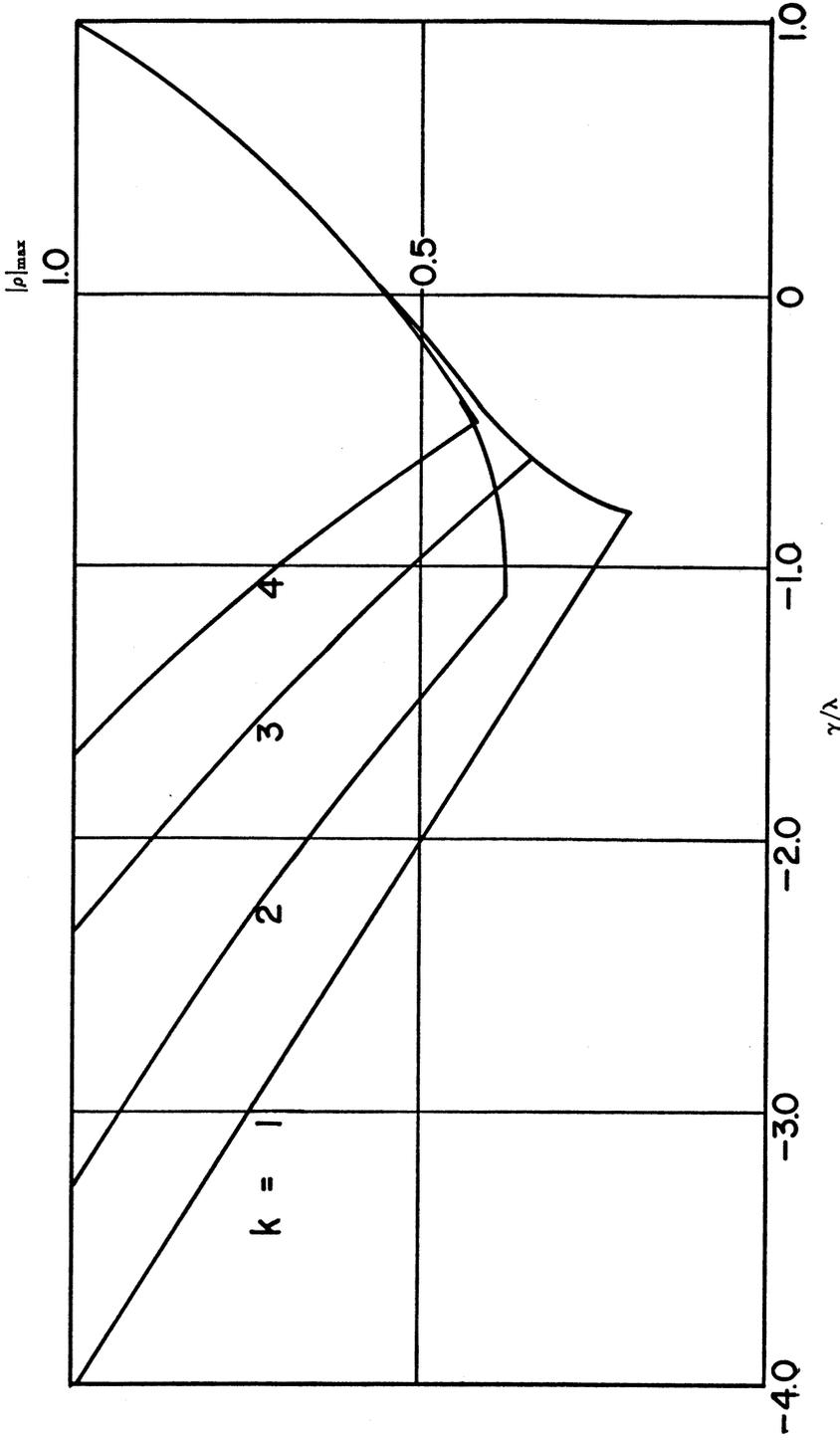


FIGURE 7. Effect of  $k$  on  $|\rho|_{\max}$  for  $\lambda h = 0.6$

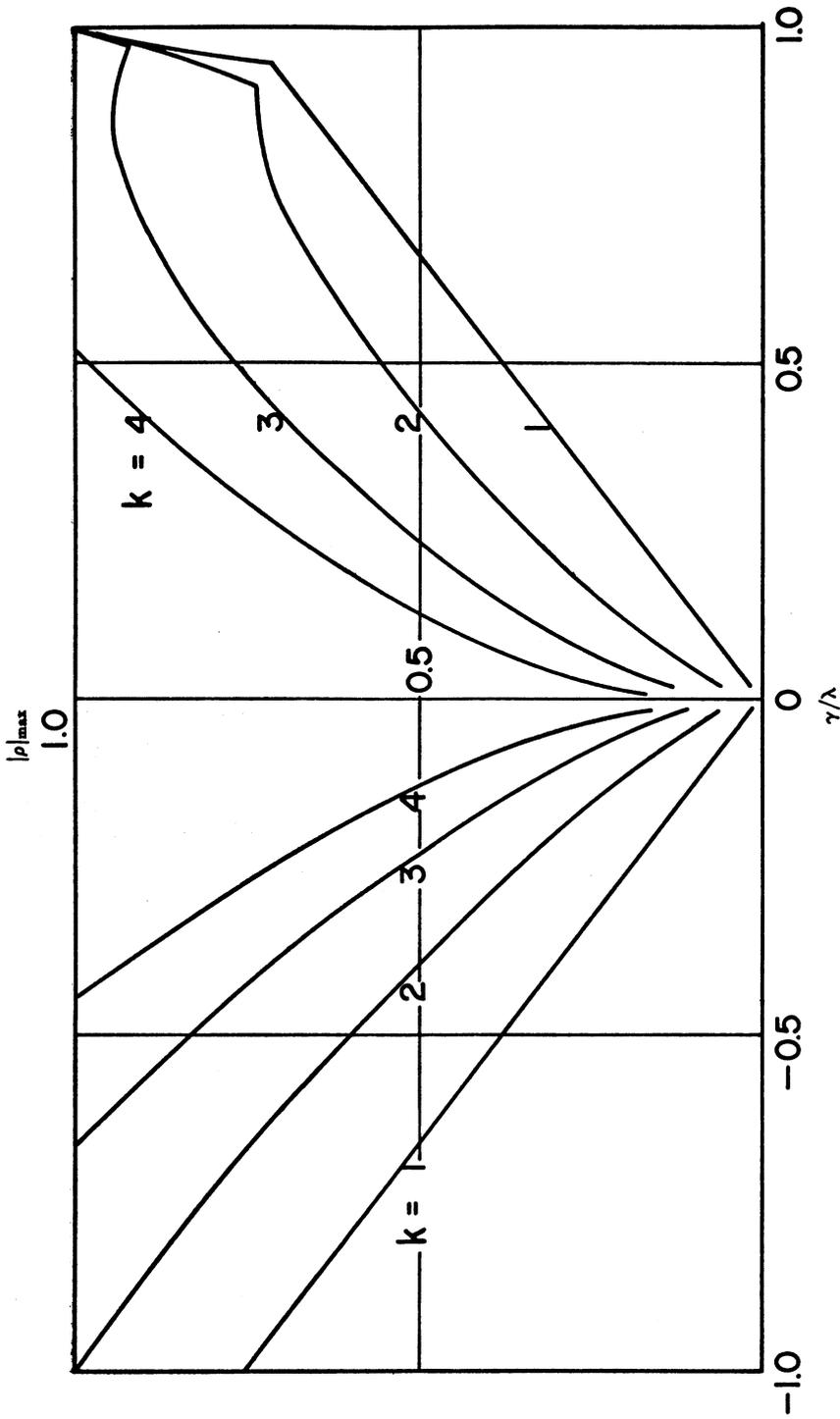


FIGURE 8. Effect of  $k$  on  $|\rho|_{\max}$  for  $\lambda h = 4.0$

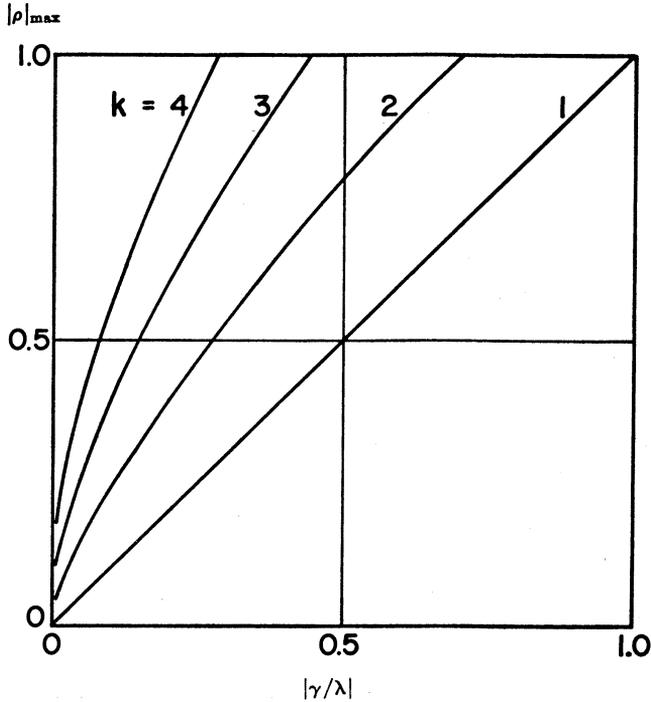


FIGURE 9. Effect of  $k$  on  $|\rho|_{\max}$  for  $\lambda h$  equal to Infinity

dition, only the absolute values of the element of  $Q_d^{-1}Q_{nd}$  are important. These elements are rational functions of  $\rho$ , which vanish at infinity; according to (29) the degree of the polynomials occurring in  $Q_d$  exceeds that of the polynomials for  $Q_{nd}$  by at least 1. We now make the assumption that zeros of all elements of  $Q_d$  lie within the unit circle. To check this, one can use the methods of [6], [7]. Then the elements of  $Q_d^{-1}Q_{nd}$  are regular functions in the region  $|\rho| \geq 1$  (including the point at infinity) and assume their maximum along the boundary  $|\rho| = 1$ . Let  $H$  be a matrix whose elements are given by

$$(33) \quad H_{ik} = \max_{\theta} [ |(Q_d^{-1}Q_{nd})_{ik}|, \rho = e^{i\theta} ].$$

We note that the diagonal elements of the matrix  $Q_d^{-1}Q_{nd}$  are zero.

We introduce an auxiliary vector  $u$  with positive components  $u_i$  and define the vector norm by

$$(34) \quad \|z\| = \max_i [(u_i)^{-1} |z_i|].$$

It follows that

$$(35) \quad |z_i| \leq u_i \|z\|.$$

In order for the definition of a matrix norm to be consistent with the given vector norm (34), one must choose

$$(36) \quad \|R\| = \max_i \left[ (u_i)^{-1} \sum_k |R_{ik}| u_k \right].$$

One has, according to the definition (34),

$$||Rz|| = \max_i \left[ (u_i)^{-1} \left| \sum_k R_{ik} z_k \right| \right] \leq \max_i \left[ (u_i)^{-1} \sum_k |R_{ik}| |z_k| \right].$$

Using (35), one arrives at the following relation, which implies the definition of the matrix norm (36):

$$||Rz|| \leq \max_i \left[ (u_i)^{-1} \sum_k |R_{ik}| u_k ||z|| \right].$$

The definition of the matrix norm (36) is now applied in (32), and (33) is substituted. This gives the following stability condition:

$$(u_i)^{-1} \sum_k H_{ik} u_k < 1.$$

This condition can be written in matrix notation

$$(37) \quad (I - H)\mathbf{u} > 0,$$

where the inequality sign is to be applied componentwise. The method is stable for step  $h$  if (37) is satisfied for vector  $\mathbf{u}$  with positive components. This requirement is not completely identical with the definition of matrices of the positive kind introduced in Section 23.1 of [8]. According to definition, the matrix  $(I - H)$  is of the positive kind if the vector  $\mathbf{u}$  has positive components whenever (37) is satisfied. We have the weaker requirement that there should exist at least one vector  $\mathbf{u}$  for which (37) is satisfied. The following theorem shows that  $(I - H)$  is actually a matrix of the positive kind.

**THEOREM.** *Let  $H = (H_{ij})$  be a matrix for which  $H_{ij} = 0$  if  $i = j$  and  $H_{ij} > 0$  if  $i \neq j$ . Then a vector  $\mathbf{u} = (u_i)$ ,  $u_i > 0$ , for which  $(I - H)\mathbf{u} > 0$  will exist if and only if the matrix  $(I - H)^{-1}$  exists and has all positive elements.*

*Remark.* The assumptions made for  $H$  are rather strong. It is likely that the proof can be carried out for irreducible matrices in nearly the same form. From a practical point of view one also should consider reducible matrices. For the purpose at hand, one can evade this problem by a continuity argument, i.e. we replace the zeros in the off-diagonal elements of  $H$  by small positive quantities and then apply the present theorem.

*Proof.* First we discuss the case that  $(I - H)^{-1}$  does not exist. We note that the matrix  $H$  satisfies the conditions of the Perron-Frobenius theorem for an irreducible matrix [9]. Accordingly, the maximum eigenvalue of  $H$  and its transpose  $H'$  is positive and the respective eigenvectors  $\mathbf{n}$  and  $\zeta$  are positive vectors.

**LEMMA 1.** *The real (and the imaginary) parts of the components of all other eigenvectors cannot have the same sign.*

The eigenvectors of  $H$  which do not belong to the maximum eigenvalue must be orthogonal to  $\zeta$  and the components of  $\zeta$  are all positive.

**LEMMA 2.** *The vectors of the null space of  $(I - H)$  can be considered as real.*

**LEMMA 3.** *A vector of the null space of  $(I - H)$  cannot have some zero components if all other components are positive.*

Assume that the  $k$ th component of a vector  $\mathbf{n}$  is zero but that all others are positive, then the  $k$ th component of  $(I - H)\mathbf{n}$  is negative, which contradicts the assumption that  $\mathbf{n}$  is a vector in the null space.

**LEMMA 4.** *If the null space of  $(I - H)$  contains one vector  $\mathbf{n}$ , whose components are*

all positive, then the null space is one dimensional.

Assume that within the null space there is a second vector  $\mathbf{n}_2$  with some positive and some negative components, then  $\mathbf{n}_3 = \mathbf{n}_1 + \alpha\mathbf{n}_2$  is also a vector of the null space. One can choose  $\alpha$  such that all components are nonnegative and at least one is zero. This contradicts Lemma 3. The same holds if  $\mathbf{n}_1$  and  $\mathbf{n}_2$  are linearly independent positive vectors.

LEMMA 5. *If  $(I - H)$  has zero for an eigenvalue, then there is no vector  $\mathbf{n}_1 > 0$  for which  $(I - H)\mathbf{n}_1 > 0$ .*

Case 1. The null space of  $(I - H)$  is one dimensional and contains one positive vector  $\mathbf{n}_1$ . Then  $\mathbf{n}_1$  is identical with the eigenvector of  $H$  which belongs to the maximum eigenvalue, for according to Lemma 1 all other eigenvectors of  $H$  have components with mixed signs. It follows that the vector adjoint to  $\mathbf{n}_1$ , to be denoted by  $\zeta_1$ , is also positive. In order for the equation  $(I - H)\mathbf{n} = \mathbf{z} > 0$  to be solvable for  $\mathbf{n}$ , one must have  $(\mathbf{z}, \zeta_1) = 0$ . But this condition cannot be satisfied for  $\zeta_1 > 0$ .

Case 2. The null space of  $(I - H)$  contains one vector  $\mathbf{n}_1$  with positive and negative components. Assume that there exists a vector  $\mathbf{n}_2 > 0$  for which  $(I - H)\mathbf{n}_2 > 0$ . One can then form a vector

$$\mathbf{n}_3 = \alpha\mathbf{n}_1 + (1 - \alpha)\mathbf{n}_2, \quad 0 < \alpha < 1,$$

for which at least one component is zero and all others are positive. Then one has on the one hand

$$(I - H)\mathbf{n}_3 = (1 - \alpha)(I - H)\mathbf{n}_2 > 0,$$

since  $(I - H)\mathbf{n}_1 = 0$ . On the other hand, assume that the  $k$ th component of  $\mathbf{n}_3$  is zero. Then one finds by direct evaluation that the  $k$ th component of  $(I - H)\mathbf{n}_3$  is negative. Because of this contradiction, the assumption that a vector  $\mathbf{n}_2 > 0$  exists is wrong.

This concludes the discussion of cases where  $(I - H)^{-1}$  does not exist.

If  $(I - H)^{-1}$  exists, then the sufficiency of the theorem is trivial. Since  $(I - H)^{-1}$  has all positive elements, an admissible vector  $\mathbf{u} > 0$  can be constructed as  $(I - H)^{-1}\mathbf{d}$ , where  $\mathbf{d} > 0$ .

The necessity is next shown by contradiction. First, assume that there exists a vector  $\mathbf{u}_I > 0$  for which

$$(38) \quad (I - H)\mathbf{u}_I > 0.$$

Next, assume that  $(I - H)^{-1}$  has some nonpositive components:

$$(39) \quad [(I - H)^{-1}]_{mn} \leq 0, \quad \text{for some suitable } (m, n).$$

Denote the  $m$ th column vector of  $(I - H)^{-1}$  by  $\mathbf{u}_{II}$ , then the components of  $\mathbf{u}_{II}$  are

$$(40) \quad (u_{II})_i = [(I - H)^{-1}]_{in}.$$

For some values of  $j$ , say  $j = k', k'', m$ ,

$$(41) \quad (u_{II})_j \leq 0.$$

We can now construct a vector

$$(42) \quad \mathbf{u}_{III} = (1 - \alpha)\mathbf{u}_I + \alpha\mathbf{u}_{II}.$$

For  $\alpha = 0$ ,  $u_{III} > 0$ ; for  $\alpha = 1$ , the components  $(u_{III})_j$ , for  $j = k', k'', m$ , are non-positive. Therefore, one can find a value of  $\alpha$ ,  $0 < \alpha \leq 1$ , such that all components of  $u_{III}$  are nonnegative and at least one is zero. Assume that

$$(43) \quad (u_{III})_{k'} = 0 \quad \text{and} \quad (u_{III})_i \geq 0 \quad \text{for } j \neq k'.$$

Based on (43), we can now compute  $(I - H)u_{III}$  by two different methods. If we write

$$(44) \quad (I - H)u_{III} = (1 - \alpha)(I - H)u_I + \alpha(I - H)u_{II}, \quad 0 < \alpha \leq 1,$$

then all the components of  $(I - H)u_{III}$  are nonnegative because of (38) and since the components  $(I - H)u_I$  are given by  $\delta_{jn}$ . In particular, one has

$$(45) \quad [(I - H)u_{III}]_{k'} \geq 0.$$

On the other hand, we can also write

$$(46) \quad [(I - H)u_{III}]_{k'} = (u_{III})_{k'} - \sum_i H_{k'i}(u_{III})_i < 0,$$

which is negative because of (43) and the definition of  $H$ . The contradiction between (45) and (46) shows that the assumption (39) is false. This completes the proof of the theorem.

Instead of determining the inverse of  $(I - H)$ , one can solve the system

$$(47) \quad (I - H)u_1 = x_1$$

with respect to  $u_1$  for some  $x_1 > 0$ . A vector  $u > 0$  for which  $(I - H)u > 0$  will exist only if  $u_1 > 0$ . The particular choice of  $x_1$  is unessential. This can be shown by replacing the vector  $u_{II}$  in the above argument by  $u$ .

If a vector  $u > 0$  satisfying (37) exists, then  $\|H\| < 1$  and (32) will hold outside of the unit circle. The explicit form of  $u$  need not be known. Sometimes, in particular if  $H$  is small, it will be easiest to find  $u$  by inspection.

It is crucial that the matrix  $H$  give a rather narrow bound for  $Q_d^{-1}Q_{nd}$ . Such a bound would be obtained if one searches along the unit circle for the maxima of the individual elements of  $Q_d^{-1}Q_{nd}$ . If  $H$  is smaller, then rougher estimations will be sufficient. One can, for instance, find upper bounds for the elements of  $Q_{nd}$  and lower bounds for  $Q_d$  along the unit circle from the coefficient of the polynomials. If the determination of the lower bound of  $Q_d$  is critical, it might be preferable to search for it along the unit circle. The following examples show the effects of different approximations for  $Q_d^{-1}Q_{nd}$ .

We consider a system of two equations in which the diagonal matrix  $\Lambda$  is given by

$$\Lambda = \begin{bmatrix} 1 & 0 \\ 0 & 100 \end{bmatrix}.$$

For the matrix  $A$ , we have assumed three different forms

$$A_1 = \begin{bmatrix} \frac{1}{2} & 1 \\ 1 & 30 \end{bmatrix}, \quad A_2 = \begin{bmatrix} \frac{1}{2} & 1 \\ 1 & 20 \end{bmatrix}, \quad A_3 = \begin{bmatrix} \frac{1}{2} & 1 \\ 1 & 10 \end{bmatrix}.$$

An exact stability analysis gives the corresponding admissible step sizes as

$$h_1 = 0.55, \quad h_2 = 3.30, \quad h_3 = 3.66.$$

If  $H$  is obtained by searching the maximum elements of  $Q_a^{-1}Q_{na}$  along the unit circle, then (47) gives the admissible step size for stability as

$$h_1 = 0.40, \quad h_2 = 2.55, \quad h_3 = 3.15.$$

These results show that the present approach gives good approximations to the actual stability limits. If  $H$  is found by bounding the elements of  $Q_a$  and  $Q_{na}$  separately, the results obtained are less satisfactory. The values of  $h$  for which stability can then be guaranteed are

$$h_1 < 0.1, \quad h_2 = 0.2, \quad h_3 = 1.31.$$

The results summarized in Fig. 5 would suggest that a method is stable for a certain  $h$  if stability has been established for a larger  $h$  but this conjecture has not been proven.

Appendix. The constant matrices  $B^P$  and  $B^C$ , for  $k = 4$ , are

$$B^P = \begin{bmatrix} 0 & -1/4 & 11/24 & -1/4 & 1/24 \\ 0 & 4/3 & -7/3 & 7/6 & -1/6 \\ 0 & -3 & 19/4 & -2 & 1/4 \\ 0 & 4 & -13/3 & 3/2 & -1/6 \\ 1 & -25/12 & 35/24 & -5/12 & 1/24 \end{bmatrix},$$

$$B^C = \begin{bmatrix} 0 & -1/12 & -1/24 & 1/12 & 1/24 \\ 0 & 1/2 & 1/6 & -1/2 & -1/6 \\ 0 & -3/2 & 1/4 & 1 & 1/4 \\ 1 & 5/6 & -5/6 & -5/6 & -1/6 \\ 0 & 1/4 & 11/24 & 1/4 & 1/24 \end{bmatrix}.$$

Let  $M \equiv \Delta h$  and  $E \equiv e^{-M}$ , then the diagonal matrices  $V_i$  and  $W_i$  for  $k = 4$  are given by

$$V_0 = M^{-5} \left[ M^4(5 - E) - M^3 \left( \frac{77}{12} - \frac{25}{12} E \right) \right. \\ \left. + M^2 \left( \frac{71}{12} - \frac{35}{12} E \right) - M \left( \frac{7}{2} - \frac{5}{2} E \right) + (1 - E) \right],$$

$$V_1 = -M^{-5} \left[ 10M^4 - M^3 \left( \frac{107}{6} - 4E \right) \right. \\ \left. + M^2 \left( \frac{59}{3} - \frac{26}{3} E \right) - M(13 - 9E) + 4(1 - E) \right],$$

$$V_2 = M^{-5} \left[ 10M^4 - M^3 \left( \frac{39}{2} - 3E \right) \right. \\ \left. + M^2 \left( \frac{49}{2} - \frac{19}{2} E \right) - M(18 - 12E) + 6(1 - E) \right],$$

$$V_3 = -M^{-5} \left[ 5M^4 - M^3 \left( \frac{61}{6} - \frac{4}{3} E \right) + M^2 \left( \frac{41}{3} - \frac{14}{3} E \right) - M(11 - 7E) + 4(1 - E) \right],$$

$$V_4 = M^{-5} \left[ M^4 - M^3 \left( \frac{25}{12} - \frac{1}{4} E \right) + M^2 \left( \frac{35}{12} - \frac{11}{12} E \right) - M \left( \frac{5}{2} - \frac{3}{2} E \right) + (1 - E) \right];$$

$$W_0 = V_4,$$

$$W_1 = -M^{-5} \left[ M^4 E - M^3 \left( 4 + \frac{5}{6} E \right) + M^2 \left( \frac{26}{3} - \frac{5}{3} E \right) - M(9 - 5E) + 4(1 - E) \right],$$

$$W_2 = M^{-5} \left[ -M^3 \left( 3 + \frac{3}{2} E \right) + M^2 \left( \frac{19}{2} - \frac{1}{2} E \right) - M(12 - 6E) + 6(1 - E) \right],$$

$$W_3 = -M^{-5} \left[ -M^3 \left( \frac{4}{3} + \frac{1}{2} E \right) + M^2 \left( \frac{14}{3} + \frac{1}{3} E \right) - M(7 - 3E) + 4(1 - E) \right],$$

$$W_4 = M^{-5} \left[ -M^3 \left( \frac{1}{4} + \frac{1}{12} E \right) + M^2 \left( \frac{11}{12} + \frac{1}{12} E \right) - M \left( \frac{3}{2} - \frac{1}{2} E \right) + (1 - E) \right].$$

We note that  $V_i$  and  $W_i$  are finite for  $M = 0$ .

Aerospace Research Laboratories  
Wright-Patterson Air Force Base  
Ohio 45433

Ohio State University Research Foundation  
Columbus, Ohio 43210

1. K. G. GUDERLEY & C. C. HSU, "A special form of Galerkin's method applied to heat transfer in plane Couette-Poiseuille flows." (In prep.)

2. C. W. GEAR, "The automatic integration of stiff ordinary differential equations. (With discussion)," *Proc. IFIP Congress Information Processing 68* (Edinburgh, 1968), vol. 1: *Mathematics, Software*, North-Holland, Amsterdam, 1969, pp. 187-193. MR 41 #4808.

3. D. A. CALAHAN, "Numerical solution of linear systems with widely separated time constants," *Proc. IEEE*, v. 55, 1967, pp. 2016-2017.

4. J. L. BLUE & H. K. GUMMEL, "Rational approximations to matrix exponential for systems of stiff differential equations," *J. Computational Physics*, v. 5, 1970, pp. 70-83. MR 40 #8267.

5. P. HENRICI, *Discrete Variable Methods in Ordinary Differential Equations*, Wiley, New York, 1962. MR 24 #B1772.

6. F. JOHN, *Lectures on Advanced Numerical Analysis*, Gordon and Breach, New York, 1967. MR 36 #4773.

7. M. MARDEN, *Geometry of Polynomials*, 2nd ed., Math. Surveys, no. 3, Amer. Math. Soc., Providence, R. I., 1966. MR 37 #1562.

8. L. COLLATZ, *Funktionalanalysis und numerische Mathematik*, Die Grundlehren der math. Wissenschaften, Band 120, Springer-Verlag, Berlin, 1964; English transl., Academic Press, New York, 1966. MR 29 #2931; MR 34 #4961.

9. R. S. VARGA, *Matrix Iterative Analysis*, Prentice-Hall, Englewood Cliffs, N. J., 1962. MR 28 #1725.

10. J. CERTAINE, "The solution of ordinary differential equations with large time constants," in *Mathematical Methods for Digital Computers*, A. Ralston and H. S. Wilf (Editors), Wiley, New York, 1960, pp. 128-132. MR 22 #8691.