

# Matricial Difference Schemes for Integrating Stiff Systems of Ordinary Differential Equations

By W. L. Miranker\*

**Abstract.** In this paper we give a description and analysis of a class of matricial difference schemes. This class of schemes is based in part on a generalization of the feature of classical numerical methods of being characterized by approximations at a single point in the complex plane. The schemes introduced here are effective for integrating stiff systems.

1. Introduction. In this paper we give a description and analysis of a class of difference schemes of matricial type. The special nature of this class of schemes makes them effective for the integration of stiff systems of differential equations.

The system

$$(1.1) \quad \dot{x} = Ax$$

of ordinary differential equations has for its solution

$$(1.2) \quad x_{n+1} = \exp(hA)x_n$$

where  $h$  is a mesh increment and  $x_n = x(nh)$ . A difference approximation to (1.1) has a solution  $u(t)$  which may be written as

$$(1.3) \quad u_{n+1} = K(hA)u_n.$$

$K(z)$  is typically a polynomial or rational function of  $z$ . The control of the error  $e_n = u_n - x_n$  depends on  $K(hA)$  being stable (boundedness of the norm of the powers of  $K(hA)$ ) and the closeness of  $K(hA)$  to  $\exp(hA)$ . Existing methods usually handle these two features in the following way;  $K(hA)$  will be close to  $\exp(hA)$  if it is close on the spectrum  $\sigma(hA)$  of  $hA$ . Since  $\sigma(A)$  has some arbitrary configuration in the complex plane, making  $K(hA)$  close to  $\exp(hA)$  is accomplished in two steps. First,  $K(z)$  is chosen close to  $\exp(z)$  in a neighborhood of  $z = 0$  (e.g.  $K(z) = 1 + z$  for the forward Euler). Then  $h$  is reduced until  $h \times \sigma(A)$  is shrunk into the neighborhood of  $z = 0$  where  $K(z)$  is close to  $\exp(z)$ . The stability of  $K(hA)$  is accomplished by making  $|K(z)| \leq 1$ , typically in a set containing  $z = 0$  (e.g. in the  $z = x + iy$  plane this set is  $x^2 + y^2 - 2x \geq 0$  for the backward Euler). Then  $h$  is reduced until  $h \times \sigma(A)$  of  $A$  is shrunk into this neighborhood.

The methods introduced here are not restricted to operate in this classical way. Rather the single point, the origin, is replaced by a set of points in the complex plane and the approximation of  $K(z)$  to  $\exp(z)$  is arranged at this set of points. To be effective, the set of points at which the approximation is to be made should have certain characteristics (to be described later) relative to the spectrum of  $A$ . Generally,

---

Received September 25, 1969, revised March 28, 1971.

AMS 1969 subject classifications. Primary 6560, 6561; Secondary 3937.

Key words and phrases. Stiff systems, matricial difference schemes.

\* This work was performed while the author was a visiting Professor of Mathematics at the Hebrew University of Jerusalem.

these characteristics can be determined at the expense of additional computations. A favorable situation for this procedure is provided by stiff systems.

A definition of a stiff system is one for which the stiffness

$$(1.4) \quad s = \max |\lambda| / \min |\lambda|$$

is a large number. The max and min are taken over the nonzero  $\lambda \in \sigma(A)$ .

Since a classical method of numerical integration when applied to a stiff system is enormously handicapped by the requirement that  $|h\lambda|$  is small for the large  $|\lambda|$  in order to gain stability and accuracy, it becomes feasible to expand the computation to gain information about  $\sigma(A)$  for exploitation by the method discussed here.

There are a number of schemes in the literature which deal with integrating stiff systems, (see references [1]–[5]). A general approach to stiffness is through the approach of  $A$ -stability, (see [7]). Here, an  $A$ -stable scheme, such as the backward Euler, is used. In the initial transitional region,  $h$  must be small to get any accuracy. However, as the solution smooths out as time increases,  $h$  is increased. The  $A$ -stability of the method guarantees not to excite the modes with large  $\lambda$  which have become quiescent, while the smallness of these modes allows accuracy to be maintained. In this connection, there are the results of G. G. Dahlquist [7] and C. W. Gear [8]. Dahlquist's result places a strong limitation on this approach, since he shows that  $A$ -stable methods of order greater than two do not exist within the class of linear multistep methods, considered by him.

In Section 2, we describe our integration schemes and give a stability and error analysis for them. In Section 3, we describe the results of calculations based on one family of our schemes. Included is a comparison with an effective classical method. In an appendix, we give an error analysis of a scheme due to W. Liniger and R. Willoughby [1]. This scheme is a scalar version of the class of schemes presented here where, moreover, the set of points in the complex plane (referred to above) consists of two points on the real line, one of which is the origin. Even for this extremely special case, enormously effective computations were performed for large nonlinear systems (see [1]).

**2. Statement of Results.** Consider the system of ordinary differential equations

$$(2.1) \quad \dot{x} = Ax, \quad t > 0.$$

Here  $x$  is an  $m$ -vector and  $A$  is an  $m \times m$  matrix. Let  $h$  denote a mesh increment and let  $x_n = x(nh)$ . Then

$$(2.2) \quad x_n = e^{Ah} x_{n-1}.$$

Now let\*\*

$$(2.3) \quad L(z) = \sum (\alpha_i + z\beta_i)e^{(r-i)z}, \quad R(z) = \sum (\gamma_i + z\delta_i)e^{(r-i)z}$$

and let

$$(2.4) \quad C(z) = L(z)[R(z)]^{-1}.$$

---

\*\* Unless otherwise specified all sums are taken from 0 to  $r$ .

For any solution  $x$  of (2.1), we have

$$(2.5) \quad [L(hA) - C(hA)R(hA)]x_{n-r} = 0$$

identically. This follows from (2.4) by inserting (2.2) and (2.3) into (2.5). This suggests the following difference scheme

$$(2.6) \quad \sum \alpha_i u_{n-i} + h \sum \beta_i A u_{n-i} - P(hA) \left[ \sum \gamma_i u_{n-i} + h \sum \delta_i A u_{n-i} \right] = 0,$$

$n = r, r+1, \dots$ , for determining the  $m$ -vector valued mesh function  $u_n$  as an approximation to  $x_n$ . Here,  $P(z)$  is an approximation to  $C(z)$  to be specified.

Let  $H$  be the shift operator

$$(2.7) \quad Hf(t) = f(t+h).$$

Let  $\mathcal{L} = \mathcal{L}(H)$  and  $\mathcal{R} = \mathcal{R}(H)$  be the operators corresponding to  $L(z)$  and  $R(z)$  respectively, i.e.

$$(2.8) \quad \mathcal{L}(H) = \sum (\alpha_i + hA\beta_i)H^{r-i}, \quad \mathcal{R}(H) = \sum (\gamma_i + hA\delta_i)H^{r-i}.$$

Then (2.5) and (2.6) can be written respectively as

$$(2.9) \quad [\mathcal{L}(H) - C(hA)\mathcal{R}(H)]x_{n-r} = 0,$$

$$(2.10) \quad [\mathcal{L}(H) - P(hA)\mathcal{R}(H)]u_{n-r} = 0.$$

By subtracting (2.9) from (2.10), the error  $e_n = u_n - x_n$  is seen to satisfy the following equation

$$(2.11) \quad [\mathcal{L}(H) - P(hA)\mathcal{R}(H)]e_{n-r} = [P(hA) - C(hA)]\mathcal{R}(H)x_{n-r}.$$

Of course,

$$(2.12) \quad \mathcal{R}(H)x_{n-r} = R(hA)x_{n-r}.$$

In order to estimate  $e_n$ , we first solve (2.11) for  $e_n$ .

Let

$$(2.13) \quad \mathcal{S}(H) = \mathcal{L}(H) - P(hA)\mathcal{R}(H) \equiv \sum s_i H^{r-i}.$$

Thus

$$(2.14) \quad s_i \equiv s_i(A) \equiv \alpha_i + hA\beta_i - P(hA)(\gamma_i + hA\delta_i).$$

Using (2.12) and (2.13), we may rewrite (2.11) as

$$(2.15) \quad \mathcal{S}(H)e_{n-r} = [P(hA) - C(hA)]R(hA)x_{n-r}.$$

Further let

$$(2.16) \quad S(z) = \sum s_i z^{r-i}.$$

Finally, define  $\sigma_j$ ,  $j = 0, 1, \dots$ , formally through

$$(2.17) \quad [z^r S(z^{-1})]^{-1} = \sum_{j=0}^{\infty} \sigma_j z^j.$$

Now multiply (2.15) by  $\sigma_{N-n}$  and sum the result over  $n$  from  $r$  to  $N$ . The left member resulting from this procedure is

$$\begin{aligned}
 \sum_{n=r}^N \sigma_{N-n} S(H) e_{n-r} &= \sum_{n=r}^N \sigma_{N-n} \sum_{i=0}^r s_i H^{r-i} e_{n-i} \\
 (2.18) \qquad \qquad \qquad &= \sigma_0 R_0 e_N + (\sigma_1 R_0 + \sigma_0 R_1) e_{N-1} \\
 &\quad + \cdots + (\sigma_{N-r} s_0 + \cdots + \sigma_{N-2r} s_r) e_r \\
 &\quad + \text{multiples of } e_0, \cdots, e_{r-1}.
 \end{aligned}$$

Using the relations (2.16) and (2.17) between the  $s_i$  and the  $\sigma_i$  simplifies (2.18) to

$$(2.19) \qquad \sum_{n=r}^N \sigma_{N-n} S(H) e_{n-r} = e_N + \text{multiples of } e_0, \cdots, e_{r-1}.$$

Thus this summing procedure solves (2.15) and gives

$$\begin{aligned}
 (2.20) \qquad e_N &= \sum_{n=r}^N \sigma_{N-n} [P(hA) - C(hA)] R(hA) x_{n-r} \\
 &\quad + \text{multiples of } e_0, \cdots, e_{r-1}.
 \end{aligned}$$

Now, to get an estimate of  $e_N$ , we need a stability and an accuracy statement. The stability statement is given by the following two lemmas.

**LEMMA 1.** *If  $\sum s_i(\lambda)z^{n-i}$  obeys the root condition for all eigenvalues  $\lambda$  of  $A$ , then the determinant of  $S(z)$  obeys the root condition.*

*Proof.* Let  $f(A) = \sum s_i(A)z^{n-i}$ . Suppose the determinant  $|f(A)|$  vanishes for a value of  $z$ , then  $|f(A) + \mu I - \mu I|$  vanishes. Then  $\mu = \mu + f(\lambda)$ , for  $\lambda$  any eigenvalue of  $A$ , or  $f(\lambda) = 0$  for that value of  $z$ . Q.E.D.

**LEMMA 2.** *Let the determinant  $|S(z)|$  of  $S(z)$  obey the root condition. If the determinant of  $s_0$  is not zero, then the matrix  $[z^r S(z^{-1})]^{-1}$  is analytic in a neighborhood of  $z = 0$ . Furthermore, the matrices  $\sigma_j$ ,  $j = 0, 1, \cdots$ , given by (2.17), have uniformly bounded norms.*

*Proof.* Since  $z^r S(z^{-1}) = \sum_{i=0}^r s_i z^i$  and  $|s_0| \neq 0$ , it is clear that  $[z^r S(z^{-1})]^{-1}$  is analytic in a neighborhood of the origin. Since  $|z^r S(z^{-1})| = z^{mr} |S(z^{-1})|$ , the root condition locates the roots of the polynomial  $|z^r S(z^{-1})|$  outside the open unit disc and those on the boundary of the unit disc are simple. (Note that the apparent root at the origin occurring from  $z^{mr}$  is annihilated by a corresponding pole of  $|S(z^{-1})|$ .) Since

$$[z^r S(z^{-1})]^{-1} = [\text{matrix of polynomials}] / |z^r S(z^{-1})|,$$

it suffices to show that the power series for the reciprocal polynomial  $|z^r S(z^{-1})|^{-1}$  has bounded coefficients, given that its only roots are outside the open unit disc, with those on the boundary being simple. Let  $mr = q$  and let

$$|z^r S(z^{-1})|^{-1} = \left[ \sum_{i=0}^q t_i z^i \right]^{-1} = \sum_{i=0}^{\infty} u_i z^i.$$

Then

$$u_n = \frac{1}{2\pi i} \oint \left( z^{n+1} \sum_{i=0}^q t_i z^i \right)^{-1} d\zeta,$$

where the contour of integration lies inside the unit disc and encircles the origin.

If we move the contour through the unit disc and out to infinity in all directions, the integral will vanish if  $q \geq 1$  and we are left with the sum of the residues. If there is a pole at  $\zeta_0$  on the unit disc, it is simple. Let the residue from it be  $\mathcal{R}_0$ . Then

$$|\mathcal{R}_0| = \left| \left( \zeta_0^{n+1} \sum_{i=0}^q j t_i \zeta_0^{i-1} \right)^{-1} \right| = \left| \sum_{i=0}^q j t_i \zeta_0^{i-1} \right|^{-1}$$

which is independent of  $n$ .

If there is a pole at  $\zeta_1$  of order  $\rho + 1$  outside the unit disc, let the residue from it be  $\mathcal{R}_1$ . Then

$$\mathcal{R}_1 = D_{\zeta_1}^{\rho} \left[ (\zeta - \zeta_1)^{\rho+1} \left( \zeta^{n+1} \sum_{i=0}^q t_i \zeta^i \right)^{-1} \right]_{\zeta=\zeta_1} = D_{\zeta_1}^{\rho} [(\zeta_1^{n+1} Q(\zeta_1))^{-1}],$$

where the polynomial  $Q$  is given by

$$Q(\zeta) = \left( \sum_{i=0}^{\infty} t_i \zeta^i \right) / (\zeta - \zeta_1)^{\rho+1}.$$

Then

$$\begin{aligned} \mathcal{R}_1 &= \sum_{j=0}^{\rho} \binom{\rho}{j} D_{\zeta_1}^j (\zeta_1^{-n}) D_{\zeta_1}^{\rho-j} (\zeta_1 Q(\zeta_1)) \\ &= \sum_{j=0}^{\rho} (-1)^j \binom{\rho}{j} n(n+1) \cdots (n+j-1) \zeta_1^{-(n+i)} D^{\rho-j} (\zeta_1 Q(\zeta_1)). \end{aligned}$$

Then

$$|\mathcal{R}_1| \leq (n + \rho)^{\rho} F / |\zeta_1|^n$$

where the factor  $F$  is independent of  $n$ . Since  $|\zeta_1| > 1$ , this residue vanishes as  $n \rightarrow \infty$ . Since there are at most a finite number of residues to be accounted for, the coefficients  $u_n$  are bounded uniformly in  $n$  and the lemma is proved.

If  $S(z)$  satisfies the hypothesis of Lemma 2, then (2.20) may be used to get

$$(2.21) \quad \|e_N\| \leq \text{const} \|[P(hA) - C(hA)]R(hA)\| \sum_{n=r}^N \|x_{n-r}\|,$$

where we have assumed that the initial errors,  $e_0, e_1, \dots, e_{r-1}$  may be neglected. If we assume that  $Nh = 1$ , this becomes

$$(2.22) \quad \|e_N\| \leq \text{const} h^{-1} \|[P(hA) - C(hA)]R(hA)\|.$$

Now we turn to the question of accuracy. We must introduce hypotheses so that we can estimate  $[P(hA) - C(hA)]R(hA)$ .

Let

$$(2.23) \quad L(z) = z^{\mu} + O(z^{\mu+1}), \quad R(z) = z^{\nu} + O(z^{\nu+1})$$

as  $z \rightarrow 0$ . This hypothesis amounts to stating that the difference scheme (2.6) in the limiting cases  $P = 0$  and  $P = \infty$  has order of accuracy  $\mu - 1$  and  $\nu - 1$ , respectively.

Let  $P(z)$  be chosen as the polynomial which has contact  $\tau_i$  with  $C(z)$  at the nonzero nodes  $hz_i, i = 1, \dots, p$ , i.e.

$$(2.24) \quad P^{(m)}(hz_i) - C^{(m)}(hz_i) = 0, \quad m = 0, \dots, \tau_i - 1.$$

Now divide the eigenvalues of  $A$  into  $p + 1$  clusters,  $k_0, \dots, k_p$ . Let  $z_0 = 0$ . All eigenvalues in  $k_i$  are closer to  $z_i$  than to all  $z_j, j \neq i, i = 0, 1, \dots, p$ . Ties are resolved randomly. Let

$$(2.25) \quad d_i = \max_{\lambda_j \in k_i} |\lambda_j - z_i|, \quad i = 0, \dots, p.$$

From (2.3) and (2.23), we have

$$(2.26) \quad |L(z)| \leq \text{const} \min(|z|, |z|^\nu), \quad |R(z)| \leq \text{const} \min(|z|, |z|^\mu)$$

for  $\text{Re } z \leq 0$ .

Now using the spectral representation theorem (see [10]) (for simplicity we treat only the case that the eigenvalues of  $A$  are distinct), we have

$$(2.27) \quad \begin{aligned} & [P(hA) - C(hA)]R(hA) \\ &= \sum_{i=0}^p \sum_{\lambda_j \in k_i} [P(h\lambda_j) - C(h\lambda_j)]R(h\lambda_j)Z_{i,j}(hA)^{***} \\ &= \sum_{\lambda_j \in k_0} [P(h\lambda_j)R(h\lambda_j) - L(h\lambda_j)]Z_{0,j}(hA) \\ &\quad + \sum_{i=1}^p \sum_{\lambda_j \in k_i} \frac{1}{\tau_i!} [h(\lambda_j - z_i)]^{\tau_i} [P^{(\tau_i)}(h\tilde{\lambda}_{i,j}) - C^{(\tau_i)}(h\tilde{\lambda}_{i,j})]R(h\lambda_j)Z_{i,j}(hA). \end{aligned}$$

Here  $\tilde{\lambda}_{i,j}$  and  $\tilde{\lambda}_{i,j}$  arise from Taylor's theorem and represent appropriate values corresponding to  $\lambda_j$  and  $z_i$ , and the functions  $P$  and  $C$ , respectively.

Then

$$(2.28) \quad \begin{aligned} ||[P(hA) - C(hA)]R(hA)|| &\leq C_1 \max(|hd_0|^\nu, |hd_0|^\mu) \\ &+ C_2 \sum_{i=1}^p \frac{1}{\tau_i!} |hd_i|^{\tau_i}. \end{aligned}$$

Inserting (2.28) into (2.22) gives

$$(2.29) \quad ||e_N|| \leq \frac{\text{const}}{h} \left[ \max(|hd_0|^\nu, |hd_0|^\mu) + \sum_{i=1}^p \frac{1}{\tau_i!} |hd_i|^{\tau_i} \right].$$

The constant will depend on many aspects of the difference scheme and conceivably could be quite large. For example, as  $\tilde{\lambda}_{i,j}$  approaches a root of  $R(z)$ , the constant  $C_2$  will grow without bound.

The derivation of (2.29) may be formulated as the following theorem.

**THEOREM.** For  $n = 0, 1, \dots, N$  with  $Nh = 1$ , let  $x_n$  be a solution of the differential equations (2.1), let  $u_n$  be a solution of the difference equation (2.6), and let  $e_n = u_n - x_n$ . If (2.6) is stable (e.g. obeys the hypothesis of Lemma 2) and if (2.6) obeys the accuracy conditions characterized by (2.23) and (2.24), then modulo the initial errors  $||e_N||$  has the bound (2.29).

*Remark.* Classical linear multistep methods correspond to the case where  $P \equiv 0$ .

---

\*\*\*  $Z_{i,j}(hA)$  are the polynomials entering into the spectral representation theorem i.e.  $Z_{i,j}$  is the polynomial of minimal degree which vanishes on the spectrum of  $A$  except that at  $\lambda_j \in k_i, Z_{i,j} = 1$ .

At another extreme, we may use (2.6) and disregard the node  $z_0$  completely. This would give a more symmetric treatment.

*Example.* A simple example of the scheme (2.6) corresponds to  $r = 1$ ,  $\alpha_0 = 1$ ,  $\alpha_1 = -1$  and  $\delta_1 = 1$ . All other coefficients are zero. We select one node, i.e.,  $p = 1$ .  $P(z)$  is taken to be the constant  $C(hz_1)$ . The difference scheme is

$$(2.30) \quad u_n - u_{n-1} = \frac{e^{hz_1} - 1}{hz_1} h\dot{u}_{n-1}.$$

For this scheme  $\mu = 1$ ,  $\nu = 1$ , and  $\tau_1 = 1$ . Thus the scheme has zeroth order accuracy at the origin and at  $z_1$ . This low accuracy scheme may be viewed as the forward Euler with the mesh increment scaled by  $(e^{hz_1} - 1)/hz_1$ .

For this scheme  $S(z) = Iz - (I + ((e^{hz_1} - 1)/z_1)A)$ . By Lemma 1, the determinant  $|S(z)|$  obeys the root condition if  $z - 1 - ((e^{hz_1} - 1)/z_1)\lambda$  does for every eigenvalue  $\lambda$  of  $A$ . This latter requirement is seen to be satisfied for any choice of  $z_1$  in an interval which itself contains the interval  $(-\infty, \lambda)$ . (We are assuming that  $\lambda < 0$ .) Thus, if  $z_1$  is chosen as any lower estimate for the spectrum of  $A$ , (2.30) will be stable. Let us choose  $z_1 = \min\{\lambda\} - d$  for some  $d \geq 0$ . To simplify ideas let us consider the special case corresponding to  $m = 2$  and to, say,  $\lambda_2 = -1$  and  $\lambda_1$ , some very large negative number. The difference scheme then becomes

$$(2.31) \quad \begin{aligned} u_n - u_{n-1} &= \frac{e^{h(\lambda_1-d)} - 1}{\lambda_1 - d} A u_{n-1} \\ &\approx \frac{1}{d - \lambda_1} A u_{n-1} \end{aligned}$$

since  $\lambda_1 \ll -1$ . This corresponds to a forward Euler with a very small mesh increment,  $1/(d - \lambda_1)$ .

Now since

$$(2.32) \quad \begin{aligned} x_n &= e^{Ah} x_{n-1} \quad \text{and} \quad u_n = \left[ I + \frac{e^{h(\lambda_1-d)} - 1}{\lambda_1 - d} A \right] u_{n-1}, \\ e_n &= \left[ I + \frac{e^{h(\lambda_1-d)} - 1}{\lambda_1 - d} A - e^{Ah} \right] e_{n-1} \\ &\equiv T(hA) e_{n-1}. \end{aligned}$$

$T(h\lambda)$  is then the difference between the exponential  $e^{h\lambda}$  and the straight line  $1 - (e^{h(\lambda_1-d)} - 1)\lambda/(\lambda_1 - d)$ . At the eigenvalue  $\lambda_1$ , we have

$$(2.33) \quad T(h\lambda_1) = \frac{d}{\lambda_1} + e^{h\lambda_1} \left[ 1 + d \left( \frac{1}{\lambda_1} - h \right) + O \left( d^2 \left( \frac{1}{\lambda_1} - h \right)^2 \right) \right] + O \left( \frac{d^2}{\lambda_1^2} \right).$$

*Remark on the Nonlinear Case.* Although we have only discussed the linear case, the numerical method which we are considering may be applied to nonlinear differential equations in a variety of ways. Since the method is itself a nonlinear method, the analysis which we have given can only be carried over in a formal way to the case of nonlinear differential equations. (This has been done in Section 5 of [9] for a related numerical method.) One method for applying the problem to the nonlinear case is as follows.

For the nonlinear differential equation

$$(2.34) \quad \dot{y} = f(y)$$

and at the mesh point  $x_n$ , we replace  $y$  by  $z = y - y_n$  to get

$$(2.35) \quad \dot{z} = J(y_n)z$$

where  $J$  is the Jacobian

$$(2.36) \quad J(y_n) = \partial f(y)/\partial y.$$

The values of  $z_{n-1}, z_{n-2}, \dots$  which are needed are chosen as  $y_{n-1} - y_n, y_{n-2} - y_n, \dots$ , respectively, and  $A$  is taken as  $J(y_n)$ . This procedure is standard in the subject and we refer to [1] for a report of calculations on a related method using this procedure for nonlinear problems.

**3. Numerical Example.** We chose

$$(3.1) \quad u_n - u_{n-1} + h[\theta \dot{u}_n + (1 - \theta) \dot{u}_{n-1}] \\ - P(hA)[u_n - u_{n-1} + h(\varphi \dot{u}_n + (1 - \varphi) \dot{u}_{n-1})] = 0$$

as an example of the class of schemes discussed above with which to perform some calculations. For  $A$  we chose the matrix

$$(3.2) \quad A = - \begin{bmatrix} 10^p & 0 & 0 \\ 0 & 10^q & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

with  $p$  and  $q$  as parameters. The function  $C(z)$  is

$$(3.3) \quad C(z) = \frac{(1 - \theta z)e^z - 1 - (1 - \theta)z}{(1 - \varphi z)e^z - 1 - (1 - \varphi)z}.$$

$P(z)$  was taken as the linear polynomial which interpolates  $C(z)$  at the two points  $-(1 + \rho/100)10^p$  and  $-(1 + \rho/100)10^q$  with  $\rho$  a parameter.

As a comparison scheme, we chose the second-order scheme

$$(3.4) \quad v_n - v_{n-1} - \frac{h}{2} [\dot{v}_n + \dot{v}_{n-1}] = 0.$$

Let  $w_n$  be the exact solution of the differential equation,  $\dot{w} = Aw$ , at the mesh point  $nh$ . Let  $E_n(u)$  be the vector whose  $i$ th component is

$$(3.5) \quad E_n^i(u) = u_n^i/w_n^i - 1.$$

Let the length of  $E_n(u)$  be  $\|E_n(u)\|$  and let

$$(3.6) \quad E_u = \sum_{n=1}^N \|E_n(u)\|/N.$$

$E_v$  is defined analogously with  $v_n^i$  replacing  $u_n^i$  in (3.5) and (3.6).  $E_u$  and  $E_v$  then are relative errors averaged over  $N$  time steps for the example scheme (3.1) and the comparison scheme (3.4), respectively. Finally, let

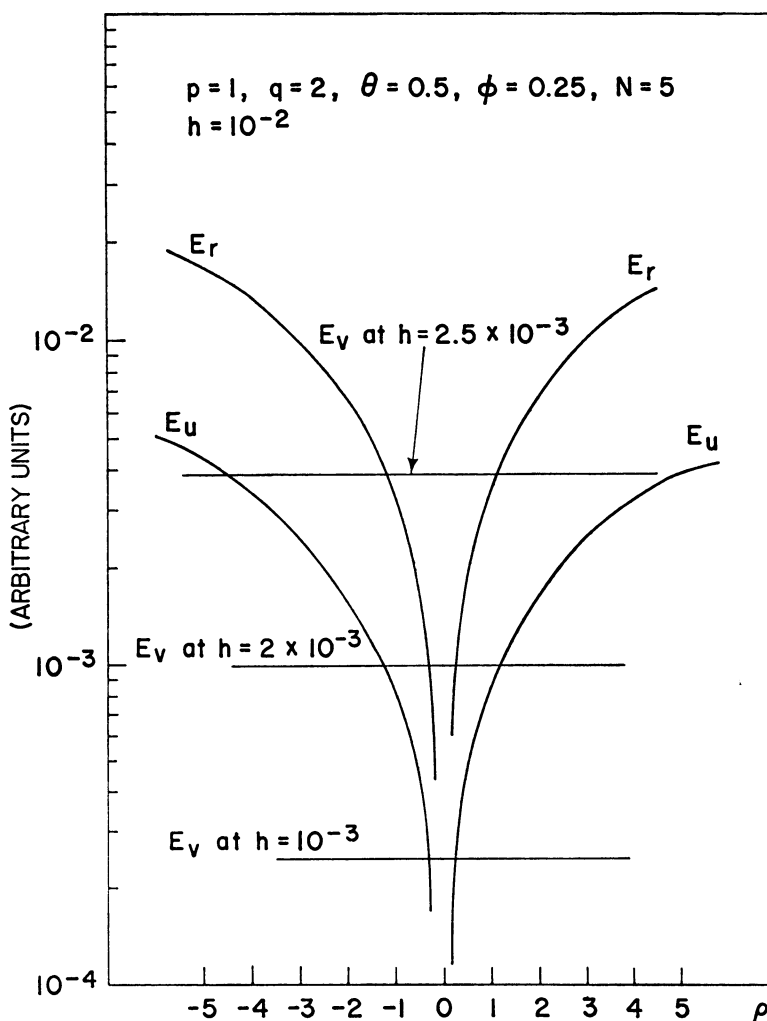
$$(3.7) \quad E_r = E_u/E_v.$$



In the following three graphs, we have plotted  $E_u$ ,  $E_v$ , and  $E_r$  for various parameter values. The initial condition was chosen to be (1, 1, 1).

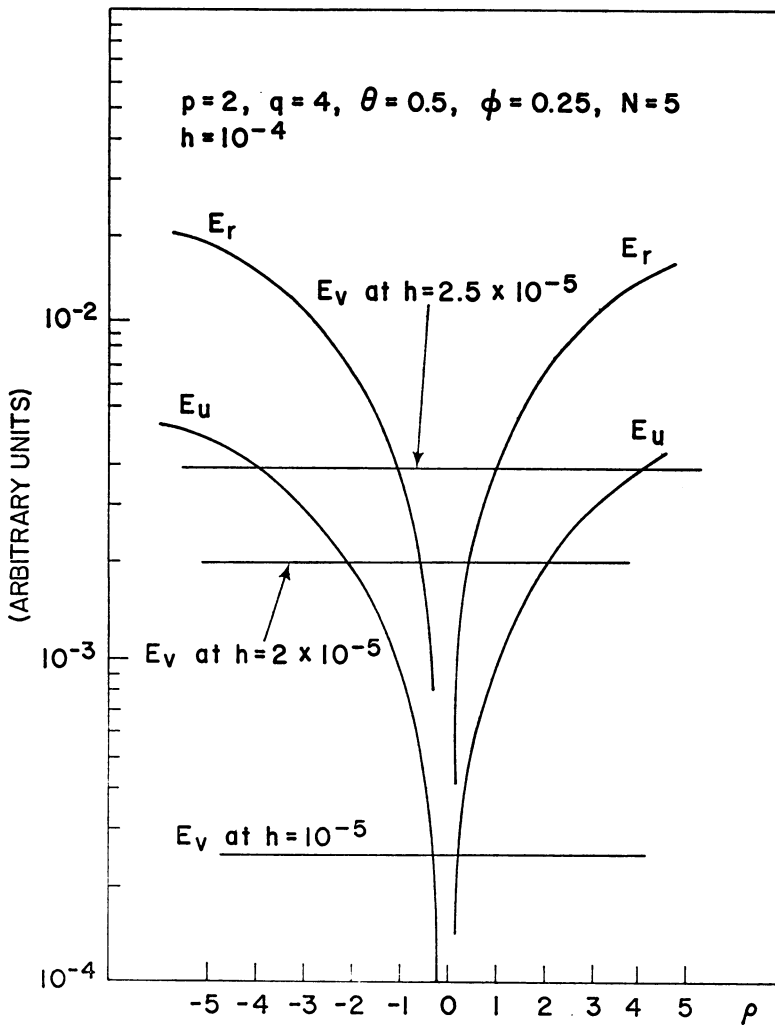
The crossings of the horizontal plot representing  $E_v$  with the curves representing  $E_u$  show with what accuracy the interpolation points must approximate the eigenvalues of  $A_i$  in order to match a calculation with the comparison formula with which a mesh increment, one order in magnitude finer, has been used.

All qualitative features of the graphs are as the theory above predicts.



#### Appendix.

*The Liniger-Willoughby Scheme.* This scheme is not in the class considered in this paper. However, it employs the idea of approximating the solution operator of the differential equation both at the origin as is usual and at one additional point using a free parameter. The scheme is simple and very effective for stiff systems with



one cluster of small and one cluster of large roots. It even works well in some other cases. The scheme is

$$(A.1) \quad u_{n+1} - u_n = \mu h \dot{u}_n + (1 - \mu) h \dot{u}_{n+1}$$

with  $\mu$  to be specified. In the linear case

$$(A.2) \quad \dot{x} = -Ax,$$

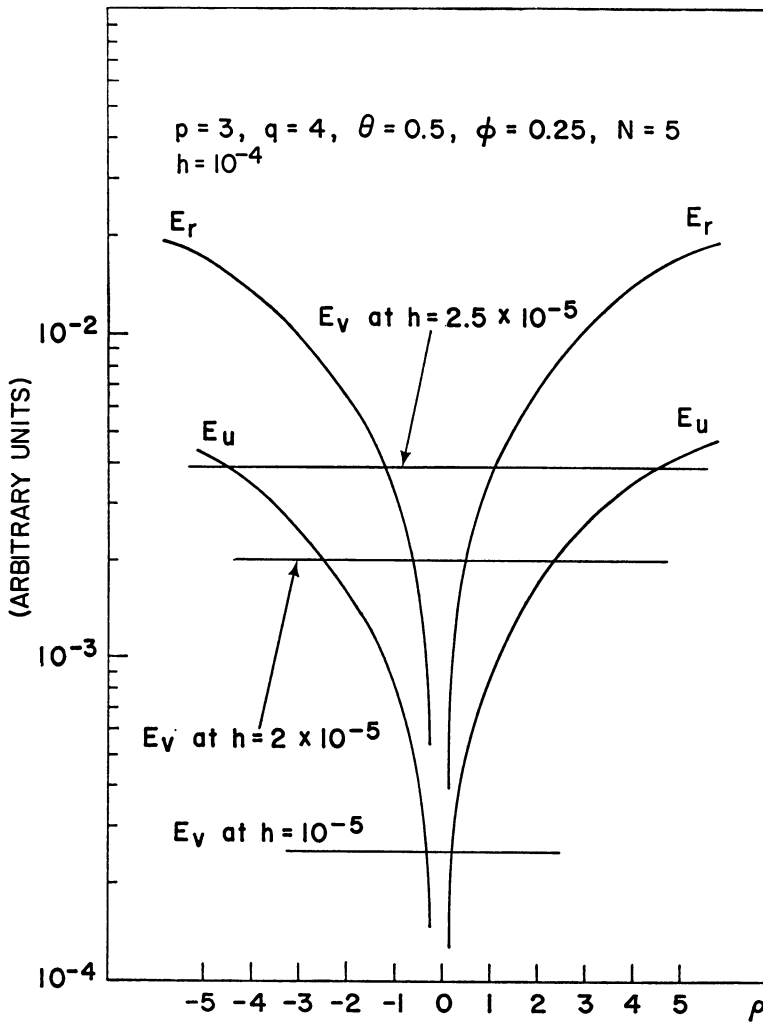
so that  $\dot{u}_n$  is replaced by  $-Au_n$  in (3.1).<sup>†</sup> The scheme becomes

$$(A.3) \quad u_{n+1} = K(hA)u_n$$

with

$$(A.4) \quad K(z) = (1 - \mu z)/(1 + (1 - \mu)z).$$

<sup>†</sup> Our analysis of this scheme makes use of the positive definiteness of  $A$ .



The error  $e_n$  obeys the equation

$$(A.5) \quad e_{n+1} = K(hA)e_n + (K(hA) - e^{-Ah})x_n.$$

Then

$$(A.6) \quad e_n = \sum_{i=0}^{n-1} K^i(hA)[K(hA) - e^{-Ah}]x_{n-i-1} + K^n(hA)e_0.$$

The function  $K(z)$  is less than one in magnitude for all real  $z > 0$  and all  $\mu$ ,  $0 < \mu < \frac{1}{2}$ . The theorem on the spectral resolution (quoted in the proof of Theorem 1 above) shows that

$$(A.7) \quad \|K^n(hA)\| < \text{const},$$

since the eigenvalues of  $A$  are strictly positive. Thus, from (A.6) and since the eigenvalues of  $A$  are strictly positive,

$$(A.8) \quad \|e_n\| \leq \text{const} \cdot [\|K(hA) - e^{-Ah}\| + \|e_0\|].$$

To estimate  $\|K(hA) - e^{-Ah}\|$ , we use the spectral resolution theorem again. We obtain the estimate

$$(A.9) \quad \|K(hA) - e^{-Ah}\| \leq \text{const} \cdot \min_{\Pi} \left[ \max_{i \in I_1} |h^2 \lambda_i^2| + \max_{i \in I_2} |r - \lambda| \right],$$

where  $I_1$  and  $I_2$  form a partition  $\Pi$  over the integers  $1, \dots, m$  and if  $i \in I_2$  and  $j \in I_2$ , then  $\lambda_i \neq \lambda_j$ . The minimum is taken over all such partitions.  $r$  is some positive scalar, fixed for convenience. Once  $r$  is fixed,  $\mu$  is chosen so that

$$(A.10) \quad K(hr) - e^{-hr} = 0.$$

The crux of the scheme is that this may be accomplished for any  $r > 0$  and by a  $\mu$  in the interval  $(0, \frac{1}{2})$ .

The estimate (A.9) then follows from the following observations.

$$(A.11) \quad |K(z) - e^{-z}| < \text{const} \cdot \min\{z^2, 1\}, \quad z \geq 0,$$

also

$$(A.12) \quad K(z) - e^{-z} = (r - z)[K'(z) + e^{-z}]$$

so that

$$(A.13) \quad |K(z) - e^{-z}| < \text{const} \cdot |r - z|, \quad z > 0.$$

The boundedness of the constant here depends on  $\mu \in (0, \frac{1}{2})$ . Since the fit of  $K(z)$  to  $e^{-z}$  at  $z = r$  is only of first order, we cannot permit repeated roots in  $I_2$  if (3.13) is to hold.

#### IBM

Thomas J. Watson Research Center  
Yorktown Heights, New York 10598

1. W. LINIGER & R. WILLOUGHBY, "Efficient numerical integration of stiff systems of ordinary differential equations," *SIAM J. Numer. Anal.*, v. 6, 1969, pp. 47-66.
2. M. E. FOWLER & R. M. WARTEN, "A numerical integration technique for ordinary differential equations with widely separated eigenvalues," *IBM J. Res. Develop.*, v. 11, 1967, pp. 537-543. MR 35 #7586.
3. C. F. CURTISS & J. O. HIRSCHFELDER, "Integration of stiff equations," *Proc. Nat. Acad. Sci. U.S.A.*, v. 38, 1952, pp. 235-243. MR 13, 873.
4. C. E. TREANOR, "A method for the numerical integration of coupled first-order differential equations with greatly different time constants," *Math. Comp.*, v. 20, 1966, pp. 39-45. MR 33 #889.
5. J. CERTAINE, "The solution of ordinary differential equations with large time constants," in *Mathematical Methods for Digital Computers*, Wiley, New York, 1960, pp. 128-132. MR 22 #8691.
6. W. L. MIRANKER & W. LINIGER, "Parallel methods for the numerical integration of ordinary differential equations," *Math. Comp.*, v. 21, 1967, pp. 303-320. MR 36 #6155.
7. G. G. DAHLQUIST, "A special stability problem for linear multistep methods," *Nordisk Tidskr. Informations-Behandling*, v. 3, 1963, pp. 27-43, MR 30 #715.
8. C. W. GEAR, *Numerical Integration of Stiff Ordinary Differential Equations*, Dept. of Computer Science, Report #221, University of Illinois, Urbana, Illinois.
9. W. L. MIRANKER, *Difference Schemes for the Integration of Stiff Systems of Ordinary Differential Equations*, IBM Research Center, Report #RC-1977, 1968.
10. F. R. GANTMACHER, *The Theory of Matrices*, GITTL, Moscow, 1953; English transl., Chelsea, New York, 1959. MR 16, 438.