

17

ON HIGHER ORDER STABLE IMPLICIT METHODS FOR SOLVING
PARABOLIC PARTIAL DIFFERENTIAL EQUATIONS

BY RICHARD S. VARGA

Reprinted from JOURNAL OF MATHEMATICS AND PHYSICS
Vol. XL, No. 3, October, 1961
Printed in U.S.A.

ON HIGHER ORDER STABLE IMPLICIT METHODS FOR SOLVING PARABOLIC PARTIAL DIFFERENTIAL EQUATIONS

BY RICHARD S. VARGA

1. Introduction. The numerical solutions of self-adjoint parabolic partial differential equations, such as

$$(1.1) \quad \frac{\partial u(x, t)}{\partial t} = \left\{ \frac{\partial}{\partial x} p(x) \frac{\partial u}{\partial x} \right\} - \sigma(x)u(x, t) + s(x),$$

$$a < x < b, \quad t > 0,$$

in one space variable where

$$(1.2) \quad u(x, 0+) = g(x), \quad a \leq x \leq b,$$

and

$$(1.3) \quad u(a, t) = \alpha, \quad u(b, t) = \beta, \quad t > 0,$$

have been considered by many authors (e.g. [5], [12], [17], [19]), and concepts such as the *stability* and *consistency* of discrete approximations, and *convergence* of the discrete approximations to the solution of (1.1) have been previously introduced. However, while most previous papers on this topic simultaneously discretized *both* space and time variables and analysed by means of Fourier series the properties of the resulting system of linear equations, our approach, like that of [8], is to discretize first only the space variables, leaving the time variable continuous, and to analyse the resulting system of ordinary differential equations by matrix methods, avoiding Fourier series arguments.

As a result of this approach, we shall introduce new numerical methods which are unconditionally stable, and show that well known numerical methods for solving (1.1), viz. the explicit method, the backward time implicit method [14], and the Crank-Nicolson implicit method [19], can be generated (Theorem 2) from particular Padé rational approximations $r_{p,q}(z)$ of e^{-z} . These results are applicable to a large class of parabolic partial differential equations which are approximated on non-uniform spatial meshes.

Considering discrete spatial meshes, we also give error estimates for the difference between the continuous time problem, and the discrete time problem used in actual numerical calculations. This consideration leads to *a priori* estimates of the discrete time mesh in order to achieve a particular accuracy. Finally, because of the monotonicity of the error of Padé rational approximations for e^{-z} , we consider also the best rational approximations for e^{-z} in the Chebyshev sense. The interesting feature of the associated matrix approximations is their accuracy for *any* size time step (§5), and as a result, offer the possibility of obtaining approximations (of limited accuracy) of (1.1) for any time t from a *single* time step.

2. The Continuous Time-Discrete Space Equations. To be explicit, we consider in particular the parabolic partial differential equation of (1.1-3), where we assume that $p(x)$ and $\sigma(x)$ are given continuous functions¹ on $a \leq x \leq b$ with

$$(2.1) \quad p(x) > 0, \quad \sigma(x) \geq 0, \quad a \leq x \leq b.$$

The given functions $s(x)$ and $g(x)$ are also assumed to be continuous on $a \leq x \leq b$. If $\{x_i\}_{i=0}^{n+1}$ is any set of spatial mesh points with $a = x_0 < x_1 < x_2 < \dots < x_n < x_{n+1} = b$, define $h_i = x_i - x_{i-1}$ for $1 \leq i \leq n + 1$. Upon integrating (1.1) from $x_i = h_i/2$ to $x_i + h_{i+1}/2$, and using central difference approximations as in [20], we obtain, denoting $u(x_i, t)$ by $u_i(t)$,

$$(2.2) \quad \left(\frac{h_i + h_{i+1}}{2}\right) \frac{d}{dt} u_i(t) = -u_i(t) \left\{ \frac{p(x_i + \frac{1}{2}h_{i+1})}{h_{i+1}} + \frac{p(x_i - \frac{1}{2}h_i)}{h_i} + \sigma_i \left(\frac{h_i + h_{i+1}}{2}\right) \right\} + \frac{p(x_i + \frac{1}{2}h_{i+1})}{h_{i+1}} u_{i+1}(t) + \frac{p(x_i - \frac{1}{2}h_i)}{h_i} u_{i-1}(t) + \left(\frac{h_i + h_{i+1}}{2}\right) s_i,$$

for $t > 0, 1 \leq i \leq n$, which, along with $u_i(0) = g_i$, defines our continuous time-discrete space approximation of (1.1). In matrix notation, this becomes

$$(2.3) \quad D \frac{d\mathbf{u}(t)}{dt} = -A\mathbf{u}(t) + \mathbf{s}, \quad t \geq 0,$$

where

$$(2.4) \quad \mathbf{u}(0) = \mathbf{g}.$$

The matrix D is a positive real diagonal $n \times n$ matrix, and, from (2.1) and (2.2), A is a real symmetric and positive definite $n \times n$ matrix. More precisely, A is a tridiagonal $n \times n$ *Stieltjes matrix* [3]. It should be pointed out that the boundary conditions of (1.3) are implicit in the vector \mathbf{s} of (2.3).

In what is to follow, we need only assume that D is a positive real diagonal $n \times n$ matrix, and that A is a positive definite Hermitian $n \times n$ matrix. Because of this, our results apply to parabolic partial differential equations in more space variables, such as

$$(2.5) \quad \frac{\partial u(\mathbf{x}, t)}{\partial t} = \text{div} \{p(\mathbf{x}) \text{grad } u(\mathbf{x}, t)\} - \sigma(\mathbf{x})u(\mathbf{x}, t) + s(\mathbf{x}),$$

for which there are spatial discretization [20] of the form (2.3) satisfying these assumptions.

We now normalize the system of ordinary differential equations of (2.3). If

¹ The derivation [20], based on integration, of the spatial difference equations of this section shows that the case where $p(x)$ and $\sigma(x)$ are piecewise continuous in $0 \leq x \leq 1$ can be treated similarly. See also [8] and [15]

$D^{\frac{1}{2}}$ is the positive real diagonal $n \times n$ matrix whose square is D , let

$$(2.6) \quad \mathbf{v}(t) \equiv D^{\frac{1}{2}}\mathbf{u}(t), \quad D^{-\frac{1}{2}}AD^{-\frac{1}{2}} \equiv B.$$

Then, (2.4) can be written as

$$(2.7) \quad \frac{d\mathbf{v}(t)}{dt} = -B\mathbf{v}(t) + \Sigma, \quad t \geq 0$$

where $\Sigma = D^{-\frac{1}{2}}\mathbf{s}$, and

$$(2.8) \quad \mathbf{v}(0) = D^{\frac{1}{2}}\mathbf{g} \equiv \mathbf{g}^*.$$

From the definition of the $n \times n$ matrix B in (2.6), we see that B is also a positive definite Hermitian matrix. Moreover, there is a 1-1 correspondence between the vectors $\mathbf{u}(t)$ and $\mathbf{v}(t)$.

That the ordinary matrix differential equation of (2.7) is *stable*² [2] from the point of view of differential equations follows immediately from the positive definite Hermitian nature of B , and the solution of (2.7), subject to the boundary condition (2.8) is given explicitly by

$$(2.9) \quad \mathbf{v}(t) = e^{-tB^*}\mathbf{g}^* + (I - e^{-tB})B^{-1}\Sigma, \quad t \geq 0,$$

or equivalently,

$$(2.10) \quad \mathbf{v}(t_0 + \Delta t) = e^{-\Delta t B} \mathbf{v}(t_0) + (I - e^{-\Delta t B})B^{-1}\Sigma$$

where $t_0 \geq 0$, $\Delta t \geq 0$, the matrix e^{-tB} being defined by the convergent matrix series $I - tB + \frac{1}{2}t^2B^2 - \dots$.

It is now interesting to show the relationship between well known numerical methods for approximating (1.1) and the continuous time-discrete space approximation of (2.7). In terms of (2.7), these well known approximations are

$$(2.11) \quad \textit{Explicit:} \quad \left(\frac{\mathbf{v}(t + \Delta t) - \mathbf{v}(t)}{\Delta t} \right) = -B\mathbf{v}(t) + \Sigma;$$

$$(2.11') \quad \textit{Backward Implicit:} \quad \left(\frac{\mathbf{v}(t + \Delta t) - \mathbf{v}(t)}{\Delta t} \right) = -B\mathbf{v}(t + \Delta t) + \Sigma;$$

$$(2.11^*) \quad \textit{Crank-Nicolson:} \quad \left(\frac{\mathbf{v}(t + \Delta t) - \mathbf{v}(t)}{\Delta t} \right) = -\frac{B}{2}(\mathbf{v}(t + \Delta t) + \mathbf{v}(t)) + \Sigma.$$

Solving for $\mathbf{v}(t + \Delta t)$ in terms of $\mathbf{v}(t)$ and $\Sigma(t)$, these expressions reduce to

$$(2.12) \quad \textit{Explicit:} \quad \mathbf{v}(t + \Delta t) = (I - \Delta t B)\mathbf{v}(t) + \Delta t \Sigma;$$

$$(2.12') \quad \textit{Backward Implicit:} \quad \mathbf{v}(t + \Delta t) = (I + \Delta t B)^{-1}\mathbf{v}(t) + \Delta t(I + \Delta t B)^{-1}\Sigma;$$

$$(2.12^*) \quad \textit{Crank-Nicolson:} \quad \mathbf{v}(t + \Delta t) = (I + \frac{1}{2}\Delta t B)^{-1}(I - \frac{1}{2}\Delta t B)\mathbf{v}(t) \\ + \Delta t(I + \frac{1}{2}\Delta t B)^{-1}\Sigma.$$

Comparing with (2.10), it is clear that these three methods each give rise to a particular rational approximation for the exponential matrix $e^{-\Delta t B}$.

² See also [8].

3. The Padé Table for e^{-z} . The Padé table [18, 22] of a function $f(z)$ analytic in a region containing the point $z = 0$ is a double entry table, such that the rational approximation for $f(z)$,

$$(3.1) \quad f(z) = \frac{n_{p,q}(z)}{d_{p,q}(z)} + O(|z|^r), \quad |z| \rightarrow 0,$$

where $n_{p,q}(z)$ and $d_{p,q}(z)$ are polynomials of degree q and p , respectively, gives the best approximation (highest power r) for $f(z)$ in the neighborhood of the origin. Except for a multiplicative factor, for fixed non-negative integers p and q the polynomials $n_{p,q}(z)$ and $d_{p,q}(z)$ are uniquely determined, and $r = p + q + 1$. For the function $f(z) = e^{-z}$, the entries of the Padé table are given explicitly [11, 18] by

$$(3.2) \quad n_{p,q}(z) = \sum_{k=0}^q \frac{(p+q-k)!q!}{(p+q)!k!(q-k)!} (-z)^k,$$

and

$$(3.2') \quad d_{p,q}(z) = \sum_{k=0}^p \frac{(p+q-k)!p!}{(p+q)!k!(p-k)!} z^k,$$

and if

$$(3.3) \quad e^{-z} \equiv \frac{n_{p,q}(z)}{d_{p,q}(z)} + R_{p,q}(z),$$

then [18]

$$(3.4) \quad R_{p,q}(z) = \frac{(-1)^{q+1}z^{p+q+1}}{(p+q)!d_{p,q}(z)} \int_0^1 e^{-z(1-u)}u^q(1-u)^p du.$$

For $z \geq 0$, we can apply the mean value theorem for integrals to the expression above, and it follows [11] that

$$(3.5) \quad R_{p,q}(z) = \frac{(-1)^{q+1}p!q!z^{p+q+1}e^{-\theta}}{(p+q)!(p+q+1)!d_{p,q}(z)}, \quad z \geq \theta \geq 0.$$

The function e^{-z} is uniformly bounded by unity for all $z \geq 0$. It is natural to ask the same question of the Padé approximations for e^{-z} .

Lemma 1. For arbitrary p and $q \geq 0$, $d_{p,q}(z) \geq 1$ for all $z \geq 0$.

Proof. The coefficients of the polynomial $d_{p,q}(z)$ are all positive real numbers, with constant coefficient unity for all $p, q \geq 0$. Thus, if $z \geq 0$, then $d_{p,q}(z) \geq 1$.

Lemma 2. $\left| \frac{n_{p,q}(z)}{d_{p,q}(z)} \right| \leq 1$ for all $z \geq 0$ if and only if $p \geq q$.

Proof. Clearly, $\left| \frac{n_{p,q}(z)}{d_{p,q}(z)} \right| = O(z^{q-p})$ for $z \rightarrow +\infty$ if $q \geq p$. Thus, the ratio being bounded for all $z \geq 0$ implies that $p \geq q$. Conversely, if $p \geq q$, then $|d_{p,q}(z)| \geq |n_{p,q}(z)|$ if

$$\frac{(p+q-k)!p!}{(p+q)!k!(p-k)!} \geq \frac{(p+q-k)!q!}{(p+q)!k!(q-k)!}, \quad 0 \leq k \leq q,$$

which is obviously true.

We remark that the error $R_{p,q}(z)$ of the Padé approximation (3.3) for e^{-z} is, from (3.5) and Lemma 1, of one sign for all $z \geq 0$ for all non-negative integers p and q .

We now consider formal Padé approximations³ for the exponential matrix e^{-tB} of the form

$$(3.6) \quad e^{-tB} \doteq [d_{p,q}(tB)]^{-1} \cdot [n_{p,q}(tB)] \equiv M_{p,q}(tB)$$

stemming naturally from the approximations of (3.3). We make the usual definition [5].

Definition 1. Let $M(t)$ be an $n \times n$ matrix whose entries are functions of the parameter t . Then $M(t_0)$ is *stable* if and only if all the eigenvalues $\lambda_i(t_0)$ of $M(t_0)$ satisfy $|\lambda_i(t_0)| \leq 1$.

For the $n \times n$ positive definite Hermitian matrix B , let $\rho(B)$ denote the *spectral radius* [10] of B , i.e. $\rho(B) = \max_{1 \leq i \leq n} |\lambda_i|$, where the λ_i 's are the (real and positive) eigenvalues of B . Moreover, let

$$(3.7) \quad \tau_{p,q} \equiv \sup \left\{ l \geq 0 \left| \left| \frac{n_{p,q}(z)}{d_{p,q}(z)} \right| \leq 1 \text{ for } 0 \leq z \leq l \right. \right\}.$$

Theorem 1. Let B be an $n \times n$ positive definite Hermitian matrix, and let

$$(3.8) \quad M_{p,q}(\Delta t B) = [d_{p,q}(\Delta t B)]^{-1} [n_{p,q}(\Delta t B)].$$

Then, $M_{p,q}(\Delta t B)$ is stable for all Δt with $0 \leq \Delta t \leq \Lambda$ if and only if

$$(3.9) \quad \Lambda \leq \tau_{p,q} / \rho(B).$$

Proof. If the eigenvalues of B are λ_i , then the eigenvalues of $d_{p,q}(\Delta t B)$ are obviously $d_{p,q}(\Delta t \lambda_i)$, and since Δt is non-negative and each λ_i is real and positive, Lemma 1 shows that the matrix $d_{p,q}(\Delta t B)$ is non-singular for any p and q , and any non-negative Δt . Thus the matrix $M_{p,q}(\Delta t B)$ is defined for all $\Delta t \geq 0$ and all non-negative integers p and q .

The eigenvalues of $M_{p,q}(\Delta t B)$ are $n_{p,q}(\Delta t \lambda_i) / d_{p,q}(\Delta t \lambda_i)$. It is clear that if $0 \leq \Delta t \leq \tau_{p,q} / \rho(B)$, the definition of $\tau_{p,q}$ shows that $M_{p,q}(\Delta t B)$ is stable for Δt in this range. Conversely, if $M_{p,q}(\Delta t B)$ is stable for all Δt with $0 \leq \Delta t \leq \Lambda$, then as $\rho(B)$ is an eigenvalue of B ,

$$\left| \frac{n_{p,q}(\Delta t \rho(B))}{d_{p,q}(\Delta t \rho(B))} \right| \leq 1$$

for all such Δt , and thus, from (3.7), $\Delta t \rho(B) \leq \Lambda \rho(B) \leq \tau_{p,q}$, which completes the proof.

From Lemma 2, we see that $\tau_{p,q} = +\infty$ if and only if $p \geq q$, which gives us

Corollary 1. Let B be an $n \times n$ positive definite Hermitian matrix. Then, $M_{p,q}(tB)$ is stable for all $t \geq 0$ if and only if $p \geq q$.

We explicitly list now particular entries from the Padé table for e^{-z} .

³ See also [13] for rational approximations of other operators useful in numerical analysis.

4. Error Analysis for the Padé Approximations. We first write (2.10) in the equivalent form

$$(4.1) \quad \mathbf{v}(t_0 + \Delta t) = B^{-1}\boldsymbol{\Sigma} + e^{-\Delta t B}\{\mathbf{v}(t_0) - B^{-1}\boldsymbol{\Sigma}\}.$$

With our definition (3.6) of the matrix Padé approximations $M_{p,q}(\Delta t B)$ of $e^{-\Delta t B}$, we correspondingly define

$$(4.2) \quad \mathbf{w}_{p,q}((j+1)\Delta t) = B^{-1}\boldsymbol{\Sigma} + M_{p,q}(\Delta t B)\{\mathbf{w}_{p,q}(j\Delta t) - B^{-1}\boldsymbol{\Sigma}\},$$

for $j \geq 0$, $\Delta t > 0$, as discrete Padé approximations of the solution vector $\mathbf{v}((j+1)\Delta t)$, where

$$(4.2') \quad \mathbf{w}_{p,q}(0) \equiv \mathbf{g}^* = \mathbf{v}(0).$$

In this form (4.2), it follows that

$$(4.3) \quad \mathbf{w}_{p,q}(j\Delta t) = B^{-1}\boldsymbol{\Sigma} + M_{p,q}^j(\Delta t B)\{\mathbf{g}^* - B^{-1}\boldsymbol{\Sigma}\}, \quad j \geq 1.$$

In order to carry out an error analysis between the continuous time-discrete space solution (4.1) and the Padé discrete time-discrete space solution (4.3), we define now the quantity

$$(4.4) \quad r_{p,q}^{(m)}(z) = \sup_{0 \leq x \leq z} |e^{-x} - M_{p,q}^m(x/m)|$$

for all non-negative integers p and q , and positive integers m . It is easy to see, using (3.5), that for fixed p, q , and z , $r_{p,q}^{(m)}(z) \rightarrow 0$ as $m \rightarrow +\infty$.

If $\|\mathbf{x}\| \equiv (\sum_{i=1}^n |x_i|^2)^{\frac{1}{2}}$ denotes the euclidean norm of the vector \mathbf{x} whose components are x_i , $1 \leq i \leq n$, then the spectral norm [10] $\|C\|$ of an arbitrary $n \times n$ complex matrix C is defined by

$$(4.5) \quad \|C\| \equiv \sup_{\mathbf{x} \neq 0} \frac{\|C\mathbf{x}\|}{\|\mathbf{x}\|}.$$

Thus,

$$(4.5') \quad \|C\mathbf{x}\| \leq \|C\| \cdot \|\mathbf{x}\|$$

for any vector \mathbf{x} . It is moreover known [10] that if C is Hermitian (or normal), then $\|C\| = \rho(B)$.

Lemma 3. Let B be an $n \times n$ positive definite Hermitian matrix. For any vector \mathbf{x} ,

$$(4.6) \quad \|(e^{-tB} - M_{p,q}^m(tB/m))\mathbf{x}\| \leq r_{p,q}^{(m)}(t\rho(B)) \cdot \|\mathbf{x}\|,$$

for all non-negative integers p, q , all positive integers m , and all $t \geq 0$.

Proof. From the discussion above, it is only necessary to show that

$$\|e^{-tB} - M_{p,q}^m(tB/m)\| \leq r_{p,q}^{(m)}(t\rho(B)).$$

But as the matrix $e^{-tB} - M_{p,q}^m(tB/m)$ is Hermitian,

$$\left\| e^{-tB} - M_{p,q}^m(tB/m) \right\| = \max_{1 \leq i \leq n} \left| e^{-t\lambda_i} - M_{p,q}^m\left(\frac{t\lambda_i}{m}\right) \right| \leq r_{p,q}^{(m)}(t\rho(B)),$$

which completes the proof.

From (4.1) and (4.3), we have

$$\mathbf{v}(m\Delta t) - \mathbf{w}_{p,q}(m\Delta t) = \{e^{-m\Delta t B} - M_{p,q}^m(\Delta t B)\}(\mathbf{g}^* - B^{-1}\boldsymbol{\Sigma}),$$

and as an immediate consequence of (4.5') and Lemma 3, we obtain

Theorem 3. Let $\mathbf{v}(t)$ be the solution of the matrix differential equation (2.7), subject to the initial condition of (2.8), and let $\mathbf{w}_{p,q}(m\Delta t)$ be defined from (4.2) and (4.2'). Then, for all non-negative integers p and q , all positive integers m , and all $\Delta t > 0$,

$$(4.7) \quad \|\mathbf{v}(m\Delta t) - \mathbf{w}_{p,q}(m\Delta t)\| \leq r_{p,q}^{(m)}(m\Delta t\rho(B)) \cdot \|\mathbf{g}^* - B^{-1}\boldsymbol{\Sigma}\|.$$

To illustrate the usefulness of higher order implicit methods, consider the particular problem (3.11) with a uniform spatial mesh $\Delta x = 0.1$, and suppose that we seek approximations (cf. [6, p. 138]) to $\mathbf{v}(T)$, where $T = 10^4$. With $\rho(B)$ bounded above in this case by $4/(\Delta x)^2 = 400$, we list estimates for the least positive integer $m_{p,q} = T/\Delta t$ for which the coefficient of $\|\mathbf{g}^* - \mathbf{h}\|$ in (4.6) is less than or equal to 0.0073, for various choices of non-negative integers p and q :

$$(4.8) \quad m_{0,1} \doteq 2.0 \times 10^6, \quad m_{1,1} \doteq 2.2 \times 10^3, \quad m_{2,2} \doteq 1.2 \times 10^3.$$

In general, it is now clear how *a priori* estimates of $m_{p,q}$, the total number of time steps, can be determined to insure a particular accuracy k between $\mathbf{v}(T)$ and $\mathbf{w}_{p,q}(T)$, where T is given. If $\tilde{\rho}$ is an estimate of the spectral radius $\rho(B)$ of the matrix B , we algebraically determine the least positive solution m of

$$(4.9) \quad r_{p,q}^m(T\tilde{\rho}) \leq \frac{k}{\|\mathbf{g}^* - B^{-1}\boldsymbol{\Sigma}\|},$$

and set $\Delta t = T/m$.

While stability considerations imply that $p \geq q$ in the Padé matrix approximations $M_{p,q}(tB)$ of (3.8) for e^{-tB} , practical considerations imply that one would examine only the diagonal entries $p = q$ of the Padé table. This follows from the fact that forming the matrix polynomial $d_{p,q}(tB)$ of degree p requires explicitly the matrices B^ν , $q \leq \nu \leq p$ which could thus also be used in the formation of $n_{p,q}(tB)$ if q were increased to p . Moreover, the matrix inversion of $d_{p,q}(tB)$ usually outweighs in computational effort the formation of $n_{p,q}(tB)$, even when $p = q$.

The diagonal Padé matrix approximations $M_{p,p}(tB)$ can also be derived as approximants of the known continued fraction expansion [22, p. 348] for e^{-z} :

$$(4.10) \quad e^{-z} = \frac{1}{1 + z \cfrac{1}{1 - z \cfrac{2}{2 + z \cfrac{3 - z}{2 + z \cfrac{5 - z}{2 - \dots}}}}},$$

and these diagonal Padé matrix approximations can also be derived from repeated differentiations of (2.7), coupled with matrix Taylor series expansions based on central differences.

We shall show in the next section why the particular constant $k = 0.0073$ was chosen in the above discussion.

5. Chebyshev Rational Approximations. We had remarked in §3 that the error $R_{p,q}(z)$ of the Padé approximation for e^{-z} is of one sign for all $z \geq 0$. While the asymptotic behavior of $R_{p,q}(z)$ in the neighborhood of $z = 0$ allows us, from Theorem 3, to obtain arbitrarily high accuracy between the continuous time and discrete time models, these Padé rational approximations nevertheless are not the best rational approximations for e^{-z} in the Chebyshev sense [1]. For a fixed non-negative integer p , consider

$$(5.1) \quad H_p \equiv \max_{z \geq 0} \left| e^{-z} - \frac{n_p(z)}{d_p(z)} \right|,$$

where $n_p(z)$ and $d_p(z)$ are any real polynomials of degree p in z , and where $d_p(z) \neq 0$ for all $z \geq 0$. Now let

$$(5.2) \quad \hat{H}_p = \min H_p$$

where the min is taken over all such real polynomials $n_p(z)$ and $d_p(z)$. The celebrated theorem of Chebyshev (see [1]) states that the rational function $Q_p(z) = n_p(z)/d_p(z)$ which minimizes H_p is uniquely determined (assuming that $Q_p(z)$ is irreducible), and is characterized by the property that $e^{-z} - Q_p(z)$, with alternate change of signs, takes on the value \hat{H}_p not less than $2p + 2$ times. Because of this required oscillation of the error, the diagonal Padé approximations for e^{-z} of §3 are evidently different from the Chebyshev rational approximations $Q_p(z)$ for e^{-z} . We now exhibit some of the rational functions $Q_p(z)$ which best approximate e^{-z} in the Chebyshev sense.⁴

$$(5.3) \quad \begin{aligned} Q_1(z) &= \frac{1 - 0.108\,196\,z}{0.937\,355 + 1.618\,932\,z}; & \hat{H}_1 &= 0.066\,83 \\ Q_2(z) &= \frac{1 - 0.189\,729\,z + 0.004\,242\,15\,z^2}{1.007\,413 + 0.674\,264\,z + 0.576\,492\,z^2}; & \hat{H}_2 &= .0073\,59. \end{aligned}$$

The particular choice of $k = 0.0073$ in §4 is now clear.

For an error analysis based on these Chebyshev rational approximations, let

$$(5.4) \quad \mathbf{w}_p(t) \equiv B^{-1}\boldsymbol{\Sigma} + Q_p(tB)\{\mathbf{g}^* - B^{-1}\boldsymbol{\Sigma}\}, \quad t \geq 0,$$

be the Chebyshev rational approximation of (4.1) corresponding to a *single* time step of length t , where

$$(5.4') \quad \mathbf{w}_p(0) \equiv \mathbf{g}^* = \mathbf{v}(0),$$

and

⁴ The author wishes to express his thanks to Mr. Fred Chapman of the Case Institute of Technology Computing Center for obtaining these results on the IBM 650.

$$(5.5) \quad Q_p(tB) \equiv [d_p(tB)]^{-1} \cdot [n_p(tB)].$$

Theorem 4. Let $\mathbf{v}(t)$ be the solution of the matrix differential equation of (2.7), subject to the initial condition of (2.8), and let $\mathbf{w}_p(t)$ be defined from (5.3). Then, for all $t \geq 0$

$$(5.6) \quad \|\mathbf{v}(t) - \mathbf{w}_p(t)\| \leq \hat{H}_p \cdot \|\mathbf{g}^* - B^{-1}\Sigma\|.$$

Proof. From (4.5'),

$$\|\mathbf{v}(t) - \mathbf{w}_p(t)\| \leq \|e^{-tB} - Q_p(tB)\| \cdot \|\mathbf{g}^* - B^{-1}\Sigma\|.$$

Using the fact that the matrix $e^{-tB} - Q_p(tB)$ is Hermitian, then

$$\|e^{-tB} - Q_p(tB)\| = \max_{1 \leq i \leq n} \left| e^{-t\lambda_i} - \frac{n_p(t\lambda_i)}{d_p(t\lambda_i)} \right| \leq \hat{H}_p,$$

which completes the proof.

Thus, in comparison with the results in (4.7), it would appear that we can obtain *in one step* approximately the same accuracy with the Chebyshev rational approximation for $p = 2$, as, say, that given by the Crank-Nicolson method, with approximately 10^3 time steps. This requires some explanation, however. If the quantity $\|\mathbf{g}^* - B^{-1}\Sigma\|$ is very large, so that $(.0073) \|\mathbf{g}^* - B^{-1}\Sigma\|$ is itself large, it means that the Chebyshev rational approximation method $p = q = 2$ with one step could give rise to approximations $\mathbf{w}_p(t)$ of $\mathbf{v}(t)$ which have unacceptably large deviations from $\mathbf{v}(t)$. In this case, use of the Padé approximations is indicated, as $r_{p,q}^{(m)}(z)$ can be made as small as possible by choosing m sufficiently large. On the other hand, if $\|\mathbf{g}^* - B^{-1}\Sigma\|$ is small, potentially large savings in digital computer time seem possible by using Chebyshev rational approximations of e^{-tB} for approximating the solution of (2.7).

6. Applications. In the numerical solution of parabolic partial differential equations in one space variable, the inversion of linear polynomials in the matrix B , corresponding to tridiagonal matrices, is carried out directly in practical applications of methods such as the backward implicit method, and the Crank-Nicolson method. But as quadratic polynomials in the matrix B are in this case only five-diagonal matrices, such direct inversions are still quite efficient. Hence, higher order methods, based either on the Padé or Chebyshev rational approximations of e^{-z} , can be used without great difficulty in one space dimension.

In two spatial dimensions, the inversion of linear polynomials in the matrix B is not generally carried out directly in large problems, but rather an iterative technique, such as the Young-Frankel successive overrelaxation method [23, 7], is used in effect to invert such polynomials. These iterative methods can also be applied to higher order methods. For example, the case $p = q = 2$, which requires the inversion of quadratic polynomials in the matrix B , can be rigorously treated by use of the iterative S2LOR [21], as well as the extension of the cyclic Chebyshev semi-iterative method to this case [9]. The methods which are contained in

this paper easily give rise to theoretical generalizations in several directions. For example, the treatment of time-varying forcing functions, as well as a certain class of hyperbolic partial differential equations, can also be examined from this point of view. We shall however leave such theoretical extensions and numerical results to a subsequent paper.

BIBLIOGRAPHY

1. N. I. ACHESER, *Theory of Approximation*, translated by Charles J. Hyman, Frederick Ungar Publishing Co., New York, 1956.
2. R. E. BELLMAN, *Stability Theory of Differential Equations*, McGraw-Hill, New York, 1953.
3. G. BIRKHOFF AND R. S. VARGA, Reactor criticability and non-negative matrices, *J. Soc. Indust. Appl. Math.*, **6** (1958): 354-377.
4. R. COURANT, K. FRIEDRICHS AND H. LEVY, Über die partiellen Differenzgleichungen der mathematischen Physik, *Math. Ann.* **100** (1928): 32-74.
5. J. DOUGLAS, On the relation between stability and convergence in the numerical solution of linear parabolic and hyperbolic differential equations, *J. Soc. Indust. Appl. Math.*, **4** (1956): 20-37.
6. E. C. DUFORT AND S. P. FRANKEL, Stability conditions in the numerical treatment of parabolic differential equations, *Math. Tables Aids to Comp.* **7** (1953): 135-152.
7. S. P. FRANKEL, Convergence rates of iterative treatments of partial differential equations, *Math. Tables Aids Comp.* **4** (1950): 65-75.
8. J. FRANKLIN, Numerical stability in digital and analog computation for diffusion problems, *J. Math. and Phys.* **37** (1959): 305-315.
9. G. H. GOLUB AND R. S. VARGA, *Chebyshev semi-iterative methods, successive overrelaxation iterative methods, and second order Richardson iterative methods*, *Numerische Mathematik* **3** (1961): 147-156, 157-168.
10. A. S. HOUSEHOLDER, The approximate solution of matrix problems, *J. Assoc. Comp. Mach.* **5** (1958): 155-169.
11. P. M. HUMMEL AND C. L. SEEBECK, A generalization of Taylor's Theorem, *Amer. Math. Monthly* **56** (1949): 243-247.
12. M. L. JUNCOSA AND D. YOUNG, On the convergence of a solution of a difference equation to a solution of the equation of diffusion, *Proc. Amer. Math. Soc.* **5** (1954): 168-174.
13. Z. KOPAL, "Operational methods in numerical analysis based on rational approximations", *On Numerical Approximation*, University of Wisconsin Press, Madison, 1959: 25-43.
14. P. LAASONEN, Über eine Methode zur Lösung der Wärmeleitungsgleichung, *Acta Math.* **81** (1949): 309-317.
15. M. LOTKIN, The numerical integration of heat conduction equations, *J. Math. and Phys.* **37** (1958): 178-187.
16. Y. L. LUKE, Rational approximations to the exponential function, *J. Assoc. Comp. Mach.* **4** (1957): 24-29.
17. G. G. O'BRIEN, M. A. HYMAN, AND S. KAPLAN, A study of the numerical solution of partial differential equations, *J. Math. and Phys.* **29** (1951): 223-251.
18. H. PADÉ, *Sur la représentation approchée d'une fonction par des fractions rationnelles*, Thesis, Ann. de l'Éc. Nor., (3), **9** (1892).
19. R. D. RICHTMYER, *Difference Methods for Initial-Value Problems*, Interscience Publishers, Inc., New York, 1957.
20. R. S. VARGA, Numerical solution of the two-group diffusion equations in x - y geometry, *IRE Trans. on Nuclear Science* NS **4** (1957): 52-62.
21. R. S. VARGA, "Factorization and normalized iterative methods", *Boundary Problems in Differential Equations*, University of Wisconsin, Madison, (1959): 121-142.

22. H. S. WALL, *Analytic Theory of Continued Fractions*, D. van Nostrand Company, Inc., Princeton, 1948.
23. D. YOUNG, Iterative methods for solving partial difference equations of elliptic type, *Trans. Amer. Math. Soc.* **76** (1954): 92-111.

CASE INSTITUTE OF TECHNOLOGY

(Received December 15, 1960)