

ON THE USE OF STABILITY REGIONS IN THE NUMERICAL ANALYSIS OF INITIAL VALUE PROBLEMS

H. W. J. LENFERINK AND M. N. SPIJKER

ABSTRACT. This paper deals with the stability analysis of one-step methods in the numerical solution of initial (-boundary) value problems for linear, ordinary, and partial differential equations. Restrictions on the stepsize are derived which guarantee the rate of error growth in these methods to be of moderate size. These restrictions are related to the stability region of the method and to numerical ranges of matrices stemming from the differential equation under consideration.

The errors in the one-step methods are measured in arbitrary norms (not necessarily generated by an inner product).

The theory is illustrated in the numerical solution of the heat equation and some other differential equations, where the error growth is measured in the maximum norm.

1. INTRODUCTION

1.1. The numerical process. In this paper we analyze the stability of the numerical process

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

Here, $h > 0$ denotes the so-called *stepsize* and A is a square matrix of order $s \geq 1$. Further, φ is a given rational function with $\varphi(0) = \varphi'(0) = 1$. We assume $\varphi(z) = P(z)/Q(z)$, where $P(z)$, $Q(z)$ are polynomials with no common zero, and write $\varphi(hA) = P(hA)Q(hA)^{-1}$ whenever the matrix $Q(hA)$ is regular. The $u_n \in \mathbb{C}^s$ are numerical approximations computed in a step-by-step fashion from (1.1) starting from a given $u_0 \in \mathbb{C}^s$.

Many numerical methods for solving *ordinary differential equations*, such as Runge-Kutta and Rosenbrock methods, result, when applied to initial value problems for linear autonomous systems, in procedures of type (1.1). Further, many numerical schemes for solving initial-boundary value problems in *partial differential equations* can be written in the form (1.1). In the latter case, s is

Received November 6, 1989; revised June 21, 1990.

1980 *Mathematics Subject Classification* (1985 Revision). Primary 65L20, 65M10.

This first author's research has been supported by the Netherlands organization for scientific research (NWO).

Offprint requests to: M. N. Spijker.

related to the discretization of the space variables, and can attain large values. For examples, we refer to §4.

1.2. Error propagation. Suppose the numerical calculations based on (1.1) were performed using a slightly perturbed starting vector, say \tilde{u}_0 , instead of u_0 . We then would obtain approximations that we denote by \tilde{u}_n .

In the stability analysis of (1.1) the crucial question is whether the difference $v_n = \tilde{u}_n - u_n$ can be bounded suitably in terms of the perturbation $v_0 = \tilde{u}_0 - u_0$. Since $v_n = \varphi(hA)\tilde{u}_{n-1} - \varphi(hA)u_{n-1} = \varphi(hA)v_{n-1}$, the stability analysis thus amounts to investigating the possible growth of vectors v_n satisfying the recurrence relation (1.1).

We shall measure the size of v_n using an arbitrary norm $|x|$ for $x = (\xi_1, \xi_2, \dots, \xi_s)^T \in \mathbb{C}^s$. This norm is *not* required to be generated by an inner product, so that our discussion will include, e.g., the important maximum norm

$$|x|_\infty = \max_{1 \leq j \leq s} |\xi_j|.$$

In this paper we focus on stability estimates of the type

$$(1.2) \quad |v_n| \leq \gamma \cdot s^p n^q |v_0| \quad \text{for } s \geq 1, n \geq 1, \text{ and } v_n \text{ satisfying (1.1).}$$

Here, γ, p, q denote nonnegative constants independent of s, n, v_0 .

1.3. Stability regions. An obvious manner to assess the stability of process (1.1) is to use the eigenvalues of the matrix $\varphi(hA)$. Denoting the *spectrum* of A by $\sigma[A]$, we see that the spectrum of $\varphi(hA)$ equals $\varphi(h\sigma[A])$. In order to guarantee (1.2), one thus arrives at the requirement that $|\varphi(h\lambda)| \leq 1$, or slightly stronger, $|\varphi(h\lambda)| < 1$, for all $\lambda \in \sigma[A]$.

Defining the *stability region* S of φ by

$$S = \{\zeta: \zeta \in \mathbb{C} \text{ with } |\varphi(\zeta)| \leq 1\},$$

the above two requirements can be cast into the form

$$(1.3.a) \quad h\sigma[A] \subset S,$$

$$(1.3.b) \quad h\sigma[A] \subset \text{int}(S),$$

respectively. Here, $\text{int}(S)$ denotes the interior of S .

In the case of the Euclidean norm and a normal matrix A one easily sees that (1.3.a) implies (1.2) with $\gamma = 1, p = 0, q = 0$. However, in more general situations, conditions (1.3) can be very unreliable. The point is that although, e.g., (1.3.b) guarantees (1.2) with $p = q = 0$, the size of γ is not under control and can be arbitrarily large. This was pointed out by, among others, Griffiths, Christie, and Mitchell [11], who displayed an instructive example where one has essentially $|v_n| \geq \alpha^n |v_0|$ ($n = 1, 2, \dots, s$), $|v_n| \leq \alpha^s |v_0|$ ($n \geq s$) with $\alpha > 1$ and arbitrary dimension $s \geq 1$. See also [15, 18, 22, 27, 28] and §§4.2, 4.3 of the present paper.

The unreliable conditions (1.3) can be converted into reliable ones, essentially by replacing $\sigma[A]$ by some appropriately chosen larger set $\tau[A]$. Here, $\tau[A]$ is associated with the matrices under consideration and satisfies $\sigma[A] \subset \tau[A] \subset \mathbb{C}$. Under such modified conditions, stability estimates (1.2), with γ nicely under control, were derived in [4, 6, 15, 20, 25, 27]. However, the conditions imposed in these references on h , A , and S are not completely satisfactory in that they cannot be fulfilled in some cases of practical interest.

1.4. Scope of the paper. This paper attempts to improve the unsatisfactory situation just mentioned. We shall derive stability estimates of type (1.2) which apply to some general situations not covered in the references cited above. We focus on modified versions of conditions (1.3), where $\sigma[A]$ is replaced by the so-called *M-numerical range* $\tau[A]$, a subset of the complex plane recently introduced in [16].

In §§2.1, 2.2 we give our basic definitions and characterizations of the *M-numerical range*. Using this concept, we review in §2.3 stability results from [4, 6, 15, 20, 25, 27]. In §2.4 we relate the *M-numerical range* to so-called circle conditions, which were basic for the stability analysis of [15, 20, 27].

Section 3 contains our main results. Section 3.1 gives an estimate (1.2) with γ nicely under control and with the optimal values $p = 0$, $q = 0$. The conditions on $h\tau[A]$ in §3.1 are rather strong. Weaker conditions are dealt with in §§3.2, 3.3.

Section 4 illustrates the stability estimates of §3. In §4.1 we apply the material from §3.1 in proving strong stability with respect to the maximum norm of finite difference methods for solving the heat equation. Sections 4.2, 4.3 contain numerical experiments pertinent to ordinary and partial differential equations, respectively.

In this paper we confine ourselves to using stability regions in deriving stability estimates for *linear, one-step* processes (1.1) with *arbitrary norms* $|\cdot|$ in (1.2). For related stability results, based on stability regions, pertinent to non-linear differential equations, multistep methods, or norms generated by an inner product, the reader may consult, e.g., [7, 8, 19, 23, 27].

2. BASIC DEFINITIONS AND A REVIEW OF KNOWN STABILITY RESULTS

2.1. Definition and elementary properties of the *M-numerical range*. For complex γ and $\rho \geq 0$ we introduce the disk

$$D[\gamma, \rho] = \{\zeta: \zeta \in \mathbb{C} \text{ with } |\zeta - \gamma| \leq \rho\}.$$

By $|\cdot|$ we denote an arbitrary norm on \mathbb{C}^s and by $\|\cdot\|$ the corresponding induced matrix norm on $\mathbb{C}^{s,s}$, defined by $\|A\| = \sup\{|Ax|: x \in \mathbb{C}^s \text{ with } |x| = 1\}$ for $s \times s$ matrices A . Let $M \geq 1$ be a given constant, and $A = (\alpha_{jk})$ an $s \times s$ matrix.

We define a disk $D[\gamma, \rho]$ to be *M-suitable for A* if

$$(2.1) \quad \|(A - \gamma I)^k\| \leq M \rho^k \quad (k = 1, 2, 3, \dots).$$

The M -numerical range of A (cf. [16]) is defined by

$$(2.2) \quad \tau_M[A] = \bigcap D[\gamma, \rho],$$

where the intersection is over all disks that are M -suitable for A .

In case $M = 1$, the set (2.2) can be seen to coincide with the so-called *algebra numerical range* (cf. [1, 2, 3]). Also, the terms Gerschgorin domain (see [26]), Hausdorff set, and field of values (see, e.g., [10]) occur in the literature to designate sets that coincide with the set $\tau_1[A]$. In the case of the maximum norm $|\cdot| = |\cdot|_\infty$, it is known (cf. [16, 21, 26]) that $\tau_1[A]$ equals the convex hull of the union of the so-called Gerschgorin disks

$$D_j[A] = \left\{ \zeta: \zeta \in \mathbb{C} \text{ with } |\zeta - \alpha_{jj}| \leq \sum_{k \neq j} |\alpha_{jk}| \right\}.$$

We thus can write

$$(2.3) \quad \tau_1[A] = \text{conv} \left\{ \bigcup_{j=1}^s D_j[A] \right\} \quad \text{if } |\cdot| = |\cdot|_\infty.$$

For general $M \geq 1$ it follows from definition (2.2) that (cf. [16])

$$(2.4.a) \quad \tau_M[A] \text{ is a compact, convex subset of the complex plane,}$$

$$(2.4.b) \quad \tau_M[\zeta_0 I + \zeta_1 A] = \zeta_0 + \zeta_1 \tau_M[A] \quad \text{for } \zeta_0, \zeta_1 \in \mathbb{C},$$

$$(2.4.c) \quad \text{conv } \sigma[A] \subset \tau_M[A],$$

$$(2.4.d) \quad \tau_N[A] \subset \tau_M[A] \subset \tau_1[A] \quad \text{for } 1 \leq M \leq N,$$

$$(2.4.e) \quad \lim_{M \rightarrow \infty} \tau_M[A] = \text{conv } \sigma[A].$$

2.2. Basic characterizations of the M -numerical range. We start with some definitions that are needed for formulating subsequent characterizations of $\tau_M[A]$.

In all of the following, V denotes an arbitrary nonempty, closed, and convex subset of \mathbb{C} . The *distance* from $\zeta \in \mathbb{C}$ to V is defined by

$$d(\zeta, V) = \inf\{|\zeta - \xi|: \xi \in V\}.$$

If ξ belongs to the *boundary* ∂V of V and

$$\text{Re}\{e^{-i\theta}(\zeta - \xi)\} \leq 0 \quad \text{for all } \zeta \in V,$$

where θ is a real constant, then θ is called a *normal direction* to V at ξ .

The M -numerical range allows the following two characterizations (cf. [16]):

$$(2.5.a) \quad \tau_M[A] \text{ is the smallest } V \text{ with the property that } (\zeta I - A) \text{ is regular and } \|(\zeta I - A)^{-k}\| \leq M \cdot [d(\zeta, V)]^{-k} \text{ for all } \zeta \notin V \text{ and } k = 1, 2, 3, \dots$$

$$(2.5.b) \quad \tau_M[A] \text{ is the smallest } V \text{ with the property that}$$

$$\|\exp[te^{-i\theta}(A - \xi I)]\| \leq M$$

for all $t \geq 0$, $\xi \in \partial V$, and normal directions θ to V at ξ .

In the following we use these characterizations to review in a coherent fashion some of the stability results to be found in the literature.

2.3. A review of some stability results from the literature. In [4] Brenner and Thomée derived important stability estimates pertinent to linear operators A in Banach spaces. Specializing their general estimates to the s -dimensional space \mathbb{C}^s , and using the characterizations (2.5), it follows that an estimate of type (1.2) holds with

$$(2.6.a) \quad p = 0, \quad q = 1/2, \quad \gamma = \gamma_0 M,$$

provided (1.3.a) is strengthened to the condition

$$(2.6.b) \quad h\tau_M[A] \subset \mathbb{C}_- \subset S.$$

Here, \mathbb{C}_- stands for the *half plane*

$$\mathbb{C}_- = \{\zeta: \zeta \in \mathbb{C} \text{ with } \operatorname{Re} \zeta \leq 0\}$$

and γ_0 is a constant depending only on φ . In a similar way one can conclude from [4, 25] that a sharper estimate (1.2) holds, with

$$(2.7.a) \quad p = 0, \quad q = 0, \quad \gamma = \gamma_0 M,$$

provided (1.3.a) is strengthened to the condition

$$(2.7.b) \quad h\tau_M[A] \subset W(\alpha) \subset \{0\} \cup \operatorname{int}(S) \quad \text{and} \quad |\varphi(\infty)| < 1.$$

Here $W(\alpha)$ stands for a *wedge*

$$W(\alpha) = \{\zeta: \zeta \in \mathbb{C} \text{ with } \zeta = 0 \text{ or } \pi \geq |\arg \zeta| \geq \pi - \alpha\}$$

with $0 \leq \alpha < \pi/2$, and γ_0 depends on φ, α only. We note that (2.7) can also be viewed as a corollary to a result by Crouzeix [6, Theorem 8].

The above conditions (2.6.b), (2.7.b) cannot be fulfilled by explicit methods (1.1) (i.e., methods where $\varphi(\zeta)$ is a polynomial). But the following two stability results are relevant also for the case of explicit methods.

In [15, 20, 27] it was shown that (1.2) holds with

$$(2.8.a) \quad p = 0, \quad q = 1/2, \quad \gamma = \gamma_0$$

under the condition

$$(2.8.b) \quad hD[-\rho, \rho] \subset S.$$

Here $\rho \geq 0$ is associated with the matrix A such that the *circle condition*

$$(2.9) \quad \|A + \rho I\| \leq \rho$$

is fulfilled. Further, γ_0 depends on $\varphi, h\rho$ only.

In [15] a sharper estimate (1.2), with

$$(2.10.a) \quad p = 0, \quad q = 0, \quad \gamma = \gamma_0,$$

was derived under a slightly stronger condition than (2.8.b), viz.

$$(2.10.b) \quad 0 < h < h_0 \quad \text{with} \quad h_0 D[-\rho, \rho] \subset S.$$

Here, ρ is as in (2.9) and γ_0 depends again on $\varphi, h\rho$ only. The proof in [15] is given only for the maximum norm $|\cdot| = |\cdot|_\infty$, but it is easily verified that the result (2.10) is still valid for an arbitrary norm $|\cdot|$ on \mathbb{C}^s .

2.4. Comparing the use of circle conditions to the use of M -numerical ranges. In §3 of this paper we shall derive stability estimates (1.2) under conditions of type

$$(2.11) \quad h\tau_M[A] \subset V.$$

Here, V will be a bounded subset of S . Similarly as (2.8.b), (2.10.b), condition (2.11) can be relevant also for the case of explicit methods.

We shall compare (2.11) with (2.8.b), where $\rho \geq 0$ is assumed to be a minimal radius with property (2.9). We review various cases.

1. Let $V = D[-r, r]$ with $r > 0$, $M = 1$, $|\cdot| = |\cdot|_\infty$. Using (2.3), we can see that (2.11) is equivalent to

$$(2.12) \quad hD[-\rho, \rho] \subset V.$$

2. Let $V = D[-r, r]$ with $r > 0$, $M = 1$. With respect to arbitrary norms $|\cdot|$, condition (2.11) is no longer equivalent to (2.12). This follows from the following counterexample, modelled after an example of J. Kraaijevanger [14].

With the choice of the Euclidean norm $|\cdot| = |\cdot|_2$ on \mathbb{C}^2 , the matrix

$$A = \begin{pmatrix} -r & 0 \\ \frac{3}{2}r & -r \end{pmatrix}$$

can be seen to satisfy

$$\tau_1[A] \subset D[-r, r], \quad \|A + rI\| > r, \quad r < \rho < \infty.$$

Hence, with $h = 1$, (2.11) holds, but (2.12) is violated.

In view of (2.2), condition (2.11) will thus, in general, be weaker than (2.12).

3. Let V not be equal to a disk, and $M = 1$. From (2.2) it follows that (2.11) is now always a *weaker* condition than (2.12). We thus cover more situations by relaxing (2.12) to (2.11).

The above makes clear that when $M = 1$ the set $\tau[A] = \tau_M[A]$, which we focus on in this paper, has fundamental advantages over the set $\tau[A] = D[-\rho, \rho]$ with ρ as above. Moreover, in view of (2.4.d), (2.4.e) we see that, by increasing M , stepsize conditions of type (2.11) can become still weaker and close to the "optimal" condition

$$h \operatorname{conv} \sigma[A] \subset V$$

(cf. (1.3.a)).

3. STABILITY ESTIMATES BASED ON THE NUMERICAL RANGE

3.1. Stability with $p = 0$, $q = 0$, for $\tau_M[A]$ within a bounded wedge. In this subsection we derive stability estimates for process (1.1) which are similar to the result (2.7). But the conditions we impose on φ are essentially weaker than in (2.7.b), and can be fulfilled, e.g., when φ is a polynomial. On the other hand, the conditions we impose on hA are stronger than in (2.7.b). Our main result will be formulated in Theorem 3.2. Its proof will rely on an application of the following lemma (with A replaced by hA).

We make the assumptions that

(3.1.a) $0 \leq \alpha < \pi/2$ and V is a compact, convex subset of $W(\alpha)$,

(3.1.b) $V \subset \text{int}(S) \cup \{0\}$,

(3.1.c) $s \geq 1$, $A \in \mathbb{C}^{s,s}$, and $M \geq 1$.

Lemma 3.1. Assume (3.1), and let $(\zeta I - A)$ be regular with

(3.2) $\|(\zeta I - A)^{-1}\| \leq Md(\zeta, V)^{-1}$ for all $\zeta \notin V$.

Then

$$\|\varphi(A)^n\| \leq \gamma_0 M \quad (n = 1, 2, 3, \dots),$$

with γ_0 depending only on φ and V .

Proof. We will express $\varphi(A)^n$ by means of an integral along a contour around V , the construction of which is possible because of (3.1.a,b).

The set $V_\delta = \{\zeta: \zeta \in \mathbb{C} \text{ and } d(\zeta, V) \leq \delta\}$ is convex for each $\delta > 0$. By using Theorem 10.4 in [24] one can see that the boundary $\partial(V_\delta)$ is equal to the range of a Lipschitz continuous, simple, positively oriented, closed curve Γ , parametrized by $\zeta = z(t)$, $0 \leq t \leq 1$.

First, let $0 \in V$, and assume, with no loss of generality, that $z(0) = \delta$. Let β be any number with $\alpha < \beta < \pi/2$. Since $\varphi(0) = \varphi'(0) = 1$ (cf. §1.1), we can choose $\delta > 0$ such that

$$V_\delta \cap W(\beta) \subset \text{int}(S) \cup \{0\},$$

φ has no poles at any $\zeta \in \mathbb{C}$ with $|\zeta| \leq \delta$.

Let the curve Γ_0 be obtained by restricting $z(\cdot)$ to $[t_0, t_1]$, where $z(t_i) \in \partial(W(\beta))$ ($i = 0, 1$). For $0 < \varepsilon \leq \delta$ we define the curves $\Gamma_1, \Gamma_2, \Gamma_3$ by

$$\begin{aligned} z_1(t) &= t \exp(i\beta) & (-|z(t_1)| \leq t \leq -\varepsilon), \\ z_2(t) &= \varepsilon \exp(it) & (-\pi + \beta \leq t \leq \pi - \beta), \\ z_3(t) &= -t \exp(-i\beta) & (\varepsilon \leq t \leq |z(t_0)|), \end{aligned}$$

respectively. By (3.2), the spectrum $\sigma[A]$ of A lies within the closed curve which consists of the segments $\Gamma_0, \Gamma_1, \Gamma_2, \Gamma_3$. As these are sufficiently smooth, we may use the Cauchy integral formula [9, p. 568] to write

(3.3)
$$\varphi(A)^n = \frac{1}{2\pi i} \sum_{j=0}^3 \int_{\Gamma_j} \varphi(\zeta)^n (\zeta I - A)^{-1} d\zeta$$

for any ε with $0 < \varepsilon \leq \delta$. We will bound the norms of the four terms on the right-hand side of (3.3) by using (3.2). First,

(3.4.a)
$$\left\| \int_{\Gamma_0} \varphi(\zeta)^n (\zeta I - A)^{-1} d\zeta \right\| \leq \left(\int_{\Gamma_0} |d\zeta| \right) M \delta^{-1}.$$

Since Γ_1 and Γ_3 lie in $\text{int}(S)$ and $\varphi(0) = \varphi'(0) = 1$, there exists a constant $L_1 > 0$ such that

$$|\varphi(\zeta)| \leq \exp(-|\zeta|L_1) \quad \text{for all } \zeta \in \Gamma_1 \cup \Gamma_3 \text{ and } 0 < \varepsilon \leq \delta.$$

We choose $\varepsilon = \delta n^{-1}$. Then we get, substituting $s = nt$,

$$(3.4.b) \quad \left\| \int_{\Gamma_1} \varphi(\zeta)^n (\zeta I - A)^{-1} d\zeta \right\| \leq M \int_{\Gamma_1} |\varphi(\zeta)|^n d(\zeta, W(\alpha))^{-1} |d\zeta| \\ \leq M \int_{\delta}^{\infty} \exp(-sL_1) (s \sin(\beta - \alpha))^{-1} ds.$$

The same bound is valid for the integral along Γ_3 in (3.3). Let $L_2 > 0$ be such that

$$|\varphi(\zeta)| \leq 1 + |\zeta|L_2 \quad \text{for all } \zeta \text{ with } |\zeta| \leq \delta.$$

We obtain

$$(3.4.c) \quad \left\| \int_{\Gamma_2} \varphi(\zeta)^n (\zeta I - A)^{-1} d\zeta \right\| \leq 2\pi M (1 + \delta n^{-1} L_2)^n [\sin(\beta - \alpha)]^{-1} \\ \leq 2\pi M \exp(\delta L_2) [\sin(\beta - \alpha)]^{-1}.$$

From (3.3), (3.4) it follows immediately that there exists a constant $\gamma_0 > 0$, which depends on φ and on V only, such that

$$\|\varphi(A)^n\| \leq \gamma_0 M \quad (\text{for } n = 1, 2, 3, \dots).$$

If $0 \notin V$, we may choose δ so small that $V_\delta \subset \text{int}(S)$ and use integration along the curve Γ specified above to obtain an integral representation even simpler than (3.3). The conclusion of the lemma follows directly from (3.2) \square

Theorem 3.2. *Assume (3.1) and $0 \in V$, $h_0 > 0$. Let the condition $h_0 \tau_M[A] \subset V$ be fulfilled. Then, for all $h \in (0, h_0]$, estimate (1.2) holds with*

$$p = 0, \quad q = 0, \quad \gamma = \gamma_0 M,$$

where γ_0 depends only on φ and V .

Proof. The conditions $h_0 \tau_M[A] \subset V$, $0 \in V$, the convexity of V , and property (2.4.b) imply that $\tau_M[hA] \subset V$ for $h \in (0, h_0]$. Characterization (2.5.a) with A replaced by hA shows that $(\zeta I - hA)$ is regular and

$$\|(\zeta I - hA)^{-1}\| \leq M d(\zeta, V)^{-1} \quad \text{for all } h \in (0, h_0] \text{ and } \zeta \notin V.$$

We may apply Lemma 3.1, with A replaced by hA . For some $\gamma_0 \geq 0$, depending only on φ and V , we thus have

$$(3.5) \quad \|\varphi(hA)^n\| \leq \gamma_0 M \quad (n = 1, 2, \dots).$$

Since the v_n satisfying (1.1) equal $v_n = \varphi(hA)^n v_0$, the proof is completed by an application of (3.5). \square

Illustrations to the above theorem will be given in §§4.1, 4.2.

3.2. Stability with $p = 0$, $q = 1$, for $\tau_M[A]$ within a bounded set. In this subsection we derive stability results for process (1.1) which, like the results (2.8), (2.10) and Theorem 3.2, are relevant to a large class of functions $\varphi(\zeta)$, including polynomials. Moreover, the assumptions we make concerning S can

be fulfilled in cases where (2.8.b), (2.10.b) are violated, and the conditions we impose on hA can be weaker than those in Theorem 3.2. On the other hand, the value $q = 1$ we arrive at is larger than in (2.8), (2.10) or Theorem 3.2.

We make the assumptions that

(3.6.a) V is a compact, convex subset of \mathbb{C} ,

(3.6.b) $V \subset S$,

(3.6.c) $s \geq 1$, $A \in \mathbb{C}^{s,s}$, and $M \geq 1$.

Lemma 3.3. Assume (3.6), and let $(\zeta I - A)$ be regular with

(3.7) $\|(\zeta I - A)^{-1}\| \leq Md(\zeta, V)^{-1}$ for all $\zeta \notin V$.

Then

$$\|\varphi(A)^n\| \leq \gamma_0 M n \quad (n = 1, 2, 3, \dots),$$

with γ_0 only depending on φ and V .

Proof. The proof rests on a representation of $\varphi(A)^n$ as a suitable contour integral. Define, for $\varepsilon > 0$,

$$V_\varepsilon = \{\zeta: \zeta \in \mathbb{C} \text{ and } d(\zeta, V) \leq \varepsilon\}.$$

In view of (3.6.a) the set V_ε is compact and convex. The boundary of V_ε is equal to the range of a Lipschitz continuous, simple, positively oriented, closed curve Γ_ε (see, e.g., [24, Theorem 10.4]).

By (3.7), the spectrum $\sigma[A]$ lies within Γ_ε . Further, in view of (3.6.b), the function φ has no poles in V_α for some $\alpha > 0$. We take $0 < \varepsilon \leq \alpha$ and start from the Cauchy integral formula [9, p. 568]

(3.8)
$$\varphi(A)^n = \frac{1}{2\pi i} \int_{\Gamma_\varepsilon} \varphi(\zeta)^n (\zeta I - A)^{-1} d\zeta \quad (n = 1, 2, 3, \dots).$$

Let $L = \max\{|\varphi'(\zeta)|: \zeta \in V_\alpha\}$. Then, using (3.6.b), we can see that

$$|\varphi(\zeta)^n| \leq (1 + \varepsilon L)^n \quad \text{whenever } \zeta \in V_\varepsilon, \quad 0 < \varepsilon \leq \alpha.$$

So, we obtain from (3.8)

(3.9)
$$\|\varphi(A)^n\| \leq \frac{1}{2\pi} \left(\int_{\Gamma_\varepsilon} |d\zeta| \right) (1 + \varepsilon L)^n \frac{M}{\varepsilon} \quad (n = 1, 2, 3, \dots).$$

The length of Γ_ε is bounded by that of Γ_α (see, e.g., [13, p. 245]), which will be denoted by P . With $\varepsilon = \alpha n^{-1}$, inequality (3.9) leads to

$$\|\varphi(A)^n\| \leq P \exp(\alpha L) M n (2\pi\alpha)^{-1} \quad (n = 1, 2, 3, \dots).$$

From this, the conclusion of the lemma follows with $\gamma_0 = P \exp(\alpha L) (2\pi\alpha)^{-1}$. \square

The next theorem follows from the above lemma by using arguments that are analogous to those used in proving Theorem 3.2 by means of Lemma 3.1.

Theorem 3.4. Assume (3.6) and $0 \in V$, $h_0 > 0$. Let the condition $h_0\tau_M[A] \subset V$ be fulfilled. Then, for all $h \in (0, h_0]$, estimate (1.2) holds with

$$p = 0, \quad q = 1, \quad \gamma = \gamma_0 M,$$

where γ_0 depends only on φ and V .

An illustration to the above theorem will be given in §4.3.

3.3. Stability with $p = 1$, $q = 0$, for $\tau_M[A]$ within a bounded set. Can one, under the general conditions on φ and V of Theorem 3.4, improve upon the value q for which the theorem is true? We will see that one can keep $q = 0$, as in Theorem 3.2. However, the value of p will go up to $p = 1$. We also have to make some further assumptions on φ and V , which are of a geometrical nature, but these will usually not form an impediment to the application of the ensuing theorem.

In addition to (3.6) we make the assumptions that

$$(3.10.a) \quad \varphi'(\zeta) \neq 0 \quad \text{for all } \zeta \in \partial V \cap \partial S,$$

(3.10.b) there exist a positive integer ν and real coefficients $\alpha_{j,k}$ (for $j \geq 0$, $k \geq 0$, $j+k \leq \nu$) with $\alpha_{j,k} \neq 0$ for some j, k satisfying $j+k = \nu$, such that the boundary ∂V of V is a subset of

$$\{\zeta: \zeta = \xi + i\eta \text{ with } \xi, \eta \in \mathbb{R} \text{ and } K(\xi, \eta) = 0\}.$$

Here

$$K(\xi, \eta) = \sum_{0 \leq j+k \leq \nu} \alpha_{j,k} \xi^j \eta^k.$$

Lemma 3.5. Assume (3.6) and (3.10). Let $(\zeta I - A)$ be regular with

$$\|(\zeta I - A)^{-1}\| \leq Md(\zeta, V)^{-1} \quad \text{for all } \zeta \notin V.$$

Then

$$\|\varphi(A)^n\| \leq \gamma_0 M s \quad (n = 1, 2, 3, \dots),$$

with γ_0 depending only on φ and V .

This lemma is an immediate consequence of the material presented in [17].

The following theorem can be proved by applying the above lemma and using arguments analogous to those used in proving Theorem 3.2.

Theorem 3.6. Assume (3.6), (3.10) and $0 \in V$, $h_0 > 0$. Let the condition $h_0\tau_M[A] \subset V$ be fulfilled. Then, for all $h \in (0, h_0]$, estimate (1.2) holds with

$$p = 1, \quad q = 0, \quad \gamma = \gamma_0 M,$$

where γ_0 depends only on φ and V .

An illustration to the above theorem will be given in §4.3.

The question arises whether the values $p = 1, q = 0$ in the above theorem can be replaced by $p = q = 0$. Unfortunately, the answer is negative. A counterexample can be constructed along the lines of [15, §6.1]. The authors have not been able to answer the question whether Theorem 3.6 is still valid for some p, q with $p + q < 1$.

4. EXAMPLES AND APPLICATIONS

4.1. **Stability estimates in the numerical solution of the heat equation.** We will apply Theorem 3.2 to derive stability estimates, with respect to the *maximum norm*, in the numerical solution of the 1-dimensional heat equation

$$\frac{\partial}{\partial t} u(x, t) = \frac{\partial^2}{\partial x^2} u(x, t) \quad (0 < x < 1; t > 0).$$

We assume homogeneous Dirichlet boundary conditions and an initial condition for u to be given. Standard space discretization with $\Delta x = (1 + s)^{-1}$ leads to an initial value problem for a system of s ordinary differential equations of type

$$\frac{d}{dt} U(t) = AU(t) \quad (t > 0).$$

Here, $U(t)$ stands for a vector in \mathbb{R}^s (unknown for $t > 0$) and A is the square tridiagonal matrix of order s with entries $-2(\Delta x)^{-2}$ on the main diagonal and $(\Delta x)^{-2}$ on the adjacent diagonals.

Consider the numerical solution of the above system by any standard one-step method (such as a Runge-Kutta method or Rosenbrock method; cf., e.g., [12]) with stepsize $\Delta t = h > 0$. One then obtains approximations $u_n \simeq U(n\Delta t)$ from a particular recurrence relation of type (1.1). In the following we study the stability of this recurrence relation.

In order to be able to apply Theorem 3.2, we first consider the numerical range of the matrix A (with respect to the maximum norm $|\cdot| = |\cdot|_\infty$ on \mathbb{C}^s). Let any α be given with $0 < \alpha < \pi/2$. It follows from [16, Theorem 3.1] that there exist $M \geq 1$ and $\lambda_0 > 0$ such that

$$\tau_M[A] \subset W(\alpha) \cap \{\zeta: \operatorname{Re} \zeta \geq -\lambda_0(\Delta x)^{-2}\}.$$

Here the quantities M and λ_0 depend on α but not on Δx . We give two applications of Theorem 3.2.

(i) Assume $\alpha \in (0, \pi/2)$ is such that, for the stability region S of the one-step method under consideration, we have

$$W(\alpha) \subset \operatorname{int}(S) \cup \{0\}.$$

Note that we do *not* assume, as in (2.7.b), that $|\varphi(\infty)| < 1$. Let $C > 0$ be any given constant, and define

$$V = W(\alpha) \cap \{\zeta: \operatorname{Re} \zeta \geq -\lambda_0 C\}.$$

Then the conditions of Theorem 3.2 are fulfilled provided $h_0(\Delta x)^{-2} \leq C$. Consequently, for any $\Delta t = h > 0$ and $\Delta x = (1 + s)^{-1}$ with

$$\Delta t(\Delta x)^{-2} \leq C,$$

there is strong stability in the following sense: Any solution v_n of our recurrence relation (1.1) satisfies $|v_n|_\infty \leq \gamma |v_0|_\infty$ with γ independent of $v_0, n, \Delta t, \Delta x$ (but possibly depending on C).

(ii) Consider any $\alpha \in (0, \pi/2)$. Since $\varphi(0) = \varphi'(0) = 1$ for the function φ under consideration, there exists $\lambda > 0$ such that (3.1.b) holds for

$$V = W(\alpha) \cap \{\zeta: \operatorname{Re} \zeta \geq -\lambda\}.$$

Hence, the conditions of Theorem 3.2 are fulfilled with $h_0 = (\Delta x)^2 \lambda \lambda_0^{-1}$. This implies again strong stability under a stepsize condition $\Delta t(\Delta x)^{-2} \leq C$ —but now with $C = \lambda \lambda_0^{-1}$, whereas in the first application (i) there was *no* restriction on C . Note that in the present application the function φ can be a polynomial, while in the first application process (1.1) was necessarily implicit.

The above conclusions do not seem to follow easily from the related material in [4, 6, 25].

4.2. A numerical illustration to Theorem 3.2. We consider the initial value problem in $\mathbb{R}^s, s \geq 2,$

$$(4.1.a) \quad \frac{d}{dt}U(t) = AU(t) \quad (t \geq 0),$$

$$(4.1.b) \quad U(0) = u_0,$$

where $A = (\alpha_{jk})$ is the nonsymmetric tridiagonal $s \times s$ matrix with

$$\begin{aligned} \alpha_{j,j-1} &= j^{-1} && (2 \leq j \leq s), \\ \alpha_{j,j} &= -(2j + 3) && (1 \leq j \leq s), \\ \alpha_{j,j+1} &= j + 1 && (1 \leq j \leq s - 1). \end{aligned}$$

We choose $h > 0$ and apply (1.1) in order to obtain numerical approximations u_n to the solution $U(nh)$ of (4.1) for $n = 1, 2, 3, \dots$. We use $\varphi(\zeta) = 1 + \zeta + c\zeta^2$ and choose c such that the stability interval along the real axis is maximal. This choice yields $c = 1/8$ and

$$\operatorname{int}(S) \cap \mathbb{R} = (-8, -4) \cup (-4, 0)$$

(cf. [12]). In the following we deal with the *maximum norm* $|\cdot| = |\cdot|_\infty$ and study the actual stability behavior of process (1.1) for two different choices of the stepsize h .

(i) Let $D = \operatorname{diag}(1!, 2!, \dots, s!)$. Then DAD^{-1} is a symmetric irreducible tridiagonal matrix. The main diagonal is the same as that of A , whereas the second diagonals contain only ones. Hence (cf., e.g., [5, p. 352]), the eigenvalues λ of A are all different from each other and satisfy

$$(-2s + 4) < \lambda < 0.$$

With $s = 40$, $h = 0.07$, we have

$$(4.2) \quad h\sigma[A] \subset (-8, -4) \cup (-4, 0) \subset \text{int}(S),$$

so that condition (1.3.b) is fulfilled here. Defining

$$(4.3.a) \quad \gamma_n = \|\varphi(hA)^n\|,$$

we see that γ_n is the smallest constant with

$$(4.3.b) \quad |v_n| \leq \gamma_n |v_0| \quad \text{for all } v_0 \in \mathbb{R}^s \text{ and } v_n \text{ satisfying (1.1).}$$

Since (1.3.b) is fulfilled, (1.2) holds for some $\gamma > 0$ and $p = q = 0$. Consequently, $\sup_{n \geq 1} \gamma_n \leq \gamma$. It is evident from the values in Table 1 that the smallest possible value of γ is quite large—from a practical point of view there is actually instability.

TABLE 1
Some values of γ_n for $h = 0.07$

n	1	4	16	64	256	1024	4096
γ_n	2.8×10^0	1.3×10^1	1.6×10^3	9.8×10^6	3.3×10^9	2.6×10^{10}	3.1×10^{10}

(ii) Define V to be the union $T_1 \cup T_2$ of the triangle

$$T_1 = \{\zeta: |\pi - \arg \zeta| \leq \pi/6 \text{ and } -27/16 \leq \text{Re } \zeta \leq 0\}$$

and the disk

$$T_2 = \{\zeta: |\zeta + 9/4| \leq 9/8\}.$$

Straightforward numerical calculations show that the general conditions (3.1.a), (3.1.b) are fulfilled in the situation at hand with $\alpha = \pi/6$. Using (2.3), we can see that also

$$h_0 \tau_M[A] \subset V$$

with $M = 1$, $h_0 = 9/(8s + 4)$. All assumptions in Theorem 3.2 being fulfilled here, we thus can conclude that the stability estimate (1.2) holds with $p = q = 0$, $\gamma = \gamma_0$ whenever $s \geq 2$ and $0 < h \leq h_0$. With $s = 40$, we arrive at $h_0 = 1/36 \simeq 0.02777$. We use $h = 0.027$ to compute some values of γ_n as defined in (4.3). They are listed in Table 2 and point to a fine stability behavior of the numerical process (1.1)—perfectly in agreement with Theorem 3.2.

TABLE 2
Some values of γ_n for $h = 0.027$

n	1	2	4	8	16	32	64
γ_n	1.2×10^0	1.0×10^0	8.4×10^{-1}	5.3×10^{-1}	2.3×10^{-1}	3.9×10^{-2}	7.4×10^{-4}

4.3. A numerical illustration to Theorems 3.4 and 3.6 with a cigar-shaped V . Theorems 3.4 and 3.6 are more general than Theorem 3.2 in that there is more freedom in choosing the set V . We will give an example where V is cigar-shaped and illustrate the practical relevance of Theorems 3.4 and 3.6 by comparing the stability results predicted by the theorems to the outcome of numerical experiments. This outcome cannot be explained by the theorems in [4, 6, 15, 20, 25, 27], where V has to be of conventional shape.

In the following we illustrate our theorems with the cigar-shaped V of type

$$V(\lambda, \rho) = \{\zeta: \zeta = \xi + \eta \text{ with } \xi \in \mathbb{R}, \eta \in \mathbb{C}, -\lambda - \rho \leq \xi \leq -\rho, |\eta| \leq \rho\},$$

where $\lambda, \rho \geq 0$ are given parameters.

Consider the initial-boundary value problem

$$\begin{aligned} \frac{\partial}{\partial t} U(x, t) &= \frac{\partial^2}{\partial x^2} U(x, t) - 200 \frac{\partial}{\partial x} U(x, t) - 137000 \cdot x \cdot U(x, t), \\ U(0, t) &= U(1, t) = 0, \quad U(x, 0) = U_0(x), \end{aligned}$$

where $0 < x < 1, t > 0$, and U_0 is a given function.

To formulate a numerical method for approximating $U(x, t)$, we choose $\Delta t = h > 0, \Delta x = (1 + s)^{-1}$, and the tridiagonal $s \times s$ matrix $A = (\alpha_{jk})$ with

$$\begin{aligned} \alpha_{j, j-1} &= (\Delta x)^{-2} + 100(\Delta x)^{-1} & (2 \leq j \leq s), \\ \alpha_{j, j} &= -2(\Delta x)^{-2} - 137000 \cdot j \cdot \Delta x & (1 \leq j \leq s), \\ \alpha_{j, j+1} &= (\Delta x)^{-2} - 100(\Delta x)^{-1} & (1 \leq j \leq s-1). \end{aligned}$$

Let u_n be computed from (1.1) starting with

$$u_0 = (U_0(\Delta x), U_0(2\Delta x), \dots, U_0(s\Delta x))^T.$$

Then the k th component of u_n approximates the true $U(x, t)$ for $x = k \cdot \Delta x, t = nh, 1 \leq k \leq s, n \geq 1$.

We use the *maximum norm* $|\cdot| = |\cdot|_\infty$ on \mathbb{R}^s . A straightforward calculation reveals that (cf. (2.3))

$$\tau_1[A] \subset V(\lambda_0, \rho_0)$$

with $\lambda_0 = 137000, \rho_0 = 2(\Delta x)^{-2}$, provided Δx is so small that $100 \cdot \Delta x \leq 1$.

Let λ, ρ be such that the stability region S of φ satisfies $V(\lambda, \rho) \subset S$. Then the condition $h_0 \tau_1[A] \subset V$ of Theorems 3.4 and 3.6 is fulfilled with $V = V(\lambda, \rho)$ if h_0 is chosen so small that

$$h_0 V(\lambda_0, \rho_0) \subset V(\lambda, \rho).$$

One easily sees that for $h_0 = \min\{\lambda/137000, \rho(\Delta x)^2/2\}$ the latter inclusion holds. Theorems 3.4 and 3.6 (provided (3.10) is fulfilled) thus apply to stepsizes h satisfying

$$(4.4) \quad h \leq \min \left\{ \frac{\lambda}{137000}, \frac{\rho \cdot (\Delta x)^2}{2} \right\}.$$

We now focus on functions φ of type

$$(4.5) \quad \varphi(\zeta) = 1 + \zeta + 0.5\zeta^2 + c\zeta^3,$$

so that the right-hand member of the stepsize restriction (4.4) depends on λ, ρ corresponding to this φ . Therefore, one might try to choose c so as to maximize the right-hand member of (4.4).

In [12, pp. 92, 93] it is shown that the set $V(\lambda, 0)$ is maximal for $c = 0.0625$ with optimal $\lambda \simeq 6.26$. However, for these values of c, λ there exists no positive ρ such that $V(\lambda, \rho) \subset S$. Modifying the value 0.0625 only slightly, one can arrive at

$$c = 0.0645 \quad \text{with } \lambda = 4.67, \quad \rho = 0.68.$$

In the following we confine our discussion first to these values.

It can be verified that (3.10.a), (3.10.b) are fulfilled with $V = V(\lambda, \rho)$. Therefore, Theorems 3.4 and 3.6 show that (1.2) holds with the maximum norm and both with $p = 0, q = 1$ and with $p = 1, q = 0$, while γ can be chosen independently of the individual $h, \Delta x$ satisfying (4.4).

For $\Delta x = 1/100, h = 3.4 \times 10^{-5}$ we have equality in (4.4), so that the error propagation should still be mild. Table 3 displays a numerical experiment, which is in agreement with this prediction.

TABLE 3
Some values of γ_n for $\Delta x = 10^{-2}, h = 3.4 \times 10^{-5}$, and $c = 0.0645$

n	1	2	4	8	16	32	64	128	256
γ_n	1.24	1.25	1.31	1.34	1.28	1.07	0.471	0.0178	0.3×10^{-6}

In this table, γ_n has a similar meaning as in §4.2.

With the same $A, \Delta x, h$ as above, but with $c = 0.0625$, we have $\lambda \simeq 6.26, \rho = 0$, so that our stepsize restriction (4.4) is violated. An easy calculation shows that the eigenvalues of A are different from each other and real with $\sigma[A] \subset (-15.7 \times 10^4, -2 \times 10^4)$. It can be verified that, with $h = 3.4 \times 10^{-5}, c = 0.0625$, the stepsize restriction (1.3.b) is fulfilled. Consequently, (1.2) holds with $p = q = 0$ for some $\gamma > 0$. Table 4 displays the actual stability behavior of the numerical process that is now under consideration.

TABLE 4
Some values of γ_n for $\Delta x = 10^{-2}, h = 3.4 \times 10^{-5}$, and $c = 0.0625$

n	1	2	8	32	128	512	2048	8192
γ_n	1.4	1.6	3.4	47	1.1×10^5	1.7×10^9	1.4×10^{11}	3.1×10^{11}

Here, γ_n has the same meaning as in the previous table, the only difference being that now $c = 0.0625$ instead of $c = 0.0645$.

The above two tables nicely illustrate the superiority of the stepsize restrictions along the lines of §3 over those based on (1.3).

BIBLIOGRAPHY

1. F. F. Bonsall and J. Duncan, *Numerical ranges of operators on normed spaces and of elements of normed algebras*, Cambridge Univ. Press, Cambridge, New York, 1971.
2. —, *Numerical ranges. II*, Cambridge Univ. Press, Cambridge, New York, 1973.
3. —, *Numerical ranges*, Studies in Functional Analysis (R. G. Bartle, ed.), Math. Assoc. Amer., 1980, pp. 1–49.
4. P. Brenner and V. Thomée, *On rational approximations of semigroups*, SIAM J. Numer. Anal. **16** (1979), 683–694.
5. R. Bulirsch and J. Stoer, *Introduction to numerical analysis*, Springer-Verlag, New York, Heidelberg, Berlin, 1980.
6. M. Crouzeix, *On multistep approximation of semigroups in Banach spaces*, J. Comput. Appl. Math. **20** (1987), 25–36.
7. M. Crouzeix and P. A. Raviart, *Approximation d'équations d'évolution linéaires par des méthodes multistep*, Etude Numérique des Grands Systèmes (Rencontres IRIA-Novosibirsk 1976), Dunod, Paris, 1976.
8. G. Dahlquist, *G-stability is equivalent to A-stability*, BIT **18** (1978), 384–401.
9. N. Dunford and J. T. Schwartz, *Linear operators. I*, Interscience, New York, London, 1958.
10. I. M. Glazman and J. I. Ljubič, *Finite-dimensional linear analysis: A systematic presentation in problem form*, MIT Press, Cambridge, MA, and London, 1974.
11. D. F. Griffiths, I. Christie, and A. R. Mitchell, *Analysis of error growth for explicit difference schemes in conduction-convection problems*, Internat. J. Numer. Methods Engrg. **15** (1980), 1075–1081.
12. P. J. Van der Houwen, *Construction of integration formulas for initial value problems*, North-Holland, Amsterdam, New York, Oxford, 1977.
13. P. J. Kelly and M. L. Weisz, *Geometry and convexity*, Wiley, New York, Chichester, Brisbane, Toronto, 1979.
14. J. F. B. M. Kraaijevanger, Private communication, 1986.
15. J. F. B. M. Kraaijevanger, H. W. J. Lenferink, and M. N. Spijker, *Stepsize restrictions for stability in the numerical solution of ordinary and partial differential equations*, J. Comput. Appl. Math. **20** (1987), 67–81.
16. H. W. J. Lenferink and M. N. Spijker, *A generalization of the numerical range of a matrix*, Linear Algebra Appl. **140** (1990), 251–266.
17. —, *On a generalization of the resolvent condition in the Kreiss matrix theorem*, Math. Comp. **57** (1991), 211–220.
18. K. W. Morton, *Stability of finite difference approximations to a diffusion-convection equation*, Internat. J. Numer. Methods Engrg. **15** (1980), 677–683.
19. O. Nevanlinna, *On the numerical integration of nonlinear initial value problems by linear multistep methods*, BIT **17** (1977), 58–71.
20. —, *Remarks on time discretization of contraction semigroups*, Report-HTKK-MAT-A225, Helsinki Univ. Techn., Inst. Math., 1984.
21. N. Nirschl and H. Schneider, *The Bauer fields of values of a matrix*, Numer. Math. **6** (1964), 355–365.
22. S. V. Parter, *Stability, convergence, and pseudo-stability of finite-difference equations for an over-determined problem*, Numer. Math. **4** (1962), 277–292.
23. S. C. Reddy and L. N. Trefethen, *Lax-stability of fully discrete spectral methods via stability regions and pseudo-eigenvalues*, Comput. Methods. Appl. Mech. Engrg. **80** (1990), 147–164.

24. R. T. Rockafellar, *Convex analysis*, Princeton Univ. Press, Princeton, NJ, 1970.
25. M. N. Le Roux, *Semidiscretization in time for parabolic problems*, *Math. Comp.* **33** (1979), 919–931.
26. G. Söderlind, *Bounds on nonlinear operators in finite-dimensional Banach spaces*, *Numer. Math.* **50** (1986), 27–44.
27. M. N. Spijker, *Stepsize restrictions for stability of one-step methods in the numerical solution of initial value problems*, *Math. Comp.* **45** (1985), 377–392.
28. L. N. Trefethen, *Lax-stability vs. eigenvalue stability of spectral methods*, *Numerical Methods in Fluid Dynamics. III* (K. W. Morton and M. J. Baines, eds.), Clarendon Press, Oxford, 1988, pp. 237–253.

DEPARTMENT OF MATHEMATICS AND COMPUTER SCIENCE, UNIVERSITY OF LEIDEN, P. O. BOX 9512, 2300 RA LEIDEN, THE NETHERLANDS