

Lecture Notes in Mathematics

Edited by A. Dold, Heidelberg and B. Eckmann, Zürich

395

Numerische Behandlung nichtlinearer Integrodifferential- und Differentialgleichungen

Vorträge einer Tagung im Mathematischen
Forschungsinstitut Oberwolfach,
2. 12. – 7. 12. 1973



Springer-Verlag
Berlin · Heidelberg · New York

TWO-SIDED ERROR ESTIMATES IN THE NUMERICAL
SOLUTION OF INITIAL VALUE PROBLEMS

by

M.N. Spijker

1. INTRODUCTION

Let u denote the solution of a finite-difference equation approximating a given differential equation. If the finite-difference equation is perturbed by a quantity w , e.g. due to round-off error, then instead of u we obtain a solution, say \tilde{u} . It is an interesting problem to find an upperbound for $||\tilde{u}-u||$ in terms of w (with $||-\|$ we denote some seminorm). Such an error estimate, which is of the form say $||\tilde{u}-u|| \leq R(w)$, is of interest only if it doesn't overestimate $||\tilde{u}-u||$ too much. In fact, we have to require that also $||\tilde{u}-u|| \geq \beta \cdot R(w)$ for some constant β which is essentially greater than zero - and if possible is approximately equal to one. Thus we arrive at the task to derive estimates of the form

$$L(w) \leq ||\tilde{u}-u|| \leq R(w)$$

where the left-hand member $L(w)$ and the right-hand member $R(w)$ are identical, except for a finite factor $\beta > 0$. An estimate of this form will be called a *two-sided error estimate*.

In this paper there will be presented a condition on the seminorm $||-\|$ which is necessary and sufficient for the existence of such a two-sided error estimate (see chapter 3). Further we shall show (in chapter 4) that in case a two-sided error estimate exists, it is possible to determine the order of magnitude of both $\sup_w ||\tilde{u}-u||$ and $\inf_w ||\tilde{u}-u||$ where w ranges over all perturbations whose seminorm $||w||$ equals an arbitrary given value $\epsilon > 0$.

We shall restrict our considerations to finite-difference methods for approximating solutions of initial value problems for ordinary and partial semi-linear differential equations. The theorems of this article will be illustrated by an application to the numerical solution of a nonlinear parabolic initial-boundary

value problem (see chapter 5). For further applications and generalizations the reader is referred to the publications [2] - [6] listed at the end of this paper.

2. NOTATIONS AND ASSUMPTIONS

2.1. The finite-difference scheme we shall deal with, is of the form

$$(1) \quad u_n = K(u_{n-1}) + h \cdot F(u_{n-1}) \quad (n=1,2,\dots,N), \quad u_0 = c,$$

and the perturbed finite-difference scheme is assumed to be of the form

$$(2) \quad \begin{aligned} \tilde{u}_n &= K(\tilde{u}_{n-1}) + h \cdot [F(\tilde{u}_{n-1}) + w_n] \quad (n=1,2,\dots,N), \\ \tilde{u}_0 &= c + w_0. \end{aligned}$$

In (1), (2) we denote by h the so-called stepsize, which belongs to some set H . It is assumed that H is contained in some interval $(0, T]$ and that $\inf H = 0$. Further u_n, \tilde{u}_n, w_n and c are vectors belonging to the normed real vectorspace \mathcal{R}_h . The norm in \mathcal{R}_h is denoted by $|a| = |a|_h$ for $a \in \mathcal{R}_h$. For $h \in H$ we use the notation $t_n = nh$ and N is the unique integer with $t_N \leq T < t_{N+1}$. Clearly $N \rightarrow \infty$ if $h \rightarrow 0$. u_n is an approximation to the true solution of some given differential equation at $t = t_n$ and c is obtained from the initial condition of the original infinitesimal problem.

An example of (1) is provided by Euler's method for solving an initial value problem for an autonomous system of ordinary differential equations which, by using vector notation, can be written in the form $\frac{d}{dt} U(t) = f(U(t))$ ($0 \leq t \leq T$), $U(0) = c$. In this example we have $F = f$ and K equals the identity. Further \mathcal{R}_h is a finite-dimensional vectorspace the dimension of which is equal to the dimension of the given system of differential equations. In this case H can be chosen equal to $H = (0, T]$ and u_n is an approximation of $U(t)$ at $t = t_n$. In this example the perturbation w_0 in (2) may represent the error occurring if the initial value c cannot be represented exactly in the computer and the error w_n ($n \geq 1$) in (2) may be due to the fact that, as a consequence of the possibly complicated structure of the function f , the value $F(\tilde{u}_{n-1}) = f(\tilde{u}_{n-1})$ is not computed exactly. For another example of (1), (2), in which K is different from the

identity operator, we refer to chapter 5.

Returning to the general scheme (1) we now shall list what conditions K and F are supposed to fulfill in the rest of this paper.

a. Assumptions on K :

1. For each $h \in H$ we denote by $K = K_h$ a linear mapping from \mathfrak{R}_h into \mathfrak{R}_h ,
2. There is a constant $\alpha > 0$ such that the norm of the n -th power of K_h satisfies $|(K_h)^n| \leq \alpha$ for all $h \in H$ and integers n with $0 \leq nh \leq T$,
3. For each $h \in H$ there is an eigenvector e_h in \mathfrak{R}_h with $|e_h| = 1$, $K(e_h) = \kappa e_h$ such that the corresponding eigenvalue $\kappa = \kappa_h$ satisfies the inequality $\kappa_h \geq 1 - \mu_1 h \geq \mu_0$ for some fixed positive constants μ_0 and μ_1 .

b. Assumption on F :

For each $h \in H$ we denote by $F = F_h$ a mapping from \mathfrak{R}_h into \mathfrak{R}_h satisfying a Lipschitz condition

$$|F(\tilde{a}) - F(a)| \leq \lambda \cdot |\tilde{a} - a|$$

where λ is an arbitrary positive constant independent of $h \in H$ and $\tilde{a}, a \in \mathfrak{R}_h$.

2.2. In order to investigate estimates for the errors $\tilde{u}_n - u_n$, resulting from the local perturbations w_n occurring in (2), it is appropriate to introduce the vectors

$$u = (u_0, u_1, \dots, u_N), \quad \tilde{u} = (\tilde{u}_0, \tilde{u}_1, \dots, \tilde{u}_N), \quad w = (w_0, w_1, \dots, w_N).$$

These vectors belong to the vectorspace X_h given by

$$X_h = \{x | x = (x_0, x_1, \dots, x_N) \text{ and each } x_n \in \mathfrak{R}_h\},$$

in which addition and multiplication with real numbers are defined coordinate-wise.

Since we want to measure the difference between \tilde{u} and u by means of a seminorm, we assume some seminorm $\|x\|_h$ to be given in X_h . It is assumed that this seminorm is *absolute*, i.e. it is required that $\|x\|_h = \|y\|_h$ for any pair of vectors

$x = (x_0, x_1, \dots, x_N)$, $y = (y_0, y_1, \dots, y_N)$ the coordinates of which satisfy

$$|x_n| = |y_n| \quad (n=0, 1, \dots, N).$$

Example 1. $\|x\|_h = |x_N|$,

Example 2. $\|x\|_h = \left(h \sum_0^N |x_n|^2 \right)^{\frac{1}{2}}$,

Example 3. $\|x\|_h = \max \{ |x_n| : 0 \leq n \leq N \}$.

The following definition formalizes the concept of a two-sided error estimate discussed in the introduction.

DEFINITION. Let γ_0 and γ_1 be positive constants and assume ϕ_h is a mapping from X_h into the set of real numbers, R . If for all $h \in H$ and all $w_h \in \mathcal{R}_h$ the relations (1), (2) imply that

$$(3) \quad \gamma_0 \cdot \phi_h[w] \leq \|\tilde{u}-u\|_h \leq \gamma_1 \cdot \phi_h[w],$$

then (3) is called a *two-sided error estimate* for the finite-difference scheme (1).

Note that (3) trivially holds with $\gamma_0 = \gamma_1 = 1$, $\phi_h[w] \equiv \|\tilde{u}-u\|_h$. It is clear that with this choice of the functional ϕ_h the error estimate is useless since the values $\phi_h[w]$ depend on w in a way which, due to the possibly complicated structure of the (nonlinear) F , is untransparent. It is in view of the existence of such trivial and simultaneously useless error estimates that we shall look for estimates of type (3) with a functional ϕ_h that is *independent of F* .

3. THE EXISTENCE OF TWO-SIDED ERROR ESTIMATES

In order to formulate our criterion for the existence of two-sided error estimates in a concise way, we introduce the *summation operator* S from X_h into X_h . For $x = (x_0, x_1, \dots, x_N) \in X_h$ we define $y = Sx$ by $y = (y_0, y_1, \dots, y_N)$ with $y_0 = 0$, $y_n = h \cdot (x_0 + x_1 + \dots + x_{n-1})$ (for $1 \leq n \leq N$).

THEOREM 1 (*The existence of two-sided error estimates*). The following propositions (i) and (ii) are equivalent:

- (i) There exists a two-sided error estimate for (1) with a functional ϕ_h independent of F ,
- (ii) There is a constant $\delta < \infty$ such that for all $h \in H$, $x \in X_h$ we have the inequality $\|Sx\|_h \leq \delta \cdot \|x\|_h$.

It should be understood that in statement (i) we deal with a *fixed* family $\{K_h\}_{h \in H}$ and a *variable* family $\{F_h\}_{h \in H}$, both families satisfying the as-

assumptions stated in chapter 2. Thus statement (i) could have been formulated alternatively as follows: (i*) There exists a family of functionals $\{\phi_h\}_{h \in H}$ such that for each family $\{F_h\}_{h \in H}$ satisfying a Lipschitz condition of the type described in chapter 2, there exist positive constants γ_0 and γ_1 such that (1), (2) always imply (3) (γ_0 and γ_1 may depend on $\{F_h\}_{h \in H}$ but not on h and w).

It is easily verified that with the seminorm $\|x\|_h$ of example 1 (chapter 2) condition (ii) of the above theorem is violated. On the other hand the seminorms of the examples 2, 3 satisfy condition (ii) with $\delta = T$.

THEOREM 2 (*The form of two-sided error estimates*). Let condition (ii) of theorem 1 be fulfilled. Then the two-sided error estimate (3) holds with

$$\begin{aligned} \gamma_0 &= (1 + \alpha \delta \lambda)^{-1}, \quad \gamma_1 = 1 + \alpha \delta \lambda \cdot \exp(\alpha \lambda T), \\ \phi_h[w] &= \|v\|_h \quad \text{where } v = (v_0, v_1, \dots, v_N) \text{ and} \\ v_n &= K^n w_0 + h \cdot (w_n + K w_{n-1} + \dots + K^{n-1} w_1). \end{aligned}$$

For the proof of the theorems 1,2 we refer to chapter 3 of [5]. Note that already in the simple case where (1) stands for Euler's method for solving an ordinary differential equation, the theorems 1,2 are nontrivial.

4. APPLICATIONS OF TWO-SIDED ERROR ESTIMATES

4.1. Two sided error estimates have applications in several different fields such as the propagation of rounding errors, the propagation of local discretization errors and in the proof of so-called equivalence theorems stating necessary and sufficient conditions for convergence (cf. [2]-[6]).

In this chapter we shall confine ourselves to showing that two-sided error estimates may be used with success in determining the orders of magnitude of the interesting functions $g_0(\epsilon, h)$, $g_1(\epsilon, h)$ defined by

$$g_0(\epsilon, h) = \inf_{\|w\|_h = \epsilon} \|\tilde{u} - u\|_h, \quad g_1(\epsilon, h) = \sup_{\|w\|_h = \epsilon} \|\tilde{u} - u\|_h$$

where $\epsilon > 0$, $h \in H$. Note that, in view of (2), \tilde{u} is uniquely determined by w , and the seminorm $\|\tilde{u} - u\|_h$ appearing in the above definitions thus indeed is a

function of the variable vector $w \in X_h$.

For the sake of simplicity we only deal with the seminorm of example 2 (section 2), i.e.

$$(4) \quad ||x||_h = \left(h \sum_0^N |x_n|^2 \right)^{\frac{1}{2}}$$

and we only consider perturbations $w = (w_0, w_1, \dots, w_N)$ the first component of which vanishes.

$$(5) \quad w_0 = 0.$$

In the subsequent sections we use the above definitions of g_0, g_1 with the restrictions expressed in (4), (5). Further we consider a fixed family $\{K_h\}_{h \in H}$ and a fixed family $\{F_h\}_{h \in H}$, both families satisfying the conditions stated in section 2.

4.2. We first give a lower bound for g_0 and an upperbound for g_1 :

$$(6) \quad (1+\alpha+\lambda h)^{-1} \cdot h \cdot \epsilon \leq g_0(\epsilon, h),$$

$$(7) \quad g_1(\epsilon, h) \leq (2\lambda)^{-1} [\exp(2\alpha\lambda T) - 1 - 2\alpha\lambda T \cdot \exp(-\alpha\lambda h)]^{\frac{1}{2}} \cdot \epsilon.$$

Proof of (6). Using (4) with $x=w$ and substituting for w_n the expression which can be obtained by subtracting (1) from (2) we have, in case $w_0=0$,

$$||w||_h = \left(h \sum_1^N h^{-2} |(\tilde{u}_n - u_n) - K(\tilde{u}_{n-1} - u_{n-1}) - h[F(\tilde{u}_{n-1}) - F(u_{n-1})]|^2 \right)^{\frac{1}{2}}.$$

In view of the assumptions on K and F we thus obtain

$$\begin{aligned} h ||w||_h &\leq \left(h \sum_1^N [|\tilde{u}_n - u_n| + (\alpha + \lambda h) |\tilde{u}_{n-1} - u_{n-1}|]^2 \right)^{\frac{1}{2}} \\ &\leq \left(h \sum_1^N |\tilde{u}_n - u_n|^2 \right)^{\frac{1}{2}} + (\alpha + \lambda h) \cdot \left(h \sum_1^N |\tilde{u}_{n-1} - u_{n-1}|^2 \right)^{\frac{1}{2}} \\ &\leq (1 + \alpha + \lambda h) ||\tilde{u} - u||_h, \end{aligned}$$

which proves (6).

Proof of (7). The upperbound (7) can be proved in a standard fashion by subtracting (1) from (2), which yields for the difference

$$d_n = \tilde{u}_n - u_n$$

the recurrence relation

$$d_n = K d_{n-1} + h[F(\tilde{u}_{n-1}) - F(u_{n-1})] + h w_n .$$

Solving for d_n we obtain (provided $\tilde{u}_0 - u_0 = w_0 = 0$) the representation

$$d_n = z_n + v_n$$

where v_n is defined in theorem 2 and $z_n = h \sum_{i=1}^{n-1} K^{n-1-i} [F(\tilde{u}_i) - F(u_i)]$

Using the assumptions on K and F stated in chapter 2 we get

$$|z_n| \leq \alpha h \sum_{i=1}^{n-1} |d_i| , \quad |v_n| \leq \alpha h \sum_{i=1}^n |w_i| .$$

Consequently

$$|d_n| \leq \alpha h \sum_{i=1}^{n-1} |d_i| + \alpha h \sum_{i=1}^n |w_i| ,$$

from which it follows by induction with respect to n that

$$|d_n| \leq \alpha h \sum_{i=1}^n (1 + \alpha h)^{n-i} |w_i| \quad (n=1, 2, \dots, N) .$$

Using the Schwarz inequality we thus obtain

$$\begin{aligned} (||d||_h)^2 &\leq h \sum_{n=1}^N \alpha^2 \left[h \sum_{i=1}^n (1 + \alpha h)^{n-i} |w_i| \right]^2 \\ &\leq \alpha^2 h \sum_{n=1}^N \left[h \sum_{i=1}^n (1 + \alpha h)^{2(n-i)} \right] \cdot \left[h \sum_{i=1}^n |w_i|^2 \right] \\ &\leq (\alpha ||w||_h)^2 \cdot h^2 \sum_{n=1}^N \sum_{i=0}^{n-1} (1 + \alpha h)^{2i} . \end{aligned}$$

Since $\sum_{n=1}^N \sum_{i=0}^{n-1} (1 + \alpha h)^{2i} \leq \sum_{n=1}^N \sum_{i=0}^{n-1} e^{2\alpha h i}$

$$= \sum_{n=1}^N (e^{2\alpha h n} - 1) / (e^{2\alpha h} - 1)$$

$$= (e^{2\alpha h N} - 1)^{-2} \cdot [e^{2\alpha h} (e^{2\alpha h N} - 1) - (e^{2\alpha h} - 1)N]$$

and $e^{2t} - 1 \geq 2t e^t$ (for $t = \alpha h$) we get

$$(||d||_h)^2 \leq (||w||_h)^2 \cdot (2\lambda)^{-2} \cdot [e^{2\alpha \lambda N h} - 1 - 2\alpha \lambda N h e^{-\alpha \lambda h}] .$$

In view of $Nh = t_N \leq T$ and the definition of d_n we thus have

$$||\tilde{u} - u||_h \leq ||w||_h \cdot (2\lambda)^{-1} \cdot [e^{2\alpha \lambda T} - 1 - 2\alpha \lambda T \cdot e^{-\alpha \lambda h}]^{\frac{1}{2}} ,$$

which completes the proof of (7).

4.3. Next we present a theorem showing that the lower and upperbound

(6), (7) are realistic, as $h \rightarrow 0$, at least in the sense that they contain the correct powers of h .

THEOREM 3. (a) For each $\epsilon > 0$ and $h \in H$ there exists a vector $w = (w_0, w_1, \dots, w_N)$ with $w_0 = 0$, $\|w\|_h = \epsilon$ such that \tilde{u}_n and u_n , computed from (2), (1), respectively satisfy

$$\|\tilde{u} - u\|_h \leq \frac{1+\alpha}{1+\mu_0} \cdot [1+\alpha\lambda T \exp(\alpha\lambda T)] \cdot h \cdot \|w\|_h.$$

(b) For each $\epsilon > 0$ and $h \in H$ there also exists a vector $w = (w_0, w_1, \dots, w_N)$ with $w_0 = 0$, $\|w\|_h = \epsilon$ such that \tilde{u}_n and u_n computed from (2), (1) satisfy

$$\|\tilde{u} - u\|_h \geq \frac{1}{\mu_1(1+\alpha\lambda T)} \cdot \left\{ \frac{1}{T} \int_0^T [1 - \exp(-\mu_1 t)]^2 dt \right\}^{\frac{1}{2}} \cdot \|w\|_h.$$

It follows from part (a) of this theorem that

$$(6^*) \quad g_0(\epsilon, h) \leq \frac{1+\alpha}{1+\mu_0} \cdot [1+\alpha\lambda T \exp(\alpha\lambda T)] \cdot h \cdot \epsilon$$

and from part (b) there follows

$$(7^*) \quad g_1(\epsilon, h) \geq \frac{1}{\mu_1(1+\alpha\lambda T)} \cdot \left\{ \frac{1}{T} \int_0^T [1 - \exp(-\mu_1 t)]^2 dt \right\}^{\frac{1}{2}} \cdot \epsilon.$$

Formula (7*) shows that (7) cannot essentially be improved since apparently it is not possible to replace (7) by an estimate of the form $g_1(\epsilon, h) \leq \gamma \cdot h^p \cdot \epsilon$ with any $p > 0$. Consequently, the order of magnitude of $g_1(\epsilon, h)$ is exactly equal to ϵ . Similarly, (6), (6*) prove that the order of magnitude of $g_0(\epsilon, h)$ is exactly equal to $h \cdot \epsilon$. In fact, according to a definition of a two-sided error estimate differing only slightly from the one given in chapter 2, one might say that (6), (6*) and (7), (7*) provide two-sided error estimates of $g_0(\epsilon, h)$ and $g_1(\epsilon, h)$, respectively.

In the above proofs of (6), (7) we have not made use of the two-sided error estimate for $\|\tilde{u} - u\|_h$ which exists in view of theorem 1. On the other hand, the following proof of theorem 3 is entirely based on this two-sided error estimate.

Proof of theorem 3. (a) Choosing $w_0 = 0$ and $w_n = [(-1)^n \cdot T^{-\frac{1}{2}} \epsilon] \cdot e_n$ (for $n=1, 2, \dots, N$) where e_n is the eigenvector whose existence was postulated in chapter 2, we have for the vectors v_n defined in theorem 2

$$v_n = (-1)^n T^{-\frac{1}{2}} \epsilon \cdot h \cdot [1 - \kappa + \kappa^2 - \dots + (-1)^{n-1} \kappa^{n-1}] \cdot e_n.$$

Consequently

$$|v_n| = T^{-\frac{1}{2}} \text{ch} \cdot [1 - (-\kappa)^n] / [1 + \kappa] .$$

Since $(\kappa)^n \leq |(K_h)^n| \leq \alpha$ and $\kappa \geq \mu_0$ (cf. chapter 2) we thus obtain

$$|v_n| \leq T^{-\frac{1}{2}} \text{ch} \cdot (1 + \alpha) / (1 + \mu_0) .$$

Hence $\|v\|_h \leq \text{ch} \cdot (1 + \alpha) / (1 + \mu_0)$. By virtue of theorem 2 we have

$$\|\tilde{u} - u\|_h \leq [1 + \alpha \lambda T \cdot \exp(\alpha \lambda T)] \cdot \|v\|_h , \text{ from which it follows that}$$

$$\|\tilde{u} - u\|_h \leq [1 + \alpha \lambda T \cdot \exp(\alpha \lambda T)] \cdot (1 + \alpha) (1 + \mu_0)^{-1} \cdot h \epsilon .$$

Since $\|w\|_h = \epsilon$ part (a) of the theorem has thus been proved.

(b) In order to prove part (b) we choose the perturbations $w_0 = 0$, $w_n = T^{-\frac{1}{2}} \epsilon \cdot e_n$ ($n=1, 2, \dots, N$) the seminorm $\|w\|_h$ of which again equals ϵ .

We have

$$\begin{aligned} |v_n| &= T^{-\frac{1}{2}} \text{ch} \cdot (1 + \kappa + \dots + \kappa^{n-1}) \\ &\geq T^{-\frac{1}{2}} \text{ch} \cdot [1 + (1 - \mu_1 h) + \dots + (1 - \mu_1 h)^{n-1}] \\ &= T^{-\frac{1}{2}} \epsilon (\mu_1)^{-1} \cdot [1 - (1 - \mu_1 h)^n] \\ &\geq T^{-\frac{1}{2}} \epsilon (\mu_1)^{-1} \cdot [1 - \exp(-\mu_1 n h)] . \end{aligned}$$

Consequently $(\|v\|_h)^2 \geq (\frac{\epsilon}{\mu_1})^2 \cdot T^{-1} \cdot h \sum_{n=1}^N [1 - \exp(-\mu_1 n h)]^2$
 $\geq (\frac{\epsilon}{\mu_1})^2 \cdot \frac{1}{T} \cdot \int_0^T [1 - \exp(-\mu_1 t)]^2 dt$. In view of theorem 2 we have

$\|\tilde{u} - u\|_h \geq (1 + \alpha \lambda T)^{-1} \cdot \|v\|_h$ from which it thus follows that part (b) of theorem 3 holds. This completes the proof of the theorem.

5. A PARABOLIC INITIAL-BOUNDARY VALUE PROBLEM

5.1. In order to illustrate the theory of the preceding chapters we consider the numerical solution by a very simple explicit finite-difference scheme of the initial-boundary value problem

$$\frac{\partial}{\partial t} U(s, t) = \frac{\partial^2}{\partial s^2} U(s, t) + f(s, U(s, t)) ,$$

$$(8) \quad U(0, t) = U(1, t) = 0 , \quad U(s, 0) = c(s) ,$$

$$0 \leq s \leq 1 , \quad 0 \leq t \leq T .$$

In (8) f and c denote given functions and f is assumed to satisfy the Lipschitz condition

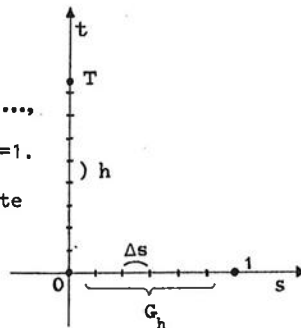
$$(9) \quad |f(s, \tilde{r}) - f(s, r)| \leq \lambda \cdot |\tilde{r} - r| ,$$

where λ is an arbitrary positive constant, uniformly for $0 \leq s \leq 1$ and $-\infty < r, \tilde{r} < \infty$. The finite-difference scheme we deal with, is of the form

$$(10) \quad \frac{1}{h} \cdot [u_n(s) - u_{n-1}(s)] = \frac{1}{(\Delta s)^2} \cdot [u_{n-1}(s - \Delta s) - 2u_{n-1}(s) + u_{n-1}(s + \Delta s)] \\ + f(s, u_{n-1}(s)) ,$$

$$u_0(s) = c(s) .$$

In (10) the variable s takes on the values $s = \Delta s, 2\Delta s, \dots, M\Delta s$ where M is an arbitrary integer ≥ 2 with $(M+1)\Delta s = 1$. Applying (10) successively with $n=1, 2, \dots, N$ we can compute approximations $u_n(s)$ of $U(s, t)$ at the points $(s, t) = (m\Delta s, nh)$ with $m=1, 2, \dots, M$ and $n=1, 2, \dots, N$.



5.2. We give some definitions which will enable us to write the finite-difference scheme (10) in the general form (1). We assume throughout that the ratio $h/(\Delta s)^2$ equals some constant σ with

$$(11) \quad 0 < \sigma \leq 1/2 .$$

The set H of stepsizes h is defined by

$$H = \{h \mid h = \sigma \cdot (\Delta s)^2, \Delta s = (M+1)^{-1} \text{ with } M=2, 3, 4, \dots\}$$

and for $h \in H$ the grid G_h on the s -axis is defined by

$$G_h = \{s \mid s = m \cdot \Delta s \text{ with } m=1, 2, \dots, M\} .$$

The vectorspace \mathcal{K}_h consists of all real functions defined on the grid G_h and the norm $|a|_h$ in \mathcal{K}_h is defined by

$$(12) \quad |a|_h = \left\{ \Delta s \sum_{m=1}^M [a(m \cdot \Delta s)]^2 \right\}^{1/2}$$

for $a \in \mathcal{K}_h$. The operators K and F mapping \mathcal{K}_h into itself are defined by

$$(Ka)(s) = \begin{cases} (1-2\sigma) \cdot a(s) + \sigma \cdot a(s + \Delta s) & (\text{for } s = \Delta s) , \\ \sigma \cdot a(s - \Delta s) + (1-2\sigma) \cdot a(s) + \sigma \cdot a(s + \Delta s) & (\text{for } s = m\Delta s, 1 < m < M) , \\ \sigma \cdot a(s - \Delta s) + (1-2\sigma) \cdot a(s) & (\text{for } s = M\Delta s) , \end{cases}$$

$$(Fa)(s) = f(s, a(s)) \quad (\text{for } s = m\Delta s, \quad 1 \leq m \leq M)$$

where a denotes an arbitrary element belonging to \mathcal{K}_h . With these definitions (1) is equivalent to (10) and (2) is equivalent to

$$(13) \quad \frac{1}{h} [\tilde{u}_n(s) - \tilde{u}_{n-1}(s)] = \frac{1}{(\Delta s)^2} [\tilde{u}_{n-1}(s - \Delta s) - 2\tilde{u}_{n-1}(s) + \tilde{u}_{n-1}(s + \Delta s)] \\ + f(s, \tilde{u}_{n-1}(s)) + w_n(s), \\ \tilde{u}_0(s) = c(s) + w_0(s)$$

where $w_n(s)$ thus represents a perturbation occurring when the approximation $\tilde{u}_n(s)$ of $U(s, nh)$ is calculated.

5.3. It is easily verified, by using (9), that the assumption on F postulated in chapter 2 is satisfied here. Further, the first assumption on K stated in chapter 2 is satisfied here as well. The second assumption on K of chapter 2 may be proved by a standard argument (cf. Rjabenki, Filippow [1]) using the inner product

$$(a, b) = \Delta s \sum_{m=1}^M a(m\Delta s) \cdot b(m\Delta s)$$

in \mathcal{K}_h . It can be verified by a straightforward calculation that the elements $e^{(1)}, e^{(2)}, \dots, e^{(M)}$ in \mathcal{K}_h given by

$$e^{(j)}(s) = \sqrt{2} \cdot \sin(sj\pi) \quad (\text{for } s \in G_h)$$

are orthonormal with respect to this inner product and that they are eigenvectors of K with eigenvalues

$$\kappa^{(j)} = 1 - 4\sigma \cdot [\sin(j\Delta s \cdot \pi/2)]^2.$$

In view of (11) we have $-1 < \kappa^{(j)} < 1$ and consequently

$$|(K_h)^n| \leq |K_h|^n = [\max_j |\kappa^{(j)}|]^n < 1, \text{ which yields the value}$$

$$(14) \quad \alpha = 1.$$

The third assumption on K of chapter 2 is fulfilled here with $e_h = e^{(1)}$,

$$\kappa_h = \kappa^{(1)} \quad \text{and}$$

$$(15) \quad \mu_0 = 1 - (\pi/3)^2 \cdot \sigma, \quad \mu_1 = \pi^2.$$

5.4. All general assumptions on K and F of chapter 2 being fulfilled

we may apply the results of chapter 4 to the finite-difference scheme (10) and the perturbed scheme (13). The seminorm $||x||_h$ (see (4)) used throughout chapter 4 yields in combination with the norm $|a|_h$ (see (12)) the formula $||x||_h = \{h\Delta s \cdot \sum_{n=0}^N \sum_{m=1}^M [x_n(m\Delta s)]^2\}^{\frac{1}{2}}$ and if $x_0 = 0$ we thus have

$$(16) \quad ||x||_h = \{h\Delta s \cdot \sum_{n=1}^N \sum_{m=1}^M [x_n(m\Delta s)]^2\}^{\frac{1}{2}}.$$

We combine the results (6), (6*), (7), (7*) of chapter 4 when applied to the situation at hand (see (14), (15)) into the following theorem.

THEOREM 4. Let u be computed from (10) and consider \tilde{u} computed from (13) with arbitrary perturbation-vector w for which $w_0(s) \equiv 0$. Then, using the notation (16), we have for any $\epsilon > 0$, $h \in H$ the inequalities

$$(17) \quad \inf_{||w||_h = \epsilon} ||\tilde{u} - u||_h \geq \frac{h\epsilon}{2 + \lambda h},$$

$$(17^*) \quad \inf_{||w||_h = \epsilon} ||\tilde{u} - u||_h \leq \frac{18}{18 - \pi^2 \sigma} \cdot [1 + \lambda T \cdot \exp(\lambda T)] \cdot h\epsilon,$$

$$(18) \quad \sup_{||w||_h = \epsilon} ||\tilde{u} - u||_h \leq \frac{1}{2\lambda} [\exp(2\lambda T) - 1 - 2\lambda T \cdot \exp(-\lambda h)]^{\frac{1}{2}} \cdot \epsilon,$$

$$(18^*) \quad \sup_{||w||_h = \epsilon} ||\tilde{u} - u||_h \geq \frac{1}{\pi^2(1 + \lambda T)} \left(\frac{1}{T} \int_0^T [1 - \exp(-\pi^2 t)]^2 dt \right)^{\frac{1}{2}} \cdot \epsilon.$$

5.5. (17*) and (18*) show that the bounds (17), (18) cannot essentially be improved (e.g. by including a factor h^p with $p \neq 0$) and that the precise orders of magnitude of $\inf ||\tilde{u} - u||_h$ and $\sup ||\tilde{u} - u||_h$ are $h\epsilon$ and ϵ , respectively. In order also to test how realistic the error bounds (17), (18) are from a numerical point of view¹⁾ a number of experiments have been performed with the initial-boundary value problems

$$\frac{\partial}{\partial t} U(s, t) = \frac{\partial^2}{\partial s^2} U(s, t) + \frac{1}{1 + [U(s, t)]^2} + \frac{1 + 2s^2(1-s)^2}{1 + s^2(1-s)^2},$$

$$(A) \quad U(0, t) = U(1, t) = 0, \quad U(s, 0) = s(1-s),$$

$$0 \leq s \leq 1, \quad 0 \leq t \leq 1$$

and

¹⁾ Here we follow a suggestion made by L. Collatz at the Conference in Oberwolfach.

$$\frac{\partial}{\partial t} U(s,t) = \frac{\partial^2}{\partial s^2} U(s,t) + \frac{U(s,t)}{100},$$

(B) $U(0,t) = U(1,t) = 0, U(s,0) = \sin(\pi s),$
 $0 \leq s \leq 1, 0 \leq t \leq 1.$

In the following tables we compare the numerical values obtained for $||\tilde{u}-u||_h$ with the theoretical bounds (17), (18). In (13) we have used the perturbations $w_n(s)$ defined by $w_0(s) = 0$ and $w_n(s) = \frac{1}{10}(-1)^n(1-\Delta s)^{-\frac{1}{2}}$ (for problem (A)),
 $w_n(s) = \frac{1}{10}(1-\Delta s)^{-\frac{1}{2}}$ (for problem (B)).

$N=h^{-1}$	Theoretical lower bound for $ \tilde{u}-u _h$ (cf.(17)), $h\epsilon/(2+h)$	$ \tilde{u}-u _h$	Theoretical upper bound for $ \tilde{u}-u _h$ (cf.(18)), $\frac{\epsilon}{2} \cdot [\exp(2) - 1 - 2 \cdot \exp(-h)]^{\frac{1}{2}}$
48	0.00103	0.00124	0.10524
192	0.00026	0.00028	0.10487
768	0.00007	0.00007	0.10478
3072	0.00002	0.00002	0.10476

Problem (A), $\lambda=1, T=1, \sigma=1/3, ||w||_h=\epsilon=0.1, w_n(m\Delta s)=(-1)^n(1-\Delta s)^{-\frac{1}{2}}\epsilon$

$N=h^{-1}$	Theoretical lower bound for $ \tilde{u}-u _h$ (cf.(17)), $h\epsilon/(2+\lambda h)$	$ \tilde{u}-u _h$	Theoretical upper bound for $ \tilde{u}-u _h$ (cf.(18)), $\frac{\epsilon}{2\lambda} \cdot [\exp(2\lambda) - 1 - 2\lambda \cdot \exp(-\lambda h)]^{\frac{1}{2}}$
32	0.00156	0.00988	0.07204
128	0.00039	0.00904	0.07122
512	0.00010	0.00870	0.07102
2048	0.00002	0.00855	0.07096

Problem (B), $\lambda=0.01, T=1, \sigma=\frac{1}{2}, ||w||_h=\epsilon=0.1, w_n(m\Delta s)=(1-\Delta s)^{-\frac{1}{2}}\epsilon$

The numerical results show that there are perturbations $w_n(s)$ for which the bounds (17), (18) don't underestimate or overestimate the actual error $\|\tilde{u}-u\|_h$ very much. In particular the lower bound (17) fits in surprisingly well with the numerical results in the first table, referring to problem (A). Thus (17), (18), when applied to the problems (A), (B), respectively appear to be realistic bounds both from a theoretical and from a more practical point of view.

Acknowledgement. I wish to thank F. Bakker and C. den Heijer who have been most helpful in performing on the IBM 370/158 of Leiden University the computations reported in section 5.5.

REFERENCES

- [1] RJABENKI, V.S., A.F. FILIPPOV, Über die Stabilität von Differenzgleichungen. VEB Deutscher Verlag der Wissenschaften, Berlin (1960).
- [2] STETTER, H.J., Analysis of discretization methods for ordinary differential equations. Springer-Verlag, Berlin etc. (1973).
- [3] SPIJKER, M.N., On the structure of error estimates for finite-difference methods, Numer. Math. 18 (1971), 73-100.
- [4] SPIJKER, M.N., Optimum error estimates for finite-difference methods, to appear in Acta Universitatis Carolinae (1974).
- [5] SPIJKER, M.N., The existence of optimum error estimates in the numerical solution of differential equations, in preparation.
- [6] STUMMEL, F., Approximation methods in analysis, Lecture Notes Series of Aarhus Universitet, Matematisk Institut (1973).

