

# Contractivity in the Numerical Solution of Initial Value Problems

M.N. Spijker

Universiteit van Leiden, Wassenaarseweg 80, Postbus 9512, NL-2300 RA Leiden,  
The Netherlands

**Summary.** Consider a linear autonomous system of ordinary differential equations with the property that the norm  $|U(t)|$  of each solution  $U(t)$  satisfies  $|U(t)| \leq |U(0)|$  ( $t \geq 0$ ). We call a numerical process for solving such a system contractive if a discrete version of this property holds for the numerical approximations. A given  $k$ -step method is said to be unconditionally contractive if for each stepsize  $h > 0$  the numerical process is contractive.

In this paper a general theory is given which yields necessary and sufficient conditions for unconditional contractivity. It turns out that unconditionally contractive methods are subject to an order barrier  $p \leq 1$ . Further the concept of a contractivity threshold is studied, which makes it possible to compare the contractivity behaviour of methods with an order  $p > 1$  as well.

Most theoretical results in this paper are formulated for differential equations in arbitrary Banach spaces. Applications are given to numerical methods for solving ordinary as well as partial differential equations.

*Subject Classifications:* AMS(MOS): 65J10, 65L20, 65M10; CR: 5.17.

## 1. Introduction

### 1.1. Formulation of Two Problems

In order to introduce the subject treated in this article we consider an initial value problem for a linear system of  $s$  ordinary differential equations

$$(1.1) \quad \frac{d}{dt} U(t) = AU(t) (t \geq 0), \quad U(0) = u_0.$$

We assume here that

$$(1.2.a) \quad s \geq 1 \text{ and } A = (a_{ij}) \text{ is a real } s \times s \text{ matrix;}$$

(1.2.b)  $|\cdot|$  is an arbitrary fixed norm on the vector space  $\mathbb{R}^s$ ;

(1.2.c) for each  $u_0 \in \mathbb{R}^s$  the solution to (1.1) satisfies  $|U(t)| \leq |u_0|$  ( $t \geq 0$ ).

Many numerical methods for the solution of initial value problems (such as linear multistep methods, Runge-Kutta methods, Rosenbrock methods) result, when applied to problem (1.1), into a procedure of type

$$(1.3) \quad P_0(hA) u_n = P_1(hA) u_{n-1} + P_2(hA) u_{n-2} + \dots + P_k(hA) u_{n-k} \\ (n = k, k + 1, k + 2, \dots).$$

Here  $k \geq 1$  is a fixed integer, and  $u_n$  are approximations to  $U(t_n)$  for  $t_n = nh$ ,  $h > 0$ ,  $n \geq k$  that are obtained from (1.3) when starting vectors  $u_0, u_1, \dots, u_{k-1} \in \mathbb{R}^s$  are given. Further  $P_0(\zeta), P_1(\zeta), \dots, P_k(\zeta)$  stand for polynomials (with real coefficients) which have no common zero.

We assume  $P_0(0) \neq 0$ , and recall that the *order* of (1.3) is the largest integer  $p$  for which  $U(t_n)$  always satisfies (1.3) up to an error  $= \mathcal{O}(h^{p+1})$  (for  $h \rightarrow 0+$ ). We assume (1.3) to be of an order  $p \geq 0$ .

*Definition 1.1.* Method (1.3) is called *unconditionally contractive* with respect to the pair  $(A, |\cdot|)$  if for all  $h > 0$  and all  $u_n (n \geq 0)$  satisfying (1.3) we have

$$(1.4) \quad |u_n| \leq \max(|u_0|, |u_1|, \dots, |u_{k-1}|) \quad (n \geq k).$$

Clearly (1.4) is a discrete version of (1.2.c), and therefore unconditional contractivity is a *natural* demand upon (1.3). It is also a *useful* demand since (1.4) implies that any (rounding) errors, which might be present in the starting vectors, are propagated in the numerical calculations no stronger than lies in the nature of the original problem (1.1).

For the case where (1.3) stands for a *linear multistep method* (i.e. all  $P_i(\zeta)$  have a degree  $\leq 1$ ) Nevanlinna and Liniger [12, p. 59] derived simple conditions on  $P_i(\zeta)$  that are necessary and sufficient for (1.3) to be unconditionally contractive with respect to all pairs  $(A, |\cdot|)$  satisfying (1.2). They arrived at a nontrivial class of linear multistep methods with this favourable contractivity property. However, all members of this class were proved to have only an order  $p \leq 1$ .

The question arises what conditions on arbitrary polynomials  $P_i(\zeta)$  are necessary and sufficient for unconditional contractivity of the general method (1.3) with respect to all pairs  $(A, |\cdot|)$  satisfying (1.2). A second natural question is whether there exist any methods (1.3) with this nice contractivity property and an order  $p > 1$ .

### 1.2. Solution of the Two Problems

Finding answers to the above two questions, as well as to related ones, is the purpose of this article. In the subsequent two theorems we formulate already our solutions to the above two problems.

We recall that a scalar function  $\psi$  is called *absolutely monotonic* on the interval  $J \subset \mathbb{R}$  if  $\psi(\zeta)$  as well as all derivatives  $\psi^{(j)}(\zeta)$  ( $j = 1, 2, 3, \dots$ ) exist in  $\mathbb{R}$  and are  $\geq 0$  (for all  $\zeta \in J$ ).

For  $1 \leq i \leq k$  we define  $\psi_i(\zeta)$  to be the rational function (without removable singularities) with  $\psi_i(\zeta) = P_i(\zeta)/P_0(\zeta)$  (for  $\zeta \in \mathbb{R}$ ,  $P_0(\zeta) \neq 0$ ).

**Theorem 1.2.** *Method (1.3) is unconditionally contractive with respect to all pairs  $(A, |\cdot|)$  satisfying (1.2) if and only if all  $\psi_i(\zeta)$  are absolutely monotonic on the interval  $(-\infty, 0]$  (here  $1 \leq i \leq k$ ).*

**Theorem 1.3.** *Suppose method (1.3) has the contractivity property of Theorem 1.2. Then its order cannot exceed  $p = 1$ .*

### 1.3. Outline of the Rest of this Paper

After a few remarks (in Sect. 1.4) in connection with the above two theorems we give in Sect. 2.1 basic definitions needed in this paper. We consider linear initial value problems in arbitrary Banach spaces, thus generalizing the situation of (1.1), (1.2).

Section 2.2 contains the principal results of Chap. 2. The Theorems 2.4, 2.5 of this section can be regarded as generalizations of the above Theorems 1.2, 1.3, respectively.

Applications of the theory of Sect. 2.2 are presented in the next Sect. 2.3. Here we deduce the Theorems 1.2, 1.3 as well as contractivity results on Runge-Kutta methods already stated (without proof) in [17].

In view of the order barrier  $p \leq 1$ , appearing in Theorem 1.3 and related theorems of Chap. 2, we turn in Chap. 3 to a framework suitable for comparing also the contractivity properties of methods with  $p > 1$ .

In Sect. 3.1 we focus on initial value problems that are less general than those considered in Chap. 2. For this restricted class of problems we give in Sect. 3.2 contractivity results valid under a condition  $0 < h < H$  on the stepsize  $h$ .

Section 3.3 contains applications of the theory of Sect. 3.2 to linear multi-step methods and general one-step methods. We describe some methods with conditional contractivity properties when applied to our restricted class of initial value problems, and with an order  $p > 1$ .

At the end of the Chaps. 2, 3 we illustrate the theories of the Sects. 2.2, 3.2 by applying them to a numerical method for solving a simple *partial differential equation* of parabolic type.

Our last Chap. 4 is of a technical nature and contains the proofs of our main Theorems 2.4, 3.3 already formulated in the Sects. 2.2, 3.2.

### 1.4. Remarks

1. The contractivity property occurring in Theorem 1.2 concerns arbitrary pairs  $(A, |\cdot|)$  satisfying (1.2). If we would consider only pairs  $(A, |\cdot|)$  satisfying

the additional condition that the norm  $|\cdot|$  is generated by some inner product in  $\mathbb{R}^s$ , the order barrier  $p \leq 1$  of Theorem 1.3 would no longer be present (see [12] or [3], [4]). On the other hand one would certainly not gain the same insight into the general error propagation of numerical methods by considering only pairs  $(A, |\cdot|)$  subject to this additional condition, nor by considering e.g. only the test equation  $\frac{d}{dt} U(t) = \lambda U(t) \ (t \geq 0)$  with  $\lambda \in \mathbb{C}$  (cf. e.g. [17], [11, pp. 258-261]).

2. Suppose the right-hand members of the inequalities in (1.2.c) and (1.4) are multiplied by an arbitrary factor  $M \geq 1$ . It can be seen, e.g. by using [10, p. 277], that the inclusion of such factors in the above would not alter the monotonicity criterion of Theorem 1.2.

3. In [1] Bolley and Crouzeix studied the positivity of numerical methods for linear initial value problems. Here they arrived also at an order barrier  $p \leq 1$  and at monotonicity criteria similar to the one of Theorem 1.2. The present paper has benefited by some of their nice results (cf. the Sects. 2.2, 4.2).

4. Theorem 1.2 would not be valid if the condition that the  $P_i(\zeta)$  have no common zero (see Sect. 1.1) would be omitted. This can be seen from the counterexample  $k = 1, P_0(\zeta) = 1 - \zeta^2, P_1(\zeta) = 1 + \zeta$ .

## 2. Unconditional Contractivity

### 2.1. Linear Autonomous Differential Equations

Throughout this paper  $\mathbb{K}$  stands consistently for the set of real or complex numbers  $\mathbb{R}$  or  $\mathbb{C}$ , respectively. Further  $X$  denotes an arbitrary Banach space over  $\mathbb{K}$ . The norm for  $x \in X$  is denoted by  $|x|$ . By  $\mathcal{B}(X)$  we denote the set of all bounded linear transformations  $T$  mapping  $X$  into itself with norm  $\|T\| = \sup \{|Tv| : v \in X, |v| = 1\} < \infty$  (see e.g. [9, Chap. 3]).

For any linear operator  $A$  we denote its domain and range by  $\mathcal{D}(A)$  and  $\mathcal{R}(A)$ , respectively and its resolvent set by

$$\rho(A) = \{\lambda \mid \lambda \in \mathbb{K}, (A - \lambda)^{-1} \text{ exists in } \mathcal{B}(X)\}.$$

In the following  $\omega$  denotes a fixed real number, and we deal with operators  $A$  satisfying condition (2.1).

(2.1)  $A$  is a linear operator with range in  $X$  and with a domain in  $X$  whose closure equals  $\overline{\mathcal{D}(A)} = X$ . Further there is some real  $\xi > \omega$  for which  $\mathcal{R}(A - \xi) = X$ .

We note that any  $A \in \mathcal{B}(X)$  automatically satisfies (2.1) (cf. [9, Chap. 3]).

Let  $u_0 \in \mathcal{D}(A)$  and consider the initial value problem

$$(2.2) \quad \frac{d}{dt} U(t) = AU(t) \ (t \geq 0), \quad U(0) = u_0.$$

We shall focus on operators  $A$  which are such that, in addition to (2.1), the following condition (I) is fulfilled.

(I) For each  $u_0 \in \mathcal{D}(A)$  problem (2.2) admits a unique solution  $U(t) \in \mathcal{D}(A)$  ( $0 \leq t < \infty$ ) and this solution satisfies  $|U(t)| \leq \exp(\omega t) \cdot |u_0|$  ( $0 \leq t < \infty$ ).

In order to state conditions on  $A$  that are equivalent to (I) but easier to check we give the subsequent definitions, where we follow closely [10, pp. 244, 245].

For  $x, y \in X, 0 \neq \tau \in \mathbb{R}$  we put

$$m_\tau[x, y] = \tau^{-1}(|x + \tau y| - |x|), \quad m_\pm[x, y] = \lim_{\tau \rightarrow 0 \pm} m_\tau[x, y],$$

and we write

$$\mu_\tau[A] = \sup m_\tau[v, Av], \quad \mu_\pm[A] = \sup m_\pm[v, Av],$$

both suprema being for all  $v \in \mathcal{D}(A)$  with  $|v| = 1$ .

We note that for the case of *finite-dimensional*  $X$  we always have

$$(2.3) \quad \mu_-[A] = \mu_+[A] = \lim_{\tau \rightarrow 0+} \tau^{-1}[\|I + \tau A\| - 1].$$

The limit in (2.3) equals the so-called *logarithmic norm* of  $A$  (see e.g. [18, p. 91]).

We consider the following two conditions (II) and (III) on  $A$ .

(II) Each real  $\xi > \omega$  belongs to  $\rho(A)$  and satisfies  $\|(A - \xi)^{-1}\| \leq (\xi - \omega)^{-1}$ .

(III)  $\mu_-[A] \leq \omega$ .

We formulate the following version of the well known *Hille-Yosida theorem*.

**Theorem 2.1.** *Let  $A$  satisfy (2.1). Then the three conditions (I), (II), (III) are equivalent to each other.*

This theorem can be proved by using the material in [10, pp. 282, 283].

*Definition 2.2.* By  $\mathcal{L}(X, \omega)$  we denote the class of all  $A$  satisfying both (2.1) and (III).

### 2.2. General $k$ -Step Methods

In this section we consider the numerical solution of the initial value problem (2.2) by the general  $k$ -step procedure

$$(2.4) \quad u_n = \psi_1(hA)u_{n-1} + \psi_2(hA)u_{n-2} + \dots + \psi_k(hA)u_{n-k} \quad (n \geq k).$$

Here  $\psi_i(\zeta)$  are rational functions which are regular at  $\zeta = 0$  and have numerators and denominators with coefficients in  $\mathbb{K}$ .

In order to define  $\psi_i(hA)$  appearing in (2.4) we consider an arbitrary linear transformation  $T$  with domain and range in  $X$ . Let  $\psi(\zeta) = p(\zeta)/q(\zeta)$  where  $p(\zeta),$

$q(\zeta)$  are polynomials with coefficients in  $\mathbb{K}$  and with no common zeros. We say that  $\psi(T)$  exists in  $\mathcal{B}(X)$  and we write  $\psi(T) = p(T)[q(T)]^{-1}$  whenever the operator  $q(T)$  is injective and  $p(T)[q(T)]^{-1} \in \mathcal{B}(X)$  (cf. [9, p. 163] or [10, p. 62]).

**Definition 2.3.** Let  $p$  be an integer  $\geq 0$ . Method (2.4) is of order  $p$  if

$$e^{k\zeta} = \psi_1(\zeta) e^{(k-1)\zeta} + \psi_2(\zeta) e^{(k-2)\zeta} + \dots + \psi_k(\zeta) + \mathcal{O}(\zeta^{q+1})$$

(for  $\zeta \rightarrow 0$ ) with  $q = p$ , and not with  $q > p$ .

The following theorem, which will be proved in Chap. 4, is basic for the rest of the present chapter. With  $\mathbb{K}_\infty^s$  we denote the  $s$ -dimensional vector-space  $\mathbb{K}^s$  equipped with the maximum norm  $|x| = |x|_\infty$  (cf. [10, p. 3]).

**Theorem 2.4.** Let  $h > 0$  and all  $\psi_i(\zeta)$  regular at  $\zeta = h\omega$ . Then the following three propositions are equivalent.

(P1) For each Banach space  $X$  (over  $\mathbb{K}$ ) and  $A \in \mathcal{L}(X, \omega)$  the operators  $\psi_i(hA)$  exist in  $\mathcal{B}(X)$  ( $i = 1, 2, \dots, k$ ) and satisfy

$$\left| \sum_{i=1}^k \psi_i(hA) v_i \right| \leq \sum_{i=1}^k \psi_i(h\omega) |v_i| \quad (\text{for all } v_i \in X).$$

(P2) For each integer  $s \geq 1$  and real (matrix)  $A \in \mathcal{L}(\mathbb{K}_\infty^s, \omega)$  such that all  $\psi_i(hA)$  exist in  $\mathcal{B}(\mathbb{K}_\infty^s)$  ( $i = 1, 2, \dots, k$ ), we have

$$\left| \sum_{i=1}^k \psi_i(hA) v_i \right|_\infty \leq \left[ \sum_{i=1}^k \psi_i(h\omega) \right] \cdot \max_{1 \leq j \leq k} |v_j|_\infty \quad (\text{for all } v_i \in \mathbb{K}_\infty^s).$$

(P3) All  $\psi_i$  ( $i = 1, 2, \dots, k$ ) are absolutely monotonic on the interval  $(-\infty, h\omega]$ .

The next theorem provides an important consequence of condition (P3) for the case where  $h\omega \geq 0$ .

**Theorem 2.5.** If all  $\psi_i$  ( $i = 1, 2, \dots, k$ ) are absolutely monotonic on  $(-\infty, 0]$  then the order of method (2.4) is not greater than 1.

*Proof.* It is easily verified that the function  $\psi$  defined by

$$\psi(\xi) = \sum_{j=1}^k \psi_j(\xi/k) \cdot \exp[(k-j)\xi/k] \quad (-\infty < \xi \leq 0)$$

is absolutely monotonic on  $(-\infty, 0]$ . Suppose the order of (2.4) exceeds 1. Then

$$\exp(\xi) = \psi(\xi) + \mathcal{O}(\xi^3) \quad (\text{for } \xi \rightarrow 0-).$$

In [1] Bolley and Crouzeix present a lemma implying that any  $\psi(\xi)$  which is absolutely monotonic on  $(-\infty, 0]$  and satisfies  $\exp(\xi) - \psi(\xi) = \mathcal{O}(\xi^3)$  (for  $\xi \rightarrow 0-$ ) must be of the form  $\psi(\xi) = \exp(\xi)$  ( $-\infty < \xi \leq 0$ ). The proof in [1] of this lemma is based on Bernstein's representation for absolutely monotonic functions  $\psi$ , which reads  $\psi(\xi) = \int_0^\infty \exp(\xi t) d\alpha(t)$  ( $-\infty < \xi \leq 0$ ). Here  $\alpha(t)$  is bounded and non-decreasing, and the integral stands for a (convergent) improper Riemann-Stieltjes integral.

Clearly, our function  $\psi$  as defined above cannot satisfy  $\psi(\xi) = \exp(\xi)$  ( $-\infty < \xi \leq 0$ ). We thus have a contradiction, which proves the theorem.  $\square$

### 2.3. Applications

2.3.1. *Proof of the Theorems 1.2 and 1.3.* The first application we shall give of the theorems from Sect. 2.2 consists in a proof of the two theorems presented in the Introduction.

1. In order to apply the results of Sect. 2.2 to the situation in Sect. 1.1, we define  $\psi_i$  as in Sect. 1.2, and  $\mathbb{K} = \mathbb{R}$ ,  $\omega = 0$ . It is easily verified that

- (2.5) (1.3) is of order  $p$  (see definition in Sect. 1.1) if and only if
- (2.4) is of order  $p$  (see Definition 2.3).

Clearly a pair  $(A, |\cdot|)$  satisfies (1.2) if and only if  $A \in \mathcal{L}(X, 0)$  where  $X = \mathbb{R}^s$  equipped with the norm  $|x|$  for  $x \in \mathbb{R}^s$ ,  $s \geq 1$ . Note also that when  $T \in \mathcal{B}(\mathbb{R}^s)$  with  $\psi_i(T)$  existing in  $\mathcal{B}(\mathbb{R}^s)$  we have

$$(2.6) \quad P_0(T) \psi_i(T) = P_i(T) \quad (i = 1, 2, \dots, k).$$

2. Suppose all rational functions  $\psi_i$  are absolutely monotonic on  $(-\infty, 0]$ . Let  $h > 0$ ,  $(A, |\cdot|)$  satisfy (1.2), and  $u_n$  be computed from (1.3). We prove (1.4).

Since proposition (P3) is valid, we have by Theorem 2.4 also (P1). Hence  $\psi_i(hA)$  exist in  $\mathcal{B}(\mathbb{R}^s)$  and the vector

$$v = \sum_{i=1}^k \psi_i(hA) u_{n-i}$$

satisfies

$$|v| \leq \sum_{i=1}^k \psi_i(0) |u_{n-i}|.$$

Since the order of (1.3) in Sect. 1.1 is assumed to be  $\geq 0$ , we obtain from (2.5) the relation

$$(2.7) \quad \psi_1(0) + \psi_2(0) + \dots + \psi_k(0) = 1.$$

It follows that  $|v| \leq \max(|u_{n-1}|, |u_{n-2}|, \dots, |u_{n-k}|)$ . Using the assumption of Sect. 1.1 that  $P_0, P_1, \dots, P_k$  have no common zero, it can be seen that  $P_0(hA)$  is regular. In view of (2.6) and the definition of  $v$  we have  $P_0(hA)(u_n - v) = 0$ . Hence  $u_n = v$  and therefore

$$|u_n| \leq \max(|u_{n-1}|, |u_{n-2}|, \dots, |u_{n-k}|).$$

This implies (1.4).

3. Suppose (1.3) is unconditionally contractive for all pairs  $(A, |\cdot|)$  satisfying (1.2). We prove (P3).

Let  $h > 0$ ,  $s \geq 1$ ,  $A \in \mathcal{L}(\mathbb{R}_\infty^s, 0)$  and suppose, as in (P2), that all  $\psi_i(hA)$  exist in  $\mathcal{B}(\mathbb{R}_\infty^s)$  ( $i = 1, 2, \dots, k$ ). Let  $v_i \in \mathbb{R}_\infty^s$  ( $i = 1, 2, \dots, k$ ) and

$$v_0 = \sum_{i=1}^k \psi_i(hA) v_i.$$

From (2.6) it follows that  $u_i = v_{k-i}$  ( $i = 0, 1, \dots, k$ ) satisfy (1.3) with  $n = k$ . Hence

$$|u_k|_\infty \leq \max(|u_0|_\infty, |u_1|_\infty, \dots, |u_{k-1}|_\infty).$$

Using (2.7) we obtain

$$|v_0|_\infty \leq \max(|v_k|_\infty, \dots, |v_2|_\infty, |v_1|_\infty) = \left[ \sum_{i=1}^k \psi_i(h\omega) \right] \cdot \max_{1 \leq j \leq k} |v_j|_\infty.$$

Hence (P2) holds and, in view of Theorem 2.4, also (P3). This completes the proof of Theorem 1.2.

4. Suppose (1.3) is unconditionally contractive for all  $(A, |\cdot|)$  satisfying (1.2). Since (P3) holds, an application of Theorem 2.5 and (2.5) shows that the order of (1.3) does not exceed 1. Theorem 1.3 has thus been proved.

2.3.2. *Linear Multistep Methods.* We consider the numerical solution of (1.1) by the linear multistep method

$$(2.8) \quad \sum_{i=0}^k \alpha_i u_{n+i} = h \sum_{i=0}^k \beta_i A u_{n+i} \quad (n=0, 1, 2, \dots).$$

Here  $\alpha_i, \beta_i$  are real coefficients with  $\alpha_k \neq 0$ ,

$$\sum_{i=0}^k \alpha_i = 0, \quad \sum_{i=0}^k i \alpha_i = \sum_{i=0}^k \beta_i = 1 \quad (\text{consistency and normalization})$$

(cf. [12, p. 462]).

It is easily verified that (2.8) is of type (1.3) with  $P_0(\zeta) = \alpha_k - \beta_k \zeta$ ,  $P_i(\zeta) = -\alpha_{k-i} + \beta_{k-i} \zeta$  ( $1 \leq i \leq k$ ) satisfying the assumptions made in Sect. 1.2. By Theorem 1.2 it follows easily that (2.8) is unconditionally contractive for all pairs  $(A, |\cdot|)$  satisfying (1.2) if and only if

$$(2.9) \quad \alpha_k > 0, \quad \beta_k > 0, \quad \alpha_i/\alpha_k \leq \beta_i/\beta_k \leq 0 \quad (i=0, 1, \dots, k-1).$$

In view of Theorem 1.3 the conditions (2.9) imply that the order of (2.8) does not exceed  $p = 1$ .

A simple example is provided by the  $\theta$ -method. Here  $k = 1$ ,  $\alpha_1 = 1$ ,  $\alpha_0 = -1$ ,  $\beta_1 = \theta$ ,  $\beta_0 = 1 - \theta$  with parameter  $\theta \in \mathbb{R}$ . Clearly (2.8) becomes equivalent to a procedure of type (1.3) with  $P_0(\zeta) = 1 - \theta \zeta$ ,  $P_1(\zeta) = 1 + (1 - \theta)\zeta$  and of type (2.4) with

$$(2.10) \quad \psi_1(\zeta) = [1 + (1 - \theta)\zeta] / [1 - \theta\zeta].$$



In view of (2.9) the  $\theta$ -method is unconditionally contractive for all pairs  $(A, |\cdot|)$  satisfying (1.2) if and only if  $\theta \geq 1$ .

In [12, pp. 58, 59] Nevanlinna and Liniger arrived at the criterion (2.9) by different means.

**2.3.3. Runge-Kutta Methods.** Any Runge-Kutta method for solving (1.1) can be written in the form

$$(2.11.a) \quad u_{n+1} = u_n + h \sum_{i=1}^m b_i A y_i \quad (n=0, 1, 2, \dots)$$

where the vectors  $y_i \in \mathbb{R}^s$  depend on  $u_n$  and satisfy

$$(2.11.b) \quad y_i = u_n + h \sum_{j=1}^m c_{ij} A y_j \quad (i=1, 2, \dots, m).$$

Here  $b_i, c_{ij}$  are real parameters with  $b_1 + b_2 + \dots + b_m = 1$ .

The order  $p$  and the concept of unconditional contractivity for (2.11) are defined in a similar way as for (1.3) in Sect. 1.1.

We define  $m \times m$  matrices  $C_i$  and polynomials  $Q_i$  ( $i=0, 1$ ) by  $C_0 = (c_{ij}), C_1 = (c_{ij} - b_j), Q_0(\zeta) = \det(I - \zeta C_0), Q_1(\zeta) = \det(I - \zeta C_1)$ . By  $\psi$  we denote the rational function, without removable singularities, such that

$$\psi(\zeta) = Q_1(\zeta)/Q_0(\zeta) \quad (\text{for } \zeta \in \mathbb{C}, Q_0(\zeta) \neq 0).$$

It can be proved that (2.11.b) admits a unique solution  $y_1, y_2, \dots, y_m$  for all  $h > 0$  and all  $A$  satisfying (1.2) if and only if

$$(2.12) \quad Q_0(\zeta) \neq 0 \quad (\text{for all complex } \zeta \text{ with } \text{Re } \zeta \leq 0).$$

**Theorem 2.6.** *Assume (2.12). Then the Runge-Kutta method (2.11) is unconditionally contractive with respect to all pairs  $(A, |\cdot|)$  satisfying (1.2) if and only if  $\psi(\zeta)$  is absolutely monotonic on  $(-\infty, 0]$ . Further if (2.11) has this contractivity property, then its order does not exceed  $p=1$ .*

*Proof.* Let  $h > 0$  and  $(A, |\cdot|)$  satisfy (1.2). It can be proved (cf. [18, pp. 132, 152]) that (2.11) is equivalent to

$$(2.13) \quad Q_0(hA) u_{n+1} = Q_1(hA) u_n \quad (n=0, 1, 2, \dots).$$

Let  $P_0(\zeta), P_1(\zeta)$  be polynomials without common zeros such that  $P_1(\zeta)/P_0(\zeta) = Q_1(\zeta)/Q_0(\zeta)$  (for all  $\zeta$  with  $Q_0(\zeta) \neq 0$ ). In view of (2.12) the matrix  $Q_0(hA)$  is regular, and therefore (2.13) is equivalent to

$$(2.14) \quad P_0(hA) u_n = P_1(hA) u_{n-1} \quad (n=1, 2, 3, \dots).$$

Since  $P_0, P_1$  satisfy the assumptions made in Sect. 1.1 with  $k=1$ , we can apply the Theorems 1.2, 1.3 to (2.14). The statements in the above theorem thus easily follow.  $\square$

2.3.4. *A Partial Differential Equation.* The material in this section provides a simple illustration to Theorem 2.4 where  $X$  is of an infinite dimension. It will also illustrate some considerations in the next chapter.

Let  $\mathbb{K} = \mathbb{R}$  and

$$X = \{x \mid x \in C[0, 1], x(0) = x(1) = 0\} \quad \text{with norm } |x| = \max_{0 \leq \xi \leq 1} |x(\xi)|.$$

The operator  $A: D \rightarrow X$  is defined by

$$(2.15) \quad \begin{aligned} (Ax)(\xi) &= a_2(\xi)x''(\xi) + a_1(\xi)x'(\xi) + a_0(\xi)x(\xi) \quad (0 \leq \xi \leq 1, x \in D), \\ D &= \{x \mid x \in C^{(2)}[0, 1], x(\xi) = a_2(\xi)x''(\xi) + a_1(\xi)x'(\xi) = 0 \text{ for } \xi = 0, 1\}. \end{aligned}$$

Here  $a_0, a_1, a_2 \in C[0, 1]$  are any given functions with  $a_2(\xi) > 0$  ( $0 \leq \xi \leq 1$ ).

It follows easily from the theory of ordinary differential equations that  $A \in \mathcal{L}(X, \omega)$  (cf. Definition 2.2) with

$$(2.16) \quad \omega = \max_{0 \leq \xi \leq 1} a_0(\xi).$$

An application of Theorem 2.1 shows that, for any  $u_0 \in D$ , the initial-boundary value problem

$$(2.17) \quad \begin{aligned} \frac{\partial}{\partial t} V(\xi, t) &= a_2(\xi) \frac{\partial^2}{\partial \xi^2} V(\xi, t) + a_1(\xi) \frac{\partial}{\partial \xi} V(\xi, t) + a_0(\xi) V(\xi, t), \\ V(\xi, 0) &= u_0(\xi), \quad V(0, t) = V(1, t) = 0 \quad (0 \leq \xi \leq 1, t \geq 0) \end{aligned}$$

has a solution, which can be written as  $V(\xi, t) = U(t)(\xi)$  ( $0 \leq \xi \leq 1, t \geq 0$ ) with  $U(t)$  as in statement (I) of Sect. 2.1.

We consider the application of procedure (2.4) with  $k=1$  and  $\psi_1$  given by (2.10). We thus arrive at the “semi-discrete” process

$$(2.18) \quad \begin{aligned} u_n(\xi) &= y(\xi) + (1 - \theta)h(Ay)(\xi), \\ y(\xi) - \theta h(Ay)(\xi) &= u_{n-1}(\xi), \quad y(0) = y(1) = 0 \quad (0 \leq \xi \leq 1, n \geq 1). \end{aligned}$$

By Theorem 2.4 we have (P3)  $\Rightarrow$  (P1). Therefore this theorem implies that for any  $u_0 \in X$ ,  $\theta \geq 1$  and  $h$  satisfying  $h\omega < 1/\theta$  the functions  $u_n(\xi)$  are uniquely determined by (2.18) and satisfy

$$\max_{\xi} |u_n(\xi)| \leq \left[ \frac{1 + (1 - \theta)h\omega}{1 - \theta h\omega} \right]^n \cdot \max_{\xi} |u_0(\xi)| \quad (n \geq 1).$$

Similar bounds can be derived in an analogous fashion for “semi-discrete” processes based on more complicated functions  $\psi_1$  (e.g.  $\psi_1 = \varphi$  given in (3.7)).

### 2.4. Remarks

1. For  $\theta=1$  process (2.18) reduces to the so-called *Rothe-method* (see [16], [14]). Further, for  $\theta=1$  the procedure (2.4) with  $k=1$ , and  $\psi_1(\zeta)$  given by

(2.10), plays an important part in the usual proof of the *Hille-Yosida theorem* (cf. e.g. [10], [9] and Sect. 2.1).

2. For procedure (2.4) with  $k=1$  very interesting upperbounds on  $|u_n|/|u_0|$  were obtained in [2]. These upperbounds are valid under much weaker assumptions on  $\psi_1(\zeta)$  than the assumption of absolute monotonicity that is required in the situation of Theorem 2.4. On the other hand, this theorem yields upperbounds which are usually sharper than those in [2] and hold through for the case of variable stepsizes  $h=h_n>0$ .

3. We have preferred the definition of  $\psi(T)$  as given in Sect. 2.2 to other definitions (see e.g. [6, p. 601]) since it is short, requires no complex variables when  $\mathbb{K}=\mathbb{R}$ , and indicates how actual (numerical) evaluations of  $\psi(T)x$  can be performed.

### 3. Conditional Contractivity

#### 3.1. A Circle Condition on the Operator $A$

We shall consider the numerical solution of the initial value problem (2.2) where the operator  $A$  satisfies some stronger condition than assumed in Chap. 2. It will turn out that for such operators the  $k$ -step procedure (2.4) can have an order  $p>1$  while it is still contractive under some condition  $0<h<H$  on the stepsize.

In all of the following  $\omega, \tau$  denote fixed real numbers with  $\tau>0, 1+\tau\omega>0$ , and we consider the following three conditions on an operator  $A\in\mathcal{B}(X)$ .

(i) For each  $u_0\in X$  the recurrence relation

$$\tau^{-1}(u_n - u_{n-1}) = Au_{n-1} \quad (n \geq 1)$$

has a solution satisfying

$$|u_n| \leq (1 + \tau\omega)^n |u_0| \quad (n \geq 1).$$

(ii)  $\|A + \tau^{-1}\| \leq \omega + \tau^{-1}$ .

(iii)  $\mu_\tau[A] \leq \omega$ .

We emphasize the apparent analogy between any of the conditions (i), (ii), (iii) and (I), (II), (III) (cf. Sect. 2.1), respectively.

We note that when  $\omega=0$ , condition (i) means that *Euler's method* with stepsize  $h=\tau$  behaves contractively (see also [8, pp. 75, 76]). Further in case  $X = \mathbb{K} = \mathbb{C}$ , condition (ii) means that  $z=A$  lies within the *circle*  $\subset \mathbb{C}$  passing through the point  $\omega$  and with center  $= -\tau^{-1}$ , while (II) means that  $z=A$  belongs to the complex half-plane  $\text{Re } z \leq \omega$ . Therefore we call (ii) a *circle condition* (see also loc. cit.). Finally an intuitive feeling of the nature of condition (iii) is obtained by discretizing the operator  $A$  of Sect. 2.3.4 as indicated in Sect. 3.3.3. While  $A$  itself satisfies (III) but violates (iii), its discretized version  $A_\delta$  fulfills both (III) and (iii) under the assumptions (3.8), (3.9).

**Theorem 3.1.** *Let  $A \in \mathcal{B}(X)$ . Then the requirements (i), (ii), (iii) are equivalent to each other and imply the properties (2.1), (I), (II), (III) listed in Sect. 2.1.*

*Proof.* It is easily verified that (i), (ii), (iii) are equivalent to each other. Since  $A \in \mathcal{B}(X)$ , it satisfies (2.1).

For  $\sigma < 0$  we have  $m_\sigma[x, y] \leq m_\tau[x, y]$  (see Sect. 2.1 and [10, p. 37]), and therefore  $m_-[x, y] \leq m_\tau[x, y]$ . It follows that  $\mu_-[A] \leq \mu[A]$ . Condition (iii) thus implies (III) and in view of Theorem 2.1 also (I), (II).  $\tau \square$

**Definition 3.2.** By  $\mathcal{L}(X, \omega, \tau)$  we denote the class of all  $A \in \mathcal{B}(X)$  satisfying (iii).

### 3.2. $k$ -Step Methods when $A$ Satisfies a Circle Condition

In this section we consider the  $k$ -step procedure (2.4) for  $A \in \mathcal{L}(X, \omega, \tau)$ . The following theorem, to be proved in Chap. 4, has some similarity to Theorem 2.4. It is basic for the present chapter.

**Theorem 3.3.** *Let  $h > 0$  and all  $\psi_i(\zeta)$  regular at  $\zeta = h\omega$ . Then the following three propositions are equivalent.*

(p-1) *For each Banachspace  $X$  (over  $\mathbb{K}$ ) and  $A \in \mathcal{L}(X, \omega, \tau)$  the operators  $\psi_i(hA)$  exist in  $\mathcal{B}(X)$  ( $1 \leq i \leq k$ ) and satisfy*

$$\left| \sum_{i=1}^k \psi_i(hA) v_i \right| \leq \sum_{i=1}^k \psi_i(h\omega) |v_i| \quad (\text{for all } v_i \in X).$$

(p-2) *For each integer  $s \geq 1$  and real (matrix)  $A \in \mathcal{L}(\mathbb{K}_\infty^s, \omega, \tau)$  such that all  $\psi_i(hA)$  exist in  $\mathcal{B}(\mathbb{K}_\infty^s)$  ( $1 \leq i \leq k$ ), we have*

$$\left| \sum_{i=1}^k \psi_i(hA) v_i \right|_\infty \leq \left[ \sum_{i=1}^k \psi_i(h\omega) \right] \cdot \max_{1 \leq j \leq k} |v_j|_\infty \quad (\text{for all } v_i \in \mathbb{K}_\infty^s).$$

(p-3) *All  $\psi_i$  ( $1 \leq i \leq k$ ) are absolutely monotonic on the interval  $[-h\tau^{-1}, h\omega]$ .*

**Definition 3.4.** A stepsize  $h_0 > 0$  is admissible if

$$|u_n| \leq \max(|u_0|, |u_1|, \dots, |u_{k-1}|) \quad (n \geq k)$$

whenever  $h = h_0$ ,  $X$  is any Banach space over  $\mathbb{K}$ ,  $A \in \mathcal{L}(X, \omega, \tau)$  and  $u_n$  satisfies (2.4). The largest number  $H \leq \infty$  with the property that each  $h_0 \in (0, H)$  is admissible, is denoted by  $H(\omega, \tau)$  and is called the *contractivity threshold* of method (2.4).

Restrictions on the stepsize  $h$  for stability or contractivity reasons are often embarrassing (in particular in the numerical solution of stiff initial value problems, see e.g. [18, 8, 12]). Therefore, the larger  $H(\omega, \tau)$  the better.

We note that thresholds quite similar to  $H(\omega, \tau)$  were introduced in [5, 17].

**Theorem 3.5.** *Let method (2.4) be of an order  $p \geq 1$ . Then*

$$H(\omega, \tau) \begin{cases} = 0 & (\text{for } \omega > 0), \\ = R\tau & (\text{for } \omega = 0), \\ \geq R\tau & (\text{for } \omega < 0). \end{cases}$$

Here

$$(3.1) \quad R = \sup \{r \mid r = 0, \text{ or } r > 0 \text{ and all } \psi_i(\zeta) \text{ are absolutely monotonic on } [-r, 0]\}.$$

*Proof.* 1. Suppose  $\omega > 0$ . Choosing  $X = \mathbb{K}$ ,  $A = \omega \in \mathcal{L}(\mathbb{K}, \omega, \tau)$ ,  $u_j = \exp(jh\omega)$  ( $0 \leq j \leq k-1$ ), we obtain by Definition 2.3 (with  $p \geq 1$ ,  $\zeta = h\omega$ ) the relation

$$u_k = \exp(kh\omega) + \mathcal{O}(h^2) \quad (\text{for } h \rightarrow 0+).$$

Consequently  $|u_k| = \gamma(h) \cdot \max(|u_0|, |u_1|, \dots, |u_{k-1}|)$  with  $\gamma(h) = 1 + h\omega + \mathcal{O}(h^2) > 1$  (for  $h \rightarrow 0+$ ). Hence  $H(\omega, \tau) = 0$ .

2. Suppose  $\omega = 0$ . We first assume  $h_0 \in (0, R\tau)$ . Since  $R > h_0 \tau^{-1}$  we have (p-3) with  $h = h_0$ , and consequently (p-1) with  $h = h_0$ . In view of  $\psi_1(0) + \psi_2(0) + \dots + \psi_k(0) = 1$ , it follows that  $h_0$  is admissible.

We next assume  $h_0 \in (R\tau, \infty)$ . With  $h = h_0$  condition (p-3) is violated. Therefore (p-2) does not hold with  $h = h_0$ , and  $h_0$  is thus not admissible.

It follows that  $H(\omega, \tau) = R\tau$ .

3. For  $\omega < 0$  we always have  $\mathcal{L}(X, \omega, \tau) \subset \mathcal{L}(X, 0, \tau)$  and consequently  $H(\omega, \tau) \geq H(0, \tau) = R\tau$ .  $\square$

In view of Theorem 3.5 we call  $R$  the *threshold-factor* of method (2.4). Clearly, the larger  $R$ , the better is the general contractivity behaviour of (2.4) in the important case  $\omega = 0$ .

### 3.3. Applications

**3.3.1. Linear Multistep Methods.** We consider the numerical solution of (2.2), with  $A \in \mathcal{L}(X, 0, \tau)$ , by the linear multistep method (2.8). We thus have a method of the general type (2.4) with

$$\psi_i(\zeta) = (-\alpha_{k-i} + \beta_{k-i}\zeta) / (\alpha_k - \beta_k\zeta) \quad (\text{for } \alpha_k - \beta_k\zeta \neq 0).$$

A straightforward calculation shows that the threshold-factor  $R$  (see (3.1)) is positive if and only if

$$(3.2) \quad \begin{aligned} \alpha_k &> 0, & \beta_k &\geq 0, & r &> 0, \\ \alpha_i &\leq 0, & \alpha_i \beta_k &\leq \beta_i \alpha_k & (0 \leq i \leq k-1). \end{aligned}$$

Here  $r$  is defined by

$$(3.3) \quad r = \min \{ -\alpha_i / \beta_i \mid 0 \leq i \leq k-1 \text{ and } \beta_i > 0 \},$$

with the convention  $\min \emptyset = \infty$ .

Further, in case (3.2) holds, the threshold-factor equals

$$(3.4) \quad R = r.$$

For the  $\theta$ -method, defined in Sect. 2.3.2, we obtain from (3.2), (3.4)

$$(3.5) \quad R = \begin{cases} 0 & (\text{for } \theta < 0), \\ (1-\theta)^{-1} & (\text{for } 0 \leq \theta < 1), \\ \infty & (\text{for } \theta \geq 1). \end{cases}$$

We note that when  $\theta = \frac{1}{2}$  we have the well-known *trapezoidal rule* or *Crank-Nicolson* method with order  $p = 2$ . This method thus has a contractive behaviour when

$$0 < h < 2\tau$$

(cf. also [1, p. 243]).

**3.3.2. One-Step Methods.** We consider the application of the general method (2.4) to  $A \in \mathcal{L}(X, 0, \tau)$  when  $k = 1$ . In view of (3.1) the threshold-factor is given by

$$(3.6) \quad R = \sup \{r \mid r = 0, \text{ or } r > 0 \text{ and } \psi_1(\zeta) \text{ absolutely monotonic on } [-r, 0]\}.$$

An example is provided by the Runge-Kutta scheme (2.11). With the assumptions and notations of Sect. 2.3.3 we have a method of type (2.4) with  $k = 1$ ,  $\psi_1(\zeta) = \psi(\zeta) = Q_1(\zeta)/Q_0(\zeta)$  (for  $\zeta \in \mathbb{C}$ ,  $Q_0(\zeta) \neq 0$ ). It thus follows that the threshold-factor of (2.11) is given by (3.6) with  $\psi_1$  replaced by  $\psi$ .

We next consider a two-parameter family of methods (2.4) with  $k = 1$  and  $\psi_1 = \varphi$  where

$$(3.7) \quad \varphi(\zeta) = [1 + (1 - 2\theta)\zeta + (\theta(\theta - 1) + \alpha)\zeta^2] \cdot [1 - \theta\zeta]^{-2}.$$

For any  $\theta, \alpha \in \mathbb{R}$  the function  $\varphi$  is a so-called restricted-denominator approximation to  $\exp(\zeta)$  with an order  $p$  satisfying  $1 \leq p \leq 3$  (see [13, Theorem 2.1]). We denote the threshold-factor by  $R(\theta, \alpha)$ .

A straightforward calculation shows that  $R(\theta, \alpha) > 0$  if and only if

$$\theta \geq 0 \quad \text{and} \quad \alpha \geq 0.$$

For  $\alpha = 0$  we have  $\varphi = \psi_1$  where  $\psi_1$  is defined in (2.10). Hence  $R(\theta, 0) = R$  with  $R$  given in (3.5).

For  $\theta = 0$  the method is explicit, and  $R(0, \alpha)$  is easily calculated. The threshold-factor  $R(0, \alpha)$  is maximal for  $\alpha = \frac{1}{4}$  with  $R(0, \frac{1}{4}) = 2$ .

Within the class of second order methods the threshold-factor can be shown [7] to be maximal for  $\theta = \frac{1}{4}$ ,  $\alpha = \frac{1}{4}$  with  $R(\frac{1}{4}, \frac{1}{4}) = 4$ .

A simple calculation using [13, Theorem 2.1] shows that there are two methods of order 3, one with  $R = 0$  and one with  $R > 0$ . The latter is obtained when  $\theta = (3 - \sqrt{3})/6$ ,  $\alpha = \sqrt{3}/6$ , with  $R(\theta, \alpha) = 1 + \sqrt{3}$ .

**3.3.3. A Partial Differential Equation.** We present a final illustration to the material of this chapter.

Let  $A$  denote the operator defined in Sect. 2.3.4. We consider the “discrete” operator  $A_\delta$  obtained from  $A$  by replacing the derivatives in (2.15) by second order, central difference approximations (cf. [11]). We deal with a uniform grid on  $[0, 1]$  with grid spacing  $\delta = (s + 1)^{-1}$  where  $s$  is an integer  $\geq 1$ . Assume  $\delta > 0$  is so small that

$$(3.8) \quad \delta \cdot \underset{\xi}{\text{Max}} |a_1(\xi)[2a_2(\xi)]^{-1}| \leq 1.$$

A simple calculation shows that  $A_\delta$  can be viewed as an element of  $\mathcal{L}(\mathbb{R}_\infty^s, \omega)$  with  $\omega$  as in (2.16). In fact,  $\mu_- [A_\delta] = \mu_\tau [A_\delta] \leq \omega$  for any  $\tau > 0$  satisfying

$$(3.9) \quad \tau \cdot \text{Max}_\xi [2\delta^{-2} a_2(\xi) - a_0(\xi)] \leq 1.$$

Applying method (2.4) with  $k=1$ ,  $\psi_1 = \varphi$  (see (3.7)) and with  $A$  replaced by  $A_\delta$  amounts to the “fully-discrete” process.

$$(I - \theta h A_\delta)^2 u_n = [I + (1 - 2\theta) h A_\delta + (\theta(\theta - 1) + \alpha)(h A_\delta)^2] u_{n-1}.$$

Here the components of  $u_n \in \mathbb{R}^s$  approximate  $V(\xi, t_n)$  (see (2.17)) at the grid points  $\xi = i\delta$  ( $i = 1, 2, \dots, s$ ) (cf. [11]).

Let  $\omega$  defined by (2.16) be nonpositive and let  $\tau$  satisfy  $\tau > 0, 1 + \omega\tau > 0$ , (3.9). Using the implication (p-3)  $\Rightarrow$  (p-1) of Theorem 3.3 one arrives at the bound

$$|u_n|_\infty \leq \varphi(h\omega)^n |u_0|_\infty \leq |u_0|_\infty \quad (n \geq 1)$$

for any stepsize  $h$  satisfying

$$0 < h < R(\theta, \alpha) \cdot \tau = H(0, \tau).$$

### 4. Proof of the Theorems 2.4, 3.3

#### 4.1. Preliminaries

We first note that it is sufficient to prove the Theorems 2.4, 3.3 for the case where  $h=1$ . The general case can be reduced to the situation where  $h=1$  by defining  $\varphi_i(\xi) = \psi_i(h\xi)$  and dealing with  $\varphi_i$  instead of  $\psi_i$ . In all of the following we therefore assume  $h=1$ .

In this section we shall prove the theorems for the case  $\mathbb{K} = \mathbb{R}$ , under the assumption that they hold when  $\mathbb{K} = \mathbb{C}$ . We shall need the following two lemmata.

**Lemma 4.1.** *Let  $T_i = (T_{imn})$  denote square matrices of order  $s \geq 1$  with entries  $\in \mathbb{K}$  ( $i = 1, 2, \dots, k$ ). Then the smallest constant  $\gamma$  with the property*

$$\left| \sum_{i=1}^k T_i v_i \right|_\infty \leq \gamma \cdot \max_{1 \leq i \leq k} |v_i|_\infty \quad (\text{for all } v_i \in \mathbb{K}^s)$$

is given by

$$\gamma = \max_{1 \leq m \leq s} \sum_{i=1}^k \sum_{n=1}^s |T_{imn}|.$$

For  $k=1$  this lemma is well-known, and its (easy) proof for  $k \geq 1$  omitted.

Let  $X$  be a real Banachspace. We denote its complexification by  $X'$  and when  $T$  is an operator with domain and range in  $X$  we define  $T'(x, y) = (Tx, Ty) \in X'$  for all  $(x, y) \in X'$  with  $x \in \mathcal{D}(T), y \in \mathcal{D}(T)$ .

**Lemma 4.2.**  $X'$  can be given a norm  $|w|$  (for  $w=(x, y)\in X'$ ) such that it is a complex Banach space with the following three properties.

(i) The isomorphism  $x \rightarrow (x, 0)$  of  $X$  into  $X'$  is an isometry.

(ii) If  $T_1, T_2, \dots, T_k$  are operators with domain  $= X$  and range  $\subset X$  and  $\gamma_1, \gamma_2, \dots, \gamma_k \in \mathbb{R}$ , then

$$\left| \sum_{i=1}^k T_i v_i \right| \leq \sum_{i=1}^k \gamma_i |v_i| \quad (\text{for all } v_i \in X)$$

if and only if

$$\left| \sum_{i=1}^k T'_i w_i \right| \leq \sum_{i=1}^k \gamma_i |w_i| \quad (\text{for all } w_i \in X').$$

(iii) Let  $\psi(\zeta) = p(\zeta)/q(\zeta)$  where  $p(\zeta), q(\zeta)$  are polynomials with real coefficients. Let  $T$  be an operator with domain and range in  $X$ . Then  $\psi(T')$  exists in  $\mathcal{B}(X')$  iff  $\psi(T)$  exists in  $\mathcal{B}(X)$ , and when they exist we have  $\psi(T') = \psi(T)'$ ,  $\|\psi(T')\| = \|\psi(T)\|$ .

This lemma is a slight extension of a theorem stated in [15, p. 6]. For its proof, which is a bit more difficult than might be expected at first, we refer to [15, pp. 6, 7].

Assume the Theorems 2.4, 3.3 hold with  $\mathbb{K} = \mathbb{C}$ . We prove that they hold also when  $\mathbb{K} = \mathbb{R}$ .

Let  $\tau \geq 0$ ,  $\omega \in \mathbb{R}$ , and  $\psi_i$  be as in Sect. 2.2 (with  $\mathbb{K} = \mathbb{R}$ ). We define  $\mathcal{L}(X, \omega, 0) = \mathcal{L}(X, \omega)$  and for  $i = 1, 2$  we denote by  $S_i(\omega, \tau, \mathbb{K})$  statement (p- $i$ ) of Theorem 3.3 when  $\tau > 0$ , and statement (Pi) of Theorem 2.4 when  $\tau = 0$ . Clearly, it is sufficient to prove the implications  $S_1(\omega, \tau, \mathbb{C}) \Rightarrow S_1(\omega, \tau, \mathbb{R})$ , and  $S_2(\omega, \tau, \mathbb{R}) \Rightarrow S_2(\omega, \tau, \mathbb{C})$ .

Using Lemma 4.2 ((i), (iii)) it can be seen that for any real Banach space  $X$  and  $A \in \mathcal{L}(X, \omega, \tau)$  we also have  $A' \in \mathcal{L}(X', \omega, \tau)$ . Hence, in view of (ii), (iii) of Lemma 4.2, we have  $S_1(\omega, \tau, \mathbb{C}) \Rightarrow S_1(\omega, \tau, \mathbb{R})$ . From Lemma 4.1 (with  $k = 1$ ) it can be seen that any real (matrix)  $A \in \mathcal{L}(\mathbb{C}_\infty^s, \omega, \tau)$  also belongs to  $\mathcal{L}(\mathbb{R}_\infty^s, \omega, \tau)$ . Applying Lemma 4.1 once more (with  $T_i = \psi_i(A)$ ) it follows that  $S_2(\omega, \tau, \mathbb{R}) \Rightarrow S_2(\omega, \tau, \mathbb{C})$ . This completes the proof.

In all of the following we assume, with no loss of generality and much gain of simplicity, that  $\mathbb{K} = \mathbb{C}$ .

### 4.2. The Proof of Theorem 3.3

1. We shall prove Theorem 3.3 by arguments that have some similarity to those in [1, pp. 239, 240].

Assume (p-3). We shall prove (p-1).

Let  $X$  be a complex Banach space and  $A \in \mathcal{L}(X, \omega, \tau)$ . From property (ii) (Sect. 3.1) we obtain  $A = \xi + T$  with  $\xi = -\tau^{-1}$ ,  $\|T\| \leq \omega - \xi$ .



Since all  $\psi_i$  are absolutely monotonic on  $[\xi, \omega]$ , the Taylor series  $\alpha_{i_0} + \alpha_{i_1} t + \alpha_{i_2} t^2 + \dots$  with  $\alpha_{i_j} = \psi_i^{(j)}(\xi)/j!$  are absolutely convergent for  $t \in \mathbb{C}$ ,  $|t| \leq \omega - \xi$ ,  $1 \leq i \leq k$ . It follows (cf. e.g. [6, p. 568]) that  $\psi_i(\xi + T)$  exist in  $\mathcal{B}(X)$  and  $\psi_i(\xi + T) = \alpha_{i_0} + \alpha_{i_1} T + \alpha_{i_2} T^2 + \dots$  (for  $1 \leq i \leq k$ ).

We thus obtain  $\|\psi_i(A)\| \leq \alpha_{i_0} + \alpha_{i_1}(\omega - \xi) + \alpha_{i_2}(\omega - \xi)^2 + \dots$  and therefore  $\|\psi_i(A)\| \leq \psi_i(\omega)$ . This proves (p-1).

2. Clearly (p-1) implies (p-2). Therefore we assume (p-2), and it remains to show that (p-3) holds.

Let  $\lambda \in [\xi, \omega)$  and  $A = \lambda + (\omega - \lambda)E$ . Here  $E$  denotes the square matrix of order  $s \geq 1$  all of whose entries  $E_{m,n}$  are zero with the exception  $E_{m,m+1} = 1$  ( $m = 1, 2, \dots, s-1$ ). For the operator norm subordinate to the maximum norm  $|\cdot|_\infty$  in  $\mathbb{C}^s$  we have  $\|A + \tau^{-1}\| = \|(\lambda - \xi) + (\omega - \lambda)E\| \leq \omega - \xi = \omega + \tau^{-1}$ , and therefore  $A \in \mathcal{L}(\mathbb{C}_\infty^s, \omega, \tau)$ .

Applying (p-2) (with  $s=1$ ) we see that all  $\psi_i(\zeta)$  are regular in a neighbourhood of  $\zeta = \lambda$ . Since the spectrum of  $E$  equals  $\{0\}$  it follows (cf. [6, p. 568]) that all  $\psi_i(A) = \psi_i(\lambda + (\omega - \lambda)E)$  exist in  $\mathcal{B}(\mathbb{C}^s)$  and

$$\psi_i(A) = \beta_{i_0} + \beta_{i_1}(\omega - \lambda)E + \beta_{i_2}(\omega - \lambda)^2 E^2 + \dots$$

with  $\beta_{ij} = \psi_i^{(j)}(\lambda)/j!$ .

It follows that  $A$  satisfies the assumptions stated in (p-2). We thus obtain, by using Lemma 4.1 with  $T_i = \psi_i(A)$ ,

$$\sum_{i=1}^k \sum_{j=0}^{s-1} |\beta_{ij}|(\omega - \lambda)^j = \max_{1 \leq m \leq s} \sum_{i=1}^k \sum_{n=m}^s |\beta_{i, n-m}|(\omega - \lambda)^{n-m} \leq \sum_{i=1}^k \psi_i(\omega)$$

for  $s = 1, 2, 3, \dots$ . Hence

$$\sum_{i=1}^k \sum_{j=0}^{\infty} |\beta_{ij}|(\omega - \lambda)^j \leq \sum_{i=1}^k \sum_{j=0}^{\infty} \beta_{ij}(\omega - \lambda)^j,$$

and therefore  $\beta_{ij} \geq 0$  ( $1 \leq i \leq k, j \geq 0$ ). This completes the proof.  $\square$

### 4.3. Proof of Theorem 2.4

1. Clearly (P1)  $\Rightarrow$  (P2). Further for any  $\tau > 0$  with  $-\tau^{-1} < \omega$  proposition (P2) implies, in view of Theorem 3.1, proposition (p-2). By virtue of Theorem 3.3, (P2) thus implies (p-3). By letting  $\tau \rightarrow 0+$  it follows that (P2)  $\Rightarrow$  (P3).

2. We assume (P3), and we shall prove (P1).

Let  $i$  be an integer with  $1 \leq i \leq k$  and write  $\psi = \psi_i$ . From (P3) it follows that  $\psi(\zeta)$  is holomorphic and satisfies  $|\psi(\zeta)| \leq \psi(\omega)$  for  $\zeta \in \mathbb{C}$ ,  $\text{Re } \zeta \leq \omega$ .

Let  $X$  be a complex Banach space and  $A \in \mathcal{L}(X, \omega)$ . It follows (see e.g. [10, pp. 279–283]) that for all  $\zeta \in \mathbb{C}$  with  $\text{Re } \zeta > \omega$  we have

$$\zeta \in \rho(A), \quad \|R(\zeta)\| \leq [(\text{Re } \zeta) - \omega]^{-1}$$

where we use the notation  $R(\zeta) = (A - \zeta)^{-1}$  for the resolvent of  $A$ .

Using the operational calculus as given in [6, pp. 599-604] it can be seen that  $\psi(A)$  exists in  $\mathcal{B}(X)$  (according to the definition given in Sect. 2.2) and satisfies

$$(4.1) \quad \psi(A) = \psi(\infty) - \frac{1}{2\pi i} \int_{\Gamma} \psi(\zeta) R(\zeta) d\zeta.$$

Here  $\Gamma$  denotes a circle with negative orientation enclosing all poles of  $\psi(\zeta)$  and lying strictly to the right of  $\zeta = \omega$ .

With  $B_\lambda$  we denote the Yosida-approximation of  $A$  given by

$$B_\lambda = -\lambda - \lambda^2 R(\lambda) \quad (\text{for } \lambda > \max(0, \omega)),$$

and with  $A_\lambda$  we denote the approximation

$$A_\lambda = B_\lambda - \varepsilon_\lambda \quad (\text{for } \lambda > \max(0, \omega))$$

where  $\varepsilon_\lambda = \omega^2(\lambda - \omega)^{-1}$ . The term  $\varepsilon_\lambda$  has been subtracted from  $B_\lambda$  here since it implies

$$A_\lambda \in \mathcal{L}(X, \omega, \tau)$$

for some  $\tau \in (0, \infty)$  (a little calculation proves that  $A_\lambda$  has property (ii) of Sect. 3.1 with  $\tau = (\lambda - \omega)(\lambda^2 - \lambda\omega + \omega^2)^{-1}$ ).

Since proposition (p-3) of Theorem 3.3 is true, we obtain by an application of this theorem to  $A_\lambda$  the inequality

$$\|\psi(A_\lambda)\| \leq \psi(\omega).$$

In view of (4.1) there follows

$$\|\psi(A)\| \leq \psi(\omega) + \alpha_\lambda$$

with

$$\alpha_\lambda = \frac{1}{2\pi} \left\| \int_{\Gamma} \psi(\zeta) C_\lambda(\zeta) d\zeta \right\|, \quad C_\lambda(\zeta) = (A - \zeta)^{-1} - (A_\lambda - \zeta)^{-1}.$$

We shall prove below that

$$(4.2) \quad \lim_{\lambda \rightarrow \infty} \|C_\lambda(\zeta)\| = 0 \quad (\text{uniformly for } \zeta \in \Gamma).$$

Consequently  $\lim_{\lambda \rightarrow \infty} \alpha_\lambda = 0$ , and  $\|\psi(A)\| \leq \psi(\omega)$ . Recalling the definition of  $\psi$  we see that (P1) holds.

3. We complete the proof by showing (4.2). Clearly

$$\begin{aligned} -(A_\lambda - \zeta)^{-1} &= [(\lambda^2 + (\lambda + \varepsilon_\lambda + \zeta)(A - \lambda))(A - \lambda)^{-1}]^{-1} \\ &= (\lambda + \varepsilon_\lambda + \zeta)^{-1} (A - \lambda)(A - \tilde{\zeta})^{-1} \end{aligned}$$

where

$$\tilde{\zeta} = (\lambda + \varepsilon_\lambda + \zeta)^{-1}(\varepsilon_\lambda + \zeta)\lambda$$

tends to  $\zeta$  for  $\lambda \rightarrow \infty$  (uniformly for  $\zeta \in \Gamma$ ).

By writing  $A - \lambda = (A - \tilde{\zeta}) + (\tilde{\zeta} - \lambda)$  we obtain

$$-(A_\lambda - \zeta)^{-1} = (\lambda + \varepsilon_\lambda + \zeta)^{-1} (I + (\tilde{\zeta} - \lambda) R(\tilde{\zeta})).$$

An application of the resolvent identity (see e.g. [10, p. 74])

$$R(\zeta) = R(\tilde{\zeta}) + (\zeta - \tilde{\zeta}) R(\zeta) R(\tilde{\zeta})$$

now yields

$$\begin{aligned} C_\lambda(\zeta) &= [1 - \lambda(\lambda + \varepsilon_\lambda + \zeta)^{-1}] R(\tilde{\zeta}) + (\lambda + \varepsilon_\lambda + \zeta)^{-1} (I + \tilde{\zeta} R(\tilde{\zeta})) \\ &\quad + (\zeta - \tilde{\zeta}) R(\zeta) R(\tilde{\zeta}). \end{aligned}$$

Let  $\delta > 0$  be such that  $\operatorname{Re} \zeta \geq \omega + \delta$  (for  $\zeta \in \Gamma$ ). Then for  $\lambda$  sufficiently large we have

$$\|R(\zeta)\| \leq \delta^{-1}, \quad \|R(\tilde{\zeta})\| \leq 2\delta^{-1} \quad (\text{for } \zeta \in \Gamma),$$

and it follows that (4.2) holds.  $\square$

*Acknowledgements.* I wish to thank A. Iserles for useful discussions in connection with the example (3.7). I am also indebted to C.B. Huijsmans for calling my attention to reference [15]. Finally I thank the referees and H. Kraaijevanger for comments on the manuscript, which have resulted in an improved presentation.

## References

1. Bolley, C., Crouzeix, M.: Conservation de la positivité lors de la discrétisation des problèmes d'évolution paraboliques. R.A.I.R.O. Analyse Numérique **12**, 237–245 (1978)
2. Brenner, P., Thomée, V.: On rational approximations of semigroups. SIAM J. Numer. Anal. **16**, 683–694 (1979)
3. Burrage, K., Butcher, J.C.: Stability criteria for implicit Runge-Kutta methods. SIAM J. Numer. Anal. **16**, 46–57 (1979)
4. Crouzeix, M.: Sur la  $B$ -stabilité des méthodes de Runge-Kutta. Numer. Math. **32**, 75–82 (1979)
5. Dahlquist, G., Jeltsch, R.: Generalized disks of contractivity for explicit and implicit Runge-Kutta methods. Report TRITA-NA-7906. Dept. Comp. Sci., Roy. Inst. of Techn., Stockholm 1979
6. Dunford, N., Schwartz, J.T.: Linear operators, Part I. New York: Interscience Publishers, Inc. 1958
7. Iserles, A.: Private communication (1980)
8. Jeltsch, R., Nevanlinna, O.: Stability of explicit time discretizations for solving initial value problems. Numer. Math. **37**, 61–69 (1981)
9. Kato, T.: Perturbation theory for linear operators. Berlin, Heidelberg, New York: Springer 1966
10. Martin, R.H.: Nonlinear operators and differential equations in Banach spaces. New York: J. Wiley and Sons 1976
11. Mitchell, A.R., Griffiths, D.F.: The finite difference method in partial differential equations. Chichester: John Wiley and Sons 1980
12. Nevanlinna, O., Liniger, W.: Contractive methods for stiff differential equations. BIT **18**, 457–474 (1978); BIT **19**, 53–72 (1979)
13. Nørsett, S.P.: Restricted Padé approximations to the exponential function. SIAM J. Numer. Anal. **15**, 1008–1029 (1978)

14. Rektorys, K.: Solution of mixed boundary value problems by the method of discretization in time. In: Numerische Behandlung von Differentialgleichungen Band 3. Albrecht, J., Collatz, L. (eds.) pp. 132-145, Basel: Birkhäuser Verlag 1981
15. Rickart, C.E.: General theory of Banach algebras. New York: Van Nostrand 1960
16. Rothe, E.: Zweidimensionale parabolische Randwertaufgaben als Grenzfall eindimensionaler Randwertaufgaben. Math. Anal. **102**, 650-670 (1930)
17. Spijker, M.N.: Contractivity of Runge-Kutta methods. In: Numerical methods for solving stiff initial value problems. Proceedings, Oberwolfach, 28.6.-4.7.1981. Dahlquist, G., Jeltsch, R. (eds.). Institut für Geometrie und Praktische Mathematik der RWTH Aachen, Bericht Nr. 9, 1981
18. Stetter, H.J.: Analysis of discretization methods for ordinary differential equations. Berlin, Heidelberg, New York: Springer 1973

Received May 11, 1982 / May 10, 1983