

# Aplikace matematiky

---

Tomáš Roubíček

Unconditional stability of difference formulas

*Aplikace matematiky*, Vol. 28 (1983), No. 2, 81–90

Persistent URL: <http://dml.cz/dmlcz/104008>

## Terms of use:

© Institute of Mathematics AS CR, 1983

Institute of Mathematics of the Czech Academy of Sciences provides access to digitized documents strictly for personal use. Each copy of any part of this document must contain these *Terms of use*.



This document has been digitized, optimized for electronic delivery and stamped with digital signature within the project *DML-CZ: The Czech Digital Mathematics Library* <http://dml.cz>

UNCONDITIONAL STABILITY OF DIFFERENCE FORMULAS

TOMÁŠ ROUBÍČEK

(Received August 27, 1980)

1. INTRODUCTION

Consider the initial-value problem

$$(1) \quad \frac{dy(t)}{dt} = Ay(t), \quad t > 0,$$

where  $A$  is a linear operator in a complex Banach space  $X$  with a domain  $D_A \subset X$ . We shall deal with the case when the problem (1) is a partial differential equation. Let us consider the numerical solution of this equation by the current finite difference method. Let a sequence of matrices  $A_n$  be given, where  $A_n$  represent the discretization of the original operator  $A$ , obtained e.g. by transcribing the problem into the difference form. The order of the matrix  $A_n$  is equal to the number of meshes of the net, say  $m(n)$ . Then  $A_n : C^{m(n)} \rightarrow C^{m(n)}$ , where  $C$  is the field of complex numbers. The numerical integration is based on a  $k$ -step formula ( $k \geq 1$ ) in the form

$$(2) \quad p_0(\Delta t \cdot A_n) u_j = p_1(\Delta t \cdot A_n) u_{j-1} + \dots + p_k(\Delta t \cdot A_n) u_{j-k},$$

$\Delta t > 0$ , where the vectors  $u$  belong to  $C^{m(n)}$  and  $p_0, \dots, p_k$  are polynomials with no common roots. Setting  $v_j = (u_j, \dots, u_{j-k+1})^T$  (here and throughout the paper, the superscript  $T$  means transposition, i.e.,  $v_j$  is a column vector) and denoting by  $C(n, \Delta t)$  the transition matrix of order  $k \cdot m(n)$ ,  $C(n, \Delta t) = [C_{i,j}]$ , where  $i, j = 1, \dots, k$  and  $C_{i,j}$  are square matrices

$$C_{i,j} = \begin{cases} (p_0(\Delta t \cdot A_n))^{-1} \cdot p_j(\Delta t \cdot A_n) & \text{for } i = 1, \\ I \text{ (unit matrix)} & \text{for } i = j + 1, \\ 0 \text{ (zero matrix)} & \text{otherwise,} \end{cases}$$

we can rewrite the formula (2) in the form  $v_j = C(n, \Delta t) v_{j-1}$ .

**Definition.** *The formula (2) will be called stable for given  $n, \Delta t$ , if the spectrum*

of the transition matrix is inside the unit circle:

$$\sigma(C(n, \Delta t)) \subset \{\lambda; |\lambda| < 1\} \quad (\sigma(\cdot) \text{ denotes the spectrum}).$$

Thus, given a sequence  $A_n$ , the stability of the formula (2) in general depends on  $n$  and  $\Delta t$ . If stability of the formula is required, then often  $\Delta t$  is bounded from above with increasing  $n$ . Then we speak about the so-called stability condition. As a rule, this condition is rather restrictive, which impedes the practicability of such formulas. If partial differential equations are involved, it is usually of special importance to have  $\Delta t$  not bounded from above. Such formulas without the stability condition are called unconditionally stable.

This paper deals with unconditional stability of difference formulas. The main results, Theorem 5.1, is deduced from a number of general assumptions, whether on the initial-value problem (1), the sequence  $A_n$  or the formula (2). The character of these assumptions is rather natural, which is of importance for the applications.

## 2. ASSUMPTIONS ON THE OPERATOR $A$

We introduce two conditions concerning the operator  $A$ .

a) (solvability of the equation  $Ay = f, f \in X$ ):

(3) The operator  $A^{-1}$  exists and is compact.

Remark. Later on we shall adopt the assumption (7) concerning the sequence  $A_n$ , which will imply that the operator  $A^{-1}$  is a uniform limit of continuous finite dimensional, and hence compact operators, so that  $A^{-1}$  itself must be compact. This is the main reason for introducing the assumption (3). By a finite dimensional operator we mean such an operator whose image is of finite dimension. The term degenerate operator is sometimes used, too.

b) (stability of the initial-value problem (1)):

Let an initial value  $y_0 \in D_A$  be given. If there is a strongly differentiable function  $y: (0, +\infty) \rightarrow D_A$  such that  $dy(t)/dt = Ay(t)$  and  $\lim_{t \rightarrow 0+} y(t) = y_0$  uniformly, then

(4) there exists a constant  $K < +\infty$  such that  $\|y(t)\| \leq K$  for  $t > 0$ , where  $\|\cdot\|$  stands for the norm in the Banach space  $X$ .

**Theorem 2.1.** *Let the conditions (3) and (4) be fulfilled. Then*

$$\sigma(A) \subset \{\lambda; \operatorname{Re} \lambda \leq 0\}.$$

*Proof.* By the assumption (3), the operator  $A$  has a compact resolvent at the point 0 and hence it has a compact resolvent at each point not belonging to its spectrum. An operator with a compact resolvent obviously has a purely point spectrum (see [4], p. 210).

Assume that  $\lambda \in \sigma(A)$ . Since the spectrum of  $A$  is a point spectrum, there exists  $y_0 \neq 0$  such that  $Ay_0 = \lambda y_0$ . Evidently  $y_0 \in D_A$  and the function  $y(t) = e^{\lambda t} y_0$  satisfies all the assumptions from the condition (4). Consequently, by this condition there exists a finite  $K$  such that  $\|y(t)\| = e^{\operatorname{Re} \lambda t} \|y_0\| < K$ . However, this is possible only if  $\operatorname{Re} \lambda \leq 0$ . Q.E.D.

Remark. By (3), the operator  $A$  has a compact resolvent and hence it is either unbounded (provided the space  $X$  is not finite dimensional) or finite dimensional. In the latter case the space  $X$  has a finite dimension as well (see [4], p. 210). The former situation occurs e.g. for partial differential equations, the latter corresponds to a system of ordinary linear differential equations. In this last case, the situation is essentially simple: the sequence  $A_n$  can be chosen constant,  $A_n = A$ , and we may skip the considerations of the next section.

### 3. ASSUMPTIONS ON THE SEQUENCE OF DISCRETIZATIONS $A_n$

In order to be able to study the convergence of the sequence of matrices  $A_n$ , we have to embed these operators in a natural way into the space  $X$  (see [1], p. 172). Until now the abstract Banach space has been sufficient, however, now it is indispensable to deal with this space in more detail. Let us point out that we have in mind the case when the problem (1) is a partial differential equation, which can be solved numerically by the finite difference method.

Let the equation (1) be given on a certain domain  $\Omega$  of the Euclidean space and let us have a net of  $m(n)$  meshes  $\Omega_n = \{x_1^n, \dots, x_{m(n)}^n\}$ ,  $\Omega_n \subset \bar{\Omega}$  (the bar denotes closure). Functions from  $X$  are defined on  $\bar{\Omega}$  as well. The generalization to vector-valued functions on  $\bar{\Omega}$  in the case of a system of partial differential equations on the domain  $\Omega$  is straightforward and we will leave it to the reader.

Let us introduce the following assumptions concerning the function space  $X$ .

Let  $x \in \bar{\Omega}$  and  $f \in X$ . Then

(5) the relation  $f \mapsto f(x)$  defines a continuous functional on  $X$ .

Further, given a finite set of points  $\{x_i\}$  different from the point  $x$ ,  $\{x_i\} \subset \bar{\Omega}$ , then

(6) there exists  $f \in X$  such that  $f$  vanishes on  $\{x_i\}$  and  $f(x) = 1$ .

The assumption (5) enables us to define, for each net  $\Omega_n$ , a continuous mapping  $\Psi_n : X \rightarrow C^{m(n)}$  by the formula  $f \mapsto (f(x_1^n), \dots, f(x_{m(n)}^n))$ . In virtue of continuity of this mapping, the kernel  $\operatorname{Ker} \Psi_n$  is closed in  $X$ . Since  $\operatorname{Ker} \Psi_n$  has a finite factor-dimension in the space  $X$ , there is a (continuous) projector onto this subspace, and even infinitely many of such projectors provided  $\operatorname{Ker} \Psi_n$  is a proper subspace of the space  $X$ . Let us fix such a projector and denote it by  $P_n : X \rightarrow \operatorname{Ker} \Psi_n$ .

In concrete problems we shall consider a "nice" projector, i.e. one with a possibly

small norm. In particular, if  $X$  is a Hilbert space, then we evidently choose an orthoprojector (its norm being one). While our considerations do not depend on this norm, in concrete problems an unsuitable choice of the sequence  $P_n$  can affect the convergence of the operators in the condition (7). In what follows, denote  $X_n = \text{Ker } P_n$  and let  $\psi_n : X_n \rightarrow C^{m(n)}$  be the restriction of the mapping  $\Psi_n$  to the subspace  $X_n$ .

**Theorem 3.1.** *The mapping  $\psi_n$  is an isomorphism, i.e.*

$$\psi_n : X_n \cong C^{m(n)}.$$

*Proof.* Let  $\psi_n f = 0$  for some  $f \in X_n = \text{Ker } P_n$ . Then  $\Psi_n f = 0$  as well and  $f \in \text{Ker } \Psi_n = \text{Im } P_n$ . Since  $X = \text{Im } P_n \oplus \text{Ker } P_n$ , we have  $f = 0$  and the mapping  $\psi_n$  is injective.

In virtue of (6) the space  $X$  contains sufficiently many functions for  $\Psi_n$  to be surjective. Consequently, for an arbitrary vector  $\varphi \in C^{m(n)}$  there exists  $f \in X$  such that  $\Psi_n f = \varphi$ . Set  $f_0 = (I - P_n)f$ . Then we have  $f_0 \in X_n$  and  $\psi_n f_0 = \Psi_n f_0 = \Psi_n(I - P_n)f = \Psi_n f = \varphi$ , and since  $\varphi$  was arbitrary,  $\psi_n$  is surjective. Q.E.D.

Thus the fulfilment of the assumptions (5) and (6) guarantees the existence of the mapping  $\psi_n^{-1}$ , which in a natural way embeds  $C^{m(n)}$  into  $X$ . The situation actually corresponds to the interpolation of a function whose values at the nodes are prescribed in advance. Now we can formulate the assumption concerning the sequence of matrices  $A_n$ :

For sufficiently large  $n$  let the matrices  $A_n$  be regular and let the sequence of operators  $\psi_n^{-1} A_n^{-1} \Psi_n : X \rightarrow X_n$  uniformly converge to  $A^{-1}$ , i.e.

$$(7) \quad \lim_{n \rightarrow \infty} \|\psi_n^{-1} A_n^{-1} \Psi_n - A^{-1}\| = 0$$

( $\|\cdot\|$  here stands for the operator norm in the Banach algebra of bounded operators in  $X$ ).

Further, let us define the extended spectrum of the operator  $A$  as usual:

$$\sigma_e(A) = \begin{cases} \sigma(A) & \text{if } A \text{ is bounded,} \\ \sigma(A) \cup \{\infty\} & \text{if } A \text{ is unbounded.} \end{cases}$$

**Theorem 3.2.** *Let the assumptions (3), (5), (6) and (7) be fulfilled. Given a neighborhood  $\Sigma$  of  $\sigma_e(A)$  (with respect to the topology of the closed complex plane, i.e. a one-point compactification of the complex plane), then there is such  $n_0$  that  $\sigma(A_n) \subset \Sigma$  for  $n \geq n_0$ .*

*Proof.* The assumption (3) implies that  $A$  is closed, and since the function  $1/z$  is holomorphic in a neighborhood of the spectrum  $\sigma(A)$ , we obtain using the operator calculus (cf. [4], p. 103) that  $\sigma(A^{-1}) = 1/\sigma_e(A)$ . Hence  $1/\Sigma$  is a neighborhood of  $\sigma(A^{-1})$ .

It is well known that the spectrum is an upper semicontinuous function of the

operator with respect to the uniform operator topology (see [4], p. 167), hence by the assumption (7) we have  $\sigma(\psi_n^{-1}A_n^{-1}\Psi_n) \subset 1/\Sigma$  for  $n$  sufficiently large.

Now we shall prove that  $\sigma(A_n^{-1}) \subset \sigma(\psi_n^{-1}A_n^{-1}\Psi_n)$ . Let  $\lambda \in \sigma(A_n^{-1})$ . The matrix  $A_n^{-1}$  has only the point spectrum, and hence there exists a nonzero  $u \in C^{m(n)}$  such that  $A_n^{-1}u = \lambda u$ . From the identities  $\psi_n^{-1}A_n^{-1}\Psi_n\psi_n^{-1}u = \psi_n^{-1}A_n^{-1}u = \lambda\psi_n^{-1}u$  we conclude that  $\psi_n^{-1}u$  is an eigenvector of the operator  $\psi_n^{-1}A_n^{-1}\Psi_n$  corresponding to the eigenvalue  $\lambda$ , hence  $\lambda \in \sigma(\psi_n^{-1}A_n^{-1}\Psi_n)$ .

Consequently,  $\sigma(A_n^{-1}) \subset 1/\Sigma$  and  $\sigma(A_n) = 1/\sigma(A_n^{-1}) \subset \Sigma$ . Q.E.D.

*Remark.* The assumption (5) in fact represents a certain smoothness requirement for the right hand side of the nonhomogeneous initial-value problem  $dy/dt - Ay = f$ . Most frequently, it is the space of functions integrable on  $\Omega$  or a space of Sobolev distribution on  $\Omega$  which are used as the space  $X$ . However, in such spaces it is meaningless to define the mapping  $\Psi_n$ . Appropriate function spaces are e.g. the space of uniformly continuous bounded functions on  $\bar{\Omega}$  or a Sobolev space with a sufficiently strong topology.

#### 4. ASSUMPTIONS ON THE INTEGRATION FORMULA

First we shall prove a theorem on the mapping of the spectrum of the matrix  $\Delta t \cdot A_n$  by the integration formula (2). Put  $\Delta t = 1$  for simplicity.

**Theorem 4.1.** *Let  $(p_0(A_n))^{-1}$  exist. Let  $\lambda \in \sigma(A_n)$  and let  $\lambda_c$  satisfy*

$$(8) \quad \lambda_c^k \cdot p_0(\lambda) = \lambda_c^{k-1} \cdot p_1(\lambda) + \dots + p_k(\lambda).$$

*Then  $\lambda_c \in \sigma(C(n, 1))$ .*

*On the other hand, if  $\lambda_c \in \sigma(C(n, 1))$ , then there exists  $\lambda \in \sigma(A_n)$  such that the identity (8) holds.*

*Proof.* Denote  $Z(\lambda_c) = \lambda_c^k \cdot p_0(A_n) - \lambda_c^{k-1} \cdot p_1(A_n) - \dots - p_k(A_n)$ . The inclusion  $\lambda_c \in \sigma(C(n, 1))$  holds if and only if there is a nonzero vector  $v$  with  $C(n, 1)v = \lambda_c v$ . Calculating the indicated product we find that  $v$  has the form  $v = (\lambda_c^{k-1}u, \dots, \lambda_c u, u)^T$ , where  $Z(\lambda_c)u = 0$ . Hence  $\lambda_c \in \sigma(C(n, 1))$  if and only if the matrix  $Z(\lambda_c)$  is not regular.

Let us express the matrix  $A_n$  in the form  $PJP^{-1}$ , where the matrix  $J$  is in the Jordan normal form, so that its diagonal is formed exactly by the eigenvalues of the matrix  $A_n$ . Then

$$Z(\lambda_c) = P[\lambda_c^k \cdot p_0(J) - \lambda_c^{k-1} \cdot p_1(J) - \dots - p_k(J)]P^{-1}.$$

Since  $J$  is in the lower triangular form, the matrix in the brackets is in the lower triangular (but generally not Jordan) form as well, and its diagonal consists exactly of the numbers  $\lambda_c^k \cdot p_0(\lambda) - \lambda_c^{k-1} \cdot p_1(\lambda) - \dots - p_k(\lambda)$ , where  $\lambda$  are the eigenvalues of the matrix  $A_n$ . Consequently, the matrix  $Z(\lambda_c)$  is not regular if and only if there is  $\lambda \in \sigma(A_n)$  satisfying the relation (8). Q.E.D.

Hence the formula (2) maps each point of the spectrum  $\sigma(A_n)$  generally onto  $k$  points of the spectrum  $\sigma(C(n, 1))$ . Define a  $k$ -valued complex function, notation  $F : C \rightarrow 2^C$ , by the relation  $F(z) = \{z_1, \dots, z_k\}$ , where  $z_i$  are exactly  $k$  roots of the equation  $z_i^k \cdot p_0(z) = z_i^{k-1} \cdot p_1(z) + \dots + p_k(z)$ . Theorem 4.1 then asserts that  $\sigma(C(n, 1)) = F(\sigma(A_n))$ , more generally

$$(9) \quad \sigma(C(n, \Delta t)) = F(\sigma(\Delta t \cdot A_n)).$$

Further, let us define a nonnegative function  $F_m = \max \{|z_1|, \dots, |z_k|\}$ .  $F_m$  is a continuous function except for those points  $z$  for which  $p_0(z) = 0$  and possibly except the point  $\infty$  (provided some of the polynomials  $p_1, \dots, p_k$  have a higher degree than  $p_0$ ).  $F_m$  is unbounded in neighborhoods of these points, and at the points we set  $F_m = +\infty$ . Now we introduce an important property of the formula (2). It will be obvious that this property depends only on the polynomials  $p_0, \dots, p_k$ .

**Definition.** The formula (2) is called *A-acceptable*, if

$$\begin{aligned} F_m(z) &\leq 1 \quad \text{for } \operatorname{Re} z \leq 0, \\ F_m(\infty) &\leq 1. \end{aligned}$$

In literature, the *A-acceptability* has been defined only for one-step formulas ( $k = 1$ ) and, moreover, consistency of the function  $F$  with the function  $e^z$  at the point 0 is required. Here we only deal stability leaving aside the problem of consistency, even if we do tacitly assume the consistency of the formula (2). The consistency implies that at the point 0 one of the values of the function  $F$  equals one, i.e.  $F_m(0) \geq 1$ . If, moreover, the formula is *A-acceptable*, then  $F_m(0) = 1$ .

It is evident from Theorem 3.2 that due to the generality of the above introduced assumptions the spectrum of the matrix  $A_n$  may lie in the right-hand halfplane as well. Then the *A-acceptability* does not guarantee the stability of the formula (2) even if the definition of stability of the formula (2) admitted eigenvalues of the transition matrix on the boundary of the unit circle. This leads us to the definition of a new notion.

**Definition.** The formula (2) is called *A<sub>n</sub>-acceptable*, if

$$\begin{aligned} F_m(z) &< 1 \quad \text{for } \operatorname{Re} z \leq 0, \quad z \neq 0, \\ F_m(\infty) &< 1. \end{aligned}$$

If the formula (2) is *A<sub>n</sub>-acceptable*, then it is *A-acceptable* as well. Further, it is evident that the roots of the polynomial  $p_0$  of an *A-acceptable* formula lie in the open right-hand halfplane.

## 5. PROOF OF UNCONDITIONAL STABILITY

The next theorem represents the result at which all the foregoing considerations have been directed.

**Theorem 5.1.** *Let  $\varepsilon > 0$ . Let the following conditions be fulfilled: The formula (2) is  $A_n$ -acceptable, the operator  $A$  satisfies (3), the initial-value problem is stable in the sense of the condition (4), the function space  $X$  satisfies the conditions (5) and (6) and the sequence of matrices  $A_n$  approximates the operator  $A$  in the sense of the condition (7).*

*Then there is  $n_0$  such that for  $n \geq n_0$  and  $\Delta t \geq \varepsilon$  the formula (2) is stable.*

*Proof.* Given  $\delta$  such that  $0 < \delta < 1/\|A^{-1}\|$ , then evidently  $\{\lambda; |\lambda| > \delta\}$  is a neighborhood of  $\sigma_e(A)$ . The set  $S = \{\lambda; \operatorname{Re} \lambda \leq 0 \text{ and } |\lambda| \geq \varepsilon\delta\} \cup \{\infty\}$  is compact with respect to the topology of the closed complex plane. Since the formula (2) is  $A_n$ -acceptable,  $F_m$  is continuous on  $S$  and satisfies  $F_m < 1$  on  $S$ . The compactness of  $S$  implies existence of such  $\delta_1$  that  $F_m \leq \delta_1 < 1$  on  $S$ . Denote  $S_1 = \{\lambda; F_m(\lambda) < 1\}$ . In virtue of continuity of the function  $F_m$  at all points at which  $F_m$  is finite we obtain that  $S_1$  is an open neighborhood of the set  $S$ .

Consider the set  $\Sigma = \{\lambda; |\lambda - \xi| > \sqrt{(\xi^2 + \varepsilon^2\delta^2)}\}$  with  $\xi > 0$ .  $\Sigma$  is the exterior of a circle with its center at  $\xi$  and with the points  $\pm i\varepsilon\delta$  lying on its boundary. In virtue of Theorem 2.1 and of the choice of  $\delta$ , the set  $\Sigma$  is a neighborhood of  $\sigma_e(\varepsilon A)$ .

We shall prove that  $F_m < 1$  on  $\Sigma$  for  $\xi$  sufficiently large. Let us put  $\Sigma = \Sigma_1 \cup \Sigma_2$  with  $\Sigma_1 = \Sigma \cap \{\lambda; \operatorname{Re} \lambda \leq 0 \text{ or } \lambda = \infty\}$ ,  $\Sigma_2 = \Sigma \cap \{\lambda; \operatorname{Re} \lambda > 0\}$ . We have  $F_m < 1$  on  $\Sigma_1$  since the formula (2) is  $A_n$ -acceptable and  $0 \notin \Sigma$ . As for  $\Sigma_2$ , for every neighborhood of the set  $S$  (naturally in the topology of the closed complex plane) we can choose  $\xi$  so large that  $\Sigma_2$  belongs to this neighborhood. In particular, there is such  $\xi < +\infty$  that  $\Sigma_2 \subset S_1$ , and hence  $F_m < 1$  on  $\Sigma_2$ .

Making use of Theorem 3.2 we conclude that there is  $n_0$  such that  $\sigma(\varepsilon A_n) \subset \Sigma$  for  $n \geq n_0$ . However,  $\sigma(\Delta t \cdot A_n) = \Delta t \varepsilon^{-1} \cdot \sigma(\varepsilon A_n) \subset \Delta t \varepsilon^{-1} \cdot \Sigma$ . Further, taking into account the geometric shape of the set  $\Sigma$  and its location with respect to the point 0 we obtain for  $\Delta t \geq \varepsilon$  the inclusion  $\Delta t \varepsilon^{-1} \Sigma \subset \Sigma$ , that is,  $\sigma(\Delta t \cdot A_n) \subset \Sigma$ . Thus  $F_m < 1$  on  $\sigma(\Delta t \cdot A_n)$ , and the relation (9) together with the definition of the function  $F_m$  complete the proof of the stability of the formula (2). Q.E.D.

## 6. EXAMPLES OF $A_n$ -ACCEPTABLE FORMULAS

**Example 1.** As the first example let us consider the usual one-step formula, which is frequently used. It is given by the polynomials

$$p_0(z) = 1 - \Theta z, \quad p_1(z) = 1 + (1 - \Theta)z.$$

It is easily seen that for  $\Theta > \frac{1}{2}$  this formula is  $A_n$ -acceptable, for  $\Theta = \frac{1}{2}$   $A$ -acceptable but not  $A_n$ -acceptable, while for  $\Theta < \frac{1}{2}$  it is neither  $A_n$ - nor  $A$ -acceptable.

**Example 2.** Another example of one-step formula, in this case of higher order, are the rational Padé approximations of the function  $e^z$ . Let nonnegative integers  $p, q$  stand for the degrees of the polynomials  $p_0$  and  $p_1$ , respectively. These polynomials

are given by the relations

$$p_0(z) = \sum_{j=0}^p \frac{(p+q-j)! p!}{(p+q)! j! (p-j)!} (-z)^j,$$

$$p_1(z) = \sum_{j=0}^q \frac{(p+q-j)! q!}{(p+q)! j! (q-j)!} z^j.$$

The formula is of the order  $p+q$ . These formulas have lately attracted considerable attention. It is evident that for  $p < q$  the formula cannot be  $A$ -acceptable. Further, it is easily seen that for  $p = q$  the formula is  $A$ -acceptable but not  $A_n$ -acceptable. The  $A$ -acceptability for  $p = q + 1$  and  $p = q + 2$  is proved in [2], and the proof easily yields that such formulas are even  $A_n$ -acceptable. On the other hand, in [3] it is proved that for  $p \geq q + 3$  and  $p \neq q + 2 \pmod{4}$  the formula is not  $A$ -acceptable (and, a fortiori, it is not  $A_n$ -acceptable). There are no results available concerning  $p = q + 2 \pmod{4}$ , however, it seems probable that in this case the formula is not  $A$ -acceptable, either.

Example 3. When applying the formula (2) we have to determine the vector  $u_j$  by solving a system of linear algebraic equations with the matrix  $p_0(\Delta t \cdot A_n)$ . For partial differential equations solved by the finite difference method the matrices  $A_n$  usually are sparse, but the matrices  $p_0(\Delta \cdot A_n)$  gradually lose this property when the degree of the polynomial  $p_0$  increases. This may be a drawback from the numerical point of view. The increase of the degree of the polynomials  $p_1, \dots, p_k$  also complicates the evaluation of the right hand side of the system of equations. The following two-step formula of the second order is very advantageous. It is given by the polynomials

$$p_0(z) = \frac{3}{2} - z, \quad p_1(z) = 2, \quad p_2(z) = -\frac{1}{2}.$$

Another advantage of the formula consists in its  $A_n$ -acceptability. In spite of the simplicity of the given formula, the proof of its  $A_n$ -acceptability is relatively difficult, since it is considerably complicated by the fact that the formula is not a one-step formula.

The values of the function  $F$  are  $z_{1,2} = (2 \pm \sqrt{(1+2z)})/(3-2z)$ . Obviously they are holomorphic functions of the argument  $z$  except for the points  $z = \frac{3}{2}$  (the root of the polynomial  $p_0$ ) and  $z = -\frac{1}{2}$  where  $z_1 = z_2$ .

Consider the set  $M = \{z; \operatorname{Re} z < 0\} \setminus \{z; z \leq -\frac{1}{2}\}$ . On this domain there exist univalent branches of the square root  $\sqrt{(1+2z)}$ . In the sequel, the symbol  $\sqrt{\quad}$  will stand for the branch of the square root with  $\sqrt{1} = 1$ . Let us distinguish the values of the function  $F : z_1 = (2 + \sqrt{(1+2z)})/(3-2z), z_2 = (2 - \sqrt{(1+2z)})/(3-2z)$ . Our aim is to show that  $|z_1| < 1, |z_2| < 1$  on the boundary of  $M$  except the point 0, where  $z_1(0) = 1$ . Since  $z_1, z_2$  are holomorphic in the domain  $M$  and continuous on the boundary of  $M$ , by virtue of the Maximum Modulus Principle this will prove the  $A_n$ -acceptability of the formula in question.

The boundary of  $M$  is the set  $\{\infty\} \cup \{z; z \leq -\frac{1}{2}\} \cup \{z; \operatorname{Re} z = 0\}$ . For the point  $\infty$  we have  $z_1(\infty) = z_2(\infty) = 0$ , and we easily find that  $|z_1(z)| = |z_2(z)| = 1/\sqrt{(3-2z)} \leq \frac{1}{2} < 1$  for  $z \leq -\frac{1}{2}$ . The situation is a little more complicated on the imaginary axis ( $\operatorname{Re} z = 0$ ), i.e. for  $z$  of the form  $z = i\omega$  with  $\omega$  real. Denote  $a + bi = \sqrt{(1 + 2i\omega)}$ ,  $a \geq 1$ ,  $b$  real. (Thus  $a^2 + b^2 = \sqrt{(1 + 4\omega^2)}$ .) Getting rid of the square root and comparing the real and imaginary parts we obtain the system of equations  $a^2 - b^2 = 1$ ,  $ab = \omega$  for  $a, b$ . The solution satisfying the inequality  $a \geq 1$  is  $a = \sqrt{(\frac{1}{2}(1 + \sqrt{(1 + 4\omega^2))})}$  and we further have

$$\begin{aligned} |2 \pm \sqrt{(1 + 2i\omega)}|^2 &= (2 \pm a)^2 + b^2 = 4 \pm 4a + a^2 + b^2 = \\ &= 4 \pm 4\sqrt{(\frac{1}{2}(1 + \sqrt{(1 + 4\omega^2))})} + \sqrt{(1 + 4\omega^2)}. \end{aligned}$$

Moreover

$$|z_1(i\omega)|^2 = (4 + 4\sqrt{(\frac{1}{2}(1 + \sqrt{(1 + 4\omega^2))})} + \sqrt{(1 + 4\omega^2))}/(9 + 4\omega^2)$$

and

$$|z_2(i\omega)|^2 = (4 - 4\sqrt{(\frac{1}{2}(1 + \sqrt{(1 + 4\omega^2))})} + \sqrt{(1 + 4\omega^2))}/(9 + 4\omega^2)$$

Now we already use the standard argument to show that  $\omega \neq 0$  fulfils the inequalities  $0 < |z_2(i\omega)|^2 < |z_1(i\omega)|^2 < 1$ . This completes the proof of the  $A_n$ -acceptability.

It is evident that this way of the discussion of the  $A_n$ -acceptability complicates very considerably with increasing  $k$ . Already for  $k = 3$  it is impracticable while for  $k > 4$  it is in principle inexecutable since there is no analytical description available for the values of the function  $F$ .

#### References

- [1] *I. Babuška, M. Práger and E. Vitásek: Numerical Processes in Differential Equations, SNTL, Prague, 1966.*
- [2] *B. L. Ehle: A-stable Methods and Padé Approximations to the Exponential. SIAM J. Math. Anal., Vol. 4 (1973), No. 4, pp. 671–680.*
- [3] *A. Iserles: On the A-acceptability of Padé Approximations. SIAM J. Math. Anal., Vol. 10 (1979), No. 5, pp. 1002–1007.*
- [4] *E. Hille and R. S. Phillips: Functional Analysis and Semi-groups. Amer. Math. Soc., Vol. 31., rev. ed., Waverly Press, Baltimore, 1957.*

#### Souhrn

### BEZPODMÍNEČNÁ STABILITA DIFERENČNÍCH FORMULÍ

TOMÁŠ ROUBÍČEK

Mějme dány parciální diferenciální rovnici evolučního typu a řešme ji přibližně metodou sítí. Článek se zabývá obecnými předpoklady jak na původní rovnici,

tak na její diskretizaci, za kterých je diferenční schema bezpodmínečně stabilní, tj. stabilní bez kriteriia stability na časový krok. Zavádí se nový pojem  $A_n$ -akceptability integrační formule a v závěru článku se uvádí několik příkladů takových formulí. Dokázané výsledky lze jednoduše aplikovat i na obyčejné diferenciální rovnice.

*Author's address:* Ing. Tomáš Roubíček, Středisko výpočetní techniky ČSAV, Pod vodárenskou věží 4, 18207 Praha 8.