

Aplikace matematiky

Petr Voňka

Stability and local error of difference methods for the solution of the ordinary differential equation of the first order

Aplikace matematiky, Vol. 17 (1972), No. 1, 18–27

Persistent URL: <http://dml.cz/dmlcz/103388>

Terms of use:

© Institute of Mathematics AS CR, 1972

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>

STABILITY AND LOCAL ERROR OF DIFFERENCE METHODS
FOR THE SOLUTION OF THE ORDINARY DIFFERENTIAL EQUATION
OF THE FIRST ORDER

PETR VOŇKA

(Received May 28, 1970)

The present paper deals with the problem of the construction of the difference formulae for the solution of the ordinary differential equation of the first order from the given characteristic polynomial (Sec. 2), with the effect of the choice of zeros of the characteristic polynomial on the local error of the difference methods for the equation $y' = Ay$, $y(0) = 1$, A being a constant (Sec. 3), and with the dependence of the error constant on the choice of zeros of the characteristic polynomial in the general case (Sec. 4).

1. INTRODUCTION

In what follows, we study some problems of the numerical solution of the differential equation

$$(1) \quad y' = f(x, y) \quad y(a) = y_0$$

by the linear multistep method (1). We assume that the assumptions of the existence theorem (1) for the solution of Eq. (1) are fulfilled. The general multistep method may be expressed in the form

$$(2) \quad \alpha_k y_{n+k} + \alpha_{k-1} y_{n+k-1} + \dots + \alpha_0 y_n = h \{ \beta_k f_{n+k} + \beta_{k-1} f_{n+k-1} + \dots + \beta_0 f_n \}$$

where k is a fixed positive integer, $f_m = f(x_m, y_m)$ ($m = 0, 1, 2, \dots$), α_i, β_i ($i = 0, 1, \dots, k$) are real constants independent of n , $\alpha_k \neq 0$, $|\alpha_0| + |\beta_0| \neq 0$.

Define the polynomials

$$(3) \quad \varrho(\xi) = \alpha_k \xi^k + \alpha_{k-1} \xi^{k-1} + \dots + \alpha_0$$

$$\sigma(\xi) = \beta_k \xi^k + \beta_{k-1} \xi^{k-1} + \dots + \beta_0$$

with the coefficients identical with those of Eq. (2). It is possible to show that the

general linear k -step method is convergent (1) if and only if the following two conditions are satisfied:

(i) All zero points of the polynomial $\varrho(\xi)$ are in their absolute value at most equal to one. Moreover, those with the absolute value equal to one are simple zeros.

(ii) $\varrho(1) = 0$, $\varrho'(1) = \sigma(1)$.

The first condition is called the condition of stability, the second one the condition of consistency. Define the differential operator associated to the relation (2) by

$$(4) \quad L(y(x), h) = \alpha_k y(x + kh) + \alpha_{k-1} y(x + (k-1)h) + \dots + \alpha_0 y(x) \\ - h\{\beta_k y'(x + kh) + \beta_{k-1} y'(x + (k-1)h) + \dots + \beta_0 y'(x)\}$$

Assuming that the differential operator (4) as a function in h has the derivatives of all orders we can expand it with respect to h :

$$(5) \quad L(y(x), h) = C_0 y(x) + C_1 h y'(x) + \dots + C_q h^q y^{(q)}(x) + \dots$$

The coefficients C_q ($q = 0, 1, 2, \dots$) are independent of the choice of $y(x)$. The differential operator (4) will be said to be of degree p if $C_0 = C_1 = \dots = C_p = 0$, $C_{p+1} \neq 0$.

The quantity

$$(6) \quad C = \frac{C_{p+1}}{\beta_0 + \beta_1 + \dots + \beta_k}$$

is the so called error constant. The constant C is invariant with respect to the multiplication of both sides of (2) by an arbitrary real number.

Suppose that an arbitrary but fixed value of k is chosen. Then it holds (1) that the degree of a stable operator (5) (i.e. the polynomial $\varrho(\xi)$ satisfies the condition of stability) cannot exceed the value $p = k + 2$. A necessary and sufficient condition for $p = k + 2$ is that k be even and all zero points of the polynomial $\varrho(\xi)$ have the absolute value equal to one, $\varrho(-1) = 0$.

2. CONSTRUCTION OF THE DIFFERENCE FORMULA (2)

Choose arbitrarily a fixed $k \geq 1$ and put $p = k + 1$. Choose a polynomial of the k -th degree $\varrho(\xi)$ so that it satisfies the condition of stability and that $\varrho(1) = 0$. Now, try to construct generally the polynomial $\sigma(\xi)$ so that the differential operator (4) may have the highest degree possible. Put

$$(7) \quad C_0 = C_1 = \dots = C_p = 0$$

Equating the coefficients at the same powers of h in the relations (4), (5) and substituting into (7) we obtain a system of p equations for p unknowns $\beta_0, \beta_1, \dots, \beta_k$:

$$(8) \quad \begin{aligned} \beta_0 + \beta_1 + \dots + \beta_k &= \alpha_1 + 2\alpha_2 + 3\alpha_3 + \dots + k\alpha_k \\ \beta_1 + 2\beta_2 + \dots + k\beta_k &= \frac{1}{2!}(\alpha_1 + 2^2\alpha_2 + 3^2\alpha_3 + \dots + k^2\alpha_k) \\ \frac{1}{(p-1)!}(\beta_1 + 2^{p-1}\beta_2 + \dots + k^{p-1}\beta_k) &= \frac{1}{p!}(\alpha_1 + 2^p\alpha_2 + 3^p\alpha_3 + \dots + k^p\alpha_k) \end{aligned}$$

The first equation in (8) corresponds to the condition $C_1 = 0$, the second one to $C_2 = 0$ etc. up to the p -th one corresponding to $C_p = 0$. The condition $C_0 = 0$ is fulfilled obviously according to the choice of $\varrho(\xi)$ ($\varrho(1) = 0$). The determinant of the system (8) is non zero for an arbitrary k . The condition of stability implies that the first component of the vector on the right hand side of the system (8) is also non zero. Hence the system (8) has for any $k \geq 1$ exactly one non zero solution. Consequently, to any polynomial $\varrho(\xi)$ of the k -th degree ($k \geq 1$) which satisfies the condition of stability and the relation $\varrho(1) = 0$, there exists exactly one polynomial $\sigma(\zeta)$ of at most the k -th degree, such that the degree of the operator (4) is at least $p = k + 1$. Moreover, the condition $\varrho(1) = 0$ guarantees the convergence of (2). The difference formula (2) whose associated operator has the degree $p = k + 2$ is called optimal.

Example 1. $k = 2$. The solution of the system

$$(9) \quad \begin{aligned} \beta_0 + \beta_1 + \beta_2 &= \alpha_1 + 2\alpha_2 \\ \beta_1 + 2\beta_2 &= \frac{1}{2}(\alpha_1 + 4\alpha_2) \\ \frac{1}{2}(\beta_1 + 4\beta_2) &= \frac{1}{6}(\alpha_1 + 8\alpha_2) \end{aligned}$$

is

$$(10) \quad \begin{aligned} \beta_0 &= \frac{5}{12}\alpha_1 + \frac{1}{3}\alpha_2 \\ \beta_1 &= \frac{2}{3}\alpha_1 + \frac{4}{3}\alpha_2 \\ \beta_2 &= -\frac{1}{12}\alpha_1 + \frac{1}{3}\alpha_2 \end{aligned}$$

The polynomial $\varrho(\xi)$ for $k = 2$ may be written in the form

$$(11) \quad \varrho(\xi) = (\xi - 1)(\xi - \lambda) \quad \lambda \in \langle -1, 1 \rangle$$

This implies

$$(12) \quad \alpha_2 = 1 \quad \alpha_1 = -(1 + \lambda) \quad \alpha_0 = \lambda$$

Substituting (12) into (10), (2) we obtain the general difference formula

$$(13) \quad y_{n+2} - (1 + \lambda)y_{n+1} + \lambda y_n = h \left(\frac{5 + \lambda}{12} f_{n+2} + \frac{2 - 2\lambda}{2} f_{n+1} - \frac{1 + 5\lambda}{12} f_n \right)$$

Substituting (10), (12) into C_3 and dividing by $\beta_0 + \beta_1 + \beta_2$ we obtain the relation for C

$$(14) \quad C = -\frac{1}{24} \frac{1 + \lambda}{1 - \lambda}$$

Substituting $\lambda = -1$ into (13) we obtain the optimal formula whose associated differential operator has the degree 4.

Example 2. $k = 4$. Analogously to Ex. 1 we solve the system (8) for $k = 4$. The solution is

$$(15) \quad \begin{aligned} \beta_0 &= \frac{224\alpha_4 + 243\alpha_3 + 232\alpha_2 + 251\alpha_1}{720} \\ \beta_1 &= \frac{512\alpha_4 + 459\alpha_3 + 496\alpha_2 + 323\alpha_1}{360} \\ \beta_2 &= \frac{64\alpha_4 + 108\alpha_3 + 32\alpha_2 - 44\alpha_1}{120} \\ \beta_3 &= \frac{512\alpha_4 + 189\alpha_3 + 16\alpha_2 + 53\alpha_1}{360} \\ \beta^4 &= \frac{224\alpha_4 - 27\alpha_3 - 8\alpha_2 - 19\alpha_1}{720} \end{aligned}$$

The constant C is evaluated also analogously to Ex. 1:

$$(16) \quad C = \frac{27\alpha_3 + 16\alpha_2 + 27\alpha_1}{1440(4\alpha_4 + 3\alpha_3 + 2\alpha_2 + \alpha_1)}$$

3. EFFECT OF THE CHOICE OF THE ZEROS OF THE POLYNOMIAL ON THE LOCAL ERRORS

The dependence of the local error of Eq. (2) on the choice of zeros of the polynomial $g(\xi)$ for fixed k will be studied in the case of the differential equation

$$(17) \quad y' = Ay \quad y(0) = 1$$

A being a real constant. The local error assumes then the form

$$(18) \quad e_n = y_n - e^{Ax_n}$$

It follows from the theory (1) of difference equations that the solution of the difference equation (2) has the form

$$(19) \quad y_n = \sum_{\mu=1}^k A_{\mu} \xi_{\mu}^n$$

where A_μ are constants and the values $\tilde{\xi}_\mu$ are the zeros of the polynomial

$$(20) \quad \tilde{\varrho}(\xi) = \varrho(\xi) - Ah \sigma(\xi)$$

Let us further assume that the zeros of $\tilde{\varrho}(\xi)$ are simple. We find easily that

$$(21) \quad \begin{aligned} \tilde{\xi}_\mu &= \xi_\mu(1 + \lambda_\mu Ah + O((Ah)^2)) \\ \lambda_\mu &= \frac{\sigma(\xi_\mu)}{\xi_\mu \varrho'(\xi_\mu)} \quad \mu = 1, 2, \dots, n \end{aligned}$$

holds. The values λ_μ are called the growth parameters. Let us number the zeros of the polynomial $\varrho(\xi)$ so that $\xi_1 = 1$. Then

$$(22) \quad \lambda_1 = \frac{\sigma(1)}{\varrho'(1)} = 1$$

The zeros of the polynomial $\varrho(\xi)$ with the absolute value equal to one may be written in the form $\xi_\mu = e^{i\varphi_\mu}$ ($\mu = 1, 2, \dots$). Then it is possible to show (1) that

$$(23) \quad \tilde{\xi}_\mu^n = e^{in\varphi_\mu}(e^{\lambda_\mu Ax} + O(h)) \quad h \rightarrow 0 \quad nh = x$$

In particular, this yields

$$(24) \quad \tilde{\xi}_1^n = e^{Ax} + O(h) \quad h \rightarrow 0$$

It can be proved (1) that the constants A_μ from (19) fulfil

$$(25) \quad \begin{aligned} \lim_{h \rightarrow 0} A_1(h) &= 1 \\ \lim_{h \rightarrow 0} A_\mu(h) &= 0 \quad \mu = 2, 3, \dots, k \\ A_\mu(h) &\neq 0 \quad \mu = 1, 2, \dots, k \end{aligned}$$

Consequently, only the first term on the right hand side of (19) approximates the solution of (17). The other terms of the sum, i.e. $\sum_{\mu=2}^k A_\mu \tilde{\xi}_\mu^n$ form the so called parasitic solution, which in case of an improper choice of the polynomial $\varrho(\xi)$ with respect to the given $A \cdot h$ may cause the results to be quite worthless. (For an example, see (1).) It is evident from (19) that, given Ah , the polynomial $\varrho(\xi)$ (and hence also the form (2)) must be chosen so that

$$(26) \quad |\tilde{\xi}_1| > |\tilde{\xi}_\mu| \quad \mu = 2, 3, \dots, k$$

Example 3. $k = 2$. For $k = 2$ the roots $\tilde{\xi}_1, \tilde{\xi}_2$ are obtained by solving the quadratic equation

$$(27) \quad a\tilde{\xi}^2 + b\tilde{\xi} + c = 0$$

Example 1 and the relation (20) yield

$$(28) \quad \begin{aligned} a &= \frac{12 - 5z - z\lambda}{12} \\ b &= \frac{-3 - 3\lambda - 2z + 2\lambda z}{3} \\ c &= \frac{12\lambda + z + 5\lambda z}{12} \end{aligned}$$

where $z = Ah$. Easy transformations lead us to

$$(29) \quad \lim_{h \rightarrow 0} \xi_{1,2}(h) = \frac{1 + \lambda \pm \sqrt{(1 - \lambda)^2}}{2} = 1 \text{ or } \lambda$$

Let us answer the question how to choose λ for given z so that $|\xi_1| > |\xi_2|$. It follows from (27), (29) that

$$(30) \quad \left| \frac{-b + \sqrt{(b^2 - 4ac)}}{2a} \right| > \left| \frac{-b - \sqrt{(b^2 - 4ac)}}{2a} \right|$$

must hold. If $a > 0$ (i.e. $z < 2$) is assumed, then a discussion of the inequality (30) shows easily that (30) is satisfied if and only if the following two conditions hold simultaneously:

$$(31) \quad \begin{aligned} b &< 0 \\ b^2 - 4ac &> 0 \end{aligned}$$

Substituting (28) into the expression $b^2 - 4ac$ and investigating the resulting quadratic function of λ for an arbitrary and fixed z , we prove easily that the second condition in (31) is fulfilled for any pair (z, λ) , $z \in (-\infty, \infty)$, $\lambda \in (-1, 1)$. The second relation (28) implies that $b < 0$ if

$$(32) \quad \lambda(2z - 3) < (2z + 3)$$

The relation (32) defines the set (z, λ) for which (if $a > 0$) $|\xi_1| > |\xi_2|$ holds. The case $a < 0$ ($z > 2$) is investigated analogously, however, it is not interesting due to the big value of z . It is evident from (32) that it is unsuitable to use the optimal formula for $Ah < 0$.

Example 4. $k = 4$. In this case it is not possible to proceed analogously to the case $k = 2$. The zeros $\xi_1, \xi_2, \xi_3, \xi_4$ of the polynomial $\tilde{q}(\xi)$ may be only estimated by means of (21). Nevertheless, let us show that in case of $Ah < 0$ it is unsuitable to use optimal formulae even for $k = 4$. For the optimal formula, the zero points of the polynomial $q(\xi)$ are

$$(33) \quad \xi_1 = 1 \quad \xi_2 = -1 \quad \xi_3 = q + i\sqrt{(1 - q^2)} \quad \xi_4 = q - i\sqrt{(1 - q^2)}$$

with $q \in (-1, 1)$. It follows from (33) that the coefficients of the polynomial $\varrho(\xi)$ are

$$(34) \quad \alpha_4 = 1 \quad \alpha_3 = -2q \quad \alpha_2 = 0 \quad \alpha_1 = 2q \quad \alpha_0 = -1$$

Substituting (34) into (15) and (21) we find

$$(35) \quad \lambda_2 = -\frac{19 + 11q}{45(1 + q)}$$

$$(36) \quad \lambda_3 = \lambda_4 = \frac{(q - 1)^2}{45(q + 1)}$$

(35) implies that the growth parameter λ_2 corresponding to the zero point $\xi_2 = -1$ is less than $-1/3$ for $k = 4$. Hence it follows from (21) that for small negative values of Ah $|\xi_2| > |\xi_1|$ holds, which is a contradiction to (26). Moreover, let us show that for $Ah > 0$ we can use the optimal formula analogously as for $k = 2$. (35) implies that for small positive Ah it is $|\xi_1| > |\xi_2|$. It holds

$$(37) \quad \tilde{\xi}_3 = (q + i\sqrt{(1 - q^2)}) \left(1 + \frac{(q - 1)^2}{45(q + 1)} Ah + O((Ah)^2) \right)$$

Hence

$$(38) \quad |\tilde{\xi}_3| \doteq 1 + \frac{(q - 1)^2}{45(1 + q)} Ah$$

It follows from (38) that for $q > -0.8$ it is $|\tilde{\xi}_1| > |\tilde{\xi}_3|$. For a suggestion concerning the choice of zero points of the polynomial $\varrho(\xi)$ for $Ah < 0$, as well as for a more detailed discussion of the problem see (2).

4. DEPENDENCE OF THE ERROR CONSTANT C ON THE CHOICE OF ZEROS OF THE POLYNOMIAL

The above considerations show that it is not always advantageous to choose the optimal difference equation for the solution of (17) or, more generally, (1). It is evident that, given k and the value Ah , we shall aim at such a choice of the zeros of the polynomial $\varrho(\xi)$ that the relations (26) hold and the absolute value of the quantity C be as small as possible. Therefore we are interested in the dependence of C on the zeros of the polynomial $\varrho(\xi)$. Before stating a general theorem, let us introduce two examples.

Example 5. $k = 2$. Example 1 implies that C is a non positive strictly monotone function of λ . Hence it is obviously profitable to choose such $\lambda = \lambda_0$ that for given Ah satisfies the relation

$$(39) \quad \lambda_0 = \min_{\lambda \in (-1, 1)} (|\tilde{\xi}_1(\lambda)| \leq |\tilde{\xi}_2(\lambda)|)$$

It is seen from (32) that $\lambda_0 = -1$ for $Ah > 0$, $\lambda_0 = (2z + 3)/(2z - 3) + \varepsilon$ for $z = Ah < 0$ where $\varepsilon > 0$.

Example 6. $k = 4$. Assume that the zeros of the polynomial $\varrho(\xi)$ are of the form

$$(40) \quad 1, b, q_1 \pm iq_2$$

where $b \in \langle -1, 1 \rangle$, $q_1^2 + q_2^2 \leq 1$, $q_1 \neq 1$, $q_2 \geq 0$. Evaluating the dependence of the coefficients of the polynomial on the zeros (40) and substituting into (16) we obtain

$$(41) \quad C = -\frac{1}{1440} \frac{bK_1 + K_2}{K_3(1 - b)}$$

where

$$(42) \quad \begin{aligned} K_1 &= 11 + 22q_1 + 27(q_1^2 + q_2^2) \\ K_2 &= 27 + 22q_1 + 11(q_1^2 + q_2^2) \\ K_3 &= 1 - 2q_1 + q_1^2 + q_2^2 \end{aligned}$$

It can be shown easily that $K_1, K_2, K_3 > 0$ for all admissible values q_1, q_2 . This implies that the constant C is non positive. (41) yields that $C = 0$ if and only if $b = -1$, $q_1^2 + q_2^2 = 1$, $q_1 \neq 1$. Further it follows from (41) that the quantity C is a strictly monotone function of b for arbitrary fixed q_1, q_2 (satisfying (40)). Choosing arbitrarily fixed values of b, q_2 from the definition domain of (40), we can show that the function $C = C(q_1)$ is strictly monotone as well. The proof of this assertion is essentially very simple, although immensely tedious and cumbersome. Actually it consists in proving that the function

$$\phi(q_1) = \frac{dC(q_1)}{dq_1}$$

has no zero in the interval $(-\sqrt{1 - q_2^2}, \sqrt{1 - q_2^2})$.

Let us now proceed to the general statement. Assume that the polynomial $\varrho(\xi)$ has r real zero points b_i ($i = 1, 2, \dots, r$) different from and s pairs of complex conjugate zero points $p_i \pm iq_i$ ($i = 1, 2, \dots, s$). Hence $k = 1 + r + 2s$. Assume that all zeros of the polynomial $\varrho(\xi)$ satisfy the stability condition. The set of these zeros denote by M .

Theorem. *The function $C = C(b_1, b_2, \dots, b_r, p_1, p_2, \dots, p_s, q_1, q_2, \dots, q_s)$ has no strict extreme on M .*

Proof. A) Assume $r \neq 0$. Then the polynomial $\varrho(\xi)$ assumes the form

$$(44) \quad \varrho(\xi) = (\xi - 1)(\xi - b_1)F$$

where F is a function independent of b_1 . In the relation (6) for C let us first consider the denominator. It is

$$(45) \quad \beta_0 + \beta_1 + \dots + \beta_k = \alpha_1 + 2\alpha_2 + \dots + k\alpha_k = \varrho'(1) = (1 - b_1) F$$

The system (8) implies

$$(46) \quad \beta_i = \sum_{j=1}^k c_{ij} \alpha_j \quad i = 0, 1, \dots, k$$

where C_{ij} are constants. On the other hand, (44) implies

$$(47) \quad \alpha_j = d_j + h_j b_1$$

Where d_j, h_j are constants independent of b_1 . Putting together (46), (47) we obtain

$$(48) \quad \beta_i = \sum_{j=1}^k C_{ij} (d_j + h_j b_1)$$

(6), (4), (5), (8) yield

$$(49) \quad C_s = \sum_{i=0}^k v_{i,s} \alpha_i + \sum_{i=0}^k w_{i,s} \beta_i$$

Where $v_{i,s}, w_{i,s}$ are constants independent of b_1 . Substituting (45), (47), (48), (49) into (6) we obtain

$$(50) \quad C = \frac{E + b_1 D}{(1 - b_1) F}$$

where the quantities E, D, F are functions independent of b_1 . This completes the proof in this case.

B) Suppose $r = 0$. Then $s \neq 0$ and the polynomial assumes the form

$$(51) \quad \varrho(\xi) = (\xi - 1) (\xi - p_1 - iq_1) (\xi - p_1 + iq_1) F$$

where F is a function independent of the values, p_1, q_1 . It holds

$$(52) \quad \begin{aligned} \varrho(\xi) &= (\xi - 1) (\xi^2 - 2p_1 \xi + p_1^2 + q_1^2) F \\ \varrho'(1) &= ((1 - p_1)^2 + q_1^2) F \end{aligned}$$

The value q_1 is from the interval $\langle 0, 1 \rangle$. Let us use the substitution

$$(53) \quad \bar{q}_1 = q_1^2$$

Substituting (52), (53) into (6) we find

$$(54) \quad C = \frac{E_1 + F_1 \bar{q}_1}{E_2 + F_2 \bar{q}_1}$$

where E_1, E_2, F_1, F_2 are functions independent of q_1 . This completes the proof.

References

- [1] *Henrici P.*: Discrete variable methods in differential equations. Wiley, New York—London 1962.
- [2] *Voňka, P.*: Stability of difference methods Thesis, Charles Univ., Fac. of Math. Phys., Dept. of Numerical Math. 1968. (Czech, unpublished.)

STABILITA A LOKÁLNÍ CHYBA DIFERENČNÍCH METOD PRO ŘEŠENÍ OBYČEJNÉ DIFERENCIÁLNÍ ROVNICE 1. ŘÁDU

PETR VOŇKA

Tato práce se zabývá problémem konstrukce diferenčních formulí pro řešení obyčejné diferenciální rovnice 1. řádu ze zadaného charakteristického polynomu (2. kapitola), vlivem volby kořenů charakteristického polynomu na lokální chybu diferenčních metod u rovnice $y' = Ay$, $y(0) = 1$, kde A je konstanta (3. kapitola) a závislostí konstanty chyby na volbě kořenů charakteristického polynomu v obecném případě (4. kapitola).

Author's address: Dr. *Petr Voňka*, Vysoká škola chemicko-technologická ČVUT, Technická 1905, Praha 6.