

Accuracy Bounds for Semidiscretizations of Hyperbolic Problems

By Rolf Jeltsch and Klaus-Günther Strack

Abstract. Bounds are given for the error constant of stable finite-difference methods for first-order hyperbolic equations in one space dimension, which use r downwind and s upwind points in the discretization of the space derivatives, and which are of optimal order $p = \min(r + s, 2r + 2, 2s)$. It is known that this order can be obtained by interpolatory methods. Examples show, however, that their error constants can be improved.

1. Introduction. We consider the linear, one-dimensional test problem

$$(1.1) \quad u_t = u_x, \quad -\infty < x < \infty, t \geq 0, \quad u(x, 0) \text{ given.}$$

We shall analyze semidiscretizations of (1.1) of the form

$$(1.2) \quad \frac{\partial u_k}{\partial t}(t) = \frac{1}{\Delta x} \sum_{j=-r}^s \alpha_j u_{k+j}(t), \quad t > 0,$$

$u_k(0)$ given for $k \in \mathbf{Z}$,

where $u_k(t)$ is an approximation of $u(k \Delta x, t)$, $k \in \mathbf{Z}$, and Δx is the steplength in space direction. The differential-difference equation (1.2) is said to be of the class $\{r, s\}$ [4]. Let

$$\|u(t)\|^2 := \Delta x \sum_{k=-\infty}^{\infty} |u_k(t)|^2.$$

The system of differential equations (1.2) is said to be *stable* if there exists an estimate

$$(1.3) \quad \|u(t)\|^2 \leq C(t) \|u(0)\|^2,$$

where $C(t)$ is a function which is bounded independently of Δx and u ; see, e.g., [1]. It is well-known that, using the Fourier transform, the stability of (1.2) is equivalent to

$$(1.4) \quad \operatorname{Re} \rho(z) \leq 0 \quad \text{for } |z| = 1,$$

where

$$(1.5) \quad \rho(z) := \sum_{j=-r}^s \alpha_j z^j,$$

Received October 18, 1983; revised July 26, 1984.

1980 *Mathematics Subject Classification.* Primary 65M20, 65M05, 65M10.

©1985 American Mathematical Society
 0025-5718/85 \$1.00 + \$.25 per page

the characteristic function of the semidiscretization used in (1.2) [1], [4], [10]. The approximation (1.2) of (1.1) has order p and error constant c_{p+1} if, for all sufficiently smooth functions $y(x)$, one has

$$(1.6) \quad \frac{1}{\Delta x} \sum_{j=-r}^s \alpha_j y(x + j\Delta x) - y'(x) \\ = c_{p+1} (\Delta x)^p y^{(p+1)}(x) + O(|\Delta x|^{p+1}), \quad \Delta x \rightarrow 0.$$

In [4] it has been shown that, for stable systems (1.2) of class $\{r, s\}$, the order cannot exceed $\min\{2r + 2, 2s, r + s\}$. For $s = 1$, this result has been given in [11], while [1] has treated the case $r = 0$. In [4] stable systems (1.2) have been given which attain the highest possible order. These systems have been called interpolatory methods, and the coefficients α_j have been given explicitly.

From (1.6) it is clear that these interpolatory methods correspond exactly to the linear multistep formulas based on differentiation as presented in [3, pp. 206–209]. We shall prove the following:

THEOREM 1. *Let the differential system (1.2) be stable and of optimal order $p = \min\{2r + 2, 2s, r + s\}$. Then the error constant c_{p+1} satisfies*

$$(1.7) \quad (-1)^{r+1} c_{p+1} > \frac{r!(r+2)!}{(2r+4)!} \quad \text{for } s \geq r + 2;$$

$$(1.8) \quad c_{p+1} = (-1)^{s-1} \frac{r!s!}{(r+s+1)!} \quad \text{if } r \leq s \leq r + 2;$$

$$(1.9) \quad (-1)^{s-1} c_{p+1} \geq \frac{(s!)^2}{(2s+1)!} \quad \text{for } r \geq s. \quad \square$$

(1.8) and (1.9) have already been shown in [4]. (1.8) and equality in (1.9) are obtained by interpolatory methods. (1.8) and (1.9) say that for a fixed s and $r \geq s$ the error constant is minimal when $r = s$. However, if one fixes r and has $s \geq r + 2$, the situation is different. The bound (1.7) is $1/(2r + 4)$ times smaller than the error constant c_{p+1} for $s = r + 2$. This suggests that for a fixed r one can possibly decrease the error constant by increasing s . In Section 3, we shall give examples of such improvements on the error constant.

To prove the main result, we introduce a comparison technique. First, we adapt the theory of order stars in Section 2. Then, in Section 3, we prove properties of the “optimal semidiscretization” and compare any stable discretization to this one.

In practice, the most convenient approximations to (1.1) are explicit, with a fixed ratio $\mu = \Delta t/\Delta x$:

$$(1.10) \quad u_{k,n+1} = \sum_{j=-r}^s a_j u_{k+j,n},$$

where $u_{k,n}$ is an approximation to $u(k\Delta x, n\Delta t)$. The order p and the error constant C_{p+1} are defined by

$$(1.11) \quad u(0, \Delta t) - \sum_{j=-r}^s a_j u(j\Delta x, 0) \\ = C_{p+1} \frac{\partial^{p+1}}{\partial x^{p+1}} u(0, 0) (\Delta x)^{p+1} + O(|\Delta x|^{p+2}),$$

where we have used $u_t = u_x$ and $\Delta t = \mu \Delta x$. Let c_{p+1} be the error constant of the semidiscrete scheme which is the “derivative” of the above fully discrete scheme; see [6, p. 783]. Then one has

$$(1.12) \quad C_{p+1} = c_{p+1}\mu + O(\mu^2), \quad \mu \rightarrow 0^+.$$

Since Theorem 1 gives bounds for c_{p+1} , one obtains bounds for C_{p+1} of the fully discrete scheme, at least asymptotically, for $\mu \rightarrow 0^+$.

2. Order Stars. Order stars have been introduced by Wanner, Hairer and Nørsett [12] to prove stability results on methods for solving ordinary differential equations. In several papers this technique and related ideas have been used to investigate stability of numerical methods for finite-dimensional systems of ordinary differential equations which originate from the semidiscretization of $u_t = u_x$, $x \in [a, b]$, $t \geq 0$, or the wave equation; see, e.g., [2], [7], [8], [9], [13]. In [4] and [6] variations of the order star technique have been developed to treat stability of semi- and full discretizations of (1.1). For a bibliography, see [5]. In this section, we shall modify the order star technique of [4] so that it can be used to compare stable semidiscretizations to the “optimal” one in Section 3.

Let

$$(2.1) \quad \rho(z) := \sum_{j=-r}^s \alpha_j z^j$$

be the characteristic function of a semidiscretization. Clearly, one has order p and error constant c_{p+1} if and only if

$$(2.2) \quad \rho(z) = \log z + c_{p+1}(z - 1)^{p+1} + O(|z - 1|^{p+2}) \quad \text{as } z \rightarrow 1$$

(see, e.g., [3, p. 227], [4, p. 56]). Since we can express the stability of a scheme, its order, and its error constant in terms of $\rho(z)$, we shall talk henceforth of rational functions only. We shall say that $\rho(z)$ is *stable* if

$$(2.3) \quad \operatorname{Re} \rho(z) \leq 0 \quad \text{for } |z| = 1,$$

and $\rho(z)$ has *order* p and *error constant* c_{p+1} if (2.2) holds. We shall also consider functions of the form

$$(2.4) \quad \varphi(z) = \frac{\tilde{\rho}(z)}{\frac{1}{2}(z + 1)} = \frac{2 \sum_{j=-r}^s \tilde{\alpha}_j z^j}{z + 1},$$

where $\tilde{\alpha}_j$ are real coefficients. The notions of order p and the error constant c_{p+1} for φ are now introduced in exactly the same way as for ρ ; one just replaces ρ in (2.2) by φ . The stability condition (2.3) is replaced by

$$(2.5) \quad \operatorname{Re} \varphi(z) = \operatorname{Re} \frac{2\tilde{\rho}(z)}{z + 1} \leq 0 \quad \text{for } |z| = 1, z \neq -1.$$

We shall prove the bound for the error constant by comparing the characteristic function ρ with a function $\varphi(z)$ of the form (2.4). This is done by considering the difference $\psi(z) := \rho(z) - \varphi(z)$. Observe that if ρ and φ are both of order at least p , then the difference $\psi(z)$ has a root of multiplicity at least $p + 1$ at $z = 1$; i.e.,

$$(2.6) \quad \psi(z) = \rho(z) - \varphi(z) = c(z - 1)^{p+1} + O(|z - 1|^{p+2}).$$

We introduce

$$(2.7) \quad S(z) := e^{\psi(z)}, \quad z \in \mathbb{C},$$

and

$$(2.8) \quad \Omega := \{z \in \mathbb{C} \mid |S(z)| > 1\}.$$

Ω is called the *order star*. Ω^c will denote the complement of Ω : i.e., $\Omega^c = \overline{\mathbb{C}} \setminus \Omega$. Since α_j and $\tilde{\alpha}_j$ are always assumed to be real, Ω and Ω^c are symmetric with respect to the real axis.

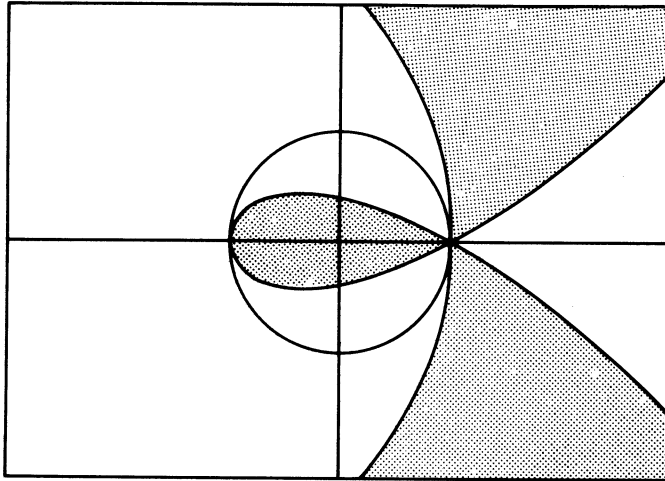


FIGURE 1(a)

Order star of $\psi_0 := (z - 1) - \frac{1}{2}(z - 1)^2 - 2(z - 1)/(z + 1)$ for $-3 \leq \operatorname{Re} z \leq 3$ and $-2 \leq \operatorname{Im} z \leq 2$

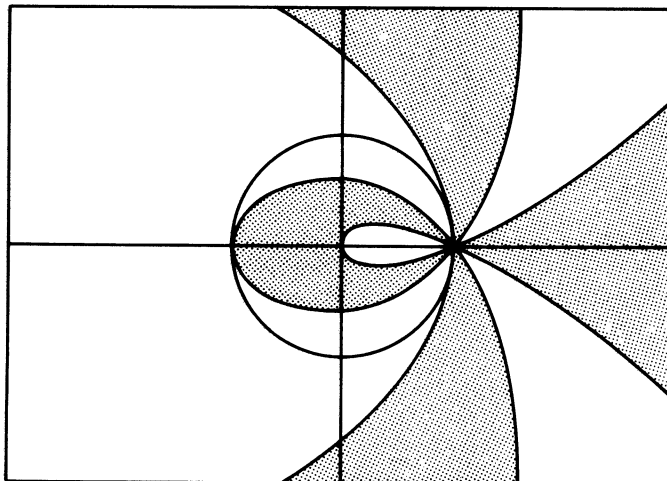


FIGURE 1(b)

Order star of $\psi_1 := -\frac{1}{4z} - \frac{5}{6} + \frac{3}{2}z - \frac{1}{2}z^2 + \frac{1}{12}z^3 - \frac{z^2 + 9z - 9 - z^{-1}}{6(z + 1)}$
for $-3 \leq \operatorname{Re} z \leq 3$ and $-2 \leq \operatorname{Im} z \leq 2$

As examples we draw the order star belonging to the interpolating function in class $\{0, 2\}$ compared with the function of type (2.4), with $r = 0, s = 2$, that has optimal error constant (see Figure 1(a)), and the order star of the interpolating function in class $\{1, 3\}$ compared with the function of type (2.4), with $r = 1, s = 3$, that has optimal error constant (see Fig. 1(b)).

Observe that in [4] the order star was formed by comparing ρ to $\log z$: i.e., by forming $\psi(z) = \rho(z) - \log z$ and thus

$$S(z) = e^{\rho(z)}/z.$$

However, we shall not need this type of order star, and thus we restrict ourselves to functions ψ of the form (2.6).

We shall need the following five lemmas, which are minor modifications of properties given in [4]. We shall therefore omit the proofs. Let D denote the unit disk $D = \{z \in \mathbb{C} \mid |z| \leq 1\}$.

LEMMA 2. *Let Ω be the order star of $\psi(z)$. Then*

$$(2.9) \quad \operatorname{Re} \psi(z) \leq 0 \quad \text{for } |z| = 1, z \neq -1$$

if and only if $\Omega \cap \partial D = \emptyset$.

LEMMA 3. *Let Ω be the order star of $\psi(z)$. Let $\psi(z)$ have a pole at z_0 of order r . Then, as z tends to z_0 , Ω consists of r sectors of angles π/r adjoining z_0 , separated by r sectors of Ω^c adjoining z_0 of the same angles π/r .*

We shall need this lemma for $z_0 = -1$, where ψ , given by (2.6), usually has a simple pole, and, for $z_0 = 0$, where ψ has a pole of order r if $\alpha_{-r} - 2\tilde{\alpha}_{-r} \neq 0$.

LEMMA 4 [4]. *Let Ω be the order star of $\psi(z)$. Then $\psi(z)$ has a root of multiplicity $p + 1$ at $z = 1$ if and only if, as z tends to one, Ω consists of $p + 1$ sectors, each of angle $\pi/(p + 1)$, separated by $p + 1$ sectors of Ω^c , each of the same angle.*

For $\psi(z)$ with (2.9) one can, in view of Lemma 2, introduce the

Definition. μ_i and μ_o denote the number of sectors of Ω , inside and outside D , respectively, approaching $z = 1$.

LEMMA 5. *For every $\psi(z)$ with (2.9) and a root of multiplicity $p + 1$ at $z = 1$, one has $\mu_i + \mu_o = p + 1$, $|\mu_i - \mu_o| \leq 1$, and $p \leq 2\mu_i$.*

Finally, we shall need the following lemma, which is part of Proposition 8 in [5].

LEMMA 6. *Let Ω be the order star of $\psi(z)$. Between any two points with $\psi(z) = 0$ that are connected by an arc of $\partial\Omega$ there is an essential singularity on this arc.*

3. Accuracy Bounds. In this section we use the tools developed in Section 2 to prove the lower bound (1.7) for the error constant given in Theorem 1. To do this let us first derive a special function $\varphi_r(z)$ of the form (2.4). Observe that by

$$(3.1) \quad \varphi(z) = \frac{\tilde{\rho}(z)}{\frac{1}{2}(z + 1)} = \log z + c_{p+1}(z - 1)^{p+1} + O(|z - 1|^{p+2}),$$

one has

$$(3.2) \quad \sum_{j=-r}^s \tilde{\alpha}_j z^{j+r} - \frac{1}{2}(z + 1)z^r \log z = c_{p+1}(z - 1)^{p+1} + O(|z - 1|^{p+2})$$

for $z \rightarrow 1$. Hence, if $p \geq r + s$, then

$$(3.3) \quad z^r \tilde{\rho}(z) = \sum_{j=1}^{r+s} b_j^r (z-1)^j,$$

where

$$(3.4) \quad \frac{1}{2}(z+1)z^r \log z = \sum_{j=1}^{\infty} b_j^r (z-1)^j.$$

In particular, for $s = r + 1$ there is a unique

$$(3.5) \quad z^r \rho_r(z) := \sum_{j=-r}^{r+1} \beta_j^r z^{j+r} := \sum_{j=1}^{2r+1} b_j^r (z-1)^j$$

and

$$(3.6) \quad \varphi_r(z) = 2\rho_r(z)/(z+1),$$

such that $p \geq 2r + 1$. Substitution of $1/z$ in (3.1) gives

$$(3.7) \quad -\varphi_r\left(\frac{1}{z}\right) = \frac{-2\sum_{j=-r}^{r+1} \beta_{1-j}^r z^j}{z+1} \\ = \log z - c_{p+1}(-1)^{p+1}(z-1)^{p+1} + O(|z-1|^{p+2}).$$

Thus $-\varphi_r(1/z)$ has order at least $2r + 1$. By the uniqueness of $\varphi_r(z)$ we have

$$(3.8) \quad \varphi_r(z) = -\varphi_r(1/z),$$

and thus

$$(3.9) \quad \beta_j^r = -\beta_{1-j}^r \quad \text{for } j = -r, -r+1, \dots, r+1.$$

We show now that $\varphi_r(z)$ has order $p = 2r + 2$ and error constant

$$(3.10) \quad \chi_r := (-1)^{r+1} \frac{r!(r+2)!}{(2r+4)!}.$$

At this point it is helpful to observe that, by (3.2), $\varphi_r(z)$ corresponds to the difference operator

$$(3.11) \quad (L_{\Delta x} y)(x) = \sum_{j=-r}^{r+1} \beta_j^r y(x+j\Delta x) - \frac{\Delta x}{2} (y'(x+\Delta x) + y'(x)) \\ = \chi_r(\Delta x)^{p+1} y^{(p+1)}(x) + O(|\Delta x|^{p+2});$$

see, for example, [3, p. 227]. Clearly, by the symmetry (3.9), the order p has to be even. Thus $p \geq 2r + 2$. To compute the error constant χ_r let $\Delta x = 1$ and

$$\hat{y}(x) = x \prod_{j=-r}^{r+1} (x-j).$$

Hence, substitution in (3.11) gives

$$(L_1 \hat{y})(0) = -\frac{1}{2}(\hat{y}'(1) + \hat{y}'(0)) = -\frac{1}{2}\hat{y}'(1) \\ = -\frac{1}{2}r!(r+1)!(-1)^r = \chi_r(2r+3)!.$$

This establishes (3.10).

Finally, we establish that φ_r satisfies (2.5). Let $|z| = 1, z \neq -1$. Hence, by (3.8),

$$\overline{\varphi_r(z)} = \varphi_r(\bar{z}) = \varphi_r(1/z) = -\varphi_r(z),$$

and therefore $\text{Re } \varphi_r(z) = 0$ for $|z| = 1, z \neq -1$. We have, therefore, established the following

PROPOSITION 7. *The function $\varphi_r(z) = 2\rho_r(z)/(z + 1)$ given by (3.5) has order $p = 2r + 2$ and error constant $\chi_r = (-1)^{r+1}r!(r + 2)!/(2r + 4)!$. The function $\varphi_r(z)$ satisfies (2.5), in fact, $\text{Re } \varphi_r(z) = 0$ for $|z| = 1, z \neq -1$.*

We give the coefficients β_j^r of $\varphi_r(z), r \geq 1$, without proof:

$$(3.12) \quad \begin{cases} \beta_1^r = \frac{r + 2}{2(r + 1)}, \\ \beta_j^r = \frac{1}{2} \frac{(-1)^j r!(r + 1)!}{(r + j)!(r + 1 - j)!j \cdot (j - 1)}, & j = 2, 3, \dots, r + 1, \\ \beta_j^r = -\beta_{1-j}^r, & j = -r, -r + 1, \dots, -1, 0. \end{cases}$$

For $r = 0$ we have the trapezoidal rule

$$\varphi_0 = 2(z - 1)/(z + 1).$$

A subset Δ_1 of $\text{Int}(\Omega^c)$ is said to be an Ω^c -component if $\partial\Delta_1 \subset \partial\Omega^c$ and Δ_1 is connected. Ω -components are defined similarly.

LEMMA 8. *Let $\psi(z)$ be given by*

$$(3.13) \quad \psi(z) = 2 \frac{\sum_{j=-r}^s \alpha_j z^j}{z + 1}.$$

For any bounded Ω -component Ω_1 one has

- (i) $\{-1\} \cap \bar{\Omega}_1 \neq \emptyset$ if $r = 0$,
- (ii) $\{0, -1\} \cap \bar{\Omega}_1 \neq \emptyset$ if $r > 0$.

The same is true for Ω^c -components.

Proof. Let $r > 0$ and let Ω_1 be a bounded Ω -component with $\{0, -1\} \cap \bar{\Omega}_1 = \emptyset$. Hence, $S(z)$ is analytic in an open set $W \supset \bar{\Omega}_1$. Hence, $|S(z)| = 1$ for $z \in \partial\Omega_1$, and $|S(z)| > 1$ for $z \in \Omega_1$. Since $S(z)$ is not constant, we have a contradiction to the maximum modulus principle. Hence, $\{0, -1\} \cap \bar{\Omega}_1 \neq \emptyset$.

For $r = 0$ observe that $z = 0$ is also a regular point of $S(z)$.

For Ω^c -components one proves the lemma in a similar fashion by considering $1/S(z)$ instead of $S(z)$. \square

For brevity we shall henceforth abbreviate “sector of Ω at $z = 1$ ”, as used in Lemma 4, by “finger”, and “sector of Ω^c at $z = 1$ ”, by “dual finger”. In the following we shall bound μ_i , i.e., the numbers of fingers in D . As a first step we have

LEMMA 9. *Let $\psi(z)$ in (3.13) satisfy (2.9).*

- (a) *An Ω -component in D has at most two fingers.*
- (b) *At most one Ω -component in D has two fingers.*

Proof. Since we work in D only, we shall not always restate this.

(a) Assume Ω_1 is an Ω -component with more than two fingers. Hence, there are at least two dual fingers in D which belong to two disjoint bounded Ω^c -components $\Delta_1,$

Δ_2 which cannot have -1 on their closures. Since Δ_1 and Δ_2 are separated by Ω_1 , one of them cannot have the origin in its closure. This is a contradiction to Lemma 8.

(b) Assume there are at least Ω -components Ω_1 and Ω_2 with two fingers each. Since Ω_1 and Ω_2 are connected and $\Omega_1 \cap \Omega_2 = \emptyset$, the fingers of Ω_1, Ω_2 cannot interlace each other when moving on a small circle around $z = 1$. Hence, between these four fingers there are at least two dual fingers in D which belong to two disjoint bounded Ω^c -components which cannot have -1 in their closures. As in the proof of part (a) this leads to a contradiction. \square

PROPOSITION 10. *Let $\psi(z) = 2 \sum_{j=-r}^s \alpha_j z^j / (z + 1)$ satisfy (2.9). Then for the number of fingers inside D one has*

$$(3.14) \quad \mu_i \leq r + 1.$$

Proof. Let $r = 0$. Hence, $z = 0$ is a regular point. Inside D we have at most one Ω -component, Ω_1 say, with $z = -1$ on its boundary. Assume now that this component has at least two fingers. Then these two fingers enclose a bounded Ω^c -component Δ_1 . Since $\bar{\Delta}_1 \cap \{-1\} = \emptyset$, this contradicts Lemma 8. Hence, $\mu_i \leq 1$. Next, we consider the case $r > 0$. Here we distinguish the following two cases:

(i) No finger in D belongs to an Ω -component which has -1 on its boundary. Then, by Lemma 3, there are at most r Ω -components in D which do not have -1 on their boundaries. Thus, by Lemma 9, one has $\mu_i \leq r + 1$.

(ii) The component Ω_1 which has -1 on its boundary has at least one finger. Since $0 \notin \text{Int } \Omega_1$, and because of the symmetry of Ω with respect to the real axis, Ω_1 has exactly two fingers. Now Ω_1 is either connected with one of the Ω -sectors at $z = 0$ and, by Lemma 9, $\mu_i \leq 2 + (r - 1) = r + 1$, or the inner boundary of Ω_1 is singularity free, which contradicts Lemma 6. \square

In the proof of the main theorem, φ_r plays an important role. It is, as we shall see, the most accurate function, although not stable itself since it is not defined at $z = -1$. We now have the tools to prove the part of Theorem 1 not covered in [4]. For convenience we restate it here as

THEOREM 11. *Let $s \geq r + 2 \geq 2$. Let the differential system (1.2) of class $\{r, s\}$ be stable and of optimal order $p = 2r + 2$. Then one has for the error constant c_{2r+3} the lower bound*

$$(3.15) \quad (-1)^{r+1} c_{2r+3} > \frac{r!(r+2)!}{(2r+4)!} = |\chi_r|.$$

Proof. By ρ we denote the characteristic polynomial of the method in consideration, having order $2r + 2$ and error constant c_{2r+3} . We now define

$$(3.16) \quad \begin{aligned} \psi(z) &:= \rho(z) - \varphi_r(z) \\ &= (c_{2r+3} - \chi_r)(z - 1)^{2r+3} + O(|z - 1|^{2r+4}). \end{aligned}$$

$\text{Re } \psi(z) = \text{Re}(\rho(z) - \varphi_r(z)) \leq 0$ for $|z| = 1, z \neq -1$, since ρ is stable and because of Proposition 7. Thus φ satisfies (2.9). We distinguish two cases.

(i) r even. Assume (3.15) is wrong. Hence, $c_{2r+3} - \chi_r \geq 0$. Proposition 10 implies that $\mu_i \leq r + 1$. Hence, by Lemma 5, $c_{2r+3} - \chi_r = 0$ is impossible. If $c_{2r+3} - \chi_r > 0$, we have that, for $\varepsilon > 0$ small enough, $1 + \varepsilon \in \Omega$ and $1 - \varepsilon \in \Omega^c$; hence, by Lemma 5 and (2.9), one has $\mu_i = r + 2$. This is a contradiction.

(ii) r odd. The assumption $c_{2r+3} = \chi_r$ yields a contradiction as before. If $\chi_r > c_{2r+3}$, we have that, for $\varepsilon > 0$ small enough, $1 + \varepsilon \in \Omega^c$ and $1 - \varepsilon \in \Omega$; hence, by Lemma 5 and (2.9), one has $\mu_i = r + 2$. This is a contradiction to Proposition 10. \square

We remark the following:

1. Since the error c_{2r+3} is bounded away from zero, we have at the same time found another way to prove the maximal order $p \leq 2r + 2$ for functions ρ in the class $\{r, s\}$. Clearly, one can give a similar proof for $p \leq 2s$.

2. Contrary to the case when s is fixed and an arbitrary number of points to the left is taken (cf. [4, Theorem 6]), we here get an improvement in the error constant of a stable method when taking more points to the right and determining coefficients α_j in a suitable way. This also reflects once more the asymmetric behavior of the advection equation.

We gain most of the possible improvement of the error constant when taking only one or two points more than in the interpolatory method with $s = r + 2$.

Next we give some examples for the improvement of the error constant. By $c(r, s)$ we denote the (absolutely) smallest error of a stable function of class $\{r, s\}$ with maximal order $p = \min\{r + s, 2(r + 1), 2s\}$.

Example 1. For order $p = 2r + 2$, one must use at least r downwind and $r + 2$ upwind points; consider, e.g., the interpolatory methods. Here we show that already adding one point more on the upwind side results in an essential improvement of the error constant. When approximating $\log z$ by stable functions of the type

$$\begin{aligned}
 (3.17) \quad \rho(z) &= \sum_{j=-r}^{r+3} \alpha_j z^j \\
 &= \sum_{j=1}^{2r+2} a_j \frac{(z-1)^j}{z^r} + a \frac{(z-1)^{2r+3}}{z^r} \quad \text{with order } p = 2r + 2,
 \end{aligned}$$

we first consider the Taylor series of $z^r \log z$ at $z = 1$, which is given by

$$z^r \log z = \sum_{j=1}^{\infty} a_j^r (z-1)^j,$$

where the a_j^r are recursively given by

$$(3.18) \quad \begin{cases} a_0^r := 0, \\ a_{j+1}^r := \frac{1}{j+1} \left(\binom{r}{j} + (r-j)a_j^r \right), & j = 0, 1, \dots, r, \\ a_{j+1}^r := \frac{r-j}{j+1} a_j^r, & j > r. \end{cases}$$

To get order $p = 2r + 2$, we must choose, analogously to (3.5), $a_j = a_j^r$ for $1 \leq j \leq 2r + 2$. Hence, the characteristic function of semidiscretizations (1.2) of class $\{r, r + 3\}$ of order $p = 2r + 2$ have the form

$$(3.19) \quad \rho(z, a) := \sum_{j=1}^{2r+2} a_j^r \frac{(z-1)^j}{z^r} + a \frac{(z-1)^{2r+3}}{z^r}.$$

$\rho(z, 0)$ belongs to the interpolatory method. It is stable and has error constant

$$(3.20) \quad c_{2r+3} = (-1)^{r+1} \frac{r!(r+2)!}{(2r+3)!} = c(r, r+2).$$

Clearly the error constant of $\rho(z, a)$ is

$$(3.21) \quad c(a) = c_{2r+3} + a.$$

In the following, we show that one has to choose

$$(3.22) \quad a = a^* := (-1)^r \frac{r!(r+2)!}{(2r+3)!} \cdot \frac{2r+3}{3r+6}$$

in order that $\rho(z, a^*)$ is stable and

$$(3.23) \quad c(a^*) = c(r, r+3),$$

i.e., $\rho(z, a)$ is unstable for all a with $|c(a)| < |c(a^*)|$. Since $\rho(z, 0)$ is stable, it is clear from (3.20) and (3.21) that one gets an improvement of the error constant compared to c_{2r+3} if

$$(3.24) \quad (-1)^r a > 0.$$

Clearly,

$$(3.25) \quad \begin{aligned} \rho(z, a) - \log z &= \rho(z, 0) - \log z + a \frac{(z-1)^{2r+3}}{z^r} \\ &= -\frac{1}{z^r} \sum_{j=2r+3}^{\infty} a_j^r (z-1)^j + a \frac{(z-1)^{2r+3}}{z^r}. \end{aligned}$$

By a simple but tedious calculation we obtain, using (3.25), (3.18), and the fact $a_{2r+3}^r = -c_{2r+3}$, that

$$(3.26) \quad \operatorname{Re} \rho(e^{i\theta}, a) = \operatorname{Re}(\rho(e^{i\theta}, a) - \log e^{i\theta}) = \Delta \theta^{2r+4} + O(\theta^{2r+6}),$$

where

$$(3.27) \quad \Delta = \left((-1)^r a - \frac{r!(r+2)!}{(2r+3)!} \cdot \frac{2r+3}{3r+6} \right) \cdot \frac{3}{2}.$$

Hence, from stability, we have the necessary condition $\Delta \leq 0$, which is equivalent to

$$(3.28) \quad (-1)^r a \leq (-1)^r a^*.$$

Since $|a^*| < |c_{2r+3}|$, we see from (3.21) that the (absolutely) smallest error constants of all $\rho(z, a)$, with (3.28), is obtained for $a = a^*$. It remains to show that $\rho(z, a^*)$ is stable. To do this, observe that

$$\begin{aligned} \operatorname{Re} \rho(e^{i\theta}, a^*) &= \operatorname{Re} \left(\sum_{j=1}^{2r+2} a_j^r \frac{(z-1)^j}{z^r} + a^* \frac{(z-1)^{2r+3}}{z^r} \right) \Bigg|_{z=e^{i\theta}} \\ &= \sum_{j=0}^{r+3} d_j \cos j\theta \end{aligned}$$

for some real coefficients d_j . With the transformation $x = \cos \theta$, $\theta \in [0, \pi]$, we obtain

$$(3.29) \quad \operatorname{Re} \rho(e^{i\theta}, a^*) = \sum_{j=0}^{r+3} d_j T_j(x) = \sum_{j=0}^{r+3} \delta_j (x-1)^j =: T(x),$$

where $T_j(x)$ are the Chebyshev polynomials. We substitute $x = \cos \theta$ in (3.29) and expand with respect to θ at $\theta = 0$. Comparing this expansion with (3.26) for $a = a^*$

yields $\delta_0 = \delta_1 = \delta_2 = \dots = \delta_{r+2} = 0$. Hence,

$$T(x) = (x - 1)^{r+3} \delta_{r+3},$$

and thus $T(x)$ does not change sign in $(-1, 1)$. Therefore, $\text{Re } \rho(e^{i\theta}, a^*)$ has constant sign for all $\theta \in [0, 2\pi]$. However,

$$\text{Re } \rho(-1, a^*) = \text{Re } \rho(-1, 0) + a^* \cdot 2^{2r+3} (-1)^{r+1} < 0,$$

since $\text{Re } \rho(-1, 0) \leq 0$ by stability of $\rho(z, 0)$. Hence, the method $\rho(z, a^*)$ is stable. We have thus proved the following

THEOREM 12. *Among all stable differential-difference equations (1.2) with r downwind and at most $r + 3$ upwind points and of optimal order $p = 2r + 2$ the one with the characteristic function $\rho(z, a^*)$ given by (3.19) and (3.22) gives the smallest error constant in absolute value. The value of the error constant is*

$$(3.30) \quad c(a^*) = (-1)^{r+1} \frac{r!(r+2)!}{(2r+3)!} \cdot \frac{r+3}{3r+6}.$$

The value of $c(a^*)$ is obtained by substitution of a^* in (3.21). We observe that

$$c(r, r+3) = \frac{r+3}{3r+6} \cdot c(r, r+2).$$

Hence, adding one upwind point results in a decrease of the error constant by at least a factor 2, and at most a factor 3. We list for $r = 0, 1, 2, 3$ the values $c(r, r+2)$, $c(r, r+3)$ as well as the bound $c(r, \infty) := \chi_r$ of the error constant for formulas with arbitrary large s .

r	p	$c(r, r+2)$	$c(r, r+3)$	$c(r, \infty)$
0	2	$-\frac{1}{3}$	$-\frac{1}{6}$	$-\frac{1}{12}$
1	4	$\frac{1}{20}$	$\frac{1}{45}$	$\frac{1}{120}$
2	6	$-\frac{1}{105}$	$-\frac{1}{252}$	$-\frac{1}{840}$
3	8	$\frac{1}{504}$	$\frac{1}{1260}$	$\frac{1}{5040}$

For the resulting functions $\rho(z, a^*)$ we can also give the coefficients as follows:

$$(3.31) \quad \begin{cases} \alpha_j = (-1)^{j+1} \frac{r!(r+1)!}{(r+j)!(r+2-j)!} \left(\frac{r+2}{j} + \frac{2r+3}{3(r+3-j)} \right), \\ \qquad \qquad \qquad -r \leq j \leq r+2, \quad j \neq 0, \\ \alpha_0 = -\left(\frac{1}{r+1} + \frac{1}{r+2} + \frac{2r+3}{(r+3)(3r+6)} \right), \\ \alpha_{r+3} = (-1)^r \frac{r!(r+2)!}{(2r+2)!(3r+6)}. \end{cases}$$

Example 2. Error constants of stable methods in the purely upwind case $r = 0$ for different values of s :

s	2	3	4	$\rightarrow \infty$
c_3	-0.3333	-0.1667	-0.1262	-0.0833
a_1	1.000	1.000	1.000	1.000
a_2	-0.500	-0.500	-0.500	-0.500
a_3	-	0.1667	0.2071	0.250
a_4	-	-	-0.0732	-0.125

The coefficients result from the transformed nonlinear system of equations in (1.4) with the side condition to “minimize” c_3 .

Institut für Geometrie und Praktische Mathematik
Templergraben 55
D-5100 Aachen, West Germany

1. B. ENGQUIST & S. OSHER, “One-sided difference approximations for nonlinear conservation laws,” *Math. Comp.*, v. 36, 1981, pp. 321–351.
2. E. HAIRER, “Unconditionally stable methods for second order differential equations,” *Numer. Math.*, v. 32, 1979, pp. 373–379.
3. P. HENRICI, *Discrete Variable Methods in Ordinary Differential Equations*, Wiley, New York, 1961.
4. A. ISERLES, “Order stars and a saturation theorem for first order hyperbolics,” *IMA J. Numer. Anal.*, v. 2, 1982, pp. 49–61.
5. A. ISERLES, *Order Stars, Approximation and Finite Differences. I: The General Theory of Order Stars*, Report DAMTP NA3, Univ. of Cambridge, 1983.
6. A. ISERLES & G. STRANG, “The optimal accuracy of difference schemes,” *Trans. Amer. Math. Soc.*, v. 277, 1983, pp. 779–803.
7. R. JELTSCH & O. NEVANLINNA, “Stability of explicit time discretizations for solving initial value problems,” *Numer. Math.*, v. 37, 1981, pp. 61–91.
8. R. JELTSCH & O. NEVANLINNA, “Stability and accuracy of time discretizations for initial value problems,” *Numer. Math.*, v. 40, 1982, pp. 245–296.
9. R. JELTSCH & O. NEVANLINNA, “Stability of semidiscretizations of hyperbolic problems,” *SIAM J. Numer. Anal.*, v. 20, 1983, pp. 1210–1218.
10. R. RICHTMYER & K. MORTON, *Difference Methods for Initial Value Problems*, 2nd ed., Wiley, New York, 1967.
11. G. STRANG, “Accurate partial difference methods. II: Nonlinear problems,” *Numer. Math.*, v. 6, 1964, pp. 37–44.
12. G. WANNER, E. HAIRER & S. P. NØRSETT, “Order stars and stability theorems,” *BIT*, v. 18, 1978, pp. 475–489.
13. G. WANNER, E. HAIRER & S. P. NØRSETT, “When I -stability implies A -stability,” *BIT*, v. 18, 1978, p. 503.