

## CONTRACTIVITY-PRESERVING IMPLICIT LINEAR MULTISTEP METHODS

H. W. J. LENFERINK

**ABSTRACT.** We investigate contractivity properties of implicit linear multistep methods in the numerical solution of ordinary differential equations. The emphasis is on nonlinear and linear systems  $\frac{d}{dt}U(t) = f(t, U(t))$ , where  $f$  satisfies a so-called circle condition in an arbitrary norm. The results for the two types of systems turn out to be closely related. We construct optimal multistep methods of given order and stepnumber, which allow the use of a maximal stepsize.

### 1. INTRODUCTION

In this paper we will study the contractivity of implicit linear multistep methods in the numerical solution of ordinary differential equations. Consider the initial value problem

$$(1.1a) \quad \frac{d}{dt}U(t) = f(t, U(t)) \quad (t \geq 0),$$

$$(1.1b) \quad U(0) = u_0,$$

where  $u_0 \in \mathbb{K}^s$  and the function  $f$  with values in  $\mathbb{K}^s$  are given ( $\mathbb{K}$  stands consistently for  $\mathbb{R}$  or  $\mathbb{C}$  in this section) and  $s \geq 1$ . To approximate the solution to (1.1), we want to use the linear multistep method

$$(1.2) \quad \begin{aligned} &u_n - h\beta_k f(t_n, u_n) \\ &= \sum_{i=0}^{k-1} (-\alpha_i u_{n-k+i} + h\beta_i f(t_{n-k+i}, u_{n-k+i})) \quad (n \geq k). \end{aligned}$$

The vector  $u_n$  is an approximation to  $U(nh)$  ( $n = 0, 1, \dots$ ),  $h > 0$  is the *stepsize*,  $k$  is the *stepnumber*, and  $u_0, \dots, u_{k-1} \in \mathbb{K}^s$  are given initial approximations. Further,  $\alpha_i$  ( $0 \leq i \leq k-1$ ) and  $\beta_i$  ( $0 \leq i \leq k$ ) are coefficients in  $\mathbb{K}$  specifying the method. The method is *explicit* if  $\beta_k = 0$ , and

---

Received January 18, 1990.

1980 *Mathematics Subject Classification* (1985 Revision). Primary 65L05, 65L20; Secondary 65M10.

This research has been supported by the Netherlands Organization for Scientific Research (NWO).

implicit if  $\beta_k \neq 0$ . The order of the method is the largest integer  $p$  such that

$$(1.3) \quad (1 - h\beta_k) \exp(kh) = \sum_{i=0}^{k-1} (-\alpha_i + h\beta_i) \exp(ih) + \mathcal{O}(h^{p+1}) \quad (\text{for } h \rightarrow 0).$$

We assume (see, e.g., [20, p. 390] or [4, 8, 10, 11, 13, 14, 17, 18, 19]) that  $f$  satisfies, for some  $\rho > 0$ , the *circle condition*

$$(1.4) \quad \|f(t, x) - f(t, y) + \rho(x - y)\| \leq \rho \|x - y\| \quad (\text{for } x, y \in \mathbb{K}^s).$$

Here,  $\|\cdot\|$  stands for an arbitrary norm in  $\mathbb{K}^s$ . This circle condition implies (cf., e.g., [20])

$$(1.5) \quad \|V(t_2) - W(t_2)\| \leq \|V(t_1) - W(t_1)\| \quad (\text{for } t_2 \geq t_1 \geq 0),$$

for any two solutions  $V$  and  $W$  of (1.1a).

In view of (1.5) it is natural to ask for a multistep method (1.2) which guarantees a similar property for the numerical approximations  $u_n$ . Such a property is also favorable with respect to the propagation of errors. Therefore, we call the  $k$ -step method (1.2) *contractive* (cf. [4, 8, 11, 13, 17, 19]) for a function  $f$ , a stepsize  $h$ , and a norm  $\|\cdot\|$  if for any two sequences  $(v_n)_{n \geq 0}$ ,  $(w_n)_{n \geq 0}$  satisfying (1.2) we have

$$(1.6) \quad \|v_n - w_n\| \leq \max_{0 \leq i \leq k-1} \|v_{n-k+i} - w_{n-k+i}\| \quad (n = k, k+1, \dots).$$

Two test equations will be considered. One of them is obtained by choosing for the function  $f$

$$(1.7) \quad f(t, x) = Ax \quad (t \geq 0, x \in \mathbb{K}^s).$$

Here,  $A$  is an  $s \times s$  matrix with coefficients in  $\mathbb{K}$  such that (1.4) holds for some  $\rho > 0$ . The contractivity of (1.2) when applied to this test problem is related to the choice of the stepsize  $h$  by Theorem 3.3 in [19]. In fact, (1.2) is contractive for all  $s \geq 1$ , each matrix  $A$ , and each norm  $\|\cdot\|$  on  $\mathbb{K}^s$  such that (1.4) is satisfied, if and only if

$$(1.8) \quad h \leq R\rho^{-1}.$$

The factor  $R$  ( $0 \leq R \leq \infty$ ), which in [8, 11, 13, 19] is called the *threshold factor*, is given by

$$(1.9) \quad \begin{aligned} R &= \inf\{-\alpha_i \beta_i^{-1} \mid 0 \leq i \leq k-1 \text{ and } \beta_i > 0\} \\ &\quad \text{if } \beta_k \geq 0 \text{ and } \alpha_i \leq 0, \alpha_i \beta_k \leq \beta_i \quad (0 \leq i \leq k-1), \\ R &= 0 \quad (\text{otherwise}). \end{aligned}$$

We define, for  $k \geq 1$ ,  $p \geq 1$ , the *optimal threshold factor*  $R_{k,p}$  by

$$(1.10) \quad R_{k,p} = \sup\{R \mid R \text{ is the threshold factor (1.9) of a } k\text{-step method (1.2) of order } p\}.$$

In the definitions of the threshold factors, we put, as usual,  $\inf \emptyset = \infty$  and  $\sup \emptyset = -\infty$ . A  $k$ -step method (1.2) of order  $p$  will be called *optimal* for (1.1), (1.7) if its threshold factor  $R$  equals  $R_{k,p}$ .

One also arrives at a stepsize restriction which involves the threshold factor (1.9) if one aims at the conservation of positivity in the numerical solution of ordinary differential equations (see [1]). This factor (1.9) is necessarily finite when the order  $p$  is larger than one (cf. [1, 15, 19]).

A threshold factor with a similar property as  $R$  given by (1.9) can be defined for methods that do not belong to the class of methods (1.2) (e.g., for explicit or implicit Runge-Kutta methods, cf. [19, p. 283]). For rather general one-step methods, the size of the corresponding optimal threshold factor was studied in [8, 11]. In [13], this was done for explicit linear multistep methods. One of the purposes of the present paper is to study the size of optimal threshold factors  $R_{k,p}$  for linear multistep methods that are implicit.

We arrive at another natural test equation if we choose  $\mathbb{K} = \mathbb{C}$  and

$$(1.11) \quad f(t, x) = a(t)x \quad (t \geq 0, x \in \mathbb{C})$$

in (1.1a). Here,  $a(\cdot)$  is a function with values in  $\mathbb{C}$  such that (1.4) holds for some  $\rho > 0$ . This scalar test equation was also dealt with in, e.g., [2, 17, 18].

Also when (1.2) is applied to the second test equation, i.e., (1.1), (1.11), the contractivity is related to the choice of the stepsize  $h$ . It was proved in [17, 18] that (1.2) is contractive for all functions  $a(\cdot)$  satisfying (1.4) for some fixed  $\rho > 0$  if and only if

$$h \leq S\rho^{-1}.$$

The constant  $S$ , which we also call a *threshold factor*, is given by

$$(1.12) \quad \begin{aligned} S &= \inf\{-\alpha_i \beta_i^{-1} \mid 0 \leq i \leq k-1 \text{ and } \beta_i > 0\} \\ &\quad \text{if } \alpha_i \leq 0 \ (0 \leq i \leq k-1), \ \beta_i \geq 0 \ (0 \leq i \leq k), \\ S &= 0 \quad (\text{otherwise}). \end{aligned}$$

In fact (cf. [17, 18, 20]), method (1.2) is contractive for any norm  $\|\cdot\|$  and any function  $f$  satisfying (1.4) if and only if the stepsize  $h$  is chosen in conformity with the latter stepsize restriction.

Similarly to (1.10), we define for  $k \geq 1, p \geq 1$  the *optimal threshold factor*

$$(1.13) \quad \begin{aligned} S_{k,p} &= \sup\{S \mid S \text{ is the threshold factor (1.12)} \\ &\quad \text{of a } k\text{-step method (1.2) of order } p\}, \end{aligned}$$

and we say that (1.2) is an *optimal*  $k$ -step method of order  $p$  for (1.1), (1.11) if the threshold factor  $S$  of this method is equal to  $S_{k,p}$ .

In the numerical solution of (stiff) differential equations, stepsize restrictions that are imposed in view of contractivity or positivity requirements can be quite restrictive and embarrassing.

In view of the above considerations, it is natural to study the size of optimal threshold factors  $R_{k,p}$  and  $S_{k,p}$  and to construct optimal linear multistep

methods. Methods that are optimal with respect to  $S$  among a subset of  $k$ -step methods of order  $p$  were found in [17], but here we shall consider optimality with respect to all methods.

Optimal methods for (1.1), (1.11) will be constructed in §2. The size of  $S_{k,p}$ , for fixed  $p$ , when  $k$  tends to infinity, is considered in §3. We will prove, for  $p = 2, 4, 6$ , that  $S_{k,p}$  is maximal for  $k = p^2/4$ . On the other hand, for  $p = 3, 5, 7$ , it turns out that  $S_{k,p}$  is a strictly increasing function of  $k$  with  $\lim_{k \rightarrow \infty} S_{k,p} = S_{(p-1)^2/4, p-1}$ .

The implicit Euler method is optimal for both test equations when  $p = 1$ . In fact, for all  $k \geq 1$ , we have  $R_{k,1} = S_{k,1} = \infty$ . This was shown in [15, 17, 19]. We will show in §4 that the equality  $R_{k,p} = S_{k,p}$  also holds for  $p = 2, 4, 6$  and  $k \geq p^2/4$ .

For  $p = 2$  and any  $k \geq 1$  we will prove that the trapezoidal rule (i.e., (1.2) with  $\beta_k = \beta_{k-1} = 1/2$ ,  $\alpha_{k-1} = -1$ ,  $\beta_i = \alpha_i = 0$  ( $0 \leq i \leq k-2$ )) has optimal threshold factors  $R$  and  $S$ . This reveals another optimal property of this method, which is already known to have optimal properties with respect to  $A$ -stability (cf. [3]).

Further, for  $1 \leq p \leq 7$ ,  $R_{k,p}$  can be only slightly larger than  $S_{k,p}$  (cf. §4).

A numerical example is presented in §5, where one of our optimal methods is compared with a  $BDF$ -formula (cf. [12]). Finally, some technical lemmas are formulated in §6.

## 2. DETERMINATION OF OPTIMAL METHODS FOR (1.1), (1.11)

**2.1. An algorithm to find optimal methods.** This and the following subsection are devoted to the construction of optimal methods (1.2) for (1.1), (1.11). By "threshold factor" etc. we mean the term which refers to this test equation.

We use the notations

$$\begin{aligned} a_i &= (1, i, i^2, \dots, i^q)^T, & b_i(x) &= (1, i+x, i^2+2ix, \dots, i^q+qxi^{q-1})^T, \\ c_i &= (i, i^2, \dots, i^q)^T, & d_i(x) &= (i+x, i^2+2ix, \dots, i^q+qix^{q-1})^T. \end{aligned}$$

Here,  $i$  is an integer and  $x$  is a real number. The length of these vectors will not always be indicated explicitly.

By  $\text{co}(X)$  we mean the convex hull (see [16]) of a subset  $X$  of a linear space. The relative interior (see [16]) of such a subset  $X$  will be denoted by  $\text{ri}(X)$ . For  $r > 0$ ,  $k \geq 1$ , and  $p \geq 1$ , we define in  $\mathbb{R}^p$

$$K^p(r) = \text{co}\{c_i, d_j(r^{-1}) \mid 0 \leq i \leq k-1, 0 \leq j \leq k\}.$$

Further,  $\#(Y)$  denotes the number of elements of a set  $Y$ . For several other notions, and Carathéodory's theorem, the reader is referred to [16].

Consider a method (1.2) with  $\beta_k \geq 0$ . From (1.3) it follows, for  $r > 0$ , that (1.2) is of order at least  $p$  if and only if the equality

$$(2.1) \quad \sum_{i=0}^{k-1} (-\alpha_i - r\beta_i)a_i + \sum_{i=0}^{k-1} r\beta_i b_i(r^{-1}) = b_k(-\beta_k)$$

holds in  $\mathbb{R}^{p+1}$ . Equivalently, we have (in  $\mathbb{R}^{p+1}$ )

$$(2.2) \quad \sum_{i=0}^{k-1} (-\alpha_i - r\beta_i)(1 + \beta_k r)^{-1} a_i + \sum_{i=0}^k r\beta_i(1 + \beta_k r)^{-1} b_i(r^{-1}) = a_k.$$

We shall use this relation most often with  $a_i, b_i$ , and  $a_k$  replaced by  $c_i, d_i$ , and  $c_k$ , respectively.

Let  $k \geq 1, p \geq 1$ , and  $r_1$  be given with  $0 < r_1 < \infty$  and  $c_k \in K^p(r_1)$ . From the definition of  $K^p(\cdot)$  it can be seen that

$$(2.3) \quad K^p(r_1) \subset K^p(r_0) \quad (\text{for all } r_0 \text{ with } 0 < r_0 \leq r_1).$$

Suppose method (1.2) has a threshold factor  $S$  with  $0 < r_0 \leq S < \infty$ . It follows from (1.12) and (2.2) that  $c_k \in K^p(S)$ . Thus, by (2.3), with  $r_1 = S$ , we have  $c_k \in K^p(r_0)$ .

Assume  $r$  is a real number with  $0 < r \leq S_{k,p}$ . We have for each  $\varepsilon \in (0, r)$  that  $c_k \in K^p(r - \varepsilon)$ , by (1.13) and the preceding conclusion (with  $r_0 = r - \varepsilon$ ). It can be seen that this implies that  $c_k \in K^p(r)$ .

Conversely, if  $c_k \in K^p(r)$  for some  $r$  with  $0 < r < \infty$ , we can use (2.2) and (1.12) to construct a  $k$ -step method (1.2) of order at least  $p$  with threshold factor  $S \geq r$ . Hence,

$$(2.4) \quad S_{k,p} \geq r \quad \text{iff} \quad c_k \in K^p(r) \quad (\text{for all } r \text{ with } 0 < r < \infty).$$

Since the implicit Euler method is known to have  $S_{1,1} = R_{1,1} = \infty$  and  $S_{k,p} < \infty, R_{k,p} < \infty$  for methods of order  $p \geq 2$ , the restriction  $r < \infty$  in (2.4) is not important.

As a consequence of (1.12), (2.2), (2.4), and Theorem 2.3 in [17], we have the following existence theorem for optimal methods.

**Theorem 2.1.** (i) *Let  $k \geq 1, p \geq 1$ . If there exists a  $k$ -step method (1.2) of order at least  $p$ , then there exists a  $k$ -step method (1.2) of order at least  $p$  that is optimal for (1.1), (1.11).*

(ii) *For all  $p \geq 1$ , there exists an integer  $k(p)$  such that  $S_{k,p} > 0$  for all  $k \geq k(p)$ .*

We now give four lemmas which enable us to devise an algorithm to compute optimal methods.

**Lemma 2.2.** *Let  $k \geq 1, p \geq 2$ , and  $r$  be given with  $S_{k,p} > 0$  and  $0 < r < \infty$ . Then  $\dim(K^p(r)) = p$ .*

*Proof.* Since  $S_{k,p} > 0$ , there exists a stable (cf. [9])  $k$ -step method of order  $p$ . Hence  $p \leq 2+k$  by Theorem 5.9 in [9]. Using [7, p. 122], and some elementary operations, it can be seen that the first  $p+1$  vectors of the sequence of vectors in  $\mathbb{R}^{p+1}$

$$b_0(r^{-1}), a_0, b_1(r^{-1}), a_1, \dots, b_k(r^{-1})$$

are linearly independent. Hence the first  $p+1$  vectors of the sequence

$$d_0(r^{-1}), c_0, d_1(r^{-1}), c_1, \dots, d_k(r^{-1})$$

are affinely independent and  $\dim(K^p(r)) = p$ .  $\square$

**Lemma 2.3.** *Let  $k \geq 1$ ,  $p \geq 2$ , and let  $S$  be a real number,  $0 < S < \infty$ . Then  $S = S_{k,p}$  if and only if  $c_k$  is an element of the boundary  $\partial K^p(S)$  of  $K^p(S)$ .*

*Proof.* Assume  $S = S_{k,p}$ . By (2.4), we have  $c_k \in K^p(S)$ . If  $c_k$  were contained in the interior  $\text{int}(K^p(S))$  of  $K^p(S)$ , there would exist  $r > S$  with  $c_k \in K^p(r)$ . However, since  $S = S_{k,p}$ , this cannot be the case, by (2.4). Hence  $c_k \in \partial K^p(S)$ .

To prove the converse, let  $c_k \in \partial K^p(S)$ . Then  $S_{k,p} \geq S$  (cf. (2.4)). Assume  $S_{k,p} > S$ . Since  $p \geq 2$ , we know that  $S_{k,p} < \infty$  (cf. §1). A contradiction can be derived as follows. Let  $G$  be a face of  $K^p(S)$  with  $\dim(G) = p-1$ ,  $c_k \in G$ . This is possible by Lemma 2.2. There exists a method of the form (1.2) with threshold factor  $S_{k,p}$  (cf. Theorem 2.1). For some  $I \subset \{0, 1, \dots, k-1\}$ ,  $J \subset \{0, 1, \dots, k\}$ , we have by (2.2) (with  $r = S_{k,p}^{-1}$ ) that

$$(2.5) \quad c_k = \sum_{i \in I} \lambda_i c_i + \sum_{j \in J} \mu_j d_j(S_{k,p}^{-1}),$$

with  $\sum_{i \in I} \lambda_i + \sum_{j \in J} \mu_j = 1$ ,  $\lambda_i > 0$  ( $i \in I$ ) and  $\mu_j > 0$  ( $j \in J$ ). The set  $G$  is a face of  $K^p(S)$ ,  $c_k \in G$ , and, by (2.3),  $K^p(S_{k,p}) \subset K^p(S)$ . Hence,

$$\{c_i, d_j(S_{k,p}^{-1}) | i \in I, j \in J\} \subset G.$$

Since  $d_j(S_{k,p}^{-1}) \in \text{ri}(\text{co}\{c_j, d_j(S_{k,p}^{-1})\})$  and  $c_j, d_j(S_{k,p}^{-1}) \in K^p(S)$ , we also have  $c_j \in G$  and  $d_j(S_{k,p}^{-1}) \in G$  for all  $j \in J$ . The convexity of  $G$  then implies that

$$(2.6) \quad \{c_i, d_j(r^{-1}) | i \in I \cup J, j \in J, r \in [S, \infty)\} \subset G.$$

Assume  $k \in J$ . Then  $Sp_0 = \text{span}\{a_i, b_j(r^{-1}) | i \in I \cup J \setminus \{k\}, j \in J \setminus \{k\}\}$  is independent of  $r$  provided  $r \in (0, \infty)$ . Further, by (2.5),

$$(1 - \mu_k)^{-1}(a_k - \mu_k b_k(S_{k,p}^{-1})) = b_k(-\mu_k(1 - \mu_k)^{-1} S_{k,p}^{-1}) \in Sp_0.$$

Consequently,  $Sp = \text{span}\{a_i, b_j(r^{-1}) | i \in I \cup J \setminus \{k\}, j \in J\}$  is independent of  $r$ , provided  $r \in (0, \infty)$ .

The last assertion also holds when  $k \notin J$ .

There exist  $r_0$  with  $0 < r_0 < \infty$  and  $I_0 \subset I \cup J \setminus \{k\}$ ,  $J_0 \subset J$  such that

$$B(r) = \{a_i, b_j(r^{-1}) | i \in I_0, j \in J_0\}$$

is a basis of  $Sp$  when  $r = r_0$ . From (2.6) we know that  $\dim(Sp) \leq p$ , so that  $\#(B(r_0)) \leq p$ . Furthermore, there exists a neighborhood  $U$  of  $r_0$  such that  $B(r)$  is a basis of  $Sp$  for all  $r \in U$ . Since  $a_k \in Sp$  (cf. (2.5)), we see that there exist indices

$$i(1) < i(2) < \dots < i(m) \quad \text{and} \quad j(1) < j(2) < \dots < j(n)$$

such that

$$\det(a_{i(1)}, a_{i(2)}, \dots, a_{i(m)}, b_{j(1)}(r^{-1}), b_{j(2)}(r^{-1}), \dots, b_{j(n)}(r^{-1}), a_k) = 0 \quad (\text{for all } r \in U).$$

Hence, the determinant must vanish for all  $r > 0$ . However, it follows from [7, p. 122] that this cannot hold.

Therefore, the assumption  $S_{k,p} > S$  cannot hold and we must have  $S_{k,p} = S$ .  $\square$

The condition stated in Lemma 2.3 is *necessary* and *sufficient* for a method (1.2) to be optimal. However, from a computational point of view, this lemma is only useful to show whether or not a *given* method is optimal. In Algorithm 2.7 below, candidates for optimal methods will be obtained by Lemma 2.5. In the proof of this lemma we need

**Lemma 2.4.** *Let  $s \geq 1$ ,  $X = \{x(1), x(2), \dots, x(n)\} \subset \mathbb{R}^s$ ,  $K = \text{co}(X)$ ,  $d$  be the dimension of  $K$ , and let  $v$  belong to the relative interior  $\text{ri}(K)$  of  $K$ . Then*

(i)  $\Lambda = \{(\lambda(1), \lambda(2), \dots, \lambda(n))^T \mid v = \sum_{i=1}^n \lambda(i)x(i), \sum_{i=1}^n \lambda(i) = 1, \lambda(i) \geq 0 (1 \leq i \leq n)\}$  is a convex set whose dimension equals  $n - d - 1$ .

(ii) For each  $i$  ( $1 \leq i \leq n$ ), there exists a set  $Y \subset X$  of affinely independent vectors such that  $x(i) \in Y$ ,  $\#(Y) \leq d + 1$ , and  $v \in \text{ri}(\text{co}(Y))$ .

*Proof.* It is easy to see that  $\Lambda$  is convex. Further, there exist numbers  $\lambda(i) > 0$  such that  $v = \sum_{i=1}^n \lambda(i)x(i)$ . By using elementary linear algebra, it can then be seen that  $\dim(\Lambda) = n - d - 1$ . This proves part (i).

To prove part (ii), let  $1 \leq i \leq n$ . If  $v = x(i)$ , the assertion holds for this value of  $i$ . Otherwise, let  $\xi > 0$  be such that  $w = v + \xi(v - x(i))$  is contained in the relative boundary (cf. [16]) of  $K$ . Apply Carathéodory's theorem to  $w$  and a face  $F$  of  $K$ ,  $F \neq K$ , with  $w \in F$ , and use  $v = (1 + \xi)^{-1}(w + \xi x(i))$ .  $\square$

In the following lemma,  $k \geq 1$  and  $p \geq 2$  are given with  $S_{k,p} > 0$ . By Lemma 2.3 we have  $c_k \in \partial K^p(S_{k,p})$ . Hence we may denote by  $F$  the smallest face of  $K^p(S_{k,p})$  such that  $c_k \in \text{ri}(F)$ . The dimension  $d$  of  $F$  satisfies  $d \leq p - 1$ . By Lemma 2.2 (with  $r = S_{k,p} < \infty$ ) there exists a  $(p - 1)$ -dimensional face  $G$  of  $K^p(S_{k,p})$  such that  $F \subset G$ . Further, we introduce  $I = \{i \mid 0 \leq i \leq k - 1 \text{ and } c_i \in F\}$  and  $J = \{j \mid 0 \leq j \leq k \text{ and } d_j(S_{k,p}^{-1}) \in F\}$ .

**Lemma 2.5.** (i) *The coefficients*

$$(-\alpha_0 - S_{k,p}\beta_0)(1 + \beta_k S_{k,p})^{-1}, \dots, (-\alpha_{k-1} - S_{k,p}\beta_{k-1})(1 + \beta_k S_{k,p})^{-1}, \\ S_{k,p}\beta_0(1 + \beta_k S_{k,p})^{-1}, \dots, S_{k,p}\beta_k(1 + \beta_k S_{k,p})^{-1}$$

*in (2.2) of the class of optimal  $k$ -step methods (1.2) of order  $p$  form a convex set in  $\mathbb{R}^{2k+1}$  of dimension  $\#(I) + \#(J) - d - 1$ .*

(ii) To each  $I_0 \subset I$ ,  $J_0 \subset J$  with  $\#(I_0) + \#(J_0) = 1$ , there exist  $I_1, J_1$ , with  $I_0 \subset I_1 \subset I$ ,  $J_0 \subset J_1 \subset J$ , such that  $Y = \{c_i, d_j(S_{k,p}^{-1}) | i \in I_1, j \in J_1\}$  is a set of at most  $d + 1$  affinely independent vectors and  $c_k \in \text{ri}(\text{co}(Y))$ .

(iii) Let  $I_1, J_1$  be as in (ii). There exists a unique optimal method (1.2) such that for the coefficients in (2.2)

$$\begin{aligned} (-\alpha_i - S_{k,p}\beta_i)(1 + \beta_k S_{k,p})^{-1} &= 0 \quad (\text{for all } i \notin I_1), \\ S_{k,p}\beta_j(1 + \beta_k S_{k,p})^{-1} &= 0 \quad (\text{for all } j \notin J_1). \end{aligned}$$

(iv) For each  $I_1, J_1$  as in (ii) there exist

$$\begin{aligned} I_2 &= \{i(1), i(2), \dots, i(m)\} \subset \{0, 1, \dots, k-1\}, \\ J_2 &= \{j(1), j(2), \dots, j(n)\} \subset \{0, 1, \dots, k\}, \end{aligned}$$

with  $I_1 \subset I_2$ ,  $J_1 \subset J_2$  and such that  $X = \{c_i, d_j(S_{k,p}^{-1}) | i \in I_2, j \in J_2\}$  is a set of  $p$  affinely independent vectors and  $G$  is contained in the affine hull  $\text{aff}(X)$  of  $X$ .

(v) If  $\#(G \cap \{c_i, d_j(S_{k,p}^{-1}) | 0 \leq i \leq k-1, 0 \leq j \leq k\}) = p$ , then the optimal  $k$ -step method (1.2) of order  $p$  is unique.

(vi) Let  $I_2, J_2$  be as in (iv). Then  $S_{k,p}$  is a simple zero of the polynomial

$$\begin{aligned} P(r) &= r^n \det(a_{i(1)}, a_{i(2)}, \dots, a_{i(m)}, \\ &\quad b_{j(1)}(r^{-1}), b_{j(2)}(r^{-1}), \dots, b_{j(n)}(r^{-1}), a_k). \end{aligned}$$

*Proof.* Each optimal  $k$ -step method (1.2) of order  $p$  corresponds in a one-to-one fashion to an expression of  $a_k$  as a convex combination (2.2) (with  $r = S_{k,p}$ ) of the vectors  $a_i, b_j(S_{k,p}^{-1})$ . This is an immediate consequence of (1.12), (1.13). Since  $F$  is a face of  $K^p(S_{k,p})$  with  $c_k \in F$ , the coefficients in (2.2) corresponding to  $i \notin I, j \notin J$  are zero. We may now use Lemma 2.4 with  $s = p$ ,  $X = \{c_i, d_j(S_{k,p}^{-1}) | i \in I, j \in J\}$ ,  $K = F$ , and  $v = c_k$ . Parts (i) and (ii) of Lemma 2.5 follow from parts (i) and (ii) of Lemma 2.4, respectively.

Part (iii) follows from part (ii) of Lemma 2.5 and formula (2.2).

Since  $F \subset G$ , part (ii) also implies part (iv).

Assume the condition of part (v) is fulfilled. Since  $\dim(G) = p - 1$  and  $F \subset G$ , we necessarily have  $\#(I) + \#(J) = \dim(F) + 1 = d + 1$ . Thus, we can apply part (i) to prove part (v).

As for part (vi), it follows from (iv) and the inclusion  $c_k \in G$  that  $P(S_{k,p}) = 0$ . Some computation, using Lemma 2.2 (with  $r = S_{k,p}$ ), shows the fact that  $G$  is a face of  $K^p(S_{k,p})$ , and [7, p. 122], that

$$(2.7) \quad \frac{dP}{dr}(S_{k,p}) \det\left(a_{i(1)}, \dots, a_{i(m)}, b_{j(1)}(S_{k,p}^{-1}), \dots, b_{j(n)}(S_{k,p}^{-1}), \begin{pmatrix} 1 \\ v \end{pmatrix}\right) < 0$$

for all  $v \in K^p(S_{k,p}) \setminus \text{aff}(X)$ . Hence  $S_{k,p}$  is a simple zero of  $P$ .  $\square$



**Corollary 2.6.** *Let  $k \geq 1$ ,  $p \geq 2$ ,  $S_{k,p} > 0$ , and let  $I$ ,  $J$ , and  $d$  be as above. An optimal  $k$ -step method of order  $p$  is unique if and only if  $\#(I) + \#(J) = d + 1$ .*

In the following algorithm, we supply integers  $k \geq 1$ ,  $p \geq 2$ . Using the algorithm, we can find all optimal  $k$ -step methods of order at least  $p$  for which there exist sets  $I_2$  and  $J_2$  as in part (iv) of Lemma 2.5. In particular, if the optimal threshold factor is positive, we obtain at least one optimal method. The sufficient condition in part (v) of the same lemma is convenient to show the uniqueness of an optimal method.

**Algorithm 2.7.** 1. Give integers  $k \geq 1$ ,  $p \geq 2$ .

2. Find a (new) pair of integer sequences  $I = \{i(1), i(2), \dots, i(m)\} \subset \{0, 1, \dots, k - 1\}$ ,  $J = \{j(1), j(2), \dots, j(n)\} \subset \{0, 1, \dots, k\}$  with  $m + n = p$ . If such a pair does not exist, then stop.

3. Compute positive roots, at least all those that are simple, of the polynomial

$$P(r) = r^n \det(a_{i(1)}, a_{i(2)}, \dots, a_{i(m)}, \\ b_{j(1)}(r^{-1}), b_{j(2)}(r^{-1}), \dots, b_{j(n)}(r^{-1}), a_k).$$

4. Verify whether for any of these roots, say  $r$ ,

- (a)  $V = \{c_i, d_j(r^{-1}) | i \in I, j \in J\}$  is a set of  $p$  affinely independent vectors in  $\mathbb{R}^p$ ,
- (b)  $K^p(r)$  lies on one side of the affine hull  $\text{aff}(V)$  of  $V$ ,
- (c)  $c_k \in K^p(r)$ .

If these conditions hold,  $r = S_{k,p}$ . Determine coefficients  $\alpha_i$  ( $0 \leq i \leq k - 1$ ) and  $\beta_j$  ( $0 \leq j \leq k$ ) of an optimal method (1.2) from equation (2.2) and

$$(-\alpha_i - r\beta_i)(1 + \beta_k r)^{-1} = 0 \quad (\text{for all } i \notin I), \\ r\beta_j(1 + \beta_k r)^{-1} = 0 \quad (\text{for all } j \notin J).$$

If, in addition,  $c_i \notin \text{aff}(V)$  (for  $i \notin I$ ),  $d_j(r^{-1}) \notin \text{aff}(V)$  (for  $j \notin J$ ), a unique optimal method has been found. Stop.

Otherwise, return to step 2.

**2.2. A survey of optimal threshold factors and optimal methods.** We can apply our algorithm with  $k \geq 1$ ,  $p \geq 2$  to find  $S_{k,p}$  and corresponding optimal methods. For order  $p = 1$ , optimal methods can be found “by hand”, using (1.12) and (2.1). Table 1 lists the positive threshold factors  $S_{k,p}$  obtained in this way for  $k \leq 20$ ,  $p \leq 8$ . We will also give some of the corresponding optimal methods.

(i)  $k \geq 1$ ,  $p = 1$ .  $S_{k,p} = \infty$  and (1.2) is an optimal method for (1.1), (1.11) whenever

$$\alpha_i \leq 0, \beta_i = 0 \quad (0 \leq i \leq k - 1); \quad \sum_{i=0}^{k-1} \alpha_i = -1, \quad \beta_k = \sum_{i=0}^{k-1} i\alpha_i + k.$$

In case  $k = 1$ , we have the so-called *implicit Euler method*.

(ii)  $k \geq 1, p = 2$ .  $S_{k,p} = 2$  and (1.2) is a unique optimal method with

$$\alpha_{k-1} = -1, \quad \beta_{k-1} = \beta_k = 1/2, \quad \alpha_i = \beta_i = 0 \quad (0 \leq i \leq k-2).$$

If  $k = 1$ , this is the so-called *trapezoidal rule*.

(iii)  $k \geq 2, p = 3$ .  $S_{k,p} = (2k-3)(k-1)^{-1}$  and (1.2) is a unique optimal method with nonzero coefficients

$$\begin{aligned} \alpha_0 &= -(2k+1)^{-1}(k-1)^{-2}, \\ \alpha_{k-1} &= -k^2(2k-3)(k-1)^{-2}(2k+1)^{-1}, \\ \beta_{k-1} &= k^2(2k+1)^{-1}(k-1)^{-1}, \\ \beta_k &= k(2k+1)^{-1}. \end{aligned}$$

(iv)  $k = 5, p = 6$ . This is a rather exceptional situation in that there exists a 1-parameter set of optimal methods, with  $S_{5,6} = 1/2$ . From the results of the application of Algorithm 2.7 one can deduce that (1.2) is optimal if and only if  $\lambda \in [0, 1]$  and

$$\begin{aligned} \alpha_0 &= \lambda \frac{-459}{3709} + (1-\lambda) \frac{-513}{5888}, & \beta_0 &= \lambda \frac{210}{3709} + (1-\lambda) \frac{135}{2944}, \\ \alpha_1 &= \lambda \frac{-125}{3709}, & \beta_1 &= \lambda \frac{250}{3709}, \\ \alpha_2 &= \lambda \frac{-1000}{3709} + (1-\lambda) \frac{-125}{368}, & \beta_2 &= \lambda \frac{2000}{3709} + (1-\lambda) \frac{375}{736}, \\ \alpha_3 &= 0, & \beta_3 &= 0, \\ \alpha_4 &= \lambda \frac{-2125}{3709} + (1-\lambda) \frac{-3375}{5888}, & \beta_4 &= \lambda \frac{4250}{3709} + (1-\lambda) \frac{3375}{2944}, \\ & & \beta_5 &= \lambda \frac{1210}{3709} + (1-\lambda) \frac{15}{46}. \end{aligned}$$

(v) For all other values of  $k$  and  $p$  in Table 1, we found a unique corresponding optimal method. The coefficients of some of them are listed in Tables 2, 3, 4.

Although, for reasons of convention, we gave the coefficients  $\alpha_i$  and  $\beta_j$  in the above enumeration, the preceding theory suggests we use, instead of the classical form (1.2), the equivalent form

$$(2.8a) \quad u_n - h\beta_k f(t_n, u_n) = \sum_{i=0}^{k-1} \gamma_i u_{n-k+i} + \sum_{j=0}^{k-1} \delta_j \left( u_{n-k+j} + hS_{k,p}^{-1} f(t_{n-k+j}, u_{n-k+j}) \right)$$

( $n = k, k+1, \dots$ ). Here we assume that  $S_{k,p} < \infty$  and

$$(2.8b) \quad \gamma_i = -\alpha_i - S_{k,p} \beta_i \quad (0 \leq i \leq k-1),$$

$$(2.8c) \quad \delta_j = S_{k,p} \beta_j \quad (0 \leq j \leq k-1).$$

TABLE 1  
Optimal threshold factors  $S_{k,p}$  for  $k$ -step methods (1.2) of order  $p$

| $p$<br>$k$ | 1        | 2 | 3      | 4      | 5      | 6      | 7      | 8      |
|------------|----------|---|--------|--------|--------|--------|--------|--------|
| 1          | $\infty$ | 2 |        |        |        |        |        |        |
| 2          | $\infty$ | 2 | 1.0000 |        |        |        |        |        |
| 3          | $\infty$ | 2 | 1.5000 | 1.0000 |        |        |        |        |
| 4          | $\infty$ | 2 | 1.6667 | 1.2432 | 0.6667 |        |        |        |
| 5          | $\infty$ | 2 | 1.7500 | 1.2432 | 0.7955 | 0.5000 |        |        |
| 6          | $\infty$ | 2 | 1.8000 | 1.2432 | 0.9294 | 0.6597 | 0.3000 |        |
| 7          | $\infty$ | 2 | 1.8333 | 1.2432 | 1.0056 | 0.7837 | 0.4676 | 0.1965 |
| 8          | $\infty$ | 2 | 1.8571 | 1.2432 | 1.0524 | 0.8683 | 0.5500 | 0.3450 |
| 9          | $\infty$ | 2 | 1.8750 | 1.2432 | 1.0837 | 0.9053 | 0.6420 | 0.4426 |
| 10         | $\infty$ | 2 | 1.8889 | 1.2432 | 1.1061 | 0.9053 | 0.6901 | 0.5328 |
| 11         | $\infty$ | 2 | 1.9000 | 1.2432 | 1.1229 | 0.9053 | 0.7329 | 0.5803 |
| 12         | $\infty$ | 2 | 1.9091 | 1.2432 | 1.1361 | 0.9053 | 0.7644 | 0.6249 |
| 13         | $\infty$ | 2 | 1.9167 | 1.2432 | 1.1466 | 0.9053 | 0.7810 | 0.6620 |
| 14         | $\infty$ | 2 | 1.9231 | 1.2432 | 1.1552 | 0.9053 | 0.7946 | 0.6920 |
| 15         | $\infty$ | 2 | 1.9286 | 1.2432 | 1.1624 | 0.9053 | 0.8057 | 0.7138 |
| 16         | $\infty$ | 2 | 1.9333 | 1.2432 | 1.1685 | 0.9053 | 0.8149 | 0.7189 |
| 17         | $\infty$ | 2 | 1.9375 | 1.2432 | 1.1737 | 0.9053 | 0.8226 | 0.7189 |
| 18         | $\infty$ | 2 | 1.9412 | 1.2432 | 1.1783 | 0.9053 | 0.8291 | 0.7189 |
| 19         | $\infty$ | 2 | 1.9444 | 1.2432 | 1.1822 | 0.9053 | 0.8346 | 0.7189 |
| 20         | $\infty$ | 2 | 1.9474 | 1.2432 | 1.1858 | 0.9053 | 0.8394 | 0.7189 |

TABLE 2  
The nonzero coefficients of optimal  $k$ -step methods of order 4

|         |                              |                            |
|---------|------------------------------|----------------------------|
| $k = 3$ | $\alpha_0 = -5/32$           | $\beta_0 = 3/32$           |
|         | $\alpha_2 = -27/32$          | $\beta_2 = 27/32$          |
|         |                              | $\beta_3 = 3/8$            |
| $k = 4$ | $\alpha_0 = -0.044656443869$ | $\beta_0 = 0.035920310223$ |
|         | $\alpha_1 = -0.032287478509$ | $\beta_1 = 0.025971083765$ |
|         | $\alpha_3 = -0.923056077622$ | $\beta_3 = 0.742478750864$ |
|         |                              | $\beta_4 = 0.394174143773$ |

These coefficients can easily be obtained from the tables with

$$S_{k,p} = -\alpha_{k-1}/\beta_{k-1},$$

as it has turned out that  $\gamma_{k-1} = 0$  for optimal methods of order  $p$  with  $2 \leq p \leq 8$ .

TABLE 3

*The nonzero coefficients of optimal k-step methods of order 5*

|          |   |  |
|----------|---|--|
| $k = 4$  | $\alpha_0 = -3/35$<br>$\alpha_1 = -8/35$<br>$\alpha_3 = -24/35$   | $\beta_1 = 12/35$<br>$\beta_3 = 36/35$<br>$\beta_4 = 12/35$  |
| $k = 5$  | $\alpha_0 = -0.024644773737$<br>$\alpha_1 = -0.049099680927$<br>$\alpha_2 = -0.166384720704$<br>$\alpha_4 = -0.759870824632$    | $\beta_1 = 0.061720193877$<br>$\beta_2 = 0.209152015371$<br>$\beta_4 = 0.955186952991$<br>$\beta_5 = 0.352588416898$       |
| $k = 6$  | $\alpha_0 = -0.005965846424$<br>$\alpha_2 = -0.062219568345$<br>$\alpha_3 = -0.107893579835$<br>$\alpha_5 = -0.823921005396$    | $\beta_2 = 0.066948740947$<br>$\beta_3 = 0.116094333637$<br>$\beta_5 = 0.886545429648$<br>$\beta_6 = 0.362686592593$       |
| $k = 7$  | $\alpha_0 = -0.001899283695$<br>$\alpha_3 = -0.063432086540$<br>$\alpha_4 = -0.079765609816$<br>$\alpha_6 = -0.854903019949$    | $\beta_3 = 0.063077153883$<br>$\beta_4 = 0.079319283337$<br>$\beta_6 = 0.850119431442$<br>$\beta_7 = 0.368707312757$       |
| $k = 8$  | $\alpha_0 = -0.000746897457$<br>$\alpha_4 = -0.062185187770$<br>$\alpha_5 = -0.065188908644$<br>$\alpha_7 = -0.871879006129$    | $\beta_4 = 0.059089849737$<br>$\beta_5 = 0.061944057008$<br>$\beta_7 = 0.828480242782$<br>$\beta_8 = 0.372647513271$       |
| $k = 9$  | $\alpha_0 = -0.000340789183$<br>$\alpha_5 = -0.060559712643$<br>$\alpha_6 = -0.056750393804$<br>$\alpha_8 = -0.882349104370$    | $\beta_5 = 0.055882324853$<br>$\beta_6 = 0.052367222428$<br>$\beta_8 = 0.814200020652$<br>$\beta_9 = 0.375456671066$       |
| $k = 10$ | $\alpha_0 = -0.000173264587$<br>$\alpha_6 = -0.059033126258$<br>$\alpha_7 = -0.051416380337$<br>$\alpha_9 = -0.889377228819$    | $\beta_6 = 0.053370030451$<br>$\beta_7 = 0.046483965160$<br>$\beta_9 = 0.804058548035$<br>$\beta_{10} = 0.377578977080$    |
| $k = 11$ | $\alpha_0 = -0.000095614309$<br>$\alpha_7 = -0.057700216566$<br>$\alpha_8 = -0.047811920852$<br>$\alpha_{10} = -0.894392248273$ | $\beta_7 = 0.051383050339$<br>$\beta_8 = 0.042577350348$<br>$\beta_{10} = 0.796471913807$<br>$\beta_{11} = 0.379248320000$ |
| $k = 12$ | $\alpha_0 = -0.000056258666$<br>$\alpha_8 = -0.056558096984$<br>$\alpha_9 = -0.045247851084$<br>$\alpha_{11} = -0.898137793267$ | $\beta_8 = 0.049784569697$<br>$\beta_9 = 0.039828864761$<br>$\beta_{11} = 0.790574753235$<br>$\beta_{12} = 0.380600650748$ |

TABLE 4  
*The nonzero coefficients of optimal k-step methods of order 6*

|         |  |  |
|---------|--|--|
| $k = 5$ | one parameter set<br>as described above  |  |
| $k = 6$ | $\alpha_0 = -0.039244965860$<br>$\alpha_2 = -0.076331008461$<br>$\alpha_3 = -0.196586017197$<br>$\alpha_5 = -0.687838008482$                                 | $\beta_0 = 0.024348556303$<br>$\beta_2 = 0.115708903958$<br>$\beta_3 = 0.298001468104$<br>$\beta_5 = 1.042682176830$<br>$\beta_6 = 0.337648783882$                               |
| $k = 7$ | $\alpha_0 = -0.013219720715$<br>$\alpha_3 = -0.084824895181$<br>$\alpha_4 = -0.141412093193$<br>$\alpha_6 = -0.760543290911$                                 | $\beta_0 = 0.010637484589$<br>$\beta_3 = 0.108232169095$<br>$\beta_4 = 0.180434500388$<br>$\beta_6 = 0.970413814127$<br>$\beta_7 = 0.346899228021$                               |
| $k = 8$ | $\alpha_0 = -0.005219538160$<br>$\alpha_4 = -0.083655994035$<br>$\alpha_5 = -0.108421088766$<br>$\alpha_7 = -0.802703379039$                                 | $\beta_0 = 0.005143493020$<br>$\beta_4 = 0.096341243207$<br>$\beta_5 = 0.124861614544$<br>$\beta_7 = 0.924422001734$<br>$\beta_8 = 0.353578574250$                               |
| $k = 9$ | $\alpha_0 = -0.001642953120$<br>$\alpha_1 = -0.001731087076$<br>$\alpha_5 = -0.081210958324$<br>$\alpha_6 = -0.096236043942$<br>$\alpha_8 = -0.819178957538$ | $\beta_0 = 0.001814860688$<br>$\beta_1 = 0.001912216389$<br>$\beta_5 = 0.089708326990$<br>$\beta_6 = 0.106305536547$<br>$\beta_8 = 0.904892335991$<br>$\beta_9 = 0.356732920744$ |

For each  $k \geq 1$  and  $p \geq 2$  with  $S_{k,p} > 0$ , there exists an optimal  $k$ -step method of order  $p$  with at most  $p$  of the coefficients  $\beta_k, \gamma_i, \delta_j$  ( $0 \leq i, j \leq k - 1$ ) not equal to zero (cf. part (iii) of Lemma 2.5). For  $k = 4, p = 5$ , the optimal method even has only four nonzero coefficients.

We found a remarkable method for  $k = 1, p = 2$  and for  $k = 4, p = 4$  and for  $k = 9, p = 6$ . In these cases, all the  $\gamma_i$  ( $0 \leq i \leq k - 1$ ) vanish. By storing the fixed linear combinations  $u_i + hS_{k,p}^{-1}f(t_i, u_i)$  (instead of the usual vectors  $u_i$  and  $f(t_i, u_i)$ ), and using these in (2.8a), we obtain an efficient implementation of these methods. Other interesting features of these methods will be exhibited in §§3 and 4.

We conjecture that also for  $k = 16, p = 8$ , the coefficients  $\gamma_i$  of the corresponding optimal method vanish. However we have no complete proof of this.

3. ASYMPTOTIC BEHAVIOR OF  $S_{k,p}$  WHEN  $k$  TENDS TO INFINITY

To a given  $k$ -step method (1.2) we *adjoin*, for  $l \geq k$ , the  $l$ -step method

$$(3.1) \quad u_n - h\beta_l^l f(t_n, u_n) = \sum_{i=0}^{l-1} (-\alpha_i^l u_{n-l+i} + h\beta_i^l f(t_{n-l+i}, u_{n-l+i})) \quad (n \geq l).$$

Here,  $\beta_i^l = \beta_{i-(l-k)}$  ( $l-k \leq i \leq l$ ),  $\alpha_i^l = \alpha_{i-(l-k)}$  ( $l-k \leq i \leq l-1$ ), and  $\beta_i^l = \alpha_i^l = 0$  ( $0 \leq i < l-k$ ).

Inspired by the manifest regularities in Table 1, we will derive two theorems concerning the form of optimal methods for fixed order  $p$  when  $k$  tends to infinity.

**Theorem 3.1.** *Let  $k \geq 1, l \geq k, p \geq 2, p$  even,  $S_{k,p} > 0$ . Let (1.2) be an optimal  $k$ -step method of order at least  $p$  such that*

$$(3.2) \quad -\alpha_i - S_{k,p}\beta_i = 0 \quad (0 \leq i \leq k-1) \quad \text{and} \quad \#\{i \mid \beta_i \neq 0\} = p.$$

*Then:*

(i) *There exist indices  $j_1 < j_2 < \dots < j_{p/2}$  such that*

$$\{j \mid \beta_j \neq 0\} = \{j_1, j_1 + 1, j_2, j_2 + 1, \dots, j_{p/2}, j_{p/2} + 1\}.$$

(ii)  $S_{l,p} = S_{k,p}$ , and the  $l$ -step method adjoined to (1.2) is an optimal  $l$ -step method of order at least  $p$ .

(iii) An optimal  $l$ -step method of order at least  $p$  is unique if and only if (1.2) is a unique optimal  $k$ -step method of order at least  $p$ .

*Proof.* Let  $j(1) < j(2) < \dots < j(p)$  be such that  $\beta_{j(m)} \neq 0$  ( $1 \leq m \leq p$ ). Since  $S_{k,p} < \infty$  for  $p \geq 2$ , we can define

$$(3.3) \quad f_n(v) = \det \left( b_{j(1)+n}(S_{k,p}^{-1}), \dots, b_{j(p)+n}(S_{k,p}^{-1}), \begin{pmatrix} 1 \\ v \end{pmatrix} \right)$$

(for all  $v \in \mathbb{R}^p$  and  $n \in \mathbb{Z}$ ).

By (2.2), (3.2), and the fact that  $\beta_{j(m)} > 0$  for  $1 \leq m \leq p$  (cf. (1.12)), the vector  $c_k$  is contained in  $\text{ri}(\text{co}\{d_{j(m)}(S_{k,p}^{-1}) \mid 1 \leq m \leq p\})$ . Using Lemmas 2.3 (with  $S = S_{k,p}$ ), 6.2, and 6.1, it can be seen that  $K^p(S_{k,p})$  lies on one side of  $\text{aff}\{d_{j(1)}(S_{k,p}^{-1}), \dots, d_{j(p)}(S_{k,p}^{-1})\}$ . It follows from Lemma 6.4 (with  $i = y = 0, x = S_{k,p}^{-1}$ ) that  $f_0(v) > 0$  for  $v = c_0$ . Therefore, we must have

$$(3.4) \quad f_0(v) \geq 0 \quad (\text{for all } v \in K^p(S_{k,p})).$$

In particular, we may choose  $v = d_j(S_{k,p}^{-1})$ ,  $0 \leq j \leq k$ . Part (i) follows by interchanging columns in (3.3) and using Lemmas 6.1 and 6.2.

To prove part (ii), we appeal to Lemma 2.3 with  $k = l$  and  $S = S_{k,p}$ . We will show that  $c_l \in \partial K^p(S_{k,p})$ .

Let  $v = c_i$  for some  $i$  with  $0 \leq i \leq l - 1$ . Then  $f_{l-k}(v) = f_0(c_{i-(l-k)})$  by Lemma 6.3. This last value is nonnegative by Lemma 6.4 and (3.4). Similarly,  $f_{l-k}(v) \geq 0$  for  $v = d_j(S_{k,p}^{-1})$  ( $0 \leq j \leq l$ ), by Lemmas 6.1 and 6.2 and by part (i). Consequently,  $f_{l-k}(v) \geq 0$  for all  $v \in K^p(S_{k,p})$ . Further,  $c_l \in K^p(S_{k,p})$  and  $f_{l-k}(c_l) = f_0(c_k) = 0$ , by (2.2) and Lemma 6.3 (part (i)).

Therefore,  $c_l \in \partial K^p(S_{k,p})$ . An application of Lemma 2.3 shows that  $S_{k,p} = S_{l,p}$ . Part (ii) now follows immediately.

Finally, part (iii) can be seen to follow from Corollary 2.6, using (3.3) and Lemmas 6.1, 6.2, 6.3, 6.4.  $\square$

**Corollary 3.2.** *Let  $p = 2, 4$ , or  $6$ , and  $k \geq p^2/4$ . There exists a unique optimal  $k$ -step method of order  $p$ . It is equal to the  $k$ -step method adjoined to the optimal  $p^2/4$ -step method of order  $p$ .*

*Proof.* The unique optimal methods of order 2, 4, 6 and stepnumber 1, 4, 9, respectively, obtained in §2.2, satisfy condition (3.2) of Theorem 3.1.  $\square$

**Theorem 3.3.** *Let  $k \geq 1$ ,  $p \geq 2$ ,  $p$  even,  $S_{k,p} > 0$ . Let (1.2) be a unique optimal  $k$ -step method of order  $p$ . Let (3.2) be satisfied and*

$$(3.5) \quad \det((0, \dots, 0, 1, 0)^T, b_{j(1)}(S_{k,p}^{-1}), \dots, b_{j(p)}(S_{k,p}^{-1}), a_k) < 0,$$

with  $j(1), \dots, j(p)$  as in the proof of Theorem 3.1. Then:

(i) *For  $l$  large enough, there exists a unique optimal  $l$ -step method (1.2) of order at least  $p + 1$ .*

(ii)  $S_{l,p+1} - S_{k,p} = \mathcal{O}(l^{-1})$  (for  $l \rightarrow \infty$ ).

If  $\alpha_i^l$  ( $0 \leq i \leq l - 1$ ),  $\beta_i^l$  ( $0 \leq i \leq l$ ) denote the coefficients as in (3.1) of the method in (i), then:

(iii) *For  $l$  large enough,*

$$(3.6a) \quad -\alpha_i^l - S_{l,p+1}\beta_i^l = 0 \quad (\text{for } 1 \leq i \leq l - 1), \quad -\alpha_0^l > 0,$$

$$(3.6b) \quad \beta_{j(m)+l-k}^l > 0 \quad (\text{for } 1 \leq m \leq p), \quad \beta_i^l = 0 \quad (i \text{ otherwise}).$$

(iv)  $-\alpha_0^l = \mathcal{O}(l^{-(p+1)})$  (for  $l \rightarrow \infty$ ),  $\beta_{j(m)+l-k}^l - \beta_{j(m)}^l = \mathcal{O}(l^{-1})$  (for  $l \rightarrow \infty$ ,  $1 \leq m \leq p$ ).

*Proof.* The proof will be given in four steps.

1. Define for  $s \geq 1$ ,  $x \geq 0$ , and  $v \in \mathbb{R}^{p+1}$

$$g(s, x, v) = \det \left( a_{-s}, b_{j(1)}(x), b_{j(2)}(x), \dots, b_{j(p)}(x), \begin{pmatrix} 1 \\ v \end{pmatrix} \right).$$

From the implicit function theorem, (2.7), Lemma 6.4 (with  $i = y = 0$ ,  $x = S_{k,p}^{-1}$ ), and (3.5), it follows that, for  $s$  large enough, there exists a unique  $x_s > S_{k,p}^{-1}$  with  $g(s, x_s, c_k) = 0$  in some neighborhood of  $S_{k,p}^{-1}$ . It can be verified that  $x_s - S_{k,p}^{-1} = \mathcal{O}(s^{-1})$  (for  $s \rightarrow \infty$ ).

2. Let  $s$  be large enough. From Lemma 6.4 (with  $i = s$ ,  $y = 0$ ,  $x = x_s$ ) and Theorem 3.1 (part (i)), we see that  $g(s, x_s, d_j(x_s)) > 0$  for  $j$  not equal to any  $j(m)$  and  $-s \leq j \leq k$ . Likewise,  $g(s, x_s, c_{-i}) > 0$  for  $1 \leq i < s$ , by Lemma 6.5 (with  $j = s$ ,  $x = x_s$ ). Finally,  $g(s, x_s, c_i) > 0$  for  $0 \leq i \leq k-1$ , since  $s$  is large enough and using Lemma 6.4 (with  $i = y = 0$ ,  $x = x_s$ ), Lemma 2.3 (with  $S = S_{k,p}$ ), the fact that (1.2) is a unique optimal method, and Corollary 2.6.

Applying, then, Lemma 6.3  $s$  times, we see that  $K^{p+1}(x_s)$  lies on one side of  $\text{aff}\{c_0, d_{j(1)+s}(x_s), d_{j(2)+s}(x_s), \dots, d_{j(p)+s}(x_s)\}$ .

3. The vectors  $a_{-s}, b_{j(1)}(x_s), \dots, b_{j(p)}(x_s)$  in  $\mathbb{R}^{p+2}$  are linearly independent (use Lemma 6.4 with  $i = s$ ,  $x = x_s$ ,  $y = 0$ ). Since  $g(s, x_s, c_k) = 0$ , there exist unique coefficients  $\gamma_s$  and  $\delta_m^s$  ( $1 \leq m \leq p$ ) such that (in  $\mathbb{R}^{p+2}$ )

$$(3.7) \quad a_k = \gamma_s a_{-s} + \sum_{m=1}^p \delta_m^s b_{j(m)}(x_s).$$

Standard perturbation analysis, (2.2), (3.2), and the result of step 1 show that  $\gamma_s = \mathcal{O}(s^{-(p+1)})$  and

$$\delta_m^s - S_{k,p} \beta_{j(m)} (1 + \beta_k S_{k,p})^{-1} = \mathcal{O}(s^{-1}) \quad (\text{for } s \rightarrow \infty, 1 \leq m \leq p).$$

Further,  $\delta_m^s$  ( $1 \leq m \leq p$ ) are positive when  $s$  is large enough. By the substitution (3.7) it is evident that  $\gamma_s > 0$  if and only if

$$(3.8) \quad \det(a_{-s}, b_{j(1)}(x_s), \dots, b_{j(p)}(x_s)) \cdot \det(a_k, b_{j(1)}(x_s), \dots, b_{j(p)}(x_s)) > 0.$$

We may use (2.7) with  $m = 0$ ,  $n = p$ ,  $v = c_0$ . Since  $x_s > S_{k,p}^{-1}$  (for  $s \rightarrow \infty$ ), we obtain from (2.7) and Lemma 6.4 that (3.8) holds (for  $s \rightarrow \infty$ ). Therefore,  $\gamma_s > 0$  for  $s$  large enough.

4. Let  $l$  be large enough. We choose  $s = l - k$ , and determine coefficients  $\alpha_i^l, \beta_j^l$  such that

$$(3.9) \quad \begin{aligned} &(-\alpha_0^l - x_{l-k}^{-1} \beta_0^l)(1 + \beta_l^l x_{l-k}^{-1})^{-1} = \gamma_{l-k}, \\ &-\alpha_i^l - x_{l-k}^{-1} \beta_i^l = 0 \quad (1 \leq i \leq l-1), \\ &x_{l-k}^{-1} \beta_{j(m)+l-k}^l (1 + \beta_l^l x_{l-k}^{-1})^{-1} = \delta_m^s \quad (1 \leq m \leq p), \\ &\beta_j^l = 0 \quad (\text{if there is no } m \text{ with } j = j(m) + l - k). \end{aligned}$$

Let (3.1) be the  $l$ -step method with these coefficients. This  $l$ -step method is of order at least  $p+1$ , by (3.7), Lemma 6.3 (part (i), applied  $l-k$  times), and (2.2) (in  $\mathbb{R}^{p+2}$  and with  $k = l$ ,  $r = x_{l-k}^{-1}$ ).



From Lemma 2.3 (with  $k = l$ ,  $p$  replaced by  $p + 1$ , and  $S = x_{l-k}^{-1}$ ), (3.7), Lemma 6.3, and the result of step 2, it follows that this is an optimal method and that  $S_{l,p+1} = x_{l-k}^{-1}$ . Part (ii) of Theorem 3.3 thus follows from step 1.

The uniqueness of this optimal method can be shown by applying Lemma 2.5 (part (v)) and using step 2 and Lemma 6.3. This gives (i).

Finally, parts (iii) and (iv) follow from step 3 and (3.9).  $\square$

The optimal methods of stepnumber 1, 4, 9 and order 2, 4, 6, respectively, satisfy the conditions of Theorem 3.3. Hence, a description of the structure of optimal  $l$ -step methods of order 3, 5, 7 is provided by this theorem for  $l$  large enough. We found that the coefficients of the optimal  $l$ -step method satisfy (3.6) for  $l \geq 2, 5, 13$ , respectively.

#### 4. OPTIMAL METHODS FOR TEST EQUATION (1.1), (1.7)

**4.1. Order relations.** In order to derive optimal methods for test equation (1.1), (1.7), we will first give convenient order relations, similar to (2.1), for test equation (1.1), (1.11).

Let (1.2) be a method with  $\beta_k \geq 0$ , and let  $r > 0$ . Then (1.2) is of order at least  $p$  if and only if (in  $\mathbb{R}^{p+1}$ )

$$(4.1a) \quad \sum_{i=0}^{k-1} (\gamma_i b_i(-\beta_k) + \delta_i b_i(r^{-1})) = b_k(-\beta_k),$$

where

$$(4.1b) \quad \gamma_i = (-\alpha_i - r\beta_i)(1 + r\beta_k)^{-1} \quad (0 \leq i \leq k - 1),$$

$$(4.1c) \quad \delta_i = r(\beta_i - \alpha_i\beta_k)(1 + r\beta_k)^{-1} \quad (0 \leq i \leq k - 1).$$

The threshold factor  $R$  given by (1.9) is at least  $r$  if and only if

$$(4.2) \quad \gamma_i \geq 0, \quad \delta_i \geq 0 \quad (0 \leq i \leq k - 1).$$

**4.2. A sufficient condition for optimality.** In this subsection we give a lemma which is helpful in verifying whether a given  $k$ -step method (1.2) of order  $p$  is optimal for (1.1), (1.7). To this end, we define the subset of  $\mathbb{R}^p$

$$K^p(r, y) = \text{co}\{v \mid v = d_i(-y) \text{ or } v = d_i(r^{-1}) \text{ for some } i \text{ with } 0 \leq i \leq k - 1\}.$$

Clearly,  $K^p(r_1, y) \subset K^p(r_0, y)$  for all  $y \geq 0$  and  $r_0, r_1$  with  $0 < r_0 \leq r_1 < \infty$ .

**Lemma 4.1.** *Let  $k \geq 1$ ,  $p \geq 2$ , and assume  $y_0 \geq 0$  and  $r > 0$  are given such that:*

$$(4.3a) \quad \text{There exist index sets } I, J \subset \{0, 1, \dots, k - 1\}, J \text{ not contained in } I, I = \{i(1), i(2), \dots, i(m)\}, J = \{j(1), j(2), \dots, j(n)\}, m + n = p - 1, \text{ with (in } \mathbb{R}^p)$$

$$d_k(-y_0) \in \text{ri}(\text{co}(\{d_i(-y_0) \mid i \in I\} \cup \{d_j(r^{-1}) \mid j \in J\})).$$

(4.3b)  $S_1 = \{d_i(-y_0) | i \in I\} \cup \{d_j(r^{-1}) | j \in J \cup \{k\}\}$  is a set of  $p$  affinely independent vectors in  $\mathbb{R}^p$ .

(4.3c)  $S_2(y) = \{d_i(-y) | i \in I \cup \{k\}\} \cup \{d_j(r^{-1}) | j \in J\}$  is a set of  $p$  affinely independent vectors in  $\mathbb{R}^p$  for all  $y \geq 0$  with  $y \neq y_0$ .

(4.3d) The set  $\{d_i(-y_0) | 0 \leq i \leq k-1, i \notin I\} \cup \{d_j(r^{-1}) | 0 \leq j \leq k-1, j \notin J\}$  does not intersect with, and lies on one side of,  $\text{aff}(S_1)$ .

(4.3e) For each  $y$  with  $y \geq 0$ ,  $y \neq y_0$ , the following property holds: the set  $\{d_i(-y) | 0 \leq i \leq k-1, i \notin I\} \cup \{d_j(r^{-1}) | 0 \leq j \leq k-1, j \notin J\}$  does not intersect with, and lies on one side of,  $\text{aff}(S_2(y))$ .

Then there exists a unique  $k$ -step method of order at least  $p$  which is optimal for (1.1), (1.7). Its threshold factor equals  $R_{k,p} = r$ . The coefficients  $\alpha_i, \beta_i$  can be found by putting

$$\gamma_i = 0, \quad \delta_j = 0 \quad (0 \leq i, j \leq k-1; i \notin I, j \notin J), \quad \beta_k = y_0$$

in (4.1) and solving for the remaining coefficients.

*Proof.* From (4.3a), (4.1), (4.2), it can be seen that there exists a  $k$ -step method (1.2) of order at least  $p$  with threshold factor  $R$  given by (1.9) with  $R \geq r$ . Assume (1.2) is such a method. By virtue of (4.1), (4.2), we have  $d_k(-\beta_k) \in K^p(r, \beta_k)$ . However, conditions (4.3c, e) imply that  $d_k(-y) \notin K^p(r, y)$  for  $y \in (0, y_0)$  or  $y \in (y_0, \infty)$ . Therefore,  $\beta_k = y_0$ .

Next we show that  $r = R_{k,p}$ . It can be verified that (4.3a, b, d) imply that  $d_k(-\beta_k) \notin K^p(r_1, y_0)$  for  $r < r_1 < \infty$ . Hence  $r = R_{k,p}$ .

Finally, we obtain from (4.3a, b, d) that  $\{i | \gamma_i \neq 0\} \subset I$ ,  $\{j | \delta_j \neq 0\} \subset J$  for the coefficients in (4.1). Conditions (4.3a, b) guarantee that there exist unique coefficients  $\alpha_i, \beta_i$  satisfying the order relations (4.1) with  $\gamma_i = 0$  ( $i \notin I$ ),  $\delta_j = 0$  ( $j \notin J$ ). Thus, we have obtained the unique  $k$ -step method (1.2) of order at least  $p$  which is optimal for (1.1), (1.7).  $\square$

**4.3. Comparison of  $R_{k,p}$  and  $S_{k,p}$ .** As an immediate consequence of definitions (1.9), (1.10), (1.12), and (1.13), we have

**Theorem 4.2.** Let  $k \geq 1$ ,  $p \geq 1$ . Then:

(i) For any method (1.2) with threshold factors  $R$  and  $S$  given by (1.9) and (1.12), one has  $R \geq S$ . If the method is explicit, then  $R = S$ .

(ii)  $R_{k,p} \geq S_{k,p}$ .

At first sight, it might be expected that  $R_{k,p} > S_{k,p}$  for given  $k \geq 1$ ,  $p \geq 2$ . Theorem 4.3 below, however, rather shows the opposite.

**Theorem 4.3.** Let  $p = 2, 4$ , or  $6$ ,  $k \geq p^2/4$ . We have:

(i)  $R_{k,p} = R_{p^2/4,p} = S_{p^2/4,p}$ .

(ii) *There exists a unique optimal  $k$ -step method (1.2) of order  $p$  for (1.1), (1.7). It is the  $k$ -step method (3.1) adjoined to the optimal  $p^2/4$ -step method of order  $p$  for (1.1), (1.11).*

*Proof.* Let  $p \in \{2, 4, 6\}$ ,  $k \geq p^2/4$ , and let (1.2) be the optimal  $k$ -step method of order  $p$  obtained in §3. Let  $j(1) < j(2) < \dots < j(p-1) < k$  be such that  $\beta_{j(i)} \neq 0$  ( $1 \leq i \leq p-1$ ). We will verify condition (4.3) of Lemma 4.1 with  $I = \{i | \alpha_i + S_{k,p}\beta_i \neq 0\} = \emptyset$ ,  $J = \{j | \beta_j \neq 0\} = \{j(i) | 1 \leq i \leq p-1\}$ , and  $r = S_{k,p}$ ,  $y_0 = \beta_k$ .

Recall that  $d_k(-\beta_k) \in \text{ri}(\text{co}\{d_{j(i)}(S_{k,p}^{-1}) | 1 \leq i \leq p-1\})$  (cf. (1.12), (2.1)). Hence (4.3a) holds.

Condition (4.3b) follows from Lemmas 6.1, 6.2. Further, the determinant of the  $p \times p$  matrix made up of the vectors  $b_{j(i)}(S_{k,p}^{-1})$  and  $b_k(-y)$  is a linear function of  $y$  and vanishes only at  $y = \beta_k$ . So also (4.3c) is satisfied.

In order to prove (4.3d), we consider the function

$$f_s(y, z) = \det(b_s(-y), b_{j(1)}(S_{k,p}^{-1}), \dots, b_{j(p-1)}(S_{k,p}^{-1}), b_k(-z)).$$

We have, by Theorem 3.1, part (i), and Lemmas 6.2 and 6.1,

$$(4.4a) \quad f_s(-S_{k,p}^{-1}, -S_{k,p}^{-1}) > 0 \quad (\text{for } 0 \leq s \leq k-1, s \notin J),$$

$$(4.4b) \quad f_s(-S_{k,p}^{-1}, -S_{k,p}^{-1}) = 0 \quad (\text{for } s \in J).$$

We also see that  $\frac{\partial}{\partial y} f_s(0, -S_{k,p}^{-1}) > 0$ . For  $s = 0, 1, \dots, k - p^2/4$ , this follows from Lemmas 6.3 and 6.4, whereas for  $s = k - p^2/4 + 1, \dots, k - 1$ , this can be shown by Lemma 6.3 and brute force computation.

The function  $f$  is linear in  $y$ . Therefore, we also have

$$(4.4c) \quad f_s(\beta_k, -S_{k,p}^{-1}) > 0 \quad (0 \leq s \leq k-1).$$

The inequalities in (4.4a, c) imply (4.3d).

Finally, we prove (4.3e). Since the function  $f$  is linear in  $z$ , it follows from (4.3a), (4.4a, b), and the positivity of  $\frac{\partial}{\partial y} f_s(0, -S_{k,p}^{-1})$  that

$$(4.5a) \quad \text{sign}(f_s(y, y)) = \text{sign}(f_s(y, -S_{k,p}^{-1})(\beta_k - y)) = \text{sign}(\beta_k - y) \\ (\text{for all } y \geq 0, y \neq \beta_k, \text{ and } 0 \leq s \leq k-1).$$

Similarly, the linearity of  $f$  in  $z$ , (4.3a), Theorem 3.1, part (i), and Lemmas 6.2 and 6.1 imply

$$(4.5b) \quad \text{sign}(f_s(-S_{k,p}^{-1}, y)) = \text{sign}(f_s(-S_{k,p}^{-1}, -S_{k,p}^{-1})(\beta_k - y)) = \text{sign}(\beta_k - y) \\ (\text{for all } y \geq 0, y \neq \beta_k, \text{ and } 0 \leq s \leq k-1, s \notin J).$$

Condition (4.3e) is satisfied because of (4.5).

Hence, condition (4.3) of Lemma 4.1 is satisfied, and (1.2) is the unique  $k$ -step method of order  $p$  which is optimal for test equation (1.1), (1.7).  $\square$

TABLE 5  
Optimal threshold factors  $R_{k,3}$  for  $2 \leq k \leq 20$

| $k$ | $R_{k,3}$ | $k$ | $R_{k,3}$ |
|-----|-----------|-----|-----------|
| 2   | 1.225     | 11  | 1.904     |
| 3   | 1.572     | 12  | 1.913     |
| 4   | 1.703     | 13  | 1.920     |
| 5   | 1.772     | 14  | 1.926     |
| 6   | 1.815     | 15  | 1.931     |
| 7   | 1.844     | 16  | 1.935     |
| 8   | 1.865     | 17  | 1.939     |
| 9   | 1.881     | 18  | 1.943     |
| 10  | 1.894     | 19  | 1.946     |
|     |           | 20  | 1.949     |

*Remark.* The unique  $k$ -step method (1.2) of order 2 which is optimal for test equation (1.1), (1.7) is the trapezoidal rule. This was already proved in [5] by using the "Greek-Roman" transformation (cf. [9, p. 230]).

In general, the optimal threshold factors  $R_{k,p}$  and  $S_{k,p}$  need not be equal. Compare, e.g., Table 1 and Table 5. The optimal methods (1.2) for (1.1), (1.7) of order 3 were found by setting  $\gamma_i = 0$  ( $1 \leq i \leq k-1$ ),  $\delta_i = 0$  ( $0 \leq i \leq k-2$ ) in (4.1), solving for the remaining coefficients, and applying Lemma 4.1.

On the other hand, for  $p = 3, 5, 7$  and  $k \geq 1$ ,  $R_{k,p}$  cannot be substantially larger than  $S_{k,p}$ , as follows from Corollary 4.4.

**Corollary 4.4.** *Let  $p = 2, 4$ , or  $6$ . We have  $R_{k,p+1} = S_{p^2/4,p} + \mathcal{O}(k^{-1})$  (for  $k \rightarrow \infty$ ).*

*Proof.* From Theorem 4.2, the definition of the threshold factor  $R$ , and Theorem 4.3, one obtains the inequality

$$S_{k,p+1} \leq R_{k,p+1} \leq R_{k,p} = S_{p^2/4,p} \quad (\text{for } k \geq p^2/4).$$

The corollary follows from this inequality and Theorem 3.3, part (ii).  $\square$

## 5. A NUMERICAL EXAMPLE

For the amplification of errors in the numerical solution  $u_n$  ( $n = 0, 1, \dots$ ) in the application of (1.2) to (1.1) it can make an essential difference whether the stepsize restriction (1.8) is satisfied or not. We wish to illustrate this by a simple numerical example.

Consider the system of ordinary differential equations

$$(5.1a) \quad \frac{d}{dt}U(t) = AU(t) + b(t) \quad (t \geq 0),$$

$$(5.1b) \quad U(0) = u_0.$$

Here,  $A = (\alpha_{i,j})$  is an  $s \times s$  matrix with nonzero coefficients  $\alpha_{i,i} = -1$  ( $1 \leq i \leq s$ ),  $\alpha_{i,i-1} = 1$  ( $2 \leq i \leq s$ ), and  $u_0 \in \mathbb{R}^s$ .

TABLE 6  
Amplification factors  $\gamma_n$

| $n$         | 1                 | 8                 | 32                   | 128                   | 512               |
|-------------|-------------------|-------------------|----------------------|-----------------------|-------------------|
| method (i)  | $1.0 \times 10^1$ | $2.0 \times 10^1$ | $6.7 \times 10^1$    | $4.8 \times 10^4$     | $1.4 \times 10^9$ |
| method (ii) | $1.0 \times 10^0$ | $1.0 \times 10^0$ | $9.6 \times 10^{-1}$ | $2.7 \times 10^{-17}$ | 0                 |

We might use method (1.2) with exact starting values  $u_0, u_1, \dots, u_{k-1}$  to obtain approximations  $u_n$  ( $n = 0, 1, \dots$ ), or with slightly perturbed starting values  $v_0, v_1, \dots, v_{k-1}$  to obtain approximations  $v_n$  ( $n = 0, 1, \dots$ ). The difference  $w_n = v_n - u_n$  ( $n = 0, 1, \dots$ ) will satisfy (1.2) with  $f(t, x) = Ax$  and initial values  $w_0, w_1, \dots, w_{k-1}$ . What is of interest to us is the maximum norm of  $w_n$ . Therefore, we consider (5.1a) with  $b(t) \equiv 0$ . Further, for  $n \geq 1$ , let the amplification factor  $\gamma_n$  be the smallest real number such that

$$\|w_{n+k-1}\|_\infty \leq \gamma_n \max\{\|w_0\|_\infty, \|w_1\|_\infty, \dots, \|w_{k-1}\|_\infty\}$$

for all  $w_0, w_1, \dots$  satisfying (1.2).

Two methods (1.2) for solving (5.1) with  $s = 40$  are compared, viz.

(i) The backward difference method of order 6 (cf. [12]) with stepsize  $h = 0.9052$ . The coefficients in (1.2) for this 6-step method are

$$\begin{aligned} \alpha_0 &= 10/147, & \alpha_1 &= -72/147, & \alpha_2 &= 225/147, & \alpha_3 &= -400/147, \\ \alpha_4 &= 450/147, & \alpha_5 &= -360/147, & \beta_i &= 0 \ (0 \leq i \leq 5), & \beta_6 &= 60/147. \end{aligned}$$

The region of stability  $S$  of this method (cf., e.g., [10]) is the set of complex numbers defined by

$$\begin{aligned} S = \{ \lambda \mid \text{the roots } \zeta_i \text{ of } \sum_{i=0}^5 \alpha_i z^i + (1 - \lambda \beta_6) z^6 = 0 \\ \text{satisfy } |\zeta_i| \leq 1, \text{ and if } |\zeta_i| = 1, \text{ then } \zeta_i \text{ is a simple root} \}. \end{aligned}$$

The eigenvalues of  $hA$  are contained in  $\text{int}(S)$ . Hence  $\max_{n \geq 1} \gamma_n$  is finite. On the other hand, the threshold factor  $R$  in (1.9) equals  $R = 0$ .

(ii) The optimal 9-step method of order 6 (cf. Table 1, Corollary 3.2, Theorem 4.3). The matrix  $A$  satisfies (1.4) with  $\rho = 1$ . So (1.8) is fulfilled with stepsize  $h = 0.9052$ . Thus (1.6) holds, and  $\gamma_n \leq 1$  for all  $n \geq 1$ .

In Table 6,  $\gamma_n$  is given for various values of  $n$  for both methods. Clearly, from a practical point of view, it is desirable that the  $\gamma_n$  be of moderate size. However, this is not the case for method (i). The table thus reveals the superiority, in the present example, of our optimal method (ii) over the well-known backward difference method (i).

### 6. TECHNICAL LEMMAS

Here we give lemmas which were repeatedly made use of in the proofs in the previous sections.

**Lemma 6.1.** *Let real coefficients  $0 < a_1 < a_2 < \dots < a_s$  and  $\alpha_1 < \alpha_2 < \dots < \alpha_s$  be given. Then the determinant of the matrix  $M = (\mu_{i,j})$ , defined by  $\mu_{i,j} = a_j^{\alpha_i}$ ,  $1 \leq i, j \leq s$ , is positive.*

For the proof we refer to [6, p. 118], where the transposed matrix  $M^T$  is considered.

**Lemma 6.2.** *Let  $s \geq 1$ ,  $x \in \mathbb{R}$ , and let the  $s \times s$  matrix  $S = (\sigma_{i,j})$  be given by  $\sigma_{i,i} = 1$  ( $1 \leq i \leq s$ ),  $\sigma_{i,i-1} = (i-1)x$  ( $2 \leq i \leq s$ ),  $\sigma_{i,j} = 0$  (otherwise). Then:*

- (i)  $Sa_i = b_i(x)$  for all  $i \in \mathbb{Z}$ .
- (ii)  $\det(Sv_1, Sv_2, \dots, Sv_s) = \det(v_1, v_2, \dots, v_s)$  (for all  $v_1, v_2, \dots, v_s \in \mathbb{R}^s$ ).

*Proof.*  $S$  is a triangular matrix with unit diagonal elements. Hence  $\det(S) = 1$ .  $\square$

**Lemma 6.3.** *Let  $s \geq 1$ , and let the  $s \times s$  matrix  $T = (\tau_{i,j})$  be given by  $\tau_{i,j} = \binom{i-1}{j-1}$  ( $j \leq i$ ),  $\tau_{i,j} = 0$  (otherwise). Then:*

- (i)  $Ta_i = a_{i+1}$  and  $Tb_i(x) = b_{i+1}(x)$  for all  $i \in \mathbb{Z}$ ,  $x \in \mathbb{R}$ .
- (ii)  $\det(Tv_1, Tv_2, \dots, Tv_s) = \det(v_1, v_2, \dots, v_s)$  (for all  $v_1, v_2, \dots, v_s \in \mathbb{R}^s$ ).

*Proof.*  $T$  is a triangular matrix with unit diagonal elements. Hence  $\det(T) = 1$ .  $\square$

**Lemma 6.4.** *Let  $i \geq 0$  and  $x, y$  be real numbers with  $x > 0$ . Let  $0 \leq j(1) < j(2) < \dots < j(s)$  be a sequence of integers. Define*

$$f(y) = \det \left( b_{-i}(-y), b_{j(1)}(x), b_{j(2)}(x), \dots, b_{j(s)}(x) \right).$$

*Then  $f(y) > 0$  and  $\frac{d}{dy}f > 0$  for all  $y \geq 0$ .*

This lemma can be proved by using Lemma 6.2, expanding the determinant along the first column, and by applying Lemma 6.1.

**Lemma 6.5.** *Let  $i, j, s$  be positive integers with  $i < j$ . Let  $0 \leq j(1) < j(2) < \dots < j(s)$  be a sequence of integers, and let  $x \in (0, \infty)$ . Then*

$$\det(a_{-j}, a_{-i}, b_{j(1)}(x), b_{j(2)}(x), \dots, b_{j(s)}(x)) > 0.$$

This lemma is closely related to Lemma 3.3 in [13].

#### ACKNOWLEDGMENT

I wish to thank Professor M. N. Spijker, Dr. K. Dekker, and the referee for carefully reading preliminary versions of this paper and making a number of valuable suggestions concerning the presentation of the material.

#### BIBLIOGRAPHY

1. C. Bolley and M. Crouzeix, *Conservation de la positivité lors de la discrétisation des problèmes d'évolution paraboliques*, RAIRO Anal. Numér. **12** (1978), 237–245.
2. J. C. Butcher, *The equivalence of algebraic stability and AN-stability*, BIT **27** (1978), 510–533.
3. G. Dahlquist, *A special stability problem for linear multistep methods*, BIT **3** (1963), 27–43.

4. G. Dahlquist and R. Jeltsch, *Generalized disks of contractivity for explicit and implicit Runge-Kutta methods*, Report TRITA-NA-7906, NADA, Roy. Inst. Techn. Stockholm.
5. K. Dekker, private communication, 1985.
6. F. R. Gantmacher, *Applications of the theory of matrices*, Interscience, New York, 1959.
7. W. Gautschi, *On inverses of Vandermonde and confluent Vandermonde matrices*, Numer. Math. **4** (1962), 117–123; Numer. Math. **5** (1963), 425–430.
8. J. A. Van de Griend and J. F. B. M. Kraaijevanger, *Absolute monotonicity of rational functions occurring in the numerical solution of initial value problems*, Numer. Math. **49** (1986), 413–424.
9. P. Henrici, *Discrete variable methods in ordinary differential equations*, Wiley, New York and London, 1962.
10. R. Jeltsch and O. Nevanlinna, *Stability of explicit time discretizations for solving initial value problems*, Numer. Math. **37** (1981), 61–69.
11. J. F. B. M. Kraaijevanger, *Absolute monotonicity of polynomials occurring in the numerical solution of initial value problems*, Numer. Math. **48** (1986), 303–322.
12. J. P. Lambert, *Computational methods in ordinary differential equations*, Wiley, London, 1973.
13. H. W. J. Lenferink, *Contractivity preserving explicit linear multistep methods*, Numer. Math. **55** (1989), 213–223.
14. O. Nevanlinna, *Remarks on time discretization of contraction semigroups*, Report HTKK-MAT-A225, Institute of Math., Helsinki University of Technology, 1984.
15. O. Nevanlinna and W. Liniger, *Contractive methods for stiff differential equations*, BIT **18** (1978), 457–474; BIT **19** (1979), 53–72.
16. R. T. Rockafellar, *Convex analysis*, Princeton Univ. Press, Princeton, N. J., 1970.
17. J. Sand, *Circle contractive linear multistep methods*, BIT **26** (1986), 114–122.
18. —, *Choices in contractivity theory*, Appl. Numer. Math. **5** (1989), 105–115.
19. M. N. Spijker, *Contractivity in the numerical solution of initial value problems*, Numer. Math. **42** (1983), 271–290.
20. R. Vanselow, *Nonlinear stability behaviour of linear multistep methods*, BIT **23** (1983), 388–396.

DEPARTMENT OF MATHEMATICS AND COMPUTER SCIENCE, LEIDEN UNIVERSITY, P. O. BOX 9512,  
2300 RA LEIDEN, THE NETHERLANDS