

A Special Class of Explicit Linear Multistep Methods as Basic Methods for the Correction in the Dominant Space Technique

By Peter Alfeld

Abstract. A class of explicit linear multistep methods is suggested as basic methods for the CDS schemes introduced in [3]. These schemes are designed for the numerical solution of certain stiff ordinary differential equations, and operate with dominant eigenvalues, and the corresponding eigenvectors, of the Jacobian. The motivation, and the stability analysis for CDS schemes assumes that the eigensystem is constant. Here methods are introduced that perform particularly well if the eigensystem is not constant. In a certain sense the methods introduced here can be considered explicit approximations to the well-known implicit backward-differentiation formulas used by Gear [6] for the stiff option of his o.d.e. solver.

1. Introduction. In [1] and [3] the CDS technique is introduced for the numerical solution of separably stiff initial-value problems of ordinary differential equations. For these systems the eigenvalues of the Jacobian can be separated into two sets, one of which dominates the other.

The CDS technique consists of taking a step from x_{n+k-1} to x_{n+k} by a conventional explicit k -step method, the basic method, and then applying a correction in the dominant space (hence CDS), i.e., the space spanned by the right eigenvectors corresponding to the dominant eigenvalues.

We will use notations and definitions as they are given in [8]. For a more detailed description and motivation of the CDS technique the reader is referred to [1] and [3].

In Section 2 the basic definitions are given. It is convenient to change slightly the notation employed in [1] and [3].

In Section 3 the effects of nonlinearity, of errors in the approximation of the dominant eigensystem, and of the change of the dominant eigensystem with the independent variable x (interprojection) are investigated.

In Section 4 explicit linear multistep methods are introduced that subdue the interprojection effects, if they are employed as basic methods in the CDS scheme.

In Section 5 numerical examples are given.

2. Correction in the Dominant Space. We consider the (in general nonlinear) initial-value problem (IVP)

Received July 24, 1978.

AMS (MOS) subject classifications (1970). Primary 65L05.

Key words and phrases. Ordinary differential equations, numerical analysis, correction in the dominant space, separably stiff systems, interprojection, backward-differentiation formulas.

$$(2.1) \quad y' = f(x, y); \quad y(a) = \eta; \quad y, f \in R^m; \quad x \in [a, b].$$

We will also consider the special case that the IVP is linear, i.e.,

$$(2.2) \quad y' = A(x)y + g(x); \quad y(a) = \eta; \quad y, f, g \in R^m; \quad x \in [a, b],$$

where A is an $m \times m$ matrix function.

In all that follows we will assume, without specifically mentioning it, that f possesses as many continuous partial derivatives as required, and also that there exists a unique solution $y = y(x)$ of (2.1).

Definition 1. Let A be an $m \times m$ matrix, and let $\lambda^{(i)}, i = 1, 2, \dots, m$, be the eigenvalues of A . A is said to be *separably stiff* if

(a) A is nondefective, and

(b) there exists a constant integer $s, 1 \leq s < m$, such that $\lambda^{(i)}, i = 1, 2, \dots, s$, are real, distinct, and negative, and

$$\min_{1 \leq i \leq s} |\lambda^{(i)}| \gg \max_{s+1 \leq i \leq m} |\lambda^{(i)}|.$$

Remark. A matrix A is nondefective if the set of right eigenvectors spans R^m . This does not imply that A is nonsingular.

Definition 2. The initial-value problem (2.1) is *separably stiff* if the matrix

$$A(x, y) = \frac{\partial}{\partial y} f(x, y)$$

is separably stiff for all $(x, y(x))$, where $x \in [a, b]$ and $y(x)$ is the exact solution of (2.1).

Remarks. 1. We are considering problems whose stiffness arises from a set (which we assume to be small) of troublesome eigenvalues $\lambda^{(i)}, i = 1, 2, \dots, s$, which are real, negative, and well-separated from the rest; note, however, that we do not require all the eigenvalues to have negative real parts.

2. We assume that $A(x, y)$ is a continuous function. Thus, $A(x, y)$ will be separably stiff for y close to $y(x)$.

Consider the Correction in the Dominant Space (CDS) scheme

$$(2.3(i)) \quad \tilde{y}_{n+k} = By_n,$$

$$(2.3(ii)) \quad y_{n+k} = \tilde{y}_{n+k} + \sum_{i=1}^s \xi_{n+k}^{(i)} c_{n+k}^{(i)}.$$

Remark. Our notation differs from that in [3]; there k is assumed to equal 1. Here we use k for convenience; it will be easier later to refer to approximations evaluated at $x_{n+j}, j = 0, 1, \dots, k$.

(2.3) applied to a separably stiff IVP (2.1) defines a sequence $y_n \approx y(x_n), n = k, k+1, k+2, \dots$, where $x_n = a + nh, h$ being a constant steplength. We assume that starting values y_0, y_1, \dots, y_{k-1} are given.

B is a symbolic notation for the application of a conventional explicit k -step method. (We will only consider the case that B is an explicit linear k -step method.)

$c_{n+k}^{(i)}$ is the right eigenvector corresponding to $\lambda_{n+k}^{(i)}$, the i th eigenvalue of $(\partial/\partial y)f(x_{n+k}, \tilde{y}_{n+k})$.

The $\xi_{n+k}^{(i)}$ are scalar correction factors that can be determined in a variety of ways, see [1], [3].

We denote the left eigenvectors corresponding to $\lambda_{n+k}^{(i)}$ by $d_{n+k}^{(i)}$, and assume the following normalization:

$$(2.4) \quad \left\{ \begin{array}{l} \text{(i)} \quad \langle c_n^{(i)}, d_n^{(i)} \rangle := (c_n^{(i)})^T d_n^{(i)} = 1, \quad i = 1, 2, \dots, m, \\ \text{(ii)} \quad \|c_n^{(i)}\| := \langle c_n^{(i)}, c_n^{(i)} \rangle^{1/2} = 1, \quad i = 1, 2, \dots, m, \\ \text{(iii)} \quad \text{The first nonvanishing component of } c_0^{(i)} \text{ is positive,} \\ \text{(iv)} \quad \text{If } {}^t c_n^{(i)} \text{ denotes the } t\text{th component of } c_n^{(i)}, \quad t = 1, 2, \dots, m, \\ \quad \text{and } \hat{t}^{(n)} c_n^{(i)} \text{ is the component of } c_n^{(i)} \text{ with maximum modulus} \\ \quad \text{(the first such component if the maximum is not unique), then} \\ \text{sgn} \{ \hat{t}^{(n)} c_{n+1}^{(i)} \} = \text{sgn} \{ \hat{t}^{(n)} c_n^{(i)} \}, \quad i = 1, 2, \dots, m, \quad n = 0, 1, 2, \dots \end{array} \right.$$

The set of right and left eigenvectors of $A(x_n, \tilde{y}_n)$ is then uniquely defined. Note that

$$(2.5) \quad \langle c_n^{(i)}, d_n^{(i)} \rangle = \delta_{ij},$$

the Kronecker delta.

Definition 3. Let $\lambda_n^{(i)} = \lambda^{(i)}(x_n, \tilde{y}_n)$, and let the right and left eigenvectors $c_n^{(i)}$ and $d_n^{(i)}$ of $A(x_n, \tilde{y}_n)$ be normalized by (2.4). The *dominant* and *subdominant eigensystems* of $A(x_n, \tilde{y}_n)$ are defined to be $\{\lambda_n^{(i)}, c_n^{(i)}, d_n^{(i)} \mid i = 1, 2, \dots, s\}$ and $\{\lambda_n^{(i)}, c_n^{(i)}, d_n^{(i)} \mid i = s + 1, s + 2, \dots, m\}$, respectively. The subspaces spanned by $\{c_n^{(i)} \mid i = 1, 2, \dots, s\}$ and $\{c_n^{(i)} \mid i = s + 1, s + 2, \dots, m\}$ are defined to be the *dominant* and *subdominant spaces* at x_n , respectively.

Remarks. 1. The CDS technique requires the explicit computation of the dominant eigensystem; see [3].

2. Our notation differs from that in [3]; there the eigensystem is denoted by $\tilde{\lambda}_n^{(i)}, \tilde{c}_n^{(i)}, \tilde{d}_n^{(i)}$.

3. It would be desirable to be able to use the dominant eigensystem evaluated at (x_n, y_n) , but this would introduce complicated implicitness.

Since we assume $A(x_n, \tilde{y}_n)$ to be nondefective, we can express any vector v as

$$v = \sum_{i=1}^m \gamma_n^{(i)} c_n^{(i)},$$

where, because of (2.5), $\gamma_n^{(i)} = \langle d_n^{(i)}, v \rangle$.

This motivates

Definition 4. Let $v \in R^m$. Then $\gamma^{(i)} = \langle d^{(i)}(x, y(x)), v \rangle$ are said to be the *components of v (at x)*. $\gamma^{(i)}, i = 1, 2, \dots, s$, are the *dominant components*, and $\gamma^{(i)}, i = s + 1, s + 2, \dots, m$, are the *subdominant components* of v (at x).

Remarks. 1. When there is no confusion, we will omit the specification "at x ".

2. In a similar manner we will talk about subdominant global and local truncation errors, etc.

It is convenient to introduce a special notation for the components of the exact solution of (2.1) and its derivative. Thus, we write

$$y(x) = \sum_{i=1}^m \phi^{(i)}(x) c^{(i)}(x, y(x)), \quad y'(x) = \sum_{i=1}^m \psi^{(i)}(x) c^{(i)}(x, y(x)),$$

where

$$\phi^{(i)}(x) = \langle d^{(i)}(x, y(x)), y(x) \rangle, \quad \psi^{(i)}(x) = \langle d^{(i)}(x, y(x)), y'(x) \rangle.$$

We finish this section by stating a consequence of Theorem 1 in [3].

THEOREM 1. *Assume the basic method in the CDS scheme (2.3) is an explicit linear k -step method with steplength h , and that the correction factors are computed by any of the methods described in [1] and [3]. Then all numerical solutions $y_n, n = 0, 1, 2, \dots$, of the separably stiff test equation*

$$y' = Ay$$

by the CDS scheme (2.3) tend to zero as n tends to infinity, provided that $h\lambda^{(i)} \in \mathcal{R}_B$, $i = s + 1, s + 2, \dots, m$, where \mathcal{R}_B is the region of absolute stability of B . \square

Remark. Thus stability requirements are imposed only by the subdominant eigenvalues. These are much less severe than those that would be imposed by the dominant eigenvalues if the explicit basic method were to be employed on its own. However, the above theorem suggests that one factor in the choice of the basic method is the size of its region of absolute stability.

3. Sensitivity and Interprojection. Let us assume that we have obtained sequences

$$\{\tilde{y}_n \mid n = 0, 1, 2, \dots\}, \quad \{y_n \mid n = 0, 1, 2, \dots\} \quad \text{and} \quad \{f_n \mid n = 0, 1, 2, \dots\},$$

where $f_n = f(x_n, y_n)$, by applying (2.3) to a separably stiff IVP (2.1).

Define for $n = 0, 1, 2, \dots, i = 1, 2, \dots, m$,

$$c_n^{*(i)} := c^{(i)}(x_n, y_n), \quad d_n^{*(i)} := d^{(i)}(x_n, y_n).$$

This is the system of eigenvectors of $A(x_n, y_n)$. Further

$$(3.1) \quad \Delta c_n^{(t)} := c_n^{*(t)} - c_n^{(t)}, \quad t = 1, 2, \dots, s.$$

(Note that Δ does denote differences but is not identical with the forward difference operator. In the sequel it will be defined for each individual case.)

The vectors given in (3.1) are the errors in the approximations to the dominant right eigenvectors at (x_n, y_n) due to the fact that we evaluate $c_n^{(t)}$ at (x_n, \tilde{y}_n) , and to round-off and truncation errors occurring in the numerical method used for computing the dominant eigensystem.

We write

$$(3.2) \quad y(x_n) = \sum_{i=1}^m \phi_n^{(i)} c_n^{*(i)}$$

and

$$(3.3) \quad y'(x_n) = \sum_{i=1}^m \psi_n^{(i)} c_n^{*(i)}$$

and define for $i = 1, 2, \dots, m$,

$$(3.4) \quad \Delta\phi_n^{(i)} := \phi_n^{(i)} - \langle d_n^{*(i)}, y_n \rangle,$$

$$(3.5) \quad \Delta\psi_n^{(i)} := \psi_n^{(i)} - \langle d_n^{*(i)}, f_n \rangle.$$

These are approximate components of the global errors of the approximations to $y(x_n)$ and $y'(x_n)$.

Let us now assume that n is fixed.

We define for $t = 1, 2, \dots, s$,

$$(3.6) \quad \Delta\xi_{n+k}^{(t)} := \langle d_{n+k}^{*(t)}, y(x_{n+k}) - \tilde{y}_{n+k} \rangle - \xi_{n+k}^{(t)}.$$

This is the difference between the ‘‘ideal’’ correction factors, which would render the dominant components of the global error zero provided $\Delta c_n^{(r)} = 0$, and the actual correction factors.

We also define

$$\Delta B y_n := y(x_{n+k}) - \tilde{y}_{n+k},$$

which is the global error in the intermediate approximation \tilde{y}_{n+k} .

From the biorthonormality of the eigensystem and (3.1), (3.2), (3.4), (3.6), it follows that the dominant components of the global error in y_{n+k} are given by

$$(3.7) \quad \Delta\phi_{n+k}^{(r)} = \Delta\xi_{n+k}^{(r)} + \sum_{t=1}^s \xi_{n+k}^{(t)} \langle d_{n+k}^{*(r)}, \Delta c_{n+k}^{(t)} \rangle, \quad r = 1, 2, \dots, s.$$

Similarly, we obtain for the subdominant components of the global error

$$(3.8) \quad \Delta\phi_{n+k}^{(r)} = \langle d_{n+k}^{*(r)}, \Delta B y_n \rangle + \sum_{t=1}^s \xi_{n+k}^{(t)} \langle d_{n+k}^{*(r)}, \Delta c_{n+k}^{(t)} \rangle,$$

$$r = s + 1, s + 2, \dots, m.$$

As an example, let us consider the case that the basic method is the explicit linear multistep method defined by

$$(3.9) \quad y_{n+k} = - \sum_{j=0}^{k-1} \alpha_j y_{n+j} + h \sum_{j=0}^{k-1} \beta_j f_{n+j}.$$

Recall that the local truncation error of (3.9) is defined by

$$(3.10) \quad l_{n+k} := y(x_{n+k}) + \sum_{j=0}^{k-1} \alpha_j y(x_{n+j}) - h \sum_{j=0}^{k-1} \beta_j y'(x_{n+j}),$$

i.e., the error that would occur in y_{n+k} , if the back values y_{n+j}, f_{n+j} ($j = 0, 1, \dots, k - 1$), were exact.

Then, from (3.2), (3.3), (3.4), (3.5), (3.9), (3.10), we obtain

$$(3.11) \quad \Delta B y_n = 1_{n+k} - \sum_{j=0}^{k-1} \sum_{i=1}^m (\alpha_j \Delta \phi_{n+j}^{(i)} - h \beta_j \Delta \psi_{n+j}^{(i)}) \tilde{c}_{n+j}^{*(i)}.$$

Defining

$$(3.12) \quad \Delta_j^* c_{n+k}^{(i)} := \tilde{c}_{n+k}^{*(i)} - \tilde{c}_{n+j}^{*(i)}$$

and substituting this into (3.11) yields

$$\begin{aligned} \langle \tilde{d}_{n+k}^{*(r)}, \Delta B y_n \rangle &= \langle \tilde{d}_{n+k}^{*(r)}, 1_{n+k} \rangle - \sum_{j=0}^{k-1} [\alpha_j \Delta \phi_{n+j}^{(r)} - h \beta_j \Delta \psi_{n+j}^{(r)}] \\ &+ \sum_{i=1}^m \sum_{j=0}^{k-1} [\alpha_j \Delta \phi_{n+j}^{(i)} - h \beta_j \Delta \psi_{n+j}^{(i)}] \langle \tilde{d}_{n+k}^{*(r)}, \tilde{\Delta}_j^* c_{n+k}^{(i)} \rangle. \end{aligned}$$

Substituting this into (3.12) we obtain for $r = s + 1, s + 2, \dots, m$,

$$(3.13) \quad \begin{aligned} \Delta \phi_{n+k}^{(r)} &= \langle \tilde{d}_{n+k}^{*(r)}, 1_{n+k} \rangle - \sum_{j=0}^{k-1} [\alpha_j \Delta \phi_{n+j}^{(r)} - h \beta_j \Delta \psi_{n+j}^{(r)}] \\ &+ \sum_{i=1}^m \sum_{j=0}^{k-1} [\alpha_j \Delta \phi_{n+j}^{(i)} - h \beta_j \Delta \psi_{n+j}^{(i)}] \langle \tilde{d}_{n+k}^{*(r)}, \tilde{\Delta}_j^* c_{n+k}^{(i)} \rangle \\ &+ \sum_{t=1}^s \xi_{n+k}^{(t)} \langle \tilde{d}_{n+k}^{*(r)}, \Delta c_{n+k}^{(t)} \rangle, \end{aligned}$$

which is the global subdominant error if the basic method is given by (3.9).

The expressions (3.7) (dominant components of the global error), and (3.8) and (3.13) (subdominant components of the global error) illustrate the influences of errors arising from various sources.

Let us discuss the individual terms:

$$(1) \quad \sum_{t=1}^s \xi_{n+k}^{(t)} \langle \tilde{d}_{n+k}^{*(r)}, \Delta c_{n+k}^{(t)} \rangle, \quad r = 1, 2, \dots, m.$$

These terms, contributing both to the dominant and to the subdominant components of the global error, are due to the impossibility of computing the dominant eigensystems exactly. Let us refer to the dependence of the global error on these terms as *sensitivity* of the CDS scheme.

In the linear case (2.2), where the dominant eigensystem depends on x only, the $\Delta c_{n+k}^{(t)}$ do not vanish because of truncation and round-off errors arising in the numerical method employed (e.g., power method). However, we can iterate the power method until we reach an accuracy limited only by the accuracy of the particular computer that is used. So $\|\Delta c_{n+k}^{(t)}\|$ will equal approximately $2^{-\tau}$, where τ is the number of bits used for the mantissa.

In the nonlinear case (2.1) we face the additional problem that the dominant eigensystem is evaluated at $(x_{n+k}, \tilde{y}_{n+k})$ instead of (x_{n+k}, y_{n+k}) .

Consider now the other factors under the sum (1). The norm of the left eigenvectors will be large if the eigenproblem associated with $A(x_{n+k}, y_{n+k})$ is ill-conditioned. Note that this is not at all related to the condition number of $A(x_{n+k}, y_{n+k})$, which, for a separably stiff system, is always large (it does not exist if $A(x_{n+k}, y_{n+k})$ is singular). In spite of this $\|d_{n+k}^{*(r)}\|$ can be small, for instance

$$\|d_{n+k}^{*(r)}\| = \|c_{n+k}^{*(r)}\| = 1,$$

if $A(x_{n+k}, y_{n+k})$ is symmetric (and thus left and right eigenvectors coincide).

For the dominant components of the error $\|d_{n+k}^{*(r)}\|$ is readily available and should be a good approximation to $\|d_{n+k}^{*(r)}\|$. For the subdominant components this is not so, however.

Since the error contributions in (1) are proportional to the correction factors $\xi_{n+k}^{(r)}$, these should be small (in modulus). It turns out that the correction factors very critically depend on the choice of the basic method, and can be excessively large; see [3].

In summary, the sensitivity of the CDS scheme depends on the conditioning of the eigenproblem, the size of the correction factors, and the errors in the right-hand eigensystem. Since the latter are very small in the linear case (and can be controlled in the nonlinear case by forming the difference $c^{(i)}(x_{n+k}, y_{n+k}) - c_{n+k}^{(i)}$ after computing y_{n+k} using (2.3)), sensitivity usually will present no serious problems, provided the eigenproblem is reasonably well-conditioned, and the basic method is sensibly chosen.

$$(2) \quad \Delta \xi_{n+k}^{(r)}, \quad r = 1, 2, \dots, s.$$

This term contributes to the dominant global error (3.7) only. Whereas the error considered under point (1) is due to not correcting exactly in the right direction, this term is due to not using exactly the right correction factor. It depends on the choice of the correction factors.

$$(3) \quad \langle d_{n+k}^{*(r)}, \Delta B y_n \rangle, \quad r = s + 1, s + 2, \dots, m.$$

These terms, contributing to the subdominant error (3.8) only, represent the $c_{n+k}^{*(r)}$ -components of the local truncation error of the basic method and of the errors due to the basic method not operating on exact back values. It plays the same role that $\Delta \xi_{n+k}^{(r)}$ does for the dominant error. Here we consider the special case that the basic method is the linear multistep method (3.9); see paragraphs (4), (5), (6) below.

$$(4) \quad \langle d_{n+k}^{*(r)}, 1_{n+k} \rangle, \quad r = s + 1, s + 2, \dots, m.$$

These terms, contributing to the subdominant error (3.13), are the $c_{n+k}^{*(r)}$ -components of the local error of the basic method. They are small if the exact solution is smooth (which will usually be the case in the steady-state region).

$$(5) \quad \sum_{j=0}^{k-1} [\alpha_j \Delta \phi_{n+j}^{(r)} - h \beta_j \Delta \psi_{n+j}^{(r)}], \quad r = s + 1, s + 2, \dots, m.$$

This term, contributing to the subdominant error (3.13), represents the error arising from the errors in the previous values that we would expect if the eigensystem were constant and known exactly, and the basic method thus operating on the components of the numerical solution individually. It does not depend on the correction (2.3(ii)) at all and is peculiar to the basic method.

$$(6) \quad \sum_{i=1}^m \sum_{j=0}^{k-1} [(\alpha_j \Delta \phi_{n+j}^{(i)} - h \beta_j \Delta \psi_{n+j}^{(i)}) \langle \tilde{d}_{n+k}^{*(r)}, \tilde{\Delta}_j c_{n+k}^{(i)} \rangle], \quad r = s + 1, s + 2, \dots, m.$$

This term contributes to the subdominant error (3.13). It arises from the fact that the eigensystem changes with x , and is peculiar to CDS schemes. It can be described as representing the partial projection of previous errors, including dominant ones, into the current subdominant space. Although it only appears in the expression (3.13) for the subdominant error, a similar process influences the dominant error. (This is masked by the concept of an “ideal” correction factor $\langle \tilde{d}_{n+k}^{*(r)}, y(x_{n+k}) - \tilde{y}_{n+k} \rangle$, cf. (3.6), that deals with all dominant errors in \tilde{y}_{n+k} , no matter what source they come from.) We refer to the feature represented by the above term as *interprojection*.

Interprojection provides for communication between dominant and subdominant components of the numerical solution, by means of the basic method. Note that it occurs for linear as well as for nonlinear problems.

Because of the ill-conditioning of $f(x, y)$, dominant errors in y_{n+j} are amplified when computing f_{n+j} ($= f(x_{n+j}, y_{n+j})$) and then, because of interprojection, partly projected into the current subdominant space. This is undesirable. One attempt to decrease the interprojection effect is to design special basic methods that will be described in the next section. Another one is to devise CDS schemes for which the dominant errors in y_{n+j} are so small that the dominant errors in f_{n+j} are still tolerable, and their projection into the current subdominant space does not contribute unacceptably to the subdominant error. (Here the first version of reduction to scalar and the gradient prediction CDS schemes are notable; see [1], [3].)

Still another, more obvious, attempt is to reduce the steplength h (and thus the $\tilde{\Delta}_j c_{n+k}^{(i)}$). But, after all, the whole point of the CDS technique is to get rid of step-length restrictions due to considerations other than accuracy requirements. So, for an implementation of the CDS technique, this strategy should be considered a last resort.

However, it is a typical feature of CDS schemes that, because of interprojection, the dominant errors are considerably smaller than the subdominant ones.

4. Minimal-Projecting Methods. In [3] explicit linear multistep methods are recommended as the best choice for the basic method in (2.3). The reason given there is that a repeated evaluation of f , such as it occurs in predictor-corrector and Runge-Kutta methods, leads to large round-off errors, because of the ill-conditioning of f . These round-off errors may even render the numerical solution meaningless.

Using a basic method that evaluates f repeatedly leads to large correction factors. The above is a practical argument for keeping the correction factors small.

There are also two more theoretical arguments in favor of small correction factors.

Firstly, we have seen (cf. (3.7) and (3.8)) that the contribution due to the sensitivity of the CDS scheme, both to the dominant and to the subdominant components of the global error are proportional to $\xi_{n+k}^{(t)}$, $t = 1, 2, \dots, s$.

Secondly, in the nonlinear case we take $c^{(t)}(x_{n+k}, \tilde{y}_{n+k})$ ($t = 1, 2, \dots, s$) as an approximation to $c^{(t)}(x_{n+k}, y_{n+k})$. This is only feasible if \tilde{y}_{n+k} is close to y_{n+k} , which is another way of saying that the correction factors are small.

In this section special explicit linear multistep methods are suggested that subdue the interprojection effect described in Section 3, and that have the additional advantage that their region of absolute stability is larger than that of an Adams-Bashforth method of the same order, for $k = 1, 2, \dots, 6$. This latter feature may be advantageous in view of Theorem 1.

Let us look again at the subdominant components of the global error given by (3.13).

The first term represents the local error, and the second the conventional error propagation in the linear multistep method. These can be expected to be small. The fourth term is due to the sensitivity of the CDS scheme and is also ignored here.

The third term represents the influence of interprojection, and may be large if $\Delta_j^* c_{n+k}^{(i)} = \bar{c}_{n+k}^{(i)} - \bar{c}_{n+j}^{(i)} \neq 0$ (cf. (3.12)), i.e., if the eigensystem is not constant.

Recall that $|\Delta\psi_{n+j}^{(i)}|$ equals approximately (exactly in the linear case) $|\lambda_{n+j}^{(i)} \Delta\phi_{n+j}^{(i)}|$, and is thus much larger than $|\Delta\phi_{n+j}^{(i)}|$ ($i = 1, 2, \dots, s$). Hence, if $|\langle \bar{d}_{n+k}^{(r)}, \Delta_j^* c_{n+k}^{(i)} \rangle|$ is sufficiently large (i.e., the eigensystem varies rapidly), then the major contribution to $\Delta\phi_n^{(r)}$ in the interprojection terms will be given by

$$(4.1) \quad \rho := h \sum_{i=1}^s \sum_{j=0}^{k-1} \beta_j \Delta\psi_{n+j}^{(i)} \langle \bar{d}_{n+k}^{(r)}, \Delta_j^* c_{n+k}^{(i)} \rangle, \quad r = s + 1, s + 2, \dots, m.$$

Assume we can write $\Delta\psi_{n+j}^{(i)}$ as a smooth function $\Delta\psi^{(i)}$ of x

$$\Delta\psi_{n+j}^{(i)} = \Delta\psi^{(i)}(x_{n+j}).$$

This is true, e.g., in the linear case (2.2) for the gradient projection scheme, described in [3], where $\Delta\psi^{(i)}(x_{n+j}) = \psi^{(i)}(x_{n+j})$. The following motivation also holds if $\Delta\psi_{n+j}^{(i)}$ can be expressed as $p(h)q(x_{n+j})$, where p is some function and q is smooth.

The equation (4.1) becomes

$$(4.2) \quad \rho(x_{n+k}) = h \sum_{i=1}^s \sum_{j=0}^{k-1} \beta_j \rho^{(i)}(x_{n+j}),$$

where

$$(4.3) \quad \rho^{(i)}(x_{n+j}) = \Delta\psi^{(i)}(x_{n+j}) \langle \bar{d}^{(r)}(x_{n+k}), c^{(i)}(x_{n+k}) - c^{(i)}(x_{n+j}) \rangle.$$

Here the eigenvectors are assumed to depend on x only (in the nonlinear case we can consider $c^{(i)}(x) = c^{(i)}(x, y(x))$).

Note that if we consider x_{n+k} fixed, $\rho^{(i)}(x_{n+j})$ defined by (4.3), and thus $\rho(x_{n+k})$ defined by (4.2), is a function of h only.

Assume $\rho^{(i)}(x)$ can be expanded about x_{n+k} . We obtain from (4.2) for arbitrarily large q

$$\begin{aligned}\rho(h) &= h \sum_{i=1}^s \sum_{j=0}^{k-1} \beta_j \left[\sum_{t=0}^q \frac{h^t}{t!} (j-k)^t \frac{d^t}{dx^t} \rho^{(i)}(x_{n+k}) + O(h^{q+1}) \right] \\ &= h \sum_{t=1}^q \sum_{i=1}^s \frac{h^t}{t!} \frac{d^t}{dx^t} \rho^{(i)}(x_{n+k}) \sum_{j=0}^{k-1} \beta_j (j-k)^t + O(h^{q+2}).\end{aligned}$$

The summation in t starts at $t = 1$ because $\rho^{(i)}(x_{n+k}) = 0$. This suggests requiring that

$$\sum_{j=0}^{k-1} \beta_j (j-k)^t = 0 \quad \text{for } t = 1, 2, \dots, q.$$

Note that this puts a restriction on the basic linear multistep method that is independent of the particular problem to be tackled.

Definition 5. The explicit linear multistep method

$$(4.4) \quad \sum_{j=0}^k \alpha_j y_{n+j} = h \sum_{j=0}^{k-1} \beta_j f_{n+j}, \quad \alpha_k = 1,$$

is said to be *nonprojecting of degree q* , if

$$(4.5) \quad \sum_{j=0}^{k-1} \beta_j (j-k)^t = 0 \quad \text{for } t = 1, 2, \dots, q.$$

The term $\sum_{j=0}^{k-1} \beta_j (j-k)^{q+1} =: \gamma \neq 0$ is called the *projection constant* of (4.4).

Remark. In the above motivation we formally proceed in a manner similar to that employed in the definition of the order and error constant of a linear multistep method. The local truncation errors (respectively the interprojection terms (4.3)) are expanded about x_n (respectively x_{n+k}), and the first p (respectively q) terms of the resulting Taylor series are required to vanish. The coefficient of the first nonzero term is called the error constant (respectively the projection constant) of the method. Of two linear multistep methods which are nonprojecting of the same degree we would normally choose the one with the smaller projection constant, provided that other features (accuracy, stability) are about comparable for both methods.

It follows at once from the definition that a consistent LMM which is nonprojecting of order $q \geq 1$ has stepnumber $k \geq 2$.

Consider now the LMM (4.4), and assume it is nonprojecting of degree q , and its order is p . Since the basic method is responsible for the accuracy of the subdominant components, where $h\lambda^{(i)}$ is small, we are interested in convergent methods (4.4).

The order of a convergent LMM cannot exceed $k + 2$, and even this value is attained only if the method is implicit and the stepnumber is even. The order of an explicit convergent linear k -step method cannot exceed k ; see [7], [8].

We have k parameters β_j ; since for consistency we have to require

$$\sum_{j=0}^{k-1} |\beta_j| \neq 0,$$

we expect to be able to satisfy $k - 1$ conditions of the form (4.5).

In order to achieve order $p = k$ we have to satisfy

$$C_0 = C_1 = \dots = C_p = 0, \quad C_{p+1} \neq 0,$$

where

$$(4.6) \quad \begin{cases} C_0 = \alpha_0 + \alpha_1 + \dots + \alpha_k, \\ C_q = \frac{1}{q!} \sum_{j=0}^h j^q \alpha_j - \frac{1}{(q-1)!} \sum_{j=0}^{k-1} j^{q-1} \beta_j \end{cases}$$

and $0^0 := 1$ (0^0 occurs in C_1); see [8, p. 23].

Taking into account the number of free parameters in the general linear multi-step method (4.4), we expect that for each $k \geq 2$ there exists a unique k th order linear k -step method, which is nonprojecting of degree $k - 1$.

Definition 6. The linear k -step method (4.4) is said to be *minimal-projecting*, if it is of order $\geq k$ and nonprojecting of degree $k - 1$.

The following theorem provides an explicit expression for the coefficients of minimal-projecting LMMs.

THEOREM 2. For $k \geq 2$ the minimal-projecting k -step method is defined by

$$(4.7) \quad \beta_j = (-1)^j \binom{k}{j}, \quad j = 0, 1, \dots, k - 1,$$

$$(4.8) \quad \alpha_j = -\beta_j / (k - j), \quad j = 0, 1, \dots, k - 1,$$

$$(4.9) \quad \alpha_k = -\sum_{j=0}^{k-1} \alpha_j \neq 0,$$

where

$$\binom{k}{j} = \frac{k!}{j!(k-j)!}, \quad \text{the binomial coefficient.}$$

Remark. The coefficients given in Theorem 2 have to be normalized so as to satisfy $\alpha_k = 1$; see (4.4).

Proof. An elementary, but very technical and tedious proof of this theorem, can be found in [2]. That proof does not add insight into the subject of this investigation and, therefore, is here omitted. An alternative proof can be organized along the following lines:

Supplementing (4.5) by the normalizing equation

$$\sum_{j=0}^{k-1} \beta_j = (-1)^{k-1}$$

yields the Vandermonde system

$$\begin{pmatrix} 1 & 1 & \dots & 1 \\ 1 & 2 & & k \\ 1 & 2^2 & & k^2 \\ \vdots & \vdots & & \vdots \\ 1 & 2^{k-1} & \dots & k^{k-1} \end{pmatrix} \begin{pmatrix} \beta_k \\ \beta_{k-1} \\ \beta_{k-2} \\ \vdots \\ \beta_0 \end{pmatrix} = \begin{pmatrix} (-1)^{k-1} \\ 0 \\ \vdots \\ 0 \end{pmatrix}.$$

This establishes the existence (and uniqueness but for normalization) and also the formula (4.7) from the well-known determinant expressions.

The existence and uniqueness of the α_j (and how to obtain them) for given β_j can be deduced from the proof of Theorem 5.7 in [7]; see also the Theorem in [1]. \square

Table 1 lists the coefficients of the minimal-projecting linear k -step methods, obtained from Theorem 2, for $k = 2, 3, \dots, 7$. For easier representation the coefficients are given as integers (which are chosen so as to have no common factor). The coefficients α_k , by which all coefficients have to be divided in order to obtain the normalized form (4.4) are printed in boxes.

TABLE 1
Minimal-projecting k -step methods for $k = 2, 3, \dots, 7$

		Projection Constant	j:	0	1	2	3	4	5	6	7
α_j	2	-1.33		1	-4	3					
			β_j		-2	4					
α_j	3	3.27		-2	9	-18	11				
			β_j		6	-18	18				
α_j	4	-11.5		3	-16	36	-48	25			
			β_j		-12	48	-72	48			
α_j	5	52.6		-12	75	-200	300	-300	137		
			β_j		60	-300	600	-600	300		
α_j	6	-293		10	-72	225	-400	450	-360	147	
			β_j		-60	360	-900	1200	-900	360	
α_j	7	1944		-60	490	-1764	3675	-4900	4410	-2940	1089
			β_j		420	-2940	8820	-14700	14700	-8820	2940

Thus, e.g., for $k = 3$ we obtain

$$3y_{n+2} - 4y_{n+1} + y_n = h(4f_{n+1} - 2f_n),$$

which is equivalent to

$$y_{n+2} - \frac{4}{3}y_{n+1} + \frac{1}{3}y_n = h\left(\frac{4}{3}f_{n+1} - \frac{2}{3}f_n\right).$$

Let us now consider a different motivation for minimal-projecting LMMs: The partial projection of back values f_{n+j} ($j < k$) could be avoided completely by using methods

$$(4.10) \quad \sum_{j=0}^k \alpha_j y_{n+j} = h\beta_k f_{n+k}.$$

These are termed backward-differentiation methods; they were first used for stiff systems by Curtiss and Hirschfelder [5] and have been much used by later authors, notably Gear, who used (4.10) for the stiff option in his general purpose integration package for ordinary differential equations [6].

(4.10) has the advantage that the zeros of the stability polynomial governing the numerical solution, tend to zero as $|h\lambda|$ tends to infinity. Thus, (4.10) possesses an infinite region of absolute stability (in the sense that the region of absolute instability is bounded).

Unfortunately, (4.10) is implicit and, hence, inferior as a basic method for CDS schemes. One idea one can consider is to choose the β_j ($j = 0, 1, \dots, k - 1$) in the LMM (4.4) such that

$$(4.11) \quad \sum_{j=0}^{k-1} \beta_j y'(x_{n+j}) = \left(\sum_{j=0}^{k-1} \beta_j \right) y'(x_{n+k}) + O(h^q)$$

for some q . (4.4) satisfying (4.11) can be considered an approximation to (4.10).

On expanding the $y'(x_{n+j})$ in (4.11) about x_{n+k} we find that (4.11) holds if and only if the corresponding LMM is nonprojecting of degree q .

Because of this property the first characteristic polynomial of the minimal-projecting k -step method is identical with that of the k -step backward-differentiation method (4.10).

This is formally stated in the following definition and theorem.

Definition 7. The polynomial $\zeta(\xi) = \sum_{j=0}^k \alpha_j \xi^j$ is said to be the *first characteristic polynomial* of the linear multistep method (4.4).

THEOREM 3. For $k \geq 2$ the first characteristic polynomial of the k -step method (4.10) is identical to that of the minimal-projecting k -step method defined in Theorem 2.

Proof. Let

$$\sum_{j=0}^k \alpha_j y_{n+j} = h\beta_k f_{n+k}$$

denote the k th order k -step method (4.10), and

$$\sum_{j=0}^k \hat{\alpha}_j y_{n+j} = h \sum_{j=0}^k \hat{\beta}_j f_{n+j},$$

the k th order minimal-projecting k -step method defined in Theorem 2.

Then, according to the definition of order (see [4]), we have for an arbitrary sufficiently smooth test function $y(x)$

$$\sum_{j=0}^k \alpha_j y(x_{n+j}) - h\beta_k y'(x_{n+k}) = O(h^{k+1}),$$

$$\sum_{j=0}^k \hat{\alpha}_j y(x_{n+j}) - h \sum_{j=0}^{k-1} \hat{\beta}_j y'(x_{n+j}) = O(h^{k+1}).$$

On subtracting and using (4.11) with $q = k$, we obtain

$$\sum_{j=0}^k (\alpha_j - \hat{\alpha}_j) y(x_{n+j}) = O(h^{k+1}).$$

Since this is true for all (sufficiently smooth) test functions $y(x)$, it follows that $\alpha_j = \hat{\alpha}_j$ for $j = 0, 1, \dots, k$, which is the statement of the theorem. \square

It is natural to require that the basic method of a CDS scheme be convergent. Since minimal-projecting LMMs are consistent, this is equivalent to them being zero-stable.

It is well known that the BDF methods (4.10) are not zero-stable for $k \geq 7$. (A numerical investigation in [2] showed that they are zero-unstable for $k = 7, 8, \dots, 17$.)

The restricted order of minimal-projecting LMMs is clearly a drawback. However, bearing in mind that the CDS scheme (2.3) is aimed at obtaining cheap numerical solutions that are stable and reasonably, but not extremely, accurate, order 6 appears to be an adequate maximum order. Note also that interprojection phenomena become stronger as the stepnumber of the basic method increases which also restricts the order of the basic method.

Let us now compare minimal-projecting methods with Adams-Bashforth methods, defined by

$$(4.12) \quad y_{n+k} = y_{n+k-1} + h \sum_{j=0}^{k-1} \beta_j f_{n+j},$$

which was first used as early as 1883 [4]. In modern terms a motivation for using (4.12) is that the zeros of the first characteristic polynomial (with the necessary exception of 1) equal zero. (4.12) probably defines the most widely used class of explicit linear multistep methods.

Both types of methods are k -step methods of order k ; the minimum stepnumber of minimal projecting methods is 2, as opposed to 1 for (4.12); (4.12) is zero-stable for all $k = 1, 2, 3, \dots$, whereas minimal-projecting methods are zero-stable for $k = 2, 3, \dots, 6$, and not zero-stable for $k = 7, 8, \dots, 17$ (see above).

Both classes of methods have the disadvantage that the β_j are numerically large and alternate in sign (see Table 2 and [7]). This can introduce and amplify round-off errors.

Further comparison is contained in Table 3. There C_{MP} and C_{AB} denote the error constants, and $(-\kappa_{MP}, 0)$ and $(-\kappa_{AB}, 0)$ are the intervals of absolute stability of the minimal-projecting and the Adams-Bashforth methods, respectively. (Recall that the interval of absolute stability is the intersection of the region of absolute stability with the real line.) For the error constants of the Adams-Bashforth methods see [8, p. 26].

TABLE 2

Comparison of Adams-Bashforth methods and minimal-projecting methods

k	C_{MP}	C_{AB}	κ_{MP}	κ_{AB}
1		0.5000		2.0000
2	0.4444	0.4167	1.3333	1.0000
3	0.4091	0.3750	0.9524	0.5000
4	0.3840	0.3486	0.7111	0.3000
5	0.3650	0.3299	0.5505	0.1633
6	0.3499	0.3156	0.4402	0.0877

It follows from Table 2 that the error constants of minimal-projecting methods are slightly larger than those of Adams-Bashforth methods, which means that the local accuracy of the latter can be expected to be slightly higher. On the other hand the stability properties of the minimal-projecting methods are much better than those of the Adams-Bashforth methods; in fact, the interval of absolute stability of the 6th-order minimal-projecting methods is almost 50% larger than that of the 4th-order Adams-Bashforth method and only slightly smaller than that of the 3rd-order Adams-Bashforth method. A heuristic explanation for this phenomenon is that minimal-projecting methods can be considered, in the sense outlined above, to be approximations to the implicit backward-differentiation methods which possess infinite regions of absolute stability.

5. Numerical Examples. In this section we compare the fourth-order Adams-Bashforth and the fourth order minimal-projecting methods as basic methods for the gradient prediction scheme described in [1]. The inverse linear multistep method used for the gradient prediction is the strongly infinite stable 4-step method given in [1].

The example problem is constructed artificially such that a function that controls the change of the eigensystem with x can be incorporated.

We consider the problem

$$(5.1) \quad y' = A(x)(z(x) - y) + z'(x), \quad w = y(0) = z(0), \quad x \in [0, 100],$$

where

$$z(x) = \left[\sin x, \sin\left(\frac{\pi x}{4}\right), e^{-x} \right]^T$$

$$A(x) = \begin{bmatrix} \alpha\eta(x) - \beta & \beta - \alpha & (\beta - \alpha)/\eta(x) \\ (\gamma - \beta)\eta(x) & \beta\eta(x) - \gamma & \beta - \gamma \\ (\alpha - \gamma)\eta^2(x) & (\gamma - \alpha)\eta(x) & \gamma\eta(x) - \alpha \end{bmatrix}.$$

The solution of (5.1) obviously is $z(x)$, independent of $A(x)$. The reasons behind the selection of the exact solution $z(x)$ were to have functions that are well-behaved over the entire interval $[0, 100]$, but not as simple as polynomials; and, although the first entries of $z(x)$ are periodic, $z(x)$ as a whole should not be periodic. Apart from this the choice of $z(x)$ is arbitrary.

The reasons for the choice of $A(x)$ is to have a matrix, whose eigensystem can be easily controlled by the choice of the parameters α , β , γ , and $\eta(x)$.

The eigensystem of $A(x)$ is given by

$$\begin{aligned} \lambda^{(1)} &= \alpha; & c^{(1)} &= \kappa(x)[1, 0, \eta(x)]^T; & d^{(1)} &= \omega(x)[\eta(x), -1, -1/\eta(x)]^T, \\ \lambda^{(2)} &= \beta; & c^{(2)} &= \kappa(x)[1, \eta(x), 0]^T; & d^{(2)} &= \omega(x)[-1, 1, 1/\eta(x)]^T, \\ \lambda^{(3)} &= \gamma; & c^{(3)} &= \kappa(x)[0, 1, -\eta(x)]^T; & d^{(3)} &= \omega(x)[\eta(x), -1, -1]^T, \end{aligned}$$

where

$$\kappa(x) = 1/\sqrt{(1 + \eta^2(x))}, \quad \omega(x) = \sqrt{1 + \eta^2(x)}/(\eta(x) - 1).$$

(In the numerical test, the dominant eigensystem was not computed from the above information, but is obtained at each step by the power method.)

The eigenvalues of $A(x)$ are chosen to be

$$\alpha = -10^6, \quad \beta = -1, \quad \gamma = -2.$$

Thus, the system (5.1) is separably stiff.

The functions $\eta(x)$, governing the change of the eigensystems are given by

$$\eta(x) = -2 + 1.5 \sin(\xi x),$$

where

$$\xi = 0.9234567, 1.1234567, \dots, 2.5234567.$$

The period of oscillation increases by steps of 0.2, the digits 234567 are chosen such as to make the period incommensurable with any occurring in the exact solution.

In Table 3 we use the following abbreviations:

AB: Adams-Bashforth method (of order 4)

MP: minimal projecting method (of order 4)

MC: = $\max_{4 \leq n \leq 1000} |\xi_n^{(1)}|$ (maximum correction factor)

MD: = $\max_{4 \leq n \leq N} |\langle d_n^{(1)}, y(x_n) - y_n \rangle|$ (maximum dominant error)

MS: = $\max_{4 \leq n \leq N} \|y(x_n) - \langle d_n^{(1)}, y(x_n) - y_n \rangle c_n^{(1)}\|$ (maximum subdominant error)

The step size employed is $h = 0.1$; the starting values are exact. We define the method to have failed if MD or MS exceeds 100.

The results are as follows:

TABLE 3
Numerical results

ξ	BASIC METHOD	MC	MD	MS
0.9234567	MP	3.10_{10}^{-4}	6.46_{10}^{-9}	9.98_{10}^{-5}
	AB	9.60_{10}^{-4}	6.46_{10}^{-9}	1.06_{10}^{-4}
1.1234567	MP	3.06_{10}^{-4}	6.34_{10}^{-9}	1.02_{10}^{-4}
	AB	Failed after 420 steps		
1.3234567	MP	2.93_{10}^{-4}	6.15_{10}^{-9}	1.01_{10}^{-4}
	AB	Failed after 213 steps		
1.5234567	MP	3.89_{10}^{-4}	6.48_{10}^{-9}	9.69_{10}^{-5}
	AB	Failed after 146 steps		
1.7234567	MP	5.02_{10}^{-3}	6.49_{10}^{-9}	9.53_{10}^{-5}
	AB	Failed after 130 steps		
1.9234567	MP	1.82_{10}^{-3}	7.20_{10}^{-9}	9.36_{10}^{-5}
	AB	Failed after 114 steps		
2.1234566	MP	4.09_{10}^{-3}	9.40_{10}^{-9}	1.84_{10}^{-4}
	AB	Failed after 104 steps		
2.3234567	MP	2.81_{10}^{-1}	4.11_{10}^{-7}	1.39_{10}^{-2}
	AB	Failed after 123 steps		
2.5234567	MP	Failed after 261 steps		
	AB	Failed after 110 steps		

Conclusions. Minimal-projecting methods appear to be more robust with respect to rapidly changing eigensystems than Adams-Bashforth methods, if employed as the basic method in CDS schemes.

They have the disadvantage that their stepnumber is restricted to $k \leq 6$. However, this is not a severe restriction, since they will be employed for separably stiff systems with significant interprojection effects in which case small stepnumbers are preferable.

The zero-stable minimal-projecting methods have comparatively large regions of absolute stability. This can be explained heuristically by considering them, in the sense outlined above, to be explicit approximations to the implicit backward-differentiation formulas.

Acknowledgements. The author wishes to acknowledge helpful remarks made by the referee, particularly those concerning the proof of Theorem 2.

Department of Mathematics
University of Utah
Salt Lake City, Utah 84112

1. P. ALFELD, "Inverse linear multistep methods for the numerical solution of initial value problems of ordinary differential equations," *Math. Comp.*, v. 33, 1979, pp. 111–124.
2. P. ALFELD, *Correction in the Dominant Space: A New Technique for the Numerical Solution of Certain Stiff Initial Value Problems*, Ph. D. Thesis, University of Dundee, 1977.
3. P. ALFELD & J. D. LAMBERT, "Corrections in the dominant space: A numerical technique for a certain class of stiff initial value problems," *Math. Comp.*, v. 31, 1977, pp. 922–938.
4. F. BASHFORTH & J. C. ADAMS, *Theories of Capillary Action*, Cambridge Univ. Press, Cambridge, 1883.
5. C. F. CURTISS & J. O. HIRSCHFELDER, "Integration of stiff systems," *Proc. Nat. Acad. Sci. U.S.A.*, v. 38, 1952, pp. 235–243.
6. C. W. GEAR, *Numerical Initial Value Problems in Ordinary Differential Equations*, Prentice-Hall, Englewood Cliffs, N. J., 1971.
7. P. HENRICI, *Discrete Variable Methods in Ordinary Differential Equations*, Wiley, New York, 1962.
8. J. D. LAMBERT, *Computational Methods in Ordinary Differential Equations*, Wiley, New York, 1973.