



Chapter 6

B-series and multivalue methods

6.1 Introduction

Multivalue and multistage methods

The history of multivalue methods parallels the history of Runge–Kutta methods. Runge–Kutta methods achieved high accuracy through the multistage approach, whereas multivalue methods obtained improvements by re-using computed information in two or more steps.

The first notable publication on multistep methods was motivated by the need for numerical results for a specific problem [3] (Bashforth and Adams, 1883) and the Adams–Bashforth method was introduced. When Adams–Moulton methods [74] (Moulton, 1926) became available, predictor–corrector methods increased in popularity and have become a dominant technique in practical computation.

So-called “stiff” problems, such as those arising from the space-discretisation of partial differential equations, brought with them special difficulties which rendered Adams methods impractical and inefficient for many problems. The backward difference methods introduced in [35] (Curtiss and Hirschfelder, 1952) were a timely response to stiffness.

Even though multistep and Runge–Kutta methods developed individually and separately, they have always had a common core. That is, they are each built up from two basic operations and nothing more: the evaluation of the function f and the calculation of linear combinations of existing vectors.

Generalizations of traditional methods

Several innovations in the 1960s illustrated that methods could exist with aspects of both classical families, although they belonged to neither of them. These new ideas included the possible use of off-step predictors as a modification to the standard linear multistep algorithms [46] (Gragg and Stetter, 1964), [43] (Gear, 1965), [10]

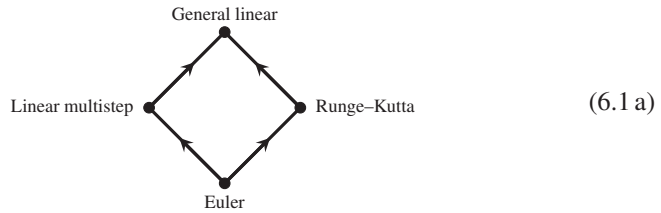
(Butcher, 1965). A second break from the strict linear multistep type of methods was the introduction of cyclic composite methods [41] (Donelson and Hansen, 1971).

Modifications to Runge–Kutta are also possible, such as cyclic composite forms of these methods and pseudo Runge–Kutta methods [27]. (Byrne and Lambert, 1966).

Methods in which stage derivatives in one step are approximated by stage derivatives already evaluated in a previous step will be referred to as “re-use” methods. An example of this is given in (6.3 h) (p. 219).

Schematic diagram of classes of numerical methods

As we have seen, multivalued and multistage methods are generalizations of the Euler method. A diagram showing the various types of methods that can be built on these ideas is given in (6.1 a):



In this diagram, \diagdown symbolizes the use of a multivalued method, rather than a one value method, whereas \diagup symbolizes a multistage rather than a one stage method. General linear methods at the top of the diagram are both multivalued and multistage.

General linear methods

In [12] (Butcher, 1966), the methods, now called General Linear Methods, were introduced with the intention that the multivalued and multistage aspects of the method should be equally balanced. These methods, using the formulation based on [6] (Burrage and Butcher, 1980), are the principal subjects of the present chapter.

In traditional and closely related methods, the quantities being approximated have a natural meaning. However, for a method written in terms of coefficient matrices, this will not always be the case. Hence, a completely fresh approach to the meaning of order is required, which will be theoretically sound and, at the same time, is practical.

Chapter outline

Section 6.2, a broad survey of linear multistep methods, is followed by Section 6.3 which attempts to motivate the need for the class of general linear methods. The formulation of these methods is presented in Section 6.4. In Section 6.5, the meaning

of order, based on the use of a starting method, to be used together with the main method, is introduced.

In Section 6.6, a general approach is discussed for determining the order of a method in terms of the B-series of the “underlying one-step method”. See [87] (Stoffer, 1993). Denote this B-series by $(\mathbf{B}_h y_0)a$. To complete the process of finding the order, an algorithm is constructed for finding the maximum p such that $a \sim E + O_{p+1}$.

6.2 Survey of linear multistep methods

It is characteristic of multistep methods that some preliminary work has to be carried out before the method can be used in its own right. In the case of a linear k -step method, $k - 1$ steps need to be performed by some other numerical method before the information is available to allow the method to be used in subsequent steps. After these $k - 1$ steps have been evaluated, approximations to the solution and the scaled derivatives at $x_i = x_0 + hi$, $i = 0, 1, \dots, k - 1$ are available. This will enable step number k and subsequent steps, to be computed using the formula

$$y_n = \sum_{i=1}^k a_i y_{n-i} + \sum_{i=0}^k b_i h f(y_{n-i}), \quad (6.2 a)$$

where a_i , $i = 1, 2, \dots, k$ and b_i , $i = 0, 1, \dots, k$ are real constants.

Note that if $b_0 \neq 0$, the method is implicit so that y_n and $h f(y_n)$ need to be evaluated together using an iterative process. Following the practice introduced in [36] (Dahlquist, 1956), a linear multistep method (6.2 a) is characterized, not directly in terms of the a_i and b_i , but in terms of two polynomials:

$$\rho(w) = w^k - a_1 w^{k-1} - \dots - a_k, \quad (6.2 b)$$

$$\sigma(w) = b_0 w^k + b_1 w^{k-1} + \dots + b_k. \quad (6.2 c)$$

It is customary to refer to the method given by (6.2 a) as the method (ρ, σ) .

Basic definitions are

Definition 6.2A The method (ρ, σ) is consistent if

$$\rho(1) = 0, \quad (6.2 d)$$

$$\rho'(1) = \sigma(1). \quad (6.2 e)$$

Definition 6.2B The method (ρ, σ) is stable if the difference equation

$$u_n = \sum_{i=1}^k a_i u_{n-i}, \quad n = k, k+1, \dots,$$

has bounded solutions for all possible choices of the initial values u_0, u_1, \dots, u_{k-1} .

The aim of these definitions is to characterize what it means for a method to be convergent. That is, if the values of y_1, \dots, y_{k-1} are determined by appropriate starting values, the value of y_n computed by the method, using stepsize H/n , the final approximation should converge to $y(x_0 + H)$ as $n \rightarrow \infty$. Informally, appropriate starting values \mathcal{S}_h^i , $i = 1, 2, \dots, k-1$, means approximations to $y(x_0 + ih)$ as $h \rightarrow 0$. The definition of convergence is given in [36] (Dahlquist, 1956) and in other references, such as [55] (Henrici, 1962), [43] (Gear, 1965), [50] (Hairer, Nørsett and Wanner, 1993), [20] (Butcher, 2016), as is the theorem relating this concept to the consistency and stability properties given in Definitions 6.2A and 6.2B. The final outcome is that consistency and stability are together necessary and sufficient for convergence.

Comments on notation

In this presentation, (6.2 a, 6.2 b, 6.2 c) correspond respectively to

$$\begin{aligned} \sum_{i=0}^k \alpha_i y_{n+i} &= \sum_{i=0}^k \beta_i h f(y_{n+i}), \\ \rho(\zeta) &= \sum_{i=0}^k \alpha_i \zeta^i, \\ \sigma(\zeta) &= \sum_{i=0}^k \beta_i \zeta^i. \end{aligned}$$

in [36]. However, the correspondence is exact only when the coefficients are scaled so that $\alpha_{k+1} = 1$. In Dahlquist's classic work, there is no such restriction.

Order conditions

Consider a single step of the method (ρ, σ) , starting from the exact initial value

$$y_k = \sum_{i=1}^k a_i y(x_{k-i}) + h \sum_{i=0}^k b_i f(y(x_{k-i})). \quad (6.2 f)$$

The error committed in this single step is $y(x_k) - y_k$, which becomes, written in B-series,

$$(\mathbf{B}_h y_0)(\mathbf{E}^k) - \sum_{i=1}^k a_i (\mathbf{B}_h y_0)(\mathbf{E}^{k-i}) - \sum_{i=0}^k b_i (\mathbf{B}_h y_0)(\mathbf{E}^{k-i} \mathbf{D}).$$

For order p , this must give an expansion for which the coefficients of $\mathbf{F}(t)$ is zero for all trees satisfying $|t| \leq p$. That is

$$\mathbf{E}^k - \sum_{i=1}^k a_i \mathbf{E}^{k-i} - \sum_{i=0}^k b_i \mathbf{E}^{k-i} \mathbf{D} = \mathcal{O}_{p+1}$$

or, what is equivalent,

$$\begin{aligned} 1 &= a_1 \mathbf{E}^{-1} + a_2 \mathbf{E}^{-2} + \cdots + a_k \mathbf{E}^{-k} \\ &\quad + b_0 \mathbf{D} + b_1 \mathbf{E}^{-1} \mathbf{D} + b_2 \mathbf{E}^{-2} \mathbf{D} + \cdots + b_k \mathbf{E}^{-k} \mathbf{D} + \mathcal{O}_{p+1}. \end{aligned} \quad (6.2 \text{ g})$$

Stability regions

We will consider the behaviour of a linear multistep method (ρ, σ) , in attempting to solve the linear problem $y' = qy$, where q is a complex scalar constant. The analysis is also applicable to $y' = Qy$, where the $N \times N$ constant matrix Q is diagonalizable. A sequence of k step values satisfies the relation

$$(1 - hqb_0)y_n = a_1 y_{n-1} + \cdots + a_k y_{n-k} + hq(b_1 y_{n-1} + \cdots + b_k y_{n-k}), \quad (6.2 \text{ h})$$

so that the sequence satisfies the difference equation with characteristic polynomial $\rho(w) - z\sigma(w)$, where $z = hq \in \mathbb{C}$. A complex number z is a “stable point” if solutions to this difference equation are bounded; the set of all stable points is the “stability region”.

A-stability

Following [37] (Dahlquist, 1963), we define

Definition 6.2C A method (ρ, σ) is A-stable if the open left-half complex plane \mathbb{C}^- is a subset of the stability region.

Significance of A-stability

If $\operatorname{Re} q < 0$ then the exact solution behaves like $\exp(qx)$. It makes sense to model problems with decaying solutions by numerical approximations which are at least bounded with increasing time. While it might be difficult to guarantee this in general, we can at least achieve it for the case of linear problems.

One-leg methods

In [38] (Dahlquist, 1976), the idea of “one-leg methods” was introduced. For these methods, the terms in (6.2 f)

$$h \sum_{i=0}^k b_i f(y(x_{k-i}))$$

are replaced by

$$h \left(\sum_{i=0}^k b_i \right) f \left(\frac{b_i}{\sum_{i=0}^k b_i} y(x_{k-i}) \right).$$

The order conditions for the linear multistep method (6.2 g) become, for the one-leg method,

$$\begin{aligned} 1 &= a_1 E^{-1} + a_2 E^{-2} + \cdots + a_k E^{-k} \\ &+ \left(\sum_{i=0}^k b_i \right) \left(\frac{b_0 + b_1 E^{-1} + b_2 E^{-2} + \cdots + b_k E^{-k}}{b_0 + b_1 + b_2 + \cdots + b_k} \right) D + O_{p+1}. \end{aligned}$$

Exercise 54 Show that the stability regions of a linear multistep method and the corresponding one-leg method are identical.

The Dahlquist barriers

The first barrier

Traditional linear multistep methods and predictor-corrector methods are governed by the Dahlquist barrier [36] (Dahlquist, 1956), quoted here without proof.

Theorem 6.2D (Dahlquist barrier) The order of a stable k -step method cannot exceed $k + 2$ (k even) or $k + 1$ (k odd).

The second Dahlquist barrier

The second barrier is concerned with the attainable order of A-stable k -step methods. The result is proved using order-stars [53] (Hairer and Wanner, 1996), or alternatively by order-arrows [20] (Butcher, 2016) and states

Theorem 6.2E (Second Dahlquist barrier) The order of an A-stable k -step method cannot exceed 2.

6.3 Motivations for general linear methods

There are several reasons why a wider formulation of numerical methods than offered by either of the traditional linear multistep or one-step schemes is appropriate. First, the barriers on what is achievable, if the constraints of the traditional methods can be overcome, and we start by considering some of these.

Breaking the Dahlquist barrier

The consequences of the first Dahlquist barrier on the order of linear k -step methods can be avoided in various ways.

Breaking the barrier using off-step points

In a number of independent contributions [46] (Gragg and Stetter, 1964), [43] (Gear, 1965), [10] (Butcher, 1965), a new approach to numerical integration was put forward in several independent projects. These can be seen as attempts to overcome limitations inherent in traditional methods by the use of off-step points.

For example, a two-predictor 2-step method which is stable and has order 5 is given by (6.3 a)–(6.3 c) below. The first predictor (6.3 a) gives an approximation $\hat{y}_{n-1/2}$ to $y(x_{n-1/2})$ from which $\hat{f}_{n-1/2} := f(\hat{y}(x_{n-1/2})) \approx y'(x_{n-1/2})$ is found. The second predictor (6.3 b) gives $\hat{y}(x_n) \approx y(x_n)$, leading to $\hat{f}_n := f(\hat{y}(x_n)) \approx y'(x_n)$. Finally, the corrector formula (6.3 c) gives $y_n \approx y(x_n)$

$$\hat{y}_{n-1/2} = y_{n-2} + \frac{9}{8}hf_{n-1} + \frac{3}{8}hf_{n-2}, \quad (6.3 \text{ a})$$

$$\hat{y}_n = \frac{28}{5}y_{n-1} - \frac{23}{5}y_{n-2} - 4hf_{n-1} - \frac{26}{15}hf_{n-2} + \frac{32}{15}h\hat{f}_{n-1/2}, \quad (6.3 \text{ b})$$

$$y_n = \frac{32}{31}y_{n-1} - \frac{1}{31}y_{n-2} + \frac{5}{31}h\hat{f}_n + \frac{4}{31}hf_{n-1} - \frac{1}{93}hf_{n-2} + \frac{64}{93}h\hat{f}_{n-1/2} \quad (6.3 \text{ c})$$

Breaking the barrier using cyclic composite methods

Cyclic composite methods were introduced in [41] (Donelson and Hansen, 1971). Consider the family of fifth order 3-step methods

$$\begin{aligned} 33y_n + (24 + 57\lambda)y_{n-1} - (57 + 24\lambda)y_{n-2} - 33\lambda y_{n-3} \\ = (10 - \lambda)hf(y_n) + (57 + 24\lambda)hf(y_{n-1}) \\ + (24 + 57\lambda)hf(y_{n-2}) - (1 - 10\lambda)hf(y_{n-3}). \end{aligned} \quad (6.3 d)$$

According to Theorem 6.2D, this method is unstable for any choice of λ . However, it can be used in a stable manner by alternating the value of λ between steps. For example $\lambda = 0$ could be used in even numbered steps and $\lambda = -\frac{361}{240}$ in odd numbered steps.

The composite method can now be written,

$$\begin{aligned} y_{2n} &= -\frac{8}{11}y_{2n-1} + \frac{19}{11}y_{2n-2} \\ &\quad + \frac{10}{33}hf(y_{2n}) + \frac{19}{11}hf(y_{2n-1}) + \frac{8}{11}hf(y_{2n-2}) - \frac{1}{33}hf(y_{2n-3}), \\ y_{2n+1} &= \frac{449}{240}y_{2n} + \frac{19}{30}y_{2n-1} - \frac{361}{240}y_{2n-2} \\ &\quad + \frac{251}{720}hf(y_{2n+1}) + \frac{19}{30}hf(y_{2n}) - \frac{449}{240}hf(y_{2n-1}) - \frac{35}{72}hf(y_{2n-2}). \end{aligned} \quad (6.3 e)$$

The stability of each of the two methods can be characterized by the companion matrices of their ρ polynomials; that is, the pair of matrices

$$M_1 = \begin{bmatrix} -\frac{8}{11} & \frac{19}{11} & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}, \quad M_2 = \begin{bmatrix} \frac{449}{240} & \frac{19}{30} & -\frac{361}{240} \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}. \quad (6.3 f)$$

Neither M_1 nor M_2 is power-bounded, a criterion equivalent to Definition 6.2B. However, for the cyclic method, stability is determined by the product

$$M_2 M_1 = \begin{bmatrix} -\frac{8}{11} & \frac{19}{11} & 0 \\ -\frac{8}{11} & \frac{19}{11} & 0 \\ 1 & 0 & 0 \end{bmatrix} =: M,$$

which is power bounded because $M^n = M^2 = \begin{bmatrix} 1 & 1 & 1 \end{bmatrix}^T \begin{bmatrix} -\frac{8}{11} & \frac{19}{11} & 0 \end{bmatrix}$, for $n = 3, 4, \dots$

Exercise 55 Show that the composite cyclic method based on (6.3 e) is stable if $\lambda = 0$ in even-numbered steps and $\lambda = \mu$, where $\mu \in (-\frac{241}{120}, -1)$ in odd-numbered steps.

Breaking the Runge–Kutta order barriers

Although explicit Runge–Kutta methods require only p stages for order $p = 1, 2, 3, 4$, it was shown in Theorem 5.5A (p. 200) that order $p \geq 5$ requires at least $p + 1$ stages. However, generalizations of Runge–Kutta methods are available to alleviate these restrictions.

Breaking the barrier by re-use of stages

The following tableau, for a 6 stage fifth order method, can be modified to become a 5 stage method in which stage number 2 is replaced by the value of stage number 4, evaluated in the *previous* step.

$$\begin{array}{c|cccccc}
 0 & 0 & & & & \\
 -\frac{1}{2} & -\frac{1}{2} & & & & \\
 \frac{1}{4} & \frac{5}{16} & -\frac{1}{16} & & & \\
 \frac{1}{2} & \frac{3}{4} & -\frac{1}{4} & & & \\
 \frac{3}{4} & \frac{15}{16} & \frac{3}{8} & \frac{3}{4} & \frac{9}{16} & \\
 1 & \frac{18}{7} & -1 & 0 & -\frac{12}{7} & \frac{8}{7} \\
 \hline
 & \frac{7}{90} & 0 & \frac{16}{45} & \frac{2}{15} & \frac{16}{45} & \frac{7}{90}
 \end{array} \quad (6.3 \text{ g})$$

Rewrite the remaining stages in step number n , after the second stage is deleted, as $Y_i^n, i = 1, 2, \dots, 5$, with a similar notation for the stage derivatives, and the method becomes

$$\begin{aligned}
 Y_1^n &= y_{n-1}, \\
 Y_2^n &= y_{n-1} + \frac{5}{16}hF_1^n - \frac{1}{16}hF_3^{n-1}, \\
 Y_3^n &= y_{n-1} + \frac{3}{4}hF_1^n - \frac{1}{4}hF_3^{n-1}, \\
 Y_4^n &= y_{n-1} - \frac{15}{16}hF_1^n + \frac{3}{8}hF_3^{n-1} + \frac{3}{4}hF_2^n + \frac{9}{16}hF_3^n, \\
 Y_5^n &= y_{n-1} + \frac{18}{7}hF_1^n - hF_3^{n-1} - \frac{12}{7}hF_3^n + \frac{8}{7}hF_4^n, \\
 y_n &= y_{n-1} + \frac{7}{90}hF_1^n + \frac{16}{45}hF_2^n + \frac{2}{15}hF_3^n + \frac{16}{45}hF_4^n + \frac{7}{90}hF_5^n.
 \end{aligned} \quad (6.3 \text{ h})$$

This method is reformulated as a general linear method in (6.4 n) (p. 225).

Breaking the barrier using effective order

Effective order, or conjugate order Section 5.7 (p. 205), is available as a means of obtaining fifth order accuracy with five stages, as long as the work expended to carry out pre- and post-processing is added to the cost.

Breaking the barrier using cyclic composite methods

Even though methods with $s = p = 5$ do not exist, it is possible to construct methods with $s = 5$ which satisfy all except one of the fifth order conditions. The following two tableaux are examples of this

0						0					
$\frac{5}{8}$	$\frac{5}{8}$					$\frac{7}{8}$	$\frac{7}{8}$				
$\frac{1}{4}$	$\frac{1}{5}$	$\frac{1}{20}$				$\frac{7}{10}$	$\frac{21}{50}$	$\frac{7}{25}$			
$\frac{7}{10}$	$\frac{1127}{1250}$	$-\frac{259}{625}$	$\frac{252}{125}$			$\frac{1}{4}$	$\frac{75}{392}$	$-\frac{11}{196}$	$\frac{45}{392}$		
1	$\frac{737}{175}$	$\frac{44}{25}$	$-\frac{32}{5}$	$\frac{10}{7}$		1	$\frac{501}{245}$	$-\frac{268}{245}$	$\frac{46}{49}$	$\frac{16}{5}$	
	$\frac{1}{14}$	0	$\frac{32}{81}$	$\frac{250}{567}$	$\frac{5}{54}$		$\frac{1}{14}$	0	$\frac{250}{567}$	$\frac{32}{81}$	$\frac{5}{54}$

If the 17 order conditions up to order 5 are tested, they are satisfied in each case, except for the single tree $t_{16} = [[[\tau^2]]]$. For this tree, the condition is $b^T A^2 c^2 = \frac{1}{60}$ but the values of the elementary weight for the two methods are $\frac{1}{960} = \frac{1}{60} - \frac{1}{64}$ and $\frac{31}{960} = \frac{1}{60} + \frac{1}{64}$, respectively. If the two methods are used cyclically, the $\pm \frac{1}{64}$, contributions to the error coefficients cancel out and fifth order is achieved after every pair of steps.

A common basis for one-step and multistep methods

The most important motivation for introducing general linear methods is that it is natural. For all step-by-step methods, some data is received at the beginning of each step and updated for output and subsequent use by the following step. The updating consists of the calculation of one or more approximations to the solution at points in or near the step; from these approximations, stage derivatives as samples of the vector field are evaluated and made available for further calculations in the step or else made available in the updating process.

6.4 Formulation of general linear methods

Following the formulation in [6] (Burrage and Butcher, 1980), we denote the data input to step number n by $y^{[n-1]}$ and the data output at the completion of the step by $y^{[n]}$. Each of these is a vector in $(\mathbb{R}^N)^r$ and is decomposed into individual components in the form

$$y^{[n-1]} = \begin{bmatrix} y_1^{[n-1]} \\ y_2^{[n-1]} \\ \vdots \\ y_r^{[n-1]} \end{bmatrix}, \quad y^{[n]} = \begin{bmatrix} y_1^{[n]} \\ y_2^{[n]} \\ \vdots \\ y_r^{[n]} \end{bmatrix}. \quad (6.4a)$$

During the computation s stages are evaluated and for each of these stages the stage derivative needs to be evaluated. These are written as vectors in $(\mathbb{R}^N)^s$ with the notation

$$Y = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_s \end{bmatrix}, \quad F = \begin{bmatrix} F_1 \\ F_2 \\ \vdots \\ F_s \end{bmatrix} := \begin{bmatrix} f(Y_1) \\ f(Y_2) \\ \vdots \\ f(Y_s) \end{bmatrix}. \quad (6.4b)$$

To express the relation between these quantities, introduce a coefficient matrix partitioned as $(s+r) \times (s+r)$:

$$\begin{bmatrix} A & U \\ B & V \end{bmatrix} := \left[\begin{array}{cccc|cccc} a_{11} & a_{12} & \cdots & a_{1s} & u_{11} & u_{12} & \cdots & u_{1r} \\ a_{21} & a_{22} & \cdots & a_{2s} & u_{21} & u_{22} & \cdots & u_{2r} \\ \vdots & \vdots & & \vdots & \vdots & \vdots & & \vdots \\ a_{s1} & a_{s2} & \cdots & a_{ss} & u_{s1} & u_{s2} & \cdots & u_{sr} \\ \hline b_{11} & b_{12} & \cdots & b_{1s} & v_{11} & v_{12} & \cdots & v_{1r} \\ b_{21} & b_{22} & \cdots & b_{2s} & v_{21} & v_{22} & \cdots & v_{2r} \\ \vdots & \vdots & & \vdots & \vdots & \vdots & & \vdots \\ b_{r1} & b_{r2} & \cdots & b_{rs} & v_{r1} & v_{r2} & \cdots & v_{rr} \end{array} \right]. \quad (6.4c)$$

The evaluation of the result consists of evaluating the stages, together with the stage derivatives,

$$Y_i = \sum_{j=1}^s h a_{ij} F_j + \sum_{j=1}^r u_{ij} y_j^{[n-1]}, \quad i = 1, 2, \dots, s, \quad (6.4d)$$

followed by the evaluation of the output values

$$y_i^{[n]} = \sum_{j=1}^s h b_{ij} F_j + \sum_{j=1}^r v_{ij} y_j^{[n-1]}, \quad i = 1, 2, \dots, r. \quad (6.4e)$$

Written more compactly, (6.4d) and (6.4e) become

$$Y = h(A \otimes \mathbf{I})F + (U \otimes \mathbf{I})y^{[n-1]}, \quad (6.4f)$$

$$y^{[n]} = h(B \otimes \mathbf{I})F + (V \otimes \mathbf{I})y^{[n-1]} \quad (6.4g)$$

or

$$\begin{bmatrix} Y \\ y^{[n]} \end{bmatrix} = \begin{bmatrix} A \otimes \mathbf{I} & U \otimes \mathbf{I} \\ B \otimes \mathbf{I} & V \otimes \mathbf{I} \end{bmatrix} \begin{bmatrix} F \\ y^{[n-1]} \end{bmatrix}. \quad (6.4h)$$

There is usually no confusion if $\otimes \mathbf{I}$ is omitted from each element in (6.4h).

Consistency, stability and convergence

Consistency and pre-consistency

The first of the consistency conditions for linear multistep methods, (6.2 d), sometimes known as “pre-consistency”, really means that there is a possibility of following a constant solution correctly. The full condition including also (6.2 e) enables linear growth to be modelled. In general linear methods we also need a condition like covariance to get the consistent behaviour that we need.

Definition 6.4A A general linear method (A, U, B, V) is pre-consistent if there exists $u \in \mathbb{R}^r$, known as the “pre-consistency vector”, such that

$$\begin{aligned} Vu &= u, \\ Uu &= \mathbf{1} := [1 \ 1 \ \dots \ 1]^T \in \mathbb{R}^s. \end{aligned}$$

Definition 6.4B A general linear method (A, U, B, V) is consistent if it is pre-consistent with pre-consistency vector u , and there exists $v \in \mathbb{R}^r$ such that

$$Bu + Vv = v + u.$$

Stability

Because there are many uses of the term “stability”, the concept considered here is sometimes referred to as “zero-stability” or “stability in the sense of Dahlquist”.

Definition 6.4C A general linear method (A, U, B, V) is stable if there exists a constant C such that

$$\|V^n\| \leq C, \quad n = 1, 2, \dots$$

For an unstable method, an error due to truncation or round-off, committed in one step of a computation, can have an impact on the overall computation which grows without bound. These informal remarks will be made more precise in the discussion of convergence.

In the meantime we can give criteria for V having bounded powers. First we remark that V satisfies Definition 6.4C if and only if the same is true for \tilde{V} defined as the Jordan canonical form of V and this is true if and only if each Jordan block J satisfies $\|J^n\| \leq C$ for all $n = 1, 2, \dots$, for some C .

Lemma 6.4D For given complex λ and positive integer m , let J be the $m \times m$ matrix

$$J = \begin{bmatrix} \lambda & 0 & 0 & \cdots & 0 \\ \mu & \lambda & 0 & \cdots & 0 \\ 0 & \mu & \lambda & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & \lambda \end{bmatrix},$$

where μ is arbitrary non-zero and does not appear in the matrix if $m = 1$. Then J has bounded powers if and only if (i) $|\lambda| < 1$ or (ii) $|\lambda| = 1$ and $m = 1$.

Proof. If $|\lambda| < 1$, choose $\mu = 1 - |\lambda|$ so that $\|J\|_\infty = 1$, and hence $\|J^n\| \leq 1$, in all cases. The necessity of $|\lambda| \leq 1$ follows from the fact that the $(1, 1)$ element of J^n is λ^n , and the necessity of $m = 1$, when $|\lambda| = 1$, follows from the fact that, if $m \geq 2$, the $(2, 1)$ element of J^n is $n\mu\lambda^{n-1}$. \square

A consequence of this result is

Theorem 6.4E A method (A, U, B, V) is stable if and only if all zeros of the minimal polynomial of V lie in the closed unit disc and those on the boundary are simple.

Convergence

In the definition of convergence below, a Lipschitz continuous problem

$$y'(x) = f(y(x)), \quad y(x_0) = y_0, \quad (6.4 \text{ i})$$

is to be solved on the interval $[x_0, \bar{x}]$ using a starting method which satisfies $y^{[0]} = uy_0 + \alpha_n$ using n steps and stepsize $h = (\bar{x} - x_0)/n$ to give a final result $y^{[n]} = uy(\bar{x}) + \beta_n$.

Definition 6.4F A pre-consistent method (A, U, B, V) is convergent if in the solution of (6.4 i) with n steps with $\|\alpha_n\| \rightarrow 0$ as $n \rightarrow \infty$, then $\|\beta_n\| \rightarrow 0$ as $n \rightarrow \infty$.

Examples of traditional methods

Example of a Runge–Kutta method

The classical Runge–Kutta method

$$\begin{array}{c|cccc}
 0 & & & & \\
 \frac{1}{2} & \frac{1}{2} & & & \\
 \frac{1}{2} & 0 & \frac{1}{2} & & \\
 1 & 0 & 0 & 1 & \\
 \hline
 & \frac{1}{6} & \frac{1}{3} & \frac{1}{3} & \frac{1}{6}
 \end{array},$$

has the representation

$$\begin{bmatrix} A & U \\ B & V \end{bmatrix} = \left[\begin{array}{cccc|c}
 0 & 0 & 0 & 0 & 1 \\
 \frac{1}{2} & 0 & 0 & 0 & 1 \\
 0 & \frac{1}{2} & 0 & 0 & 1 \\
 0 & 0 & 1 & 0 & 1 \\
 \hline
 \frac{1}{6} & \frac{1}{3} & \frac{1}{3} & \frac{1}{6} & 1
 \end{array} \right].$$

This is not a unique representation of this method. An alternative is to compute, in step number n , the scaled derivative of the output result and export this as an additional output. The method now has $r = 2$ and a starting step, consisting of evaluating $hf(y_0)$ to serve as the second input in the following step. The modified method now becomes

$$\begin{bmatrix} A & U \\ B & V \end{bmatrix} = \left[\begin{array}{cccc|cc}
 0 & 0 & 0 & 0 & 1 & \frac{1}{2} \\
 \frac{1}{2} & 0 & 0 & 0 & 1 & 0 \\
 0 & 1 & 0 & 0 & 1 & 0 \\
 \frac{1}{3} & \frac{1}{3} & \frac{1}{6} & 0 & 1 & \frac{1}{6} \\
 \hline
 \frac{1}{3} & \frac{1}{3} & \frac{1}{6} & 0 & 1 & \frac{1}{6} \\
 0 & 0 & 0 & 1 & 0 & 0
 \end{array} \right]. \quad (6.4j)$$

Although this method does not have any special advantages, it points the way to methods for which the output at the end of step n contains approximations to each of $y(x_n)$, $hy'(x_n)$ and $\frac{1}{2}h^2y''(x_n)$. A method based on this generalisation is analysed in Section 6.5 (p. 235).

Examples of linear multistep methods

The Adams–Bashforth method of order 2 is given by

$$y_n = y_{n-1} + \frac{3}{2}hf(y_{n-1}) - \frac{1}{2}hf(y_{n-2}). \quad (6.4k)$$

This has the representation

$$\begin{bmatrix} A & U \\ B & V \end{bmatrix} = \left[\begin{array}{c|cccc} 0 & 1 & \frac{3}{2} & -\frac{1}{2} & \\ \hline 0 & 1 & \frac{3}{2} & -\frac{1}{2} & \\ 1 & 0 & 0 & 0 & \\ 0 & 0 & 1 & 0 & \end{array} \right].$$

Similarly, the third order method in the same family is

$$y_n = y_{n-1} + \frac{23}{12}hf(y_{n-1}) - \frac{4}{3}hf(y_{n-2}) + \frac{5}{12}hf(y_{n-3}), \quad (6.4l)$$

with representation

$$\begin{bmatrix} A & U \\ B & V \end{bmatrix} = \left[\begin{array}{c|ccccc} 0 & 1 & \frac{23}{12} & -\frac{4}{3} & \frac{5}{12} \\ \hline 0 & 1 & \frac{23}{12} & -\frac{4}{3} & \frac{5}{12} \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{array} \right]. \quad (6.4m)$$

Examples of non-traditional methods

Examples of re-use methods

It is easy to adapt the pattern in (6.4j) to more complicated methods, such as (6.3g) (p. 219), but with one of the stage derivatives re-used in a later step, as in (6.3h) (p. 219). The general linear representation becomes

$$\begin{bmatrix} A & U \\ B & V \end{bmatrix} = \left[\begin{array}{ccccc|cc} 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ \frac{5}{16} & 0 & 0 & 0 & 0 & 1 & -\frac{1}{16} \\ \frac{3}{4} & 0 & 0 & 0 & 0 & 1 & -\frac{1}{4} \\ -\frac{15}{16} & \frac{3}{4} & \frac{9}{16} & 0 & 0 & 1 & \frac{3}{8} \\ \frac{18}{7} & 0 & -\frac{12}{7} & \frac{8}{7} & 0 & 1 & -1 \\ \hline \frac{7}{90} & \frac{16}{45} & \frac{2}{15} & \frac{16}{45} & \frac{7}{90} & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 \end{array} \right]. \quad (6.4n)$$

Another example of a re-use method was given in (1.6d) (p. 32). For this method, the second order Runge–Kutta method

$$\begin{array}{c|ccc}
 0 & & & \\
 \frac{1}{2} & \frac{1}{2} & & \\
 1 & 0 & 1 & \\
 \hline
 & \frac{1}{6} & \frac{2}{3} & \frac{1}{6}
 \end{array},$$

applied in step number n , is modified by adding to the second stage an approximation to $\frac{1}{8}h^2y''(x_{n-1})$, given by $\frac{1}{4}hy'(x_{n-2}) - \frac{3}{4}hy'(x_{n-3/2}) + \frac{1}{2}hy'(x_{n-1})$, previously computed in step number $n-1$. The purpose of the starting method (1.6 e) (p. 33) is to provide an approximation to $\frac{1}{8}h^2y''(x_0)$, to use in step number 1.

Example of an off-step points method

The method given by (6.3 a) (p. 217) has a representation as an $rs = 43$ general linear method. The input quantities and the stage values will be written as

$$\begin{aligned}
 y_1^{[n-1]} &= y_{n-1}, \\
 y_2^{[n-1]} &= y_{n-2}, \\
 y_3^{[n-1]} &= hy'_{n-1}, \\
 y_4^{[n-1]} &= hy'_{n-2}, \\
 Y_1 &= \hat{y}_{n-1/2}, \\
 Y_2 &= \hat{y}_n, \\
 Y_3 &= y_n.
 \end{aligned}$$

With this notation the method can be written

$$\begin{bmatrix} A & U \\ B & V \end{bmatrix} = \left[\begin{array}{ccc|cccc}
 0 & 0 & 0 & 0 & 1 & \frac{9}{8} & \frac{3}{8} \\
 \frac{32}{15} & 0 & 0 & \frac{28}{5} & -\frac{23}{5} & -4 & -\frac{26}{15} \\
 \frac{64}{93} & \frac{5}{31} & 0 & \frac{32}{31} & -\frac{1}{31} & \frac{4}{31} & -\frac{1}{93} \\
 \hline
 \frac{64}{93} & \frac{5}{31} & 0 & \frac{32}{31} & -\frac{1}{31} & \frac{4}{31} & -\frac{1}{93} \\
 0 & 0 & 0 & 1 & 0 & 0 & 0 \\
 0 & 0 & 1 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 1 & 0
 \end{array} \right].$$

Example of a cyclic composite method

The method (6.3 e) (p. 218) carries the approximations over two steps and, hence, it will be convenient to rescale so that h is replaced by $h/2$ and renumber in half steps. It will then be convenient to substitute $y_{n-1/2}$ from the first equation into the second.

$$\begin{aligned}
y_{n-1/2} &= -\frac{8}{11}y_{n-1} + \frac{19}{11}y_{n-3/2} \\
&\quad + \frac{5}{33}hf(y_{n-1/2}) + \frac{19}{22}hf(y_{n-1}) + \frac{4}{11}hf(y_{n-3/2}) - \frac{1}{66}hf(y_{n-2}), \\
y_n &= \frac{449}{240}y_{n-1/2} + \frac{19}{30}y_{n-1} - \frac{361}{240}y_{n-3/2} \\
&\quad + \frac{251}{1440}hf(y_n) + \frac{19}{60}hf(y_{n-1/2}) - \frac{449}{480}hf(y_{n-1}) - \frac{35}{144}hf(y_{n-3/2}) \\
&= -\frac{8}{11}y_{n-1} + \frac{19}{11}y_{n-3/2} + \frac{4753}{7920}hf(y_{n-1/2}) + \frac{251}{1440}hf(y_n) \\
&\quad + \frac{449}{660}hf(y_{n-1}) + \frac{3463}{7920}hf(y_{n-3/2}) - \frac{449}{15840}hf(y_{n-2}).
\end{aligned}$$

To recast the method in general linear notation, write

$$\begin{aligned}
y_1^{[n-1]} &= y_{n-1}, & y_2^{[n-1]} &= y_{n-3/2}, & y_3^{[n-1]} &= hf(y_{n-1}), \\
y_4^{[n-1]} &= hf(y_{n-3/2}), & y_5^{[n-1]} &= hf(y_{n-2}), \\
Y_1 = y_2^{[n]} &= y_{n-1/2}, & Y_2 = y_1^{[n]} &= y_n, & hF_1 &= hf(y_{n-1/2}), & hF_2 &= hf(y_n),
\end{aligned}$$

and the method becomes

$$\begin{aligned}
Y_1 &= \frac{5}{33}hF_1 - \frac{8}{11}y_1^{[n-1]} + \frac{19}{11}y_2^{[n-1]} + \frac{19}{22}y_3^{[n-1]} + \frac{4}{11}y_4^{[n-1]} - \frac{1}{66}y_5^{[n-1]}, \\
Y_2 &= \frac{4753}{7920}hF_1 + \frac{251}{1440}hF_2 - \frac{8}{11}y_1^{[n-1]} + \frac{19}{11}y_2^{[n-1]} + \frac{449}{660}y_3^{[n-1]} + \frac{3463}{7920}y_4^{[n-1]} - \frac{449}{15840}y_5^{[n-1]}
\end{aligned}$$

or, using coefficient tableaux,

$$\begin{bmatrix} A & U \\ B & V \end{bmatrix} = \left[\begin{array}{cc|cccc} \frac{5}{33} & 0 & -\frac{8}{11} & \frac{19}{11} & \frac{19}{22} & \frac{4}{11} & -\frac{1}{66} \\ \frac{4753}{7920} & \frac{251}{1440} & -\frac{8}{11} & \frac{19}{11} & \frac{449}{660} & \frac{3463}{7920} & -\frac{449}{15840} \\ \hline 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 \end{array} \right].$$

Example of an Almost Runge–Kutta method

The following “ARK” method, introduced in [18] (Butcher, 1997), is intended to re-use past information in a special way, which makes its behaviour very similar to that of a classical Runge–Kutta method:

$$\begin{bmatrix} A & U \\ B & V \end{bmatrix} = \left[\begin{array}{cccc|ccc} 0 & 0 & 0 & 0 & 1 & 1 & \frac{1}{2} \\ \frac{1}{16} & 0 & 0 & 0 & 1 & \frac{7}{16} & \frac{1}{16} \\ -\frac{1}{16} & 1 & 0 & 0 & 1 & -\frac{7}{16} & -\frac{5}{16} \\ \frac{1}{6} & \frac{1}{3} & \frac{1}{3} & 0 & 1 & \frac{1}{6} & 0 \\ \hline \frac{1}{6} & \frac{1}{3} & \frac{1}{3} & 0 & 1 & \frac{1}{6} & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ -1 & \frac{4}{3} & -\frac{4}{3} & 2 & 0 & -1 & 0 \end{array} \right].$$

Transformations

Let T denote a non-singular $r \times r$ matrix. If the quantities evaluated in step number n are replaced by independent linear combinations

$$\hat{y}^{[n]} = (T^{-1} \otimes I)y^{[n]}, \quad y^{[n]} = (T \otimes I)\hat{y}^{[n]},$$

then (6.4 h) transforms to

$$\begin{bmatrix} Y \\ \hat{y}^{[n]} \end{bmatrix} = \begin{bmatrix} A \otimes I & (UT) \otimes I \\ (T^{-1}B) \otimes I & (T^{-1}VT) \otimes I \end{bmatrix} \begin{bmatrix} F \\ \hat{y}^{[n-1]} \end{bmatrix}.$$

That is, the transformation is

$$\begin{bmatrix} A & U \\ B & V \end{bmatrix} \mapsto \begin{bmatrix} A & UT \\ T^{-1}B & T^{-1}VT \end{bmatrix}.$$

In working with a specific method, it is sometimes convenient to transform it to a representation that helps with the understanding of the method or makes its analysis more convenient. For example, a diagonal form of the matrix V might be preferable, if this is possible.

Transformation of the Adams–Bashforth methods

Using

$$T = \begin{bmatrix} 1 & 0 & \frac{1}{2} \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix},$$

the method (6.4) transforms to

$$\left[\begin{array}{cc|cc} A & UT \\ T^{-1}B & T^{-1}VT \end{array} \right] = \left[\begin{array}{c|ccc} 0 & 1 & \frac{3}{2} & 0 \\ \hline 0 & 1 & 1 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{array} \right].$$

We see in this example that the method has been reduced to one with $r = 2$ because the third input is not used in generating the first or second output. Hence, the method can be written

$$\left[\begin{array}{cc|c} A & U \\ B & V \end{array} \right] = \left[\begin{array}{c|cc} 0 & 1 & \frac{3}{2} \\ \hline 0 & 1 & 1 \\ 1 & 0 & 0 \end{array} \right]. \quad (6.4 \text{ o})$$

By carrying out a further transformation, it is possible to diagonalize V :

$$\left[\begin{array}{cc|c} A & U \\ B & V \end{array} \right] \mapsto \left[\begin{array}{cc|cc} A & UT \\ T^{-1}B & T^{-1}VT \end{array} \right] = \left[\begin{array}{c|ccc} 0 & 1 & \frac{1}{2} \\ \hline 1 & 1 & 0 \\ 1 & 0 & 0 \end{array} \right]. \quad (6.4 \text{ p})$$

Exercise 56 Find the transformation required to convert the representation (6.4 o) to (6.4 p).

The order 3 Adams–Bashforth method (6.4 m) transforms, using

$$T = \begin{bmatrix} 1 & 0 & 0 & -\frac{5}{12} \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix},$$

to

$$\left[\begin{array}{cc|ccc} A & UT \\ T^{-1}B & T^{-1}VT \end{array} \right] = \left[\begin{array}{c|ccccc} 0 & 1 & \frac{2}{3} & -\frac{4}{3} & 0 \\ \hline 0 & 1 & \frac{2}{3} & -\frac{11}{12} & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{array} \right].$$

This simple transformation has converted this method so that r becomes 3, instead of 4, because the zeros in the last column indicate that the value of $\tilde{y}_4^{[n-1]}$ is not needed in the computation of step number n . Hence, we could represent the method as

$$\left[\begin{array}{c|ccc} 0 & 1 & \frac{2}{3} & -\frac{4}{3} \\ \hline 0 & 1 & \frac{2}{3} & -\frac{11}{12} \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{array} \right].$$

The reduction of r , by the use of a transformation, in the general linear representation of an Adams method, was considered in [22] (Butcher, Hill, 2006).

Backward differences

In the pre-computer days, the tedious process of solving differential equations by hand calculations was prone to error, and checks to recognize when something had gone wrong, were valuable. The well-established practice of compiling difference tables, to check for errors and to facilitate interpolation, can be incorporated into Adams–Bashforth methods, by using a modified formulation. For example, in the third order case, we could use first and second order differences of the f_n values instead of the f_n values themselves. That is, (6.4 l) could be rewritten as

$$y_n = y_{n-1} + hf(y_{n-1}) + \frac{1}{2}h(\nabla f)(y_{n-1}) + \frac{5}{12}h(\nabla^2 f)(y_{n-1}),$$

where

$$(\nabla f)(y_{n-1}) := f(y_{n-1}) - f(y_{n-2}),$$

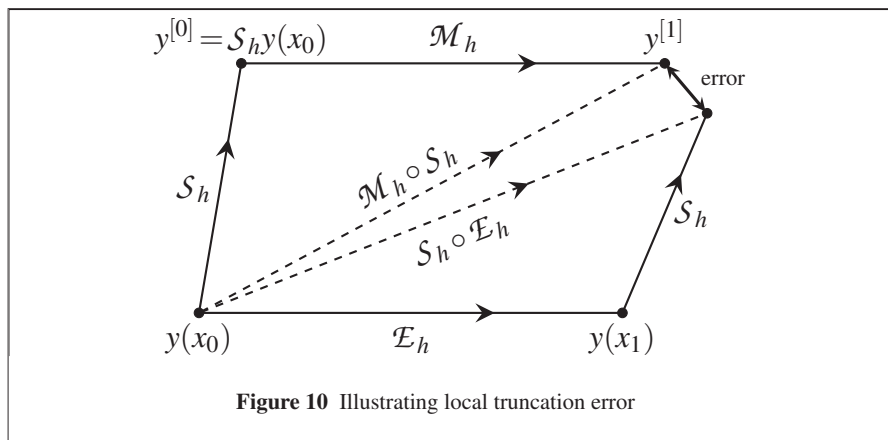
$$(\nabla^2 f)(y_{n-1}) := f(y_{n-1}) - 2f(y_{n-2}) + f(y_{n-3}).$$

To rewrite the general linear formulation (6.4 m) (p. 225), so that $y^{[n-1]}$ transforms to a vector with components y_{n-1} , $hf(y_{n-1})$, $h(\nabla f)(y_{n-1})$, $h(\nabla^2 f)(y_{n-1})$, transform according to

$$T = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 1 & -1 & 0 \\ 0 & 1 & -2 & 1 \end{bmatrix}, \quad \begin{bmatrix} A & UT \\ T^{-1}B & T^{-1}VT \end{bmatrix} = \begin{bmatrix} 0 & 1 & 1 & \frac{1}{2} & \frac{5}{12} \\ 0 & 1 & 1 & \frac{1}{2} & \frac{5}{12} \\ 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & -1 & 0 & 0 \\ 1 & 0 & -1 & -1 & 0 \end{bmatrix}.$$

The Nordsieck representation

For Adams–Bashforth methods in general, and for the third order method in particular, it was proposed in [76] (Nordsieck, 1962) and [44] (Gear, 1967) to use approximations to $y(x_{n-1})$ and to the scaled derivatives derivatives $hy'(x_{n-1})$, $\frac{1}{2!}h^2y''(x_{n-1})$, $\frac{1}{3!}h^3y^{(3)}(x_{n-1})$, \dots , rather than to the data in the original formulation, as input to step n . This has the advantage that a change of step-size can be carried out by a simple



rescaling by powers of h . In the general linear methods formulation the change to the new format is accomplished by the transformation:

$$T = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 1 & -2 & 3 \\ 0 & 1 & -4 & 12 \end{bmatrix}, \quad \begin{bmatrix} A & UT \\ T^{-1}B & T^{-1}VT \end{bmatrix} = \left[\begin{array}{c|cccc} 0 & 1 & 1 & 1 & 1 \\ \hline 0 & 1 & 1 & 1 & 1 \\ 1 & 0 & 0 & 0 & 0 \\ \frac{3}{4} & 0 & -\frac{3}{4} & -\frac{1}{2} & \frac{3}{4} \\ \frac{1}{6} & 0 & -\frac{1}{6} & -\frac{1}{3} & \frac{1}{2} \end{array} \right].$$

The idea of combining the features of Runge–Kutta methods with various other types of numerical integrators is an old one and examples of these “mixed” methods abound. We refer back to Section 1.6 (p. 28) for a first introduction to mixed or “general linear methods”.

6.5 Order of general linear methods

A general linear method operates, in each step, on r input values and, at the end of the step, exports the same number of approximations for use as input by the following step. In many cases the input vectors will have a natural interpretation, but we might need to avoid making such an assumption *a priori*. Let S_h denote the mapping from a given initial value y_0 and the input to the first step, $y^{[0]}$. Also let M_h be the mapping which moves the solution forward through a single step $y^{[0]} \mapsto y^{[1]}$ and E_h the flow of the solution through a time step h . That is, E_h is the mapping $y(x_0) \mapsto y(x_0 + h)$. The basic idea is based on Figure 10.

The “error” in this figure is assumed to be of order p . That is,

$$\text{error} := M_h(S_h(y_0)) - S_h(E_h(y_0)) = \mathcal{O}(h^{p+1}).$$

To convert this figure to a B-series formulation, introduce $\eta \in \mathbf{B}^s$ to represent the vector of stage values and $\eta D \in (\mathbf{B}^0)^s$ to represent the vector of stage derivatives and $\zeta \in \mathbf{B}^{*r}$ to represent the starting method. For order p , these quantities are connected by

$$\begin{aligned}\eta &= A\eta D + U\zeta, \\ E\zeta &= B\eta D + V\zeta + O_{p+1}.\end{aligned}$$

Starting and finishing methods

In addition to the starting method \mathcal{S}_h , we introduce a “finishing method” \mathcal{F}_h , which acts as a one sided inverse of \mathcal{S}_h so that $\mathcal{F}_h \circ \mathcal{S}_h = \text{id}$. The practical role of \mathcal{F}_h is to generate approximations to the solution to a given initial value problem after any desired number of steps. From a theoretical point of view, it gives an alternative interpretation of truncation error, in which Figure 11 is substituted for Figure 10.

An example of order analysis

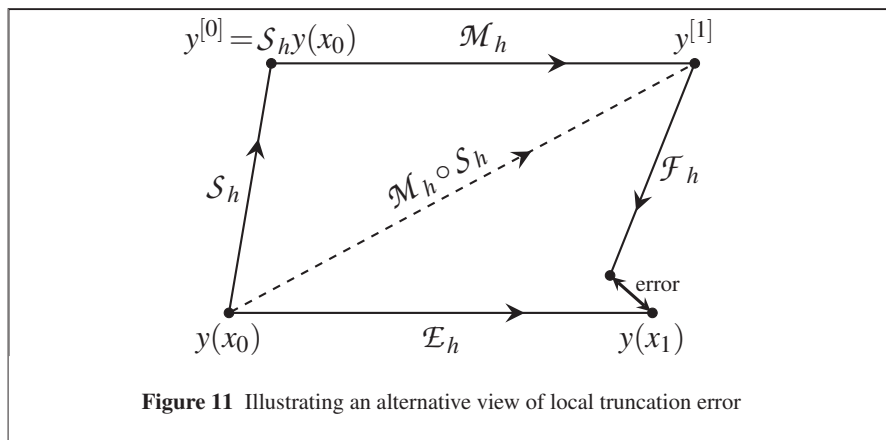
As an example of this analysis, consider the method (6.4 n) (p. 225). The conditions for order 5 are

$$\begin{aligned}\eta_1 &= \zeta_1, \\ \eta_2 &= \zeta_1 - \frac{1}{16}\zeta_2 + \frac{5}{16}\eta_1 D, \\ \eta_3 &= \zeta_1 - \frac{1}{4}\zeta_2 + \frac{3}{4}\eta_1 D, \\ \eta_4 &= \zeta_1 + \frac{3}{8}\zeta_2 - \frac{15}{16}\eta_1 D + \frac{3}{4}\eta_2 D + \frac{9}{16}\eta_3 D, \\ \eta_5 &= \zeta_1 - \zeta_2 + \frac{18}{7}\eta_1 D - \frac{12}{7}\eta_3 D + \frac{8}{7}\eta_4 D, \\ E\zeta_1 &= \zeta_1 + \frac{7}{90}\eta_1 D + \frac{16}{45}\eta_2 D + \frac{2}{15}\eta_3 D + \frac{16}{45}\eta_4 D + \frac{7}{90}\eta_5 D, \\ E\zeta_2 &= \eta_3 D.\end{aligned}$$

In terms of the tableau matrices, this is

$$\left[\begin{array}{cc} A & U \\ B & V \end{array} \right] = \left[\begin{array}{ccccc|cc} 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ \frac{5}{16} & 0 & 0 & 0 & 0 & 1 & -\frac{1}{16} \\ \frac{3}{4} & 0 & 0 & 0 & 0 & 1 & -\frac{1}{4} \\ -\frac{15}{16} & \frac{3}{4} & \frac{9}{16} & 0 & 0 & 1 & \frac{3}{8} \\ \frac{18}{7} & 0 & -\frac{12}{7} & \frac{8}{7} & 0 & 1 & -1 \\ \hline \frac{7}{90} & \frac{16}{45} & \frac{2}{15} & \frac{16}{45} & \frac{7}{90} & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 \end{array} \right].$$

The aim will be to obtain order 5 with the trivial starting method for the first component. That is, $\zeta_1 = 1$. We will, as a first step, attempt to interconnect the first three stages, with the assumption that



$$\zeta_2(\emptyset) = 0,$$

$$\zeta_2(\mathbf{t}_i) = \theta_i, \quad i = 1, 2, \dots, 8.$$

The values of θ_i need to satisfy

$$\zeta_2 = E^{-1}\eta_3 D + O_5, \quad (6.5 \text{ a})$$

as we see from the second row of B . It follows also from this row that $\theta_1 = 1$. The calculations will be shown in a tabular fashion.

	\emptyset	\cdot	\mathbf{t}	\mathbf{v}	\mathbf{i}	\mathbf{v}	\mathbf{v}	\mathbf{Y}	\mathbf{i}
η_1	1	0	0	0	0	0	0	0	0
$\eta_1 D$	0	1	0	0	0	0	0	0	0
η_2	1	$\frac{1}{4}$	$-\frac{1}{16}\theta_2$	$-\frac{1}{16}\theta_3$	$-\frac{1}{16}\theta_4$	$-\frac{1}{16}\theta_5$	$-\frac{1}{16}\theta_6$	$-\frac{1}{16}\theta_7$	$-\frac{1}{16}\theta_8$
$\eta_2 D$	0	1	$\frac{1}{4}$	$\frac{1}{16}$	$-\frac{1}{16}\theta_2$	$\frac{1}{64}$	$-\frac{1}{64}\theta_2$	$-\frac{1}{16}\theta_3$	$-\frac{1}{16}\theta_4$
η_3	1	$\frac{1}{2}$	$-\frac{1}{4}\theta_2$	$-\frac{1}{4}\theta_3$	$-\frac{1}{4}\theta_4$	$-\frac{1}{4}\theta_5$	$-\frac{1}{4}\theta_6$	$-\frac{1}{4}\theta_7$	$-\frac{1}{4}\theta_8$
$\eta_3 D$	0	1	$\frac{1}{2}$	$\frac{1}{4}$	$-\frac{1}{4}\theta_2$	$\frac{1}{8}$	$-\frac{1}{8}\theta_2$	$-\frac{1}{4}\theta_3$	$-\frac{1}{4}\theta_4$
$E^{-1}\eta_3 D$	0	1	$-\frac{1}{2}$	$\frac{1}{4}$	$-\frac{1}{4}\theta_2$	$-\frac{1}{8}$	$\frac{1}{8}\theta_2$	$\frac{1}{6} + \frac{1}{2}\theta_2 - \frac{1}{4}\theta_3$	$\frac{1}{12} + \frac{1}{4}\theta_2 - \frac{1}{4}\theta_4$

(6.5 b)

From (6.5 a), and the last row of (6.5 b), we find that $\theta_2 = -\frac{1}{2}$, $\theta_3 = \frac{1}{4}$, $\theta_4 = \frac{1}{8}$, $\theta_5 = -\frac{1}{8}$, $\theta_6 = -\frac{1}{16}$, $\theta_7 = -\frac{7}{48}$, $\theta_8 = -\frac{7}{96}$. We can now rewrite (6.5 b) with some additional data added.

	\emptyset	\cdot	\mathfrak{t}	\mathfrak{v}	\mathfrak{t}	\mathfrak{v}	\mathfrak{t}	\mathfrak{v}	\mathfrak{t}
ζ_2	0	1	$-\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{8}$	$-\frac{1}{8}$	$-\frac{1}{16}$	$-\frac{7}{96}$	$-\frac{7}{96}$
η_1	1	0	0	0	0	0	0	0	0
$\eta_1 D$	0	1	0	0	0	0	0	0	0
η_2	1	$\frac{1}{4}$	$\frac{1}{32}$	$-\frac{1}{64}$	$-\frac{1}{128}$	$\frac{1}{128}$	$\frac{1}{256}$	$\frac{7}{768}$	$\frac{7}{1536}$
$\eta_2 D$	0	1	$\frac{1}{4}$	$\frac{1}{16}$	$\frac{1}{32}$	$\frac{1}{64}$	$\frac{1}{128}$	$-\frac{1}{64}$	$-\frac{1}{128}$
η_3	1	$\frac{1}{2}$	$\frac{1}{8}$	$-\frac{1}{16}$	$-\frac{1}{32}$	$\frac{1}{32}$	$\frac{1}{64}$	$\frac{7}{192}$	$\frac{7}{384}$
$\eta_3 D$	0	1	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{8}$	$\frac{1}{16}$	$-\frac{1}{16}$	$-\frac{1}{32}$
η_4	1	$\frac{3}{4}$	$\frac{9}{32}$	$\frac{9}{32}$	$\frac{9}{64}$	$\frac{9}{256}$	$\frac{9}{512}$	$-\frac{13}{128}$	$-\frac{13}{256}$
$\eta_4 D$	0	1	$\frac{3}{4}$	$\frac{9}{16}$	$\frac{9}{32}$	$\frac{27}{64}$	$\frac{27}{128}$	$\frac{9}{32}$	$\frac{9}{64}$
η_5	1	1	$\frac{1}{2}$	$-\frac{1}{28}$	$-\frac{1}{56}$	$\frac{11}{28}$	$\frac{11}{56}$	$\frac{193}{336}$	$\frac{193}{672}$
$\eta_5 D$	0	1	1	1	$\frac{1}{2}$	1	$\frac{1}{2}$	$-\frac{1}{28}$	$-\frac{1}{56}$
$1 + \sum_{i=1}^5 b_i \eta_i D$	1	1	$\frac{1}{2}$	$\frac{1}{3}$	$\frac{1}{6}$	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{12}$	$\frac{1}{24}$
E	1	1	$\frac{1}{2}$	$\frac{1}{3}$	$\frac{1}{6}$	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{12}$	$\frac{1}{24}$
$E^{-1} \eta_3 D$	0	1	$-\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{8}$	$-\frac{1}{8}$	$-\frac{1}{16}$	$-\frac{7}{48}$	$-\frac{7}{96}$

(6.5 c)

and the fourth order conditions are verified. To explore the order conditions for order 5, we note that the $C(2)$ conditions are satisfied because $\xi(\emptyset) = 1$ and $\xi(\mathfrak{t}) = \frac{1}{2}\xi(\cdot)^2$, for $\xi = \eta_i$, $i = 1, 2, 3, 4, 5$. Hence we need only consider the trees $\cdot, \mathfrak{t}, \mathfrak{v}, \mathfrak{v}, \mathfrak{Y}, \mathfrak{v}, \mathfrak{Y}, \mathfrak{Y}, \mathfrak{Y}$. Reconstructing the information for these trees, but only as far as it involves ζ_1 , we obtain

	\emptyset	\cdot	\mathfrak{t}	\mathfrak{v}	\mathfrak{v}	\mathfrak{Y}	\mathfrak{v}	\mathfrak{Y}	\mathfrak{Y}
ζ_1	1	0	0	0	0	0	0	0	0
ζ_2	0	1	$-\frac{1}{2}$	$\frac{1}{4}$	$-\frac{1}{8}$	$-\frac{7}{48}$			
η_1	1	0	0	0	0	0			
$\eta_1 D$	0	1	0	0	0	0	0	0	0
η_2	1	$\frac{1}{4}$	$\frac{1}{32}$	$-\frac{1}{64}$	$\frac{1}{128}$	$\frac{7}{768}$			
$\eta_2 D$	0	1	$\frac{1}{4}$	$\frac{1}{16}$	$\frac{1}{64}$	$-\frac{1}{64}$	$\frac{1}{256}$	$-\frac{1}{256}$	$\frac{1}{128}$
η_3	1	$\frac{1}{2}$	$\frac{1}{8}$	$-\frac{1}{16}$	$\frac{1}{32}$	$\frac{7}{192}$			
$\eta_3 D$	0	1	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{8}$	$-\frac{1}{16}$	$\frac{1}{16}$	$-\frac{1}{32}$	$\frac{1}{32}$
η_4	1	$\frac{3}{4}$	$\frac{9}{32}$	$\frac{9}{32}$	$\frac{9}{256}$	$-\frac{13}{128}$			
$\eta_4 D$	0	1	$\frac{3}{4}$	$\frac{9}{16}$	$\frac{27}{64}$	$\frac{9}{32}$	$\frac{81}{256}$	$\frac{27}{256}$	$\frac{9}{512}$
η_5	1	1	$\frac{1}{2}$	$-\frac{1}{28}$	$\frac{11}{28}$	$\frac{193}{336}$			
$\eta_5 D$	0	1	1	1	1	$-\frac{1}{28}$	1	$-\frac{1}{28}$	$\frac{11}{28}$
$1 + \sum_{i=1}^5 b_i \eta_i D$	1	1	$\frac{1}{2}$	$\frac{1}{3}$	$\frac{1}{4}$	$\frac{1}{12}$	$\frac{1}{5}$	$\frac{1}{15}$	$\frac{1}{20}$
E	1	1	$\frac{1}{2}$	$\frac{1}{3}$	$\frac{1}{4}$	$\frac{1}{12}$	$\frac{1}{5}$	$\frac{1}{15}$	$\frac{1}{20}$

(6.5 d)

Starting method for ζ_2

Strictly speaking, the value of ζ_2 used in (6.5 d) does not satisfy the order 5 condition $E\zeta_2 = \eta_3 D + O_6$, but it is satisfactory to use a starting method based on ζ_2 without causing a loss of global order. Hence we will derive a generalized four stage Runge–Kutta method

$$\begin{array}{c|c} c & A \\ \hline 0 & b^T \end{array}, \quad (6.5 \text{ e})$$

satisfying

$$\Phi(t) = \theta(t), \quad |t| \leq 4. \quad (6.5 \text{ f})$$

Exercise 57 Show that for the four stage method (6.5 e). satisfying (6.5 f), $c_4 = -\frac{1}{2}$.

A possible solution of these conditions gives the starting method

$$\begin{array}{c|cccc} -\frac{7}{12} & -\frac{7}{12} & & & \\ -\frac{7}{6} & 0 & -\frac{7}{6} & & \\ -\frac{1}{2} & -\frac{11}{28} & 0 & -\frac{3}{28} & \\ \hline 0 & 0 & 0 & 0 & 1 \end{array}.$$

A generalization of a classical Runge–Kutta method

The idea of adding additional inputs to a Runge–Kutta method was foreshadowed in Section 6.4 (p. 225). As an example, we will use the ansatz (6.5 g), in an attempt to construct a $pqrs = 4333$ method, where $a_{21}, \dots, u_{22}, \dots, b_{31}, \dots, v_{32}$ are to be determined.

$$\begin{bmatrix} A & U \\ B & V \end{bmatrix} = \left[\begin{array}{ccc|ccc} 0 & 0 & 0 & 1 & \frac{1}{2} & \frac{1}{8} \\ a_{21} & 0 & 0 & 1 & u_{22} & u_{23} \\ \frac{2}{3} & \frac{1}{6} & 0 & 1 & \frac{1}{6} & 0 \\ \hline \frac{2}{3} & \frac{1}{6} & 0 & 1 & \frac{1}{6} & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ b_{31} & b_{32} & b_{33} & 0 & v_{32} & 0 \end{array} \right].$$

Use a similar notation to (6.5 d) but because we are deriving a $pqrs$ method with $q = 3$, we need only use a single third order tree \mathbf{v} and a single fourth order tree \mathbf{v} .

	\emptyset	\bullet	\mathbf{i}	\mathbf{v}	\mathbf{v}
ζ_1	1	0	0	0	0
ζ_2	0	1	0	0	0
ζ_3	0	0	1	$\frac{1}{3}$	θ
η_1	1	$\frac{1}{2}$	$\frac{1}{8}$	$\frac{1}{24}$	
$\eta_1 D$	0	1	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{8}$
$a_{21}\eta_1 D + \sum_{i=1}^2 u_{2i}\zeta_i$	1	$a_{21} + u_{22}$	$\frac{1}{2}a_{21} + u_{23}$	$\frac{1}{4}a_{21} + \frac{1}{3}u_{23}$	
η_2	1	1	$\frac{1}{2}$	$\frac{1}{3}$	
$\eta_2 D$	0	1	1	1	1
η_3	1	1	$\frac{1}{2}$	$\frac{1}{3}$	
$\eta_3 D$	0	1	1	1	1
$\sum_{i=1}^3 b_{3i}\eta_i D + v_{32}\zeta_2$	0	$b_{31} + b_{32} + b_{33} + v_{32}$	$\frac{1}{2}b_{31} + b_{32} + b_{33}$	$\frac{1}{4}b_{31} + b_{32} + b_{33}$	$\frac{1}{8}b_{31} + b_{32} + b_{33}$
$E\zeta_3$	0	0	1	$\frac{7}{3}$	$\theta + 4$

From $a_{21}\eta_1 D + \sum_{i=1}^2 u_{2i}\zeta_i = \eta_2 + O_4$, it is found that $a_{21} = 2$, $u_{21} = -1$, $u_{22} = -\frac{1}{2}$ and from $\sum_{i=1}^3 b_{3i}\eta_i D + v_{32}\zeta_2 = E\zeta_3 + O_5$, $\theta = -1$, $b_{31} = -\frac{16}{3}$, $b_{32} + b_{33} = \frac{11}{3}$, $v_{32} = \frac{5}{3}$. Making the arbitrary choice $b_{33} = 0$, we arrive at the 4333 method:

$$\begin{bmatrix} A & U \\ B & V \end{bmatrix} = \left[\begin{array}{ccc|ccc} 0 & 0 & 0 & 1 & \frac{1}{2} & \frac{1}{8} \\ 2 & 0 & 0 & 1 & -1 & -\frac{1}{2} \\ \hline \frac{2}{3} & \frac{1}{6} & 0 & 1 & \frac{1}{6} & 0 \\ \frac{2}{3} & \frac{1}{6} & 0 & 1 & \frac{1}{6} & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ -\frac{16}{3} & \frac{11}{3} & 0 & 0 & \frac{5}{3} & 0 \end{array} \right]. \quad (6.5g)$$

An example of method derivation

We will construct a method based on a specific design

$$\left[\begin{array}{cccc|cccc} 0 & 0 & 0 & 0 & 1 & u_1 & & \\ a_{21} & 0 & 0 & 0 & 1 & u_2 & & \\ a_{31} & a_{32} & 0 & 0 & 1 & u_3 & & \\ a_{41} & a_{42} & a_{43} & 0 & 1 & u_4 & & \\ \hline b_1 & b_2 & b_3 & b_4 & 1 & 0 & & \\ \beta_1 & \beta_2 & \beta_3 & \beta_4 & 0 & 0 & & \end{array} \right].$$

Our aim will be to obtain order 5, with the additional assumptions that the first component of the starting method is the identity mapping and that each of the stages has internal order 3.

1. First choose c_1, c_2, c_3, c_4 so that an order 5 quadrature formula $\int_0^1 \phi(x) dx \approx \sum_{i=1}^4 b_i \phi(c_i)$ is possible. The choice actually made was

$$c^\tau = \left[0 \quad \frac{3}{4} \quad \frac{3}{10} \quad 1 \right], \quad b^\tau = \left[\frac{5}{54} \quad \frac{32}{81} \quad \frac{250}{567} \quad \frac{1}{14} \right].$$

2. To obtain stage order 3 for the first stage, choose $u_1 = 0$. Without loss of generality choose $u_2 = 1$. Also choose $a_{21} = \frac{3}{4}$; this will have the effect of forcing $\sum_{i=1}^4 \beta_i$ to equal zero in the following item.
3. Solve for $u_3, a_{31}, a_{32}, a_{41}, a_{42}, a_{43}, \beta_1, \beta_2, \beta_3$, to ensure that the stage order is 3, with u_4, β_4 free parameters. Note that the conditions for these stage orders are equivalent to the requirements that the following quadrature formulae are exact for ϕ any polynomial of degree less than 3:

$$\begin{aligned} \int_0^{c_2} \phi(x) dx &\approx \beta_1 \phi(-1) + \beta_2 \phi(c_2 - 1) \\ &\quad + \beta_3 \phi(c_3 - 1) + (\beta_4 + a_{21}) \phi(0), \\ \int_0^{c_3} \phi(x) dx &\approx u_3 \beta_1 \phi(-1) + u_3 \beta_2 \phi(c_2 - 1) + u_3 \beta_3 \phi(c_3 - 1) \\ &\quad + (u_3 \beta_3 + a_{31}) \phi(0) + a_{32} \phi(c_2), \\ \int_0^1 \phi(x) dx &\approx u_4 \beta_1 \phi(-1) + u_4 \beta_2 \phi(c_2 - 1) + u_4 \beta_3 \phi(c_3 - 1) \\ &\quad + (u_4 \beta_3 + a_{41}) \phi(0) + a_{42} \phi(c_2) + a_{43} \phi(c_3). \end{aligned}$$

4. Impose an additional condition so that the remaining order condition is satisfied. This remaining condition is

$$b^\tau (Ac^3 + \beta^\tau (c-1)^3 u) = \frac{1}{60},$$

corresponding to t_{14} . This gives

$$(u_4 + \frac{208}{27})(\beta_4 - \frac{27}{7}) = 0.$$

We consider two cases:

- (i) $u_4 = -\frac{208}{27}$.
 - (ii) $\beta_4 = \frac{27}{7}$.
5. The remaining parameter is chosen so that $\beta^\tau u = 0$. This condition is considered advantageous because, at least for small h , a small perturbation of $y_2^{[n-1]}$ will have a small effect on $y_2^{[n]}$. The solutions are: Case (i): $\beta_4 = \frac{27}{448}$. Case (ii): $u_4 = \frac{691}{540}$.

The methods found using these steps are

$$\begin{aligned}
 \text{(i)} \quad \begin{bmatrix} A & U \\ B & V \end{bmatrix} &= \left[\begin{array}{cccc|cc} 0 & 0 & 0 & 0 & 1 & 0 \\ \frac{3}{4} & 0 & 0 & 0 & 1 & 1 \\ \frac{93}{250} - \frac{9}{125} & 0 & 0 & 0 & 1 & \frac{44}{125} \\ -\frac{139}{27} & \frac{148}{81} & \frac{350}{81} & 0 & 1 & -\frac{208}{27} \\ \hline \frac{5}{54} & \frac{32}{81} & \frac{250}{567} & \frac{1}{14} & 1 & 0 \\ \frac{113}{64} & \frac{41}{24} & -\frac{2375}{672} & \frac{27}{448} & 0 & 0 \end{array} \right], \\
 \text{(ii)} \quad \begin{bmatrix} A & U \\ B & V \end{bmatrix} &= \left[\begin{array}{cccc|cc} 0 & 0 & 0 & 0 & 1 & 0 \\ \frac{3}{4} & 0 & 0 & 0 & 1 & 1 \\ \frac{93}{250} - \frac{9}{125} & 0 & 0 & 0 & 1 & \frac{44}{125} \\ \frac{8881}{8640} & \frac{1069}{3240} - \frac{1855}{5184} & 0 & 0 & 1 & \frac{691}{540} \\ \hline \frac{5}{54} & \frac{32}{81} & \frac{250}{567} & \frac{1}{14} & 1 & 0 \\ -\frac{19}{16} & -\frac{37}{6} & \frac{1175}{336} & \frac{27}{7} & 0 & 0 \end{array} \right]. \quad (6.5h)
 \end{aligned}$$

Constructing a starting method

We will confine our attention to Case (ii), given by (6.5 h).

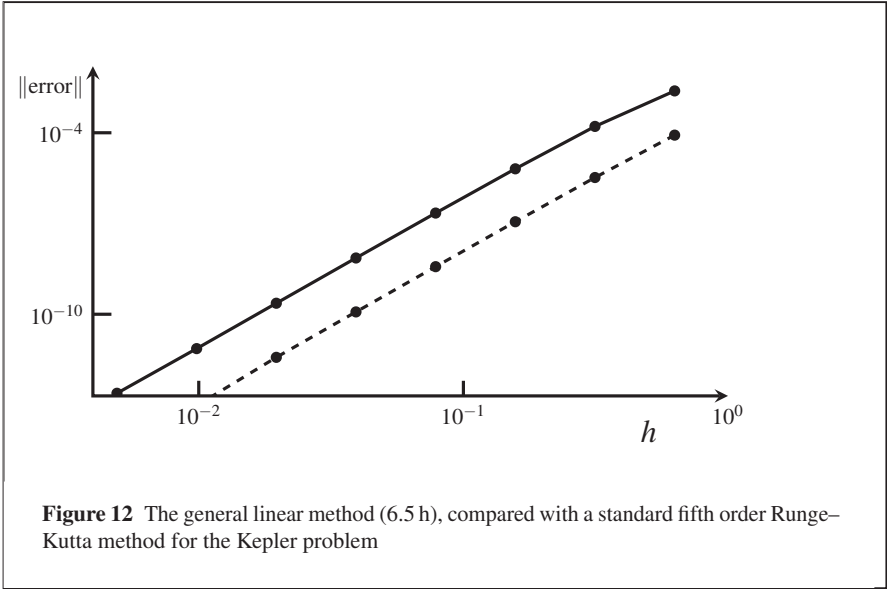
To get the right-hand sides for the starting order conditions, we need a Runge–Kutta method (with y_0 deleted)

$$\begin{array}{c|ccc}
 0 & & & \\
 \tilde{c}_2 & \tilde{a}_{21} & & \\
 \tilde{c}_3 & \tilde{a}_{31} & \tilde{a}_{32} & \\
 \hline
 0 & \tilde{b}_1 & \tilde{b}_2 & \tilde{b}_3
 \end{array},$$

with

$$\begin{aligned}
 \tilde{b}_1 + \tilde{b}_2 + \tilde{b}_3 &= \beta^\tau 1 = 0, \\
 \tilde{b}_2 \tilde{c}_2 + \tilde{b}_3 \tilde{c}_3 &= \beta^\tau (c-1) = \frac{9}{32}, \\
 \tilde{b}_2 \tilde{c}_2^2 + \tilde{b}_3 \tilde{c}_3^2 &= \beta^\tau (c-1)^2 = \frac{9}{64}, \\
 \tilde{b}_3 \tilde{a}_{32} \tilde{c}_2 &= \frac{1}{2} \beta^\tau (c-1)^2 = \frac{9}{128},
 \end{aligned}$$

with possible solution



0			
$\frac{1}{4}$	$\frac{1}{4}$		
$\frac{1}{2}$	0	$\frac{1}{2}$	
0	$-\frac{9}{16}$	0	$\frac{9}{16}$

Recursive evaluation of starting methods

Let R_h be a given starting method for the non-principal values. Calculate $y^{[0]}$. $\hat{y}^{[0]} = R_h y_0$. Use the method to find $y^{[1]}$. Evaluate $\hat{y}^{[1]}$. Evaluate $R_h y_1^{[1]}$. Then evaluate $(I - V)^{-1}(\hat{y}^{[1]} - R_h y_1^{[1]})$. Add this to $\hat{y}^{[0]}$ to get R_h^+ .

As a numerical test for (6.5 h), the Kepler problem with eccentricity $e = 0.1$ was solved over a half period. The results compared with a standard fifth order Runge–Kutta method (shown with dashed lines) are presented in Figure 12.

6.6 An algorithm for determining order

In this section we will confine our attention to methods which can be written using partitioned matrices:

$$\begin{bmatrix} A & U \\ B & V \end{bmatrix} = \left[\begin{array}{c|cc} A & \mathbf{1} & \widehat{U} \\ \hline b^\top & 1 & \mathbf{0}^\top \\ \hline \widehat{B} & \mathbf{0} & \widehat{V} \end{array} \right], \quad (6.6a)$$

where $1 \notin \sigma(\widehat{V})$.

The invariant subspace and the underlying one-step method

Following [64] (Kirchgraber, 1986), and [87] (Stoffer, 1993), we consider an approach to order of accuracy in which the “error” in Figure 10 (p. 231) is eliminated but, at the same time, \mathcal{E}_h is replaced by an approximation with \mathcal{S}_h replaced by the mapping \mathcal{S}_h^* . We now have

$$\mathcal{M}_h(\mathcal{S}_h^*(y_0)) = \mathcal{S}_h^*(\mathcal{E}_h^*(y_0)), \quad (6.6b)$$

Associated with the underlying one step method is “the invariant subspace”, see [87]. In this brief introduction to these important concepts, our aim will be limited to finding B-series for \mathcal{S}_h^* and \mathcal{E}_h^* .

Transformations

Let $\mathcal{T}_h : \mathbb{R}^N \rightarrow \mathbb{R}^N$ be a central mapping then

$$\mathcal{M}_h(\widetilde{\mathcal{S}}_h(\widetilde{y}_0)) = \widetilde{\mathcal{S}}_h(\widetilde{\mathcal{E}}_h(\widetilde{y}_0)),$$

where $\widetilde{y}_0 = \mathcal{T}_h^{-1}y_0$ and

$$\begin{aligned} \mathcal{S}_h^* &= \widetilde{\mathcal{S}}_h \circ \mathcal{T}_h, & \widetilde{\mathcal{S}}_h &= \mathcal{S}_h^* \circ \mathcal{T}_h^{-1}, \\ \mathcal{E}_h^* &= \mathcal{T}_h^{-1} \circ \widetilde{\mathcal{E}}_h \circ \mathcal{T}_h, & \widetilde{\mathcal{E}}_h &= \mathcal{T}_h \circ \mathcal{E}_h^* \circ \mathcal{T}_h^{-1}. \end{aligned}$$

This can be verified by substitution into (6.6b).

If \mathcal{S}_h^* is partitioned as

$$\mathcal{S}_h^* = \begin{bmatrix} (\mathcal{S}_I^*)_h \\ \widehat{\mathcal{S}}_h^* \end{bmatrix},$$

a convenient choice of \mathcal{T}_h is

$$\mathcal{T}_h = (\mathcal{S}_I^*)_h^{-1}$$

because the first component of $\widetilde{\mathcal{S}}_h$ will be the identity mapping and the analysis is simplified.

The quantities involved in the calculation of a single step, with input $y^{[0]}$, and output $y^{[1]}$, followed by the corresponding B-series quantities, are

$$\begin{aligned}
y^{[0]} &= \begin{bmatrix} (\mathcal{S}_I^*)_h(y_0) \\ \widehat{\mathcal{S}}^*_{*h}(y_0) \end{bmatrix}, & y^{[1]} &= \begin{bmatrix} (\mathcal{E}^*_{*h}(\mathcal{S}_I^*)_h)(y_0) \\ (\mathcal{E}^*_{*h}\widehat{\mathcal{S}}^*_{*h})(y_0) \end{bmatrix}, & Y, & hF, \\
\zeta &= \begin{bmatrix} \zeta_1 \\ \widehat{\zeta} \end{bmatrix}, & E^*\zeta &= \begin{bmatrix} E^*\zeta_1 \\ E^*\widehat{\zeta} \end{bmatrix}, & \eta^*, & \eta^*D,
\end{aligned}$$

The calculation of the single step uses the formulae

$$\begin{aligned}
Y &= hAF + \mathbf{1}y_1^{[0]} + \widehat{U}\widehat{y}^{[0]}, \\
y_1^{[1]} &= hb^T F + y_1^{[0]}, \\
\widehat{y}^{[1]} &= hB^T F + \widehat{V}\widehat{y}^{[0]},
\end{aligned}$$

corresponding to the B-series relations,

$$\begin{aligned}
\eta^* &= A(\eta^*D) + \mathbf{1}\zeta_1 + \widehat{U}\widehat{\zeta}, \\
E^*\zeta_1 &= b^T(\eta^*D) + \zeta_1, \\
E^*\widehat{\zeta} &= \widehat{B}(\eta^*D) + \widehat{V}\widehat{\zeta}.
\end{aligned} \tag{6.6 c}$$

Substitute $\eta = \zeta_1\eta$, $\zeta^* = \zeta_1\xi$ and it follows that

$$\begin{aligned}
\eta &= A(\eta D) + \mathbf{1} + \widehat{U}\xi, \\
\widetilde{E} &= b^T(\eta D) + 1, \\
\widetilde{E}\xi &= \widehat{B}(\eta D) + \widehat{V}\xi,
\end{aligned}$$

where $\widetilde{E} = \xi^{-1}E^*\xi$.

To obtain tree-by-tree formulae for η , \widetilde{E} and ξ , start with $\eta(\varnothing) = \mathbf{1}$, $\widetilde{E}(\varnothing) = 1$, $\xi(\varnothing) = \mathbf{0}$. Then for $t = [t_1 t_2 \cdots t_m]$, we find

$$\begin{aligned}
\eta(t) &= A(\eta D)(t) + \widehat{U}\xi(t), \\
\widetilde{E}(t) &= b^T(\eta D)(t), \\
\xi(t) &= (I - \widehat{V})^{-1}(\widehat{B}(\eta D)(t) - \sum_{t' < t} \widetilde{E}(t \setminus t')\xi(t')),
\end{aligned} \tag{6.6 d}$$

with $(\eta D)(t) = \prod_{i=1}^m \eta(t_i)$ and with exponentiations taking precedence over other operations. The details of these calculations, to order 5, are shown in Table 18 (p. 242).

Test for conjugacy

Having calculated \widetilde{E} , the final step in the test for order p is to determine if ξ exists such that

Table 18 Details of (6.6 d)					
n	\mathbf{t}_n	$ \mathbf{t}_n $	$\tilde{\mathbf{E}}_n$	ξ_n	η_n
0	\emptyset	0	1	$\mathbf{0}$	$\mathbf{1}$
1	\bullet	1	$b^T \mathbf{1}$	$(I - \widehat{V})^{-1}(\widehat{B}\mathbf{1})$	$A\mathbf{1} + \widehat{U}\xi_1$
2	$\mathbf{1}$	2	$b^T \eta_1$	$(I - \widehat{V})^{-1}(\widehat{B}\eta_1 - \mathbf{E}_1 \xi_1)$	$A\eta_1 + \widehat{U}\xi_2$
3	\mathbf{v}	3	$b^T \eta_1^2$	$(I - \widehat{V})^{-1}(\widehat{B}\eta_1^2 - \tilde{\mathbf{E}}_1^2 \xi_1 - 2\tilde{\mathbf{E}}_1 \xi_2)$	$A\eta_1^2 + \widehat{U}\xi_3$
4	$\mathbf{1}$	3	$b^T \eta_2$	$(I - \widehat{V})^{-1}(\widehat{B}\eta_2 - \tilde{\mathbf{E}}_2 \xi_1 - \tilde{\mathbf{E}}_1 \xi_2)$	$A\eta_2 + \widehat{U}\xi_4$
5	\mathbf{v}	4	$b^T \eta_1^3$	$(I - \widehat{V})^{-1}(\widehat{B}\eta_1^3 - \tilde{\mathbf{E}}_1^3 \xi_1 - 3\tilde{\mathbf{E}}_1^2 \xi_2 - 3\tilde{\mathbf{E}}_1 \xi_3)$	$A\eta_1^3 + \widehat{U}\xi_5$
6	\mathbf{v}	4	$b^T \eta_1 \eta_2$	$(I - \widehat{V})^{-1}(\widehat{B}\eta_1 \eta_2 - \tilde{\mathbf{E}}_1 \tilde{\mathbf{E}}_2 \xi_1 - (\tilde{\mathbf{E}}_1^2 + \tilde{\mathbf{E}}_2) \xi_2 - \tilde{\mathbf{E}}_1 \xi_3 - \tilde{\mathbf{E}}_1 \xi_4)$	$A\eta_1 \eta_2 + \widehat{U}\xi_6$
7	\mathbf{Y}	4	$b^T \eta_3$	$(I - \widehat{V})^{-1}(\widehat{B}\eta_3 - \tilde{\mathbf{E}}_3 \xi_1 - \tilde{\mathbf{E}}_1^2 \xi_2 - 2\tilde{\mathbf{E}}_1 \xi_4)$	$A\eta_3 + \widehat{U}\xi_7$
8	$\mathbf{1}$	4	$b^T \eta_4$	$(I - \widehat{V})^{-1}(\widehat{B}\eta_4 - \tilde{\mathbf{E}}_4 \xi_1 - \tilde{\mathbf{E}}_2 \xi_2 - \tilde{\mathbf{E}}_1 \xi_4)$	$A\eta_4 + \widehat{U}\xi_8$
9	\mathbf{v}	5	$b^T \eta_1^4$	$(I - \widehat{V})^{-1}(\widehat{B}\eta_1^4 - \tilde{\mathbf{E}}_1^4 \xi_1 - 4\tilde{\mathbf{E}}_2^3 \xi_2 - 6\tilde{\mathbf{E}}_1^2 \xi_3 - 4\tilde{\mathbf{E}}_1 \xi_5)$	$A\eta_1^4 + \widehat{U}\xi_9$
10	\mathbf{v}	5	$b^T \eta_1^2 \eta_2$	$(I - \widehat{V})^{-1}(\widehat{B}\eta_1^2 \eta_2 - \tilde{\mathbf{E}}_1^2 \tilde{\mathbf{E}}_2 \xi_1 - (\tilde{\mathbf{E}}_1^3 + 2\tilde{\mathbf{E}}_1^2) \xi_2 - (2\tilde{\mathbf{E}}_1^2 + \tilde{\mathbf{E}}_2) \xi_3 - \tilde{\mathbf{E}}_1^2 \xi_4 - \tilde{\mathbf{E}}_1 \xi_5 - 2\tilde{\mathbf{E}}_1 \xi_6)$	$A\eta_1^2 \eta_2 + \widehat{U}\xi_{10}$
11	\mathbf{v}	5	$b^T \eta_1 \eta_3$	$(I - \widehat{V})^{-1}(\widehat{B}\eta_1 \eta_3 - \tilde{\mathbf{E}}_1 \tilde{\mathbf{E}}_3 \xi_1 - (\tilde{\mathbf{E}}_1^3 + \tilde{\mathbf{E}}_3) \xi_2 - \tilde{\mathbf{E}}_1^2 \xi_3 - 2\tilde{\mathbf{E}}_1^2 \xi_4 - 2\tilde{\mathbf{E}}_1 \xi_6 - \tilde{\mathbf{E}}_1 \xi_7)$	$A\eta_1 \eta_3 + \widehat{U}\xi_{11}$
12	\mathbf{v}	5	$b^T \eta_1 \eta_4$	$(I - \widehat{V})^{-1}(\widehat{B}\eta_1 \eta_4 - \tilde{\mathbf{E}}_1 \tilde{\mathbf{E}}_4 \xi_1 - (\tilde{\mathbf{E}}_1 \tilde{\mathbf{E}}_2 + \tilde{\mathbf{E}}_4) \xi_2 - \tilde{\mathbf{E}}_2 \xi_3 - \tilde{\mathbf{E}}_1^2 \xi_4 - \tilde{\mathbf{E}}_1 \xi_6 - \tilde{\mathbf{E}}_1 \xi_8)$	$A\eta_1 \eta_4 + \widehat{U}\xi_{12}$
13	\mathbf{v}	5	$b^T \eta_2^2$	$(I - \widehat{V})^{-1}(\widehat{B}\eta_2^2 - \tilde{\mathbf{E}}_2^2 \xi_1 - 2\tilde{\mathbf{E}}_1 \tilde{\mathbf{E}}_2 \xi_2 - \tilde{\mathbf{E}}_1^2 \xi_3 - 2\tilde{\mathbf{E}}_1 \xi_4 - 2\tilde{\mathbf{E}}_1 \xi_6)$	$A\eta_2^2 + \widehat{U}\xi_{13}$
14	\mathbf{Y}	5	$b^T \eta_5$	$(I - \widehat{V})^{-1}(\widehat{B}\eta_5 - \tilde{\mathbf{E}}_5 \xi_1 - \tilde{\mathbf{E}}_1^3 \xi_2 - 3\tilde{\mathbf{E}}_1^2 \xi_4 - 3\tilde{\mathbf{E}}_1 \xi_7)$	$A\eta_5 + \widehat{U}\xi_{14}$
15	\mathbf{Y}	5	$b^T \eta_6$	$(I - \widehat{V})^{-1}(\widehat{B}\eta_6 - \tilde{\mathbf{E}}_6 \xi_1 - \tilde{\mathbf{E}}_1 \tilde{\mathbf{E}}_2 \xi_2 - (\tilde{\mathbf{E}}_1^2 + \tilde{\mathbf{E}}_2) \xi_4 - \tilde{\mathbf{E}}_1 \xi_7 - \tilde{\mathbf{E}}_1 \xi_8)$	$A\eta_6 + \widehat{U}\xi_{15}$
16	\mathbf{Y}	5	$b^T \eta_7$	$(I - \widehat{V})^{-1}(\widehat{B}\eta_7 - \tilde{\mathbf{E}}_7 \xi_1 - 4\tilde{\mathbf{E}}_3 \xi_2 - \tilde{\mathbf{E}}_1^2 \xi_4 - 2\tilde{\mathbf{E}}_1 \xi_8)$	$A\eta_7 + \widehat{U}\xi_{16}$
17	$\mathbf{1}$	5	$b^T \eta_8$	$(I - \widehat{V})^{-1}(\widehat{B}\eta_8 - \tilde{\mathbf{E}}_8 \xi_1 - \tilde{\mathbf{E}}_4 \xi_2 - \tilde{\mathbf{E}}_2 \xi_4 - \tilde{\mathbf{E}}_1 \xi_8)$	$A\eta_8 + \widehat{U}\xi_{17}$

$$\tilde{E} = \xi^{-1} E \xi,$$

to order p ; that is, that $(\xi \tilde{E})(t) = (E \xi)(t)$ for all t such that $|t| \leq p$. Write $\tilde{E} = a$, $E = b$, $\xi = x$.

Test that x exists such that $xa = bx + \mathcal{O}_{p+1}$

Given two members $a, b, \in \mathbf{B}$ we will consider the question “For a given integer p , does there exist $x \in \mathbf{B}$ such that $xax^{-1} = b + \mathcal{O}_{p+1}$?” To avoid uninteresting complications, we will assume that $b(\tau) \neq 0$.

Conjugacy to order 4

Using the notation based on $x_i = x(\mathbf{t}_i)$, we can write down the conditions $(xa)_i - x_i = (bx)_i - x_i$ for $i \leq 4$ (where the term x_i is subtracted from each side for convenience):

$$\begin{aligned} a_1 &= b_1, \\ x_1 a_1 + a_2 &= b_2 + b_1 x_1, \\ x_1 a_2 + x_2 a_1 + a_4 &= b_4 + b_2 x_1 + b_1 x_2, \\ a_2 x_1^2 + a_1 x_3 + 2a_4 x_1 + a_7 &= b_7 + b_1^2 x_2 + 2b_1 x_4 + b_3 x_1, \\ a_1 x_4 + a_2 x_2 + a_4 x_1 + a_8 &= b_8 + b_1 x_4 + b_2 x_2 + b_4 x_1. \end{aligned} \quad (6.6 \text{ e})$$

$$\begin{aligned} x_1^2 a_1 + 2x_1 a_2 + a_3 &= b_3 + b_1^2 x_1 + 2b_1 x_2, \\ a_1 x_1^3 + 3a_2 x_1^2 + 3a_3 x_1 + a_5 &= b_5 + b_1^3 x_1 + 3b_1^2 x_2 + 3b_1 x_3 + x_5, \\ x_1 a_1 x_2 + a_2 (x_1^2 + x_2) + x_1 a_3 + a_4 x_1 + a_6 &= b_6 + b_1 x_1 b_2 + x_2 (b_1^2 + b_2) + b_1 x_3 + x_4 b_1. \end{aligned} \quad (6.6 \text{ f})$$

From the equations (6.6 e), solve for a_1, a_2, a_4, a_7, a_8 ; and for the equations in (6.6 f), solve for x_2, x_3, x_4 , to obtain

$$\begin{aligned} x_2 &= \frac{-b_1^2 x_1 + b_1 x_1^2 + 2b_2 x_1 + a_3 - b_3}{2b_1}, \\ x_3 &= \frac{b_1^3 x_1 - 3b_1^2 x_1^2 - 2b_1 x_1^3 - 6b_1 b_2 x_1 + 6b_2 x_1^2 - 3a_3 b_1 + 6a_3 x_1 + 3b_1 b_3 + 2a_5 - 2b_5}{6b_1}, \\ x_4 &= \frac{2b_1^3 x_1 - 3b_1^2 x_1^2 + b_1 x_1^3 - 6b_1 b_2 x_1 + 6b_2 x_1^2 + 3a_3 x_1 - 3b_3 x_1 + 6b_4 x_1 - 2a_5 + 6a_6 + 2b_5 - 6b_6}{6b_1}. \end{aligned} \quad (6.6 \text{ g})$$

From these equations, we can summarize the conditions for conjugacy to order $p \leq 4$:

$$\begin{aligned} a_1 &= b_1, \\ a_2 &= b_2, \\ a_4 &= b_4, \\ -a_1 a_3 + a_5 - 2a_6 + a_7 &= b_7 - b_1 b_3 + b_5 - 2b_6, \\ a_8 &= b_8. \end{aligned} \quad (6.6 \text{ h})$$

The role of x_1

In the derivation of the order 4 conjugacy conditions, the transformation parameters $x(\mathbf{t})$, $|\mathbf{t}| \leq 3$, were allowed to take on arbitrary values. However, $x(\tau) = x_1$ does not seem to be constrained in any way even though (6.6 h) might have been expected to depend on x_1 . It would have been a simpler computation to take this quantity to be

zero so that

$$\begin{aligned} x_1 &= 0, \\ x_2 &= \frac{a_3 - b_3}{2b_1}, \\ x_3 &= \frac{b_3 - a_3}{2} + \frac{a_5 - b_5}{3b_1}, \\ x_4 &= \frac{a_6 - b_6}{b_1} + \frac{b_5 - a_5}{3b_1}. \end{aligned}$$

The fact that the conjugate order conditions do not depend on x_1 , for any order, follows from the observation that if $xa \sim bx$ then $x'a \sim bx'$, where $x' = bx$ and, if the order conditions are polynomials in x_1 then the value of these polynomials are unchanged if x_1 is replaced by $x'_1 = x_1 + b_1$, which is impossible if $b_1 \neq 0$.

Order by order conjugation

We now return to the case of arbitrary p . For $|t| = 1$, we have the single necessary condition $x(\tau) + a(\tau) = b(\tau) + x(\tau)$, implying

$$a(\tau) = b(\tau).$$

Before moving on to higher orders, we remark that the value of $x(\tau)$ is irrelevant because $a = x^{-1}bx$ is equivalent to $a = \hat{x}^{-1}b\hat{x}$, where $\hat{x} = b^\theta x$, for any real θ , so that $\hat{x}(\tau) = x(\tau) + \theta$.

For $t = [\tau]$, we have

$$x(t) + x(\tau) + a(t) = b(t) + x(\tau) + x(t),$$

so that

$$a([\tau]) = b([\tau]).$$

For $p > 2$, we will assume that the order is already known to be at least equal to $p - 1$.

Divide the trees of order p into two subsets

$$S_1 \quad t = [\tau^n t_1 \cdots t_m], \quad n \geq 1, \quad \tau \notin \{t_1, \dots, t_m\}$$

S_2 the remaining trees of order p .

For $t \in S_1$, let $t' = [\tau^{n-1} t_1 \cdots t_m]$. Then the order condition becomes

$$nb(\tau)x(t') + \sum_{t'' \leq t, |t''| < |t|} b(t \setminus t'')x(t'') = \sum_{t'' \leq t} x(t \setminus t'')a(t'').$$

This defines $x(t')$ for all t' of order $p - 1$.

For $t \in S_2$, using the known values of x for all orders less than p , we obtain an order condition for each t in this set.

Algorithm 7 Determine whether two central B-series for which $\tau \nrightarrow 0$ are conjugate

```

Input:     $a, b \in \mathbb{B}, a = b + \mathcal{O}_p$ 
Output:  if  $\exists x(a' := xax^{-1} = b + \mathcal{O}_{p+1})$  then [true,  $a', x$ ] else [false]
%
%     $\tilde{T}_p := T_p^\equiv \setminus (T_{p-1}^\equiv * \tau)$ 
%     $\mu(t) := m$ , where  $t = [\tau^{m-1}t_1t_2 \cdots t_k], |t_i| > 1, i = 1, 2, \dots, k$ 
%
1  $x \leftarrow 1 + \mathcal{O}_p$ 
2 for  $t \in T_{p-1}^\equiv$  do
3    $x(t) \leftarrow (a(t * \tau) - b(t * \tau)) / \mu(t)b(\tau)$ 
4 end for
5  $TEST \leftarrow \text{true}$ 
6 for  $t \in T_p^\equiv \setminus (T_{p-1}^\equiv * \tau)$  do
7   if  $(xa)(t) \neq (bx)(t)$  then
8      $TEST \leftarrow \text{false}$ 
9   end if
10 end for
11 if  $TEST$  then
12   return [true,  $xax^{-1}$ ]
13 else
14   return [false]
15 end if

```

Summary of Chapter 6 and the way forward

Summary

Multivalued methods have a long history in the form of linear multistep methods. In this chapter, an amalgam of multivalued and multistage (Runge–Kutta) methods is considered as a family of method, in its own right and given the name “general linear methods”.

After a review of linear multistep methods, the prototypical multivalued methods, it is shown by example that new methods flow from these by allowing multiple vector field calculations. Similarly, Runge–Kutta methods, the prototypical one-step methods, are also simply examples of known, and not so well known, multistage multivalued methods.

The insight provided by this wide range of example methods underlines the use of the natural and highly flexible general linear formulation. The fundamental questions of consistency, order, and convergence, take on a natural and straightforward meaning in the general context.

The theory of order for these methods is an important application of B-series analysis. This is closely related to the existence of the underlying one-step method and the theory of invariant subspaces.

Throughout the chapter new methods are introduced and analysed.

The way forward

The B-series approach is adapted to structure-preserving algorithms, as exemplified in the cases of symplectic preservation and energy-preservation, in Chapter 7.

Teaching and study notes

The following is a sample of the many publications on multivalued and general linear methods.

Burrage K. and Butcher, J.C. *Non-linear stability of a general class of differential equation methods* (1980) [6]

Butcher, J.C. *On the convergence of numerical solutions to ordinary differential equations* (1966) [12]

Butcher, J.C. *General linear methods* (2006) [19]

D'Ambrosio, R. and Hairer, E. *Long-term stability of multi-value methods for ordinary differential equations* (2014) [40]

Hairer, E. and Wanner, G. *On the Butcher group and general multi-value methods* (1974) [52]

Jackiewicz, Z. *General Linear Methods for Ordinary Differential Equations* (2009) [63]

Stoffer, D. *General linear methods: connection to one step methods and invariant curves* (1993) [87]

Projects

Project 22 Read up about predictor-corrector methods.

Project 23 Find criteria for the matrix V in a method (A, U, B, V) to be stable.

Project 24 Read up about the underlying one-step method and invariant sub-spaces, starting with [87].

Project 25 Learn about “Peer methods”, as a special class of general linear methods.