

(Preprint for)

J. ÁLVAREZ AND J. ROJO

*An improved class of generalized  
Runge-Kutta methods for stiff problems*  
Technical Report, (2000)

# An improved class of generalized Runge-Kutta methods for stiff problems

Jorge Alvarez<sup>a,1</sup> and Jesús Rojo<sup>a,2</sup>

<sup>a</sup> *Departamento de Matemática Aplicada a la Ingeniería, E.T.S. de Ingenieros Industriales, Universidad de Valladolid, E-47011 Valladolid, Spain*  
*E-mail: joralv@wmatem.eis.uva.es, jesroj@wmatem.eis.uva.es*

---

## Abstract

A new family of explicit  $p$ -stage methods for the numerical integration of scalar autonomous ODEs is proposed. These methods can be seen as a generalization of the explicit  $p$ -stage Runge-Kutta ones, while providing better order and stability results. In fact, we show that it is possible to obtain A-stable and L-stable formulas of order three and five, with only two and three evaluations per step respectively, and without losing the explicitness of the formulas. It is also possible to generalize the methods to get formulas for some non-autonomous scalar ODEs and systems. We obtain linearly implicit A-stable methods which do not require Jacobian evaluations. Some numerical examples are discussed in order to show the good performance of the new schemes.

*Keywords:* Generalized Runge-Kutta Methods, Stiff ODEs, Linear Stability, Numerical Experiments.

---

## 1 Introduction.

During the last decades there has been a considerable amount of research on methods for numerical integration of stiff systems of ODEs, usually looking for better stability properties. Nearly all such methods are implicit in character.

---

<sup>1</sup> This work was partially supported by Consejería de Educación y Cultura de la Junta de Castilla y León and Unión Europea (Fondo Social Europeo) under project VA19/00B and by Junta de Castilla y León under project VA11/99.

<sup>2</sup> This work was partially supported by Consejería de Educación y Cultura de la Junta de Castilla y León and Unión Europea (Fondo Social Europeo) under project VA19/00B and by Spanish DGES under project PB98-0359.

The most widely used algorithms are those based on linear multistep formulas like the BDF methods (see e.g. [21]), because they are very efficient for general stiff problems. However, from a result of Dahlquist it is known that no linear multistep method of order greater than two can be A-stable [11], and so these formulas are not suited to some stiff problems (for example, those with Jacobians whose eigenvalues have large imaginary parts).

Implicit Runge-Kutta formulas [9,6] have been widely used because of their excellent stability properties (such as A-stability, L-stability and B-stability), but the need for solving nonlinear algebraic equations at each step makes these formulas generally too costly when considering some huge systems of ODEs. In fact, if the dimension of the differential system is  $N$ , then it is well known that an  $s$ -stage fully implicit Runge-Kutta method involves the solution of an  $s \cdot N$ -dimensional nonlinear algebraic system. In order to retain the good stability properties of the method, this algebraic system has to be solved in general by a Newton-type method. The Newton-Raphson method involves the evaluation of  $s$  Jacobian matrices and also the solution of a  $s \cdot N$ -dimensional linear system (whose LU factorization requires approximately  $s^3 N^3/3$  multiplicative operations) for each iteration. Even if a simplified (modified) Newton iteration is used by replacing the  $s$  Jacobian matrices involved in Newton's method by a matrix  $J$  equal to  $f'$  evaluated at some point (losing the quadratic convergence property) so that the cost is reduced, each step requires approximately  $s^3 N^3/3$  multiplicative operations (to perform the LU factorization).

To reduce the amount of computational effort required to solve the nonlinear equations (by Newton-type iterations) when integrating with a fully implicit Runge-Kutta method, some classes of Runge-Kutta formulas have been developed. With the class of diagonally implicit Runge-Kutta (DIRK) methods (also called semi-implicit or semi-explicit Runge-Kutta methods) the algebraic cost of the LU factorization is reduced to  $s N^3/3$  multiplicative operations. By considering the class of singly diagonally implicit Runge-Kutta (SDIRK) methods and also the class of singly implicit Runge-Kutta (SIRK) methods, it is possible to reduce even more the algebraic cost of the LU factorization to  $N^3/3$  multiplicative operations (see e.g. [15,1,9] for more details).

In an attempt to overcome some of the handicaps mentioned previously, the class of formulas appropriate for solving stiff problems has been enlarged by introducing the so called multivalued (or general linear) methods, which combine linear multistep and Runge-Kutta methods (see e.g. [7,8,14,5,10]).

Many other attempts have been made in order to reduce the computation cost per step by considering linearly implicit methods, in this way eliminating the need for solving nonlinear systems which, as pointed before, usually are solved by Newton-type iteration (and hence require additional function evaluations for every iteration at every step). Such formulas have the computational ad-

vantage that it is necessary to solve only linear systems of algebraic equations at each step.

Among the many different RK-like methods of this type we have the Rosenbrock methods [24] and the ROW-methods (also called Rosenbrock-Wanner methods and modified Rosenbrock methods) [22,19,20]. These formulas, however, require the exact Jacobian at every step. Therefore the computations are costly when the Jacobian matrix is expensive to evaluate. For this reason, extensions of Rosenbrock methods have been considered in which the exact Jacobian is fixed for some number of steps so that the computation cost is reduced (see e.g. [31,32,17,29]). Moreover, Rosenbrock-type methods in which the exact Jacobian is no longer needed have been considered. The so called W-methods [28], the MROW-methods [33] and the generalized Runge-Kutta methods [30] (see also [12]) fall into this class. For an excellent survey of some of these methods the reader is referred to [15].

Special schemes for special problems have been also developed during the last years. Without being exhaustive, we have the symplectic methods for the integration of Hamiltonian problems (see e.g. [25]) and special methods for problems with oscillatory solutions (see e.g. [27]).

In our recent papers [3,4], examples are shown of explicit and linearly implicit two-stage methods of order three for the numerical integration of scalar autonomous ODEs, which do not require Jacobian evaluations, some of them being as well A-stable and L-stable. Some comparisons with Runge-Kutta methods as well as numerical experiments are also reported. In [2] we describe how to construct from a given function  $R$  a one-parameter family of explicit (or linearly implicit) two-stage methods, having  $R$  as the associated stability function, and illustrate this fact by obtaining a two-stage third order formula whose associated stability function is given by  $e^z$ .

Our first aim in the present paper is to generalize the new class of explicit (linearly implicit)  $p$ -stage methods for the numerical integration of scalar autonomous ODEs, which do not require Jacobian evaluations. These methods can be seen as a generalization of the explicit Runge-Kutta methods providing better order and stability results with the same number of stages. In fact, from Butcher's theory we know that an  $p$ -stage explicit Runge-Kutta method cannot have order greater than  $p$ . Moreover, the stability function of such methods is a polynomial, and so none of them is A-stable. We will show that it is possible to obtain A-stable explicit (linearly implicit) formulas for scalar autonomous problems of order three and five, with only two and three stages respectively, from our class of methods.

By a further generalization of our schemes, we show that it is possible to obtain A-stable linearly implicit formulas for some non-autonomous scalar ODEs and

systems, which do not require Jacobian evaluations.

Finally we illustrate the efficiency of those schemes by carrying out some numerical experiments.

## 2 The new family of methods.

We begin by considering the scalar autonomous initial value problem

$$y'(x) = f(y(x)), \quad y(x_0) = y_0. \quad (1)$$

For this problem, let us consider the family of explicit  $p$ -stage methods defined by

$$y_{n+1} = y_n + hF_{p+1}(k_1, k_2, \dots, k_p), \quad (2)$$

where the stages are given by

$$\begin{aligned} k_1 &= f(y_n) \\ k_2 &= f(y_n + hF_2(k_1)) \\ k_3 &= f(y_n + hF_3(k_1, k_2)) \\ &\vdots \\ k_p &= f(y_n + hF_p(k_1, k_2, \dots, k_{p-1})), \end{aligned} \quad (3)$$

and for each  $i$  with  $2 \leq i \leq p + 1$ ,  $F_i$  is any homogeneous function of degree one, that is

$$F_i(\alpha x_1, \alpha x_2, \dots, \alpha x_{i-1}) = \alpha F_i(x_1, x_2, \dots, x_{i-1}), \quad (4)$$

holds for each  $\alpha \in \mathbb{R}$  and  $(x_1, x_2, \dots, x_{i-1})$  in a subset of  $\mathbb{R}^{i-1}$ .

The family of methods we have just defined, may be shown to be a generalization of the explicit  $p$ -stage Runge-Kutta methods for problem (1). In fact, taking  $F_i(x_1, x_2, \dots, x_{i-1}) = a_{i1}x_1 + a_{i2}x_2 + \dots + a_{ii-1}x_{i-1}$  ( $2 \leq i \leq p$ ) in (3) and  $F_{p+1}(x_1, x_2, \dots, x_p) = b_1x_1 + b_2x_2 + \dots + b_px_p$  in (2), we get all the explicit  $p$ -stage Runge-Kutta formulas as a subfamily of our class.

Moreover, our methods can be seen as generalized explicit Runge-Kutta methods whose coefficients depend on the stages and are no longer constants. To show this it is enough to note that from the fact that  $F_i$  is an homogeneous function of degree one, by using Euler's theorem for homogeneous functions, we obtain

$$F_i(k_1, k_2, \dots, k_{i-1}) = \sum_{j=1}^{i-1} \frac{\partial F_i}{\partial x_j}(k_1, k_2, \dots, k_{i-1}) k_j, \quad 2 \leq i \leq p + 1,$$

from which it is clearly enough to take the Butcher array

$$\begin{array}{c|cccc}
 & a_{21} & & & \\
 & a_{31} & a_{32} & & \\
 & \vdots & \vdots & \ddots & \\
 & a_{p1} & a_{p2} & \cdots & a_{pp-1} \\
 \hline
 & b_1 & b_2 & \cdots & b_{p-1} & b_p
 \end{array}$$

where now the parameters  $a_{ij}$  and  $b_i$  are given in terms of the stages  $k_i$  and the functions  $F_i$  through the relations

$$\begin{aligned}
 a_{ij} &= \frac{\partial F_i}{\partial x_j}(k_1, k_2, \dots, k_{i-1}) \quad (1 \leq j < i \leq p), \\
 b_i &= \frac{\partial F_{p+1}}{\partial x_i}(k_1, k_2, \dots, k_p) \quad (1 \leq i \leq p).
 \end{aligned}$$

Our new  $p$ -stage methods are difficult to study because of the nonlinearities that can arise from the homogeneous functions in (4). For reasons that will be clear later when studying the consistency, convergence and linear stability properties of the methods, we will make some assumptions in order to simplify our analysis by giving a better formulation of our initial definitions.

We will assume in what follows that for each  $i$  with  $2 \leq i \leq p+1$  all partial derivatives of functions  $F_i$  up to order  $q \geq 1$  exist and are continuous in a neighbourhood of the point  $(1, 1, \dots, 1) \in \mathbb{R}^{i-1}$ . Therefore the values given by

$$c_i = F_i(1, 1, \dots, 1) \quad 2 \leq i \leq p+1, \quad (5)$$

exist.

As will become clear in the following argument, these  $c_i$  play a similar role to those associated with the classical Runge-Kutta methods. In fact, we have that the stages  $k_i$  can be seen as approximations to  $y'(x_n + c_i h)$ . Now we give the promised better formulation of the methods.

### 3 A useful formulation.

In order to simplify the study of the preceding family of methods, we introduce the terms

$$s_i = \frac{k_i - k_1}{k_1} = \frac{k_i}{k_1} - 1, \quad 2 \leq i \leq p, \quad (6)$$

where the stages  $k_i$  are given by (3). From considerations that will be clear later, we take  $s_i = 0$  when  $k_1 = 0$  in (6).

It is a simple task to show that  $s_i = O(h)$ , and we will exploit this property in order to simplify our study. Moreover, the terms  $s_i$  can be seen as approximations to  $c_i h f_y(y_n)$ . In fact  $s_i = c_i h f_y(y_n) + O(h^2)$  (here we assume that  $f$  has a sufficient number of bounded derivatives), and so we can obtain approximations to the Jacobian  $f_y(y_n)$  by taking  $s_i/(c_i h)$  (when  $c_i \neq 0$ ).

It is easy to show recursively that the stages  $k_i$  (with  $2 \leq i \leq p$ ) can be obtained from  $k_1$  and  $s_j$  with  $2 \leq j \leq i-1$  (see e.g. (8) below). Therefore, in terms of  $k_1$  and  $s_i$ , any method of the preceding family takes the form

$$y_{n+1} = y_n + h k_1 G_{p+1}(s_2, s_3, \dots, s_p), \quad (7)$$

where the  $s_i$  are given by (6) in terms of the stages

$$\begin{aligned} k_1 &= f(y_n) \\ k_2 &= f(y_n + h k_1 G_2) \\ k_3 &= f(y_n + h k_1 G_3(s_2)) \\ &\vdots \\ k_p &= f(y_n + h k_1 G_p(s_2, s_3, \dots, s_{p-1})), \end{aligned} \quad (8)$$

and the functions  $G_i$  can be obtained from the homogeneous functions  $F_i$  through the relations

$$\begin{aligned} G_i(s_2, s_3, \dots, s_{i-1}) &= \frac{1}{k_1} F_i(k_1, k_2, \dots, k_{i-1}) = F_i\left(1, \frac{k_2}{k_1}, \dots, \frac{k_{i-1}}{k_1}\right) \\ &= F_i(1, 1 + s_2, \dots, 1 + s_{i-1}), \quad 2 \leq i \leq p+1, \end{aligned} \quad (9)$$

Note that when  $i = 2$  we have that  $G_2 = (1/k_1)F_2(k_1) = F_2(1) = c_2$  holds from (5). Now by our previous assumptions and relations (9) it follows that for each  $i$  with  $2 \leq i \leq p+1$ , all partial derivatives of the functions  $G_i$  up to order  $q \geq 1$  exist and are continuous in a neighbourhood of the point  $(0, 0, \dots, 0) \in \mathbb{R}^{i-2}$ , and therefore  $c_i = G_i(0, 0, \dots, 0)$  holds. From the above relations, we can write any method (2) in the form given in (7).

It is also possible to obtain from a given method in terms of  $k_1$  and  $s_i$  the associated expression in terms of the  $k_i$  by using the relations

$$\begin{aligned} F_i(k_1, \dots, k_{i-1}) &= k_1 G_i(s_2, \dots, s_{i-1}) \\ &= k_1 G_i\left(\frac{k_2}{k_1} - 1, \dots, \frac{k_{i-1}}{k_1} - 1\right), \quad 2 \leq i \leq p+1. \end{aligned}$$

For  $i = 2$  we have  $F_2(k_1) = k_1 G_2$ .

#### 4 Consistency and order of the methods.

In what follows we will assume that  $f : \mathbb{R} \rightarrow \mathbb{R}$  is Lipschitz in  $\mathbb{R}$ , i. e. there exists a Lipschitz constant  $L$  such that

$$|f(y) - f(y^*)| \leq L|y - y^*|,$$

for every  $y, y^* \in \mathbb{R}$ . With the previous assumptions it is well know that for any  $y_0 \in \mathbb{R}$ , there exists a unique solution  $y(x)$  of problem (1) (throughout any interval  $[x_0, b]$ ), where  $y(x)$  is continuous and differentiable.

We will investigate the consistency of any  $p$ -stage method given by (7). When doing so, we will assume the existence (and continuity) of  $y'(x)$  in  $[x_0, b]$ , but not necessarily that of higher derivatives.

Using Henrici's notation for one step methods, our methods can be expressed as  $y_{n+1} = y_n + h \Phi(x_n, y_n, h)$ , with the increment function  $\Phi$  (not depending explicitly on  $x$  because (1) is autonomous) given by

$$\Phi(x, y, h) = k_1 G_{p+1}(s_2, s_3, \dots, s_p),$$

through (6–9) (with  $y$  in place of  $y_n$ ).

As is usual with other one-step methods, we define the local truncation error  $T(x, h)$  of any formula of our family to be

$$T(x, h) = y(x + h) - y(x) - h k_1 G_{p+1}(s_2, s_3, \dots, s_p), \quad x \in [x_0, b], \quad (10)$$

where  $h > 0$ , the stages  $k_i$  are given by (8) with the exact solution of (1)  $y(x)$  in place of  $y_n$ , the terms  $s_i$  by (6) and the functions  $G_i$  through (9).

**Definition 1** *Method (7) is said to be consistent (with (1) satisfying our previous assumptions) if, in the limit as  $h \rightarrow 0$  we have that  $T(x, h)/h \rightarrow 0$  uniformly for  $x \in [x_0, b]$ .*

It is known that consistency is a necessary condition for convergence, and therefore we want to establish the condition that method (7) must satisfy if we require it to be consistent. In order to do so, we need the following lemma:

**Lemma 2** *Suppose that for each  $j$  with  $2 \leq j \leq p+1$ , the functions  $G_j$  in (9) are continuous in a neighbourhood  $W_{j-2}$  of the point  $(0, 0, \dots, 0) \in \mathbb{R}^{j-2}$ . Let  $s_i$  ( $2 \leq i \leq p$ ) be given through (6) in terms of the stages  $k_i$  in (8) with  $y$  in place of  $y_n$ . Then for every  $\epsilon > 0$  there exist  $h_1 > 0$  such that*

$$\begin{aligned} |s_i| &\leq hL(|c_i| + \epsilon), \quad 2 \leq i \leq p, \\ |G_j(s_2, s_3, \dots, s_{j-1})| &\leq |c_j| + \epsilon, \quad 2 \leq j \leq p+1, \end{aligned}$$

hold for every  $0 < h \leq h_1$  and  $y \in \mathbb{R}$ .

Recall that as we have pointed out before,  $c_i = G_i(0, 0, \dots, 0)$ .

**Proof.** The proof follows from the Lipschitz character of  $f$ , and from the fact that the functions  $G_j$  are continuous in a neighbourhood  $W_{j-2}$  of the point  $(0, 0, \dots, 0) \in \mathbb{R}^{j-2}$ , by a recursive procedure.

From the continuity of all  $G_j$  in  $W_{j-2}$  we have that for any given  $\epsilon > 0$  there exist  $\delta > 0$  such that for  $\max_{1 \leq i \leq p-1} |x_i| \leq \delta$  we have that  $|G_j(x_1, x_2, \dots, x_{j-2})| \leq |c_j| + \epsilon$  (for  $3 \leq j \leq p+1$ ). When  $j = 2$ ,  $G_2$  is constant, and therefore  $|G_2| = |c_2| \leq |c_2| + \epsilon$  also holds.

Now by taking  $h_1 = \delta/L(\epsilon + \max_{2 \leq i \leq p} |c_i|)$  we obtain recursively (from the Lipschitz character of  $f$ ) that

$$\begin{aligned} |s_i| &= \frac{|k_i - k_1|}{|k_1|} = \frac{|f(y + hk_1 G_i(s_2, s_3, \dots, s_{i-1})) - f(y)|}{|k_1|} \\ &\leq hL|G_i(s_2, s_3, \dots, s_{i-1})| \leq hL(|c_i| + \epsilon), \end{aligned} \quad (11)$$

$$|G_{i+1}(s_2, s_3, \dots, s_i)| \leq |c_{i+1}| + \epsilon,$$

for  $2 \leq i \leq p$ ,  $0 < h \leq h_1$  and  $y \in \mathbb{R}$ .

When  $k_1 = 0$  the result also holds from our previous assumption that  $s_i = 0$  in that case. This assumption is easily justified by using that function  $f$  is continuous as follows. If  $k_1 = f(y) = 0$  we consider  $(y^{(n)})_{n \in \mathbb{N}}$  with  $y^{(n)} \rightarrow y$  such that for every  $n \in \mathbb{N}$   $f(y^{(n)}) \neq 0$ . Obviously (11) is satisfied with this  $y^{(n)}$  and so we have that in the limit case when  $y^{(n)} \rightarrow y$  it must also hold.  $\square$

From the above lemma it is clear why  $s_i = O(h)$ . Now we have

**Theorem 3** *Method (7) is consistent with (1) (under the above assumptions) iff*

$$G_{p+1}(0, 0, \dots, 0) = c_{p+1} = 1. \quad (12)$$

**Proof.** For any given  $x \in [x_0, b]$  and  $h > 0$  we get from the mean value theorem that

$$y(x+h) - y(x) = hy'(\alpha_x),$$

where  $\alpha_x \in (x, x + h)$ . It then follows from (10) that

$$\frac{T(x, h)}{h} = y'(\alpha_x) - k_1 G_{p+1}(s_2, s_3, \dots, s_p),$$

Now in the limit as  $h \rightarrow 0$  we have that  $y'(\alpha_x) \rightarrow y'(x)$  uniformly for  $x \in [x_0, b]$ . From the previous lemma we easily get that as  $h \rightarrow 0$ ,  $s_i \rightarrow 0$  for each  $i$  (with  $2 \leq i \leq p$ ). Finally, from the continuity of the function  $G_{p+1}$ , and taking into account that  $y'(x) = f(y(x))$  and  $k_1 = f(y(x))$ , we get

$$\lim_{h \rightarrow 0} \frac{T(x, h)}{h} = (1 - G_{p+1}(0, 0, \dots, 0))f(y(x)),$$

from which we obtain the consistency condition (12).  $\square$

Note at this point that using Henrici's notation for one step methods, the consistency condition reads  $\Phi(y, 0) = f(y)$ . Obviously this consistency condition takes the form (12).

Now we define the consistency of order  $q$  in the usual way, that is,

**Definition 4** *method (7) is said to be consistent (with the differential equation (1)) of order  $q$ , if  $q$  is the largest integer such that there exists  $N \geq 0$  and  $h_0 > 0$  with  $\sup_{x_0 \leq x \leq b} |T(x, h)| \leq N h^{q+1}$  for all  $h \in (0, h_0]$ .*

If all the partial derivatives of  $f(y)$  up to order  $q$  exist (and are continuous), then consistency follows from the consistency of order  $q \geq 1$ .

## 5 Convergence of the methods.

Now we will study the convergence of the methods, by considering the so called global error (the error of the computed solution after several steps). We will consider the stepsize fixed in order to simplify our study, but all results remain valid with little changes in other cases. Our task is now to estimate the global error, and this can be done in two different ways:

- by considering the local errors and how they are transported along the exact solution curves.
- by studying how the local truncation errors are transported along the numerical solutions.

For more details see e.g. [13,9]. The first approach is perhaps easier, and can yield sharp results when sharp estimates of error propagation for the exact

solutions are known. However we will follow the second approach that is the preferred one because it generalizes easily to multistep methods, and can be an important tool for the existence of asymptotic expansions.

As is usual when considering one-step methods, we need to show that the increment function  $\Phi(y, h) = k_1 G_{p+1}(s_2, s_3, \dots, s_p)$  given by (6–9) (with  $y$  in place of  $y_n$ ) satisfies a Lipschitz condition in  $y$ , for  $h$  small enough. Even though this property is nearly automatic for most of the one-step methods from the Lipschitz condition that satisfies function  $f$ , for our methods it is not as easy as we will see in what follows.

**Lemma 5** *If all partial derivatives of the functions  $G_i$  up to order  $q \geq 1$  exist and are continuous in a convex neighbourhood  $W_{i-2}$  of the point  $(0, 0, \dots, 0) \in \mathbb{R}^{i-2}$ , then there exist a constant  $\Lambda$  such that for  $0 < h \leq h_0$*

$$|\Phi(y, h) - \Phi(y^*, h)| \leq \Lambda |y - y^*|, \quad (13)$$

holds for every  $y, y^* \in \mathbb{R}$ .

**Proof.** We begin noting that

$$\begin{aligned} k_1 G_i(s_2, \dots, s_{i-1}) - k_1^* G_i(s_2^*, \dots, s_{i-1}^*) \\ = \frac{k_1 - k_1^*}{2} (G_i(s_2, \dots, s_{i-1}) + G_i(s_2^*, \dots, s_{i-1}^*)) \\ + \frac{k_1 + k_1^*}{2} (G_i(s_2, \dots, s_{i-1}) - G_i(s_2^*, \dots, s_{i-1}^*)), \end{aligned}$$

holds for  $2 \leq i \leq p+1$  and  $y, y^* \in \mathbb{R}$ . Here  $s_i$  and  $k_i$  (respectively  $s_i^*$  and  $k_i^*$ ) are given through (6–9) with  $y$  (respectively  $y^*$ ) in place of  $y_n$ . The functions  $G_i$  are given as usually by (9).

Now we must prove some inequalities that will be useful later. Since  $f$  is Lipschitz, we have

$$\begin{aligned} |k_1 - k_1^*| &\leq L |y - y^*|, \\ |k_i - k_i^*| &\leq L |y - y^*| + hL |k_1 G_i(s_2, \dots, s_{i-1}) - k_1^* G_i(s_2^*, \dots, s_{i-1}^*)|, \end{aligned} \quad (14)$$

for  $2 \leq i \leq p$ .

From the previous lemma we have that for given  $y, y^* \in \mathbb{R}$  there exist  $h_2 > 0$  such that  $(s_2, \dots, s_i)$  and  $(s_2^*, \dots, s_i^*)$  belong to  $W_{i-1}$  when  $0 < h \leq h_2$ . Therefore, by the mean value theorem we obtain for  $0 < h \leq h_2$  and for each  $i$  with  $2 \leq i \leq p+1$

$$G_i(s_2, \dots, s_{i-1}) - G_i(s_2^*, \dots, s_{i-1}^*) = \sum_{j=1}^{i-2} (s_{j+1} - s_{j+1}^*) \frac{\partial G_i}{\partial x_j}(\xi_i), \quad (15)$$

where  $\xi_i$  is an internal point of the line segment in  $W_{i-2}$  joining  $(s_2, \dots, s_{i-1})$  to  $(s_2^*, \dots, s_{i-1}^*)$ .

As in the preceding lemma, by using the continuity of the first partial derivatives of functions  $G_i$  in  $W_{i-2}$ , we obtain that for any given  $\epsilon > 0$  there exists  $h_3 > 0$  such that for  $0 < h \leq h_3$  and any given  $y \in \mathbb{R}$  we have

$$\left| \frac{\partial G_i}{\partial x_j}(s_2, \dots, s_{i-1}) \right| \leq |c_{i,j}| + \epsilon, \quad 2 \leq i \leq p+1, \quad 1 \leq j \leq i-2,$$

where  $c_{i,j} = \frac{\partial G_i}{\partial x_j}(0, \dots, 0)$ . It is easily seen that taking  $h_4 = \min(h_2, h_3)$  we have that for  $0 < h \leq h_4$  relation (15) holds with

$$\left| \frac{\partial G_i}{\partial x_j}(\xi_i) \right| \leq |c_{i,j}| + \epsilon.$$

We will also need the following obvious identity

$$(k_1 + k_1^*)(s_i - s_i^*) = (2 + s_i + s_i^*)(k_1^* - k_1) + 2(k_i - k_i^*),$$

from which we obtain by using the previous lemma and (14)

$$|(k_1 + k_1^*)(s_i - s_i^*)| \leq 2(2 + hL(|c_i| + \epsilon))L|y - y^*| + 2hL\nu_i,$$

where

$$\nu_i = |k_1 G_i(s_2, \dots, s_{i-1}) - k_1^* G_i(s_2^*, \dots, s_{i-1}^*)|, \quad 2 \leq i \leq p+1.$$

Note that  $\nu_{p+1} = |\Phi(y, h) - \Phi(y^*, h)|$ . Now from above inequalities we obtain the following recursion formula:

$$\begin{aligned} \nu_i \leq & \left( |c_i| + \epsilon + \sum_{j=1}^{i-2} (|c_{i,j}| + \epsilon)(2 + hL(|c_{j+1}| + \epsilon)) \right) L|y - y^*| \\ & + hL \sum_{j=1}^{i-2} (|c_{i,j}| + \epsilon) \nu_{j+1}, \end{aligned}$$

for every  $i$  with  $2 \leq i \leq p+1$ . Now we define  $\nu = (\nu_2, \nu_3, \dots, \nu_{p+1})^T$ ,  $C^\epsilon = (\epsilon + |c_2|, \epsilon + |c_3|, \dots, \epsilon + |c_{p+1}|)^T$  and

$$C_{\partial}^\epsilon = \begin{pmatrix} 0 & 0 & 0 & \cdots & 0 & 0 & 0 \\ \epsilon + |c_{3,1}| & 0 & 0 & \cdots & 0 & 0 & 0 \\ \epsilon + |c_{4,1}| & \epsilon + |c_{4,2}| & 0 & \cdots & 0 & 0 & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots & \vdots & \vdots \\ \epsilon + |c_{p-1,1}| & \epsilon + |c_{p-1,2}| & \epsilon + |c_{p-1,3}| & \cdots & 0 & 0 & 0 \\ \epsilon + |c_{p,1}| & \epsilon + |c_{p,2}| & \epsilon + |c_{p,3}| & \cdots & \epsilon + |c_{p,p-2}| & 0 & 0 \\ \epsilon + |c_{p+1,1}| & \epsilon + |c_{p+1,2}| & \epsilon + |c_{p+1,3}| & \cdots & \epsilon + |c_{p+1,p-2}| & \epsilon + |c_{p+1,p-1}| & 0 \end{pmatrix}$$

in terms of which (16) takes the form

$$\nu \leq (C^\epsilon + C_{\partial}^\epsilon(2\mathbb{1} + hLC^\epsilon))L|y - y^*| + hLC_{\partial}^\epsilon\nu,$$

where the inequalities are considered componentwise and  $\mathbb{1}$  is given by  $\mathbb{1} = (1, 1, \dots, 1)^T$ . It follows

$$(I - hLC_{\partial}^\epsilon)\nu \leq (C^\epsilon + C_{\partial}^\epsilon(2\mathbb{1} + hLC^\epsilon))L|y - y^*|.$$

Taking  $h$  small enough so that  $hL\|C_{\partial}^\epsilon\| < 1$  holds ( $\|\cdot\|$  is any fixed matrix norm), it follows that  $(I - hLC_{\partial}^\epsilon)^{-1}$  exists with all entries nonnegative. Therefore we obtain

$$\nu \leq (I - hLC_{\partial}^\epsilon)^{-1}(C^\epsilon + C_{\partial}^\epsilon(2\mathbb{1} + hLC^\epsilon))L|y - y^*|,$$

from which

$$\nu_{p+1} \leq e_p(I - hLC_{\partial}^\epsilon)^{-1}(C^\epsilon + C_{\partial}^\epsilon(2\mathbb{1} + hLC^\epsilon))L|y - y^*|,$$

holds, with  $e_p = (0, \dots, 0, 1) \in \mathbb{R}^p$ . Finally, taking into account that  $\nu_{p+1} = |\Phi(y, h) - \Phi(y^*, h)|$  the result follows, with the Lipschitz constant given by

$$\Lambda = e_p(I - hLC_{\partial}^\epsilon)^{-1}(C^\epsilon + C_{\partial}^\epsilon(2\mathbb{1} + hLC^\epsilon))L,$$

for  $h$  small enough and every  $y, y^* \in \mathbb{R}$ .  $\square$

**Definition 6** *Method ( $\gamma$ ) is said to be convergent if for all  $x \in [x_0, b]$ , in the limit as  $h \rightarrow 0$  (with  $x_n = x$  kept fixed) we have that  $y_n \rightarrow y(x)$ .*

Now we can prove the convergence of the methods.

**Theorem 7** *Under the above assumptions, any consistent method (7) is convergent.*

**Proof.** From the Lipschitz condition (13) satisfied by the increment function  $\Phi$ , it follows that

$$|y_{i+1} - y_{i+1}^*| \leq (1 + h\Lambda) |y_i - y_i^*| \leq \exp(h\Lambda) |y_i - y_i^*|,$$

for  $h$  small enough.

We consider the numerical solution for a point  $x > x_0$ , obtained by the step-by-step procedure (7), with constant step size  $h = x_{i+1} - x_i$  ( $0 \leq i \leq n-1$ ) and  $x_n = x$ . Now our task is to estimate the global error  $E = |y(x_n) - y_n|$  at any fixed point  $x > x_0$ , following the second approach. We have that

$$E = |y(x_n) - y_n| \leq \sum_{i=1}^n |y_{n,i} - y_{n,i-1}|,$$

holds, where  $y_{n,i}$  denotes the approximation obtained by carrying out  $n-i$  steps ( $y_{j,i}$   $j = i+1, i+2, \dots, n$ ) with method (7), using the exact value  $y_{i,i} = y(x_i)$  as the approximation at  $x_i$ . Note that with our notations  $y_{n,n} = y(x_n)$  and  $y_{n,0} = y_n$ . Now we have to bound the individual terms in the right hand side of (5). For small enough  $h$  we obtain from (5) the inequalities

$$\begin{aligned} |y_{n,i} - y_{n,i-1}| &\leq (1 + h\Lambda)^{n-i} |y_{i,i} - y_{i,i-1}| \\ &\leq \exp((n-i)h\Lambda) |y_{i,i} - y_{i,i-1}|, \quad 1 \leq i \leq n. \end{aligned}$$

Note that from  $x_i = x_0 + ih$  we have that  $(n-i)h = x_n - x_i$  in the above inequality. Note also that  $l_i = y_{i,i} - y_{i,i-1} = T(x_{i-1}, h)$  is the local truncation error associated with the  $i$ -th integration step (at  $x_i$ ). Therefore, it follows from (5) and (16), taking  $l = \max_{1 \leq i \leq n} |l_i|$  that

$$\begin{aligned} E = |y(x_n) - y_n| &\leq l \frac{\exp(\Lambda(x_n - x_0)) - 1}{\exp(\Lambda h) - 1} \\ &\leq \frac{l}{\Lambda h} (\exp(\Lambda(x_n - x_0)) - 1). \end{aligned}$$

Finally, when  $h \rightarrow 0$  (with  $nh = x - x_0$  fixed) we have from the consistency assumption that  $l/h \rightarrow 0$ , and convergence is clear from the above inequality.  $\square$

**Definition 8** *Method (7) is said to be convergent of order  $q$  if for all  $x_n \in [x_0, b]$  (with  $x_n = x$  kept fixed) and all  $h \in (0, h_0]$  there is a constant  $M \geq 0$  such that  $|y(x_n) - y_n| \leq M h^q$ .*

Now we can see that any consistent method of order  $q$  is convergent of order  $q$ .

**Theorem 9** *Under the above assumptions, any consistent method (7) of order  $q$  is convergent of order  $q$ .*

**Proof.** The proof easily follows from (16) taking into account that now we have that  $l \leq N h^{q+1}$  for small enough  $h$  and therefore (16) now reads

$$\begin{aligned} E = |y(x_n) - y_n| &\leq N h^{q+1} \frac{\exp(\Lambda(x_n - x_0)) - 1}{\exp(\Lambda h) - 1} \\ &\leq \frac{N}{\Lambda} h^q (\exp(\Lambda(x_n - x_0)) - 1), \end{aligned}$$

and it is enough to take  $M = (N/\Lambda) (\exp(\Lambda(b - x_0)) - 1)$ .  $\square$

Note that errors introduced at each step due to rounding, have not been considered in our preceding analysis. We only comment at this point that, as usual, for a given arithmetical accuracy there is a minimum step size  $h$  below which rounding errors will produce inaccuracies larger than those due to truncation errors.

## 6 Methods of polynomial type.

Now, for every fixed  $p$ , we restrict our attention to the family of methods given by (2-3), where now all the  $F_i$  ( $2 \leq i \leq p+1$ ) are assumed to be homogeneous functions (of degree one) of the special form

$$F_i(x_1, x_2, \dots, x_{i-1}) = \sum_{j_2+j_3+\dots+j_{i-1}=0}^{r_i} A_{j_2 j_3 \dots j_{i-1}} x_1 \left(\frac{x_2}{x_1}\right)^{j_2} \left(\frac{x_3}{x_1}\right)^{j_3} \dots \left(\frac{x_{i-1}}{x_1}\right)^{j_{i-1}}, \quad (16)$$

with all  $r_i$  being nonnegative integers.

The above family of methods, still contains all explicit  $p$ -stage Runge-Kutta methods. In fact, taking  $r_i = 1$  in (16) we get from (2) and (3) all explicit  $p$ -stage Runge-Kutta methods.

As we will see in what follows, we restrict our attention to the above family of methods because in terms of  $s_i$  the associated functions  $G_i$  are of polynomial type. This greatly simplifies the study of the order conditions. Moreover, order conditions for the general methods can be easily obtained from the order conditions for methods of polynomial type, by considering the Taylor expansion of the functions  $G_i$  in terms of the  $s_i$ .

From now on, we will assume that all quantities  $c_i = F_i(1, 1, \dots, 1)$  in (16) are different from zero. Even though we lose a bit of generality with this assumption, we obtain many advantages that will be clear later. In fact, it can be seen that for a given number of stages  $p$  (at least for  $p = 2, 3, 4$ ), the highest order is attained only when all  $c_i \neq 0$ , and so, all interesting methods (from an order point of view) are considered.

With the above assumption we now change our definition of the  $s_i$  in (6), hoping that this will not confuse the reader. We define

$$s_i = \frac{k_i - k_1}{c_i k_1}, \quad 2 \leq i \leq p, \quad (17)$$

with the stages  $k_i$  given by (3). Now it is easy to see that  $s_i = hf_y(y_n) + O(h^2)$ , and this will be useful later when looking for methods with good linear stability properties, since it simplifies the study of the order conditions.

In terms of  $k_1$  and  $s_i$ , the preceding method takes the form (7–8), but now the functions  $G_i$  are given in terms of the new  $s_i$  and the functions  $F_i$  in (16) through the relations

$$\begin{aligned} G_i(s_2, s_3, \dots, s_{i-1}) &= \frac{1}{k_1} F_i(k_1, k_2, \dots, k_{i-1}) = F_i\left(1, \frac{k_2}{k_1}, \dots, \frac{k_{i-1}}{k_1}\right) \\ &= F_i(1, 1 + c_2 s_2, \dots, 1 + c_{i-1} s_{i-1}), \quad 2 \leq i \leq p+1. \end{aligned} \quad (18)$$

It can be seen that the functions  $G_i$  take the form

$$G_i(z_2, z_3, \dots, z_{i-1}) = c_i \left( 1 + \sum_{j_2+j_3+\dots+j_{i-1}=1}^{r_i} a_{j_2 j_3 \dots j_{i-1}} z_2^{j_2} z_3^{j_3} \dots z_{i-1}^{j_{i-1}} \right). \quad (19)$$

Note that for this special class of methods, as we have pointed out before, all  $G_i$  are polynomial functions of the  $s_i$ . Therefore, in what follows we will refer to this formulas as methods of polynomial type.

Now we study the attainable order, in terms of the number of stages, of the methods given by (7), (8), (17) and (19). From now on, we will suppose that  $f$  is smooth enough in order that all the derivatives which occur when considering the Taylor expansion of the local truncation error make sense.

## 7 Two-stage methods of polynomial type. Attainable order.

To show that we can obtain explicit two-stage methods of polynomial type from (7–8) (together with (17) and (19)) of order three, it is enough to consider the Taylor expansion of the associated local truncation error. Note that when

doing so, only the parameters  $c_2$ ,  $c_3$  and  $a_i$  with  $1 \leq i \leq 2$  will appear in the order conditions, that is, it suffices to take in (19)  $r_3 = 2$  (the other parameters can be arbitrarily chosen). This easily follows from the fact that  $s_2 = O(h)$ , and therefore  $s_2^k = O(h^k)$  for any given  $k \in \mathbb{N}$ . We obtain the following order conditions

$$c_3 = 1, \quad (20)$$

$$c_3 a_1 = 1/2, \quad (21)$$

$$c_3 c_2 a_1 = 1/3, \quad (22)$$

$$c_3 a_2 = 1/6, \quad (23)$$

from which the general form of a third order two-stage method of polynomial type is given by

$$y_{n+1} = y_n + h k_1 G_3(s_2),$$

where  $k_1$ ,  $k_2$  and  $s_2$  are given by

$$k_1 = f(y_n), \quad k_2 = f\left(y_n + \frac{2}{3} h k_1\right), \quad s_2 = \frac{3(k_2 - k_1)}{2k_1}, \quad (24)$$

and  $G_3$  takes the form

$$G_3(s_2) = 1 + \frac{1}{2} s_2 + \frac{1}{6} s_2^2 + \sum_{i=3}^{r_3} a_i s_2^i, \quad (25)$$

that is, it is enough to take  $c_3 = 1$  (consistency condition),  $c_2 = 2/3$ ,  $a_1 = 1/2$  and  $a_2 = 1/6$ . Parameters  $a_i$  with  $i \geq 3$  can be arbitrarily chosen.

The additional order conditions (now with  $r_3 = 3$  in (19)) for order four are given by

$$c_3 c_2^2 a_1 = 1/4, \quad (26)$$

$$c_3 c_2 a_2 = 1/6, \quad (27)$$

$$c_3 a_3 = 1/24, \quad (28)$$

and it is easy to check that no two-stage method of polynomial type has order greater than three (conditions (26) and (27) cannot be satisfied). However, condition (28) can be satisfied by taking  $a_3 = 1/24$ , obtaining in this way third order methods that minimize the principal part of the local truncation error. Note at this point that the other terms of the principal part of the local truncation error (those associated with conditions (26) and (27)) are the same for any method of order three.

## 8 Three-stage methods of polynomial type. Attainable order.

When considering three-stage methods of polynomial type from (7–8) together with (17) and (19), it is possible to get a family of fifth order formulas, depending on many free parameters. As in the two-stage case, only some of this parameters will appear in the order conditions, due to the fact that both  $s_2$  and  $s_3$  are  $O(h)$ . However, the resulting order conditions are still too cumbersome; hence we define a new term  $\tilde{s}_3 = s_3 - s_2$  to make our study easier. Note that from our previous considerations it is clear that  $\tilde{s}_3 = O(h^2)$ . Now we can describe any three-stage method of polynomial type by

$$y_{n+1} = y_n + hk_1 \tilde{G}_4(s_2, \tilde{s}_3),$$

with  $\tilde{s}_3 = s_3 - s_2$ , and where  $s_2$  and  $s_3$  are given through (17) in terms of the stages  $k_i$  ( $1 \leq i \leq 3$ )

$$k_1 = f(y_n), \quad k_2 = f(y_n + hk_1 G_2), \quad k_3 = f(y_n + hk_1 G_3(s_2)).$$

Now functions  $G_2$ ,  $G_3$  and  $\tilde{G}_4$  are given by

$$\begin{aligned} G_2 &= c_2, \\ G_3(s_2) &= c_3 \left( 1 + \sum_{i=1}^{r_3} a_i s_2^i \right), \\ \tilde{G}_4(s_2, \tilde{s}_3) &= c_4 \left( 1 + \sum_{i+2j=1}^{\tilde{r}_4} a_{ij} s_2^i \tilde{s}_3^j \right), \end{aligned} \quad (29)$$

and it is easily seen when looking for formulas of order five, that it is enough to consider the parameters in (29) with  $r_3 = 3$  and  $\tilde{r}_4 = 4$ . This means that the order conditions (for order five) are completely determined in terms of the parameters  $c_i$  ( $2 \leq i \leq 4$ ),  $a_i$  ( $1 \leq i \leq 3$ ), and  $a_{ij}$  ( $1 \leq i + 2j \leq 4$ ). This follows from the fact that  $s_2 = O(h)$  and  $\tilde{s}_3 = O(h^2)$ .

Now, any three-stage methods of polynomial type must satisfy the following order conditions in order to be of order five

$$c_4 = 1, \quad (30)$$

$$c_4 a_{10} = 1/2, \quad (31)$$

$$c_4 (c_2(a_{10} - a_{01}) + c_3 a_{01}) = 1/3, \quad (32)$$

$$c_4 (a_{20} + a_1 a_{01}) = 1/6, \quad (33)$$

$$c_4 (c_2^2(a_{10} - a_{01}) + c_3^2 a_{01}) = 1/4, \quad (34)$$

$$c_4 (c_2(2a_{20} - a_{11} + a_1 a_{01}) + c_3(a_{11} + 2a_1 a_{01})) = 1/3, \quad (35)$$

$$c_4 (a_{30} + a_1 a_{11} + a_2 a_{01}) = 1/24, \quad (36)$$

$$c_4 (c_2^3(a_{10} - a_{01}) + c_3^3 a_{01}) = 1/5, \quad (37)$$

$$c_4 \left( c_2^2(2a_{20} - a_{11} + a_1 a_{01}) + c_3^2(a_{11} + 3a_1 a_{01}) \right) = 7/20, \quad (38)$$

$$c_4 \left( c_2^2(a_{20} - a_{11} + a_{02}) + c_2 c_3(a_{11} - 2a_{02} + 2a_1 a_{01}) + c_3^2 a_{02} \right) = 2/15, \quad (39)$$

$$c_4 \left( c_2(3a_{30} - a_{21} + 2a_1(a_{11} - a_{02}) + 2a_2 a_{01}) \right. \\ \left. + c_3(a_{21} + 2a_1(a_{11} + a_{02}) + (a_1^2 + 2a_2)a_{01}) \right) = 11/60, \quad (40)$$

$$c_4 \left( a_{40} + a_1 a_{21} + a_1^2 a_{02} + a_2 a_{11} + a_3 a_{01} \right) = 1/120. \quad (41)$$

It is easy to check that the above system has many solutions. In fact, we get two doubly infinite families of solutions, taking  $a_2$  and  $a_3$  as free parameters. The remaining coefficients can be computed as follows:

*Step 1* Obviously, we have  $c_4$  from (30). Coefficients  $a_{10}$ ,  $a_{01}$ ,  $c_2$  and  $c_3$ , can be chosen such that (31), (32), (34) and (37) are satisfied. We get a pair of solutions in this way.

*Step 2* Now, for each solution of the previous step, we can obtain  $a_1$ ,  $a_{20}$ ,  $a_{11}$  and  $a_{02}$ , by solving the linear system given by (33), (35), (38) and (39).

*Step 3* Finally, we obtain  $a_{40}$ ,  $a_{30}$  and  $a_{21}$  by solving the remaining linear system (36), (40) and (41) (in terms of the parameters  $a_2$  and  $a_3$ ).

The result is

$$c_2 = \frac{6 \mp \sqrt{6}}{10}, \quad c_3 = \frac{6 \pm \sqrt{6}}{10}, \quad c_4 = 1, \quad a_1 = \frac{-3 \pm 2\sqrt{6}}{5}, \quad a_{10} = \frac{1}{2}, \\ a_{01} = \frac{9 \pm \sqrt{6}}{36}, \quad a_{20} = \frac{3 \mp \sqrt{6}}{12}, \quad a_{11} = \frac{3 \pm 2\sqrt{6}}{72}, \quad a_{02} = \frac{1 \pm 4\sqrt{6}}{72}, \\ a_{30} = \frac{-9 \mp \sqrt{6}}{36} a_2, \quad a_{21} = \frac{-(3 + 11a_2) \pm (1 - 4a_2)\sqrt{6}}{24}, \\ a_{40} = \frac{-3(3 - 10a_2 + 30a_3) \pm (3 + 20a_2 - 10a_3)\sqrt{6}}{360}. \quad (42)$$

The other parameters may be arbitrarily chosen.

For order six we must add to (30–41) the following conditions (now it is enough to take  $r_3 = 4$  and  $\tilde{r}_4 = 5$  in (29))

$$c_4 \left( c_2^4(a_{10} - a_{01}) + c_3^4 a_{01} \right) = 1/6, \quad (43)$$

$$c_4 \left( c_2^3(2a_{20} - a_{11} + a_1 a_{01}) + c_3^3(a_{11} + 4a_1 a_{01}) \right) = 11/30, \quad (44)$$

$$c_4 \left( 2c_2^3(a_{20} - a_{11} + a_{02}) + c_2^2 c_3(a_{11} - 2a_{02} + 2a_1 a_{01}) \right. \\ \left. + c_2 c_3^2(a_{11} - 2a_{02} + 3a_1 a_{01}) + 2c_3^3 a_{02} \right) = 1/4, \quad (45)$$

$$c_4 \left( c_2^2(3a_{30} - a_{21} + 2a_1(a_{11} - a_{02}) + 2a_2 a_{01}) \right)$$

$$+c_3^2(a_{21} + a_1(3a_{11} + 2a_{02}) + 3(a_1^2 + a_2)a_{01}) = 4/15, \quad (46)$$

$$\begin{aligned} c_4 \left( c_2^2(3a_{30} - 2a_{21} + a_{12} + a_1(a_{11} - 2a_{02}) + a_2a_{01}) \right. \\ \left. + 2c_2c_3(a_{21} - a_{12} + a_1(2a_{11} - a_{02}) + (a_1^2 + 2a_2)a_{01}) \right. \\ \left. + c_3^2(a_{12} + 4a_1a_{02}) \right) = 17/90, \quad (47) \end{aligned}$$

$$\begin{aligned} c_4 \left( c_2(4a_{40} - a_{31} + a_1(3a_{21} - 2a_{12}) + 2(a_1^2 - a_2)a_{02}) \right. \\ \left. + 3a_2a_{11} + 3a_3a_{01} \right) + c_3(a_{31} + 2(a_1(a_{21} + a_{12}) \\ + (2a_1^2 + a_2)a_{02} + (a_1a_2 + a_3)a_{01}) + (a_1^2 + 2a_2)a_{11}) = 13/180, \quad (48) \end{aligned}$$

$$c_4(a_{50} + a_1(a_{31} + a_1a_{12} + 2a_2a_{02}) + a_2a_{21} + a_3a_{11} + a_4a_{01}) = 1/720. \quad (49)$$

Now from (42) and conditions (43–49), it is not difficult to see that order six cannot be attained with only three stages. In fact, conditions (43–45) cannot be satisfied, and the coefficients of the associated terms in the principal part of the local truncation error are the same for any method of order five. However, conditions (46–49) can be satisfied by taking

$$\begin{aligned} a_2 &= \frac{-519 \pm 226\sqrt{6}}{300}, \quad a_{12} = \frac{-103 \pm 42\sqrt{6}}{864}, \\ a_{31} &= \frac{-3(119 + 660a_3) \pm 5(25 - 144a_3)\sqrt{6}}{4320}, \\ a_{50} &= \frac{(22597 + 2400a_3 - 7200a_4) \mp 8(1143 - 200a_3 + 100a_4)\sqrt{6}}{28800}, \quad (50) \end{aligned}$$

obtaining in this way three-stage methods of order five that minimize the principal part of the local truncation error. Note that now  $a_{30}$ ,  $a_{21}$  and  $a_{40}$  are given in (42) by

$$\begin{aligned} a_{30} &= \frac{221 \mp 101\sqrt{6}}{720}, \quad a_{21} = \frac{-123 \mp 22\sqrt{6}}{1440}, \\ a_{40} &= \frac{(59 - 180a_3) \mp 2(9 + 10a_3)\sqrt{6}}{720}, \quad (51) \end{aligned}$$

## 9 Methods of rational type.

In the last three sections we have only considered methods of polynomial type. Now we will consider methods of rational type, that is, methods given by (2–3), where now all the  $F_i$  (with  $2 \leq i \leq p + 1$ ) are supposed to be homogeneous functions (of degree one) of rational type. More precisely, we will consider functions  $F_i$  given in terms of the quotient of two homogeneous polynomials  $\tilde{N}_i(x_1, x_2, \dots, x_{i-1})$  and  $\tilde{D}_i(x_1, x_2, \dots, x_{i-1})$  with degrees  $r_i + 1$  and  $r_i$  respectively, for some nonnegative integer  $r_i$ , that is

$$\begin{aligned}
F_i(x_1, x_2, \dots, x_{i-1}) &= \frac{\tilde{N}_i(x_1, x_2, \dots, x_{i-1})}{\tilde{D}_i(x_1, x_2, \dots, x_{i-1})} \\
&= \frac{\sum_{j_1+j_2+\dots+j_{i-1}=r_i+1} N_{j_1 j_2 \dots j_{i-1}} x_1^{j_1} x_2^{j_2} \dots x_{i-1}^{j_{i-1}}}{\sum_{j_1+j_2+\dots+j_{i-1}=r_i} D_{j_1 j_2 \dots j_{i-1}} x_1^{j_1} x_2^{j_2} \dots x_{i-1}^{j_{i-1}}}. \quad (52)
\end{aligned}$$

It is not difficult to see that the new family of methods contains all the previous methods of polynomial type as a subfamily. Now, assuming as before that all quantities  $c_i$  are different from zero, we get from (17) and relation (18) that, in terms of  $k_1$  and  $s_i$ , any method takes the form (7–8), with the functions  $G_i$  given by

$$G_i(z_2, z_3, \dots, z_{i-1}) = c_i \left( \frac{1 + \sum_{j_2+j_3+\dots+j_{i-1}=1}^{n_i^*} n_{j_2 j_3 \dots j_{i-1}} z_2^{j_2} z_3^{j_3} \dots z_{i-1}^{j_{i-1}}}{1 + \sum_{j_2+j_3+\dots+j_{i-1}=1}^{d_i^*} d_{j_2 j_3 \dots j_{i-1}} z_2^{j_2} z_3^{j_3} \dots z_{i-1}^{j_{i-1}}} \right). \quad (53)$$

Moreover,  $n_i^*$  and  $d_i^*$  can be obtained from the functions  $F_i$  (in (52)) of the associated method, by means of

$$\begin{aligned}
n_i^* &= \max\{j_2 + j_3 + \dots + j_{i-1} / N_{j_1 j_2 \dots j_{i-1}} \neq 0\} \\
d_i^* &= \max\{j_2 + j_3 + \dots + j_{i-1} / D_{j_1 j_2 \dots j_{i-1}} \neq 0\},
\end{aligned}$$

and therefore, are always lower or equal than  $r_i + 1$  and  $r_i$  respectively.

Now it is a simple task to obtain the order conditions for the rational methods, from the order conditions for the methods of polynomial type. All we need is to consider the Taylor's expansion of the functions  $G_i(s_2, \dots, s_{i-1})$  in (53) (as functions of the  $s_j$  with  $2 \leq j \leq i-1$ ), and then compare with the associated expansion of a method of polynomial type with the same number of stages. To show this, we will obtain the order conditions for the two and three-stage methods of rational type from the order conditions of the corresponding methods of polynomial type. Note at this point that it is also possible to obtain the order conditions for general methods (that is, for arbitrarily given  $G_i$ ) in the same way, because only expansions of the  $G_i$  in terms of the  $s_j$  are involved.

## 10 Two-stage methods of rational type.

To obtain all the explicit two-stage methods (of rational type) of order three, it suffices to note that from (24) and (25) we have that  $G_2 = c_2 = 2/3$  and

also

$$G_3(s_2) = c_3 \left( \frac{1 + \sum_{i=1}^{n_3^*} n_i s_2^i}{1 + \sum_{i=1}^{d_3^*} d_i s_2^i} \right) = 1 + \frac{1}{2}s_2 + \frac{1}{6}s_2^2 + O(s_2^3), \quad (54)$$

must hold. Obviously it is enough to consider  $n_3^* = d_3^* = 2$  in (54), obtaining the following conditions

$$c_2 = 2/3, \quad (55)$$

$$c_3 = 1, \quad (56)$$

$$n_1 = 1/2 + d_1, \quad (57)$$

$$n_2 = 1/6 + (1/2)d_1 + d_2. \quad (58)$$

If we want to minimize the principal part of the local truncation error, all we need is to take  $n_3^* = d_3^* = 3$  in (54) and expand to higher order the second term, obtaining

$$n_3 = 1/24 + (1/6)d_1 + (1/2)d_2 + d_3. \quad (59)$$

The general form of a third order two-stage method of our family is given by

$$y_{n+1} = y_n + hk_1 G_3(s_2),$$

where  $k_1$ ,  $k_2$  and  $s_2$  are given by

$$k_1 = f(y_n), \quad k_2 = f\left(y_n + \frac{2}{3}hk_1\right), \quad s_2 = \frac{3(k_2 - k_1)}{2k_1},$$

and  $G_3$  takes the form

$$G_3(s_2) = \frac{1 + \frac{1+2d_1}{2}s_2 + \frac{1+3d_1+6d_2}{6}s_2^2 + \sum_{i=3}^{n_3^*} n_i s_2^i}{1 + d_1 s_2 + d_2 s_2^2 + \sum_{i=3}^{d_3^*} d_i s_2^i}. \quad (60)$$

Some third order methods with special properties (such as A-stability, L-stability, order four when applied to linear problems, etc) have been developed in [3,4], where also some numerical experiments can be found. However, when extending our methods in order to apply them to some special systems of ODE's, it is desirable to consider formulas in which the denominator in (60) is given in the form  $(1 - as_2)^\alpha$  ( $\alpha \in \mathbb{N}$ ) so that only one LU-decomposition per step is needed. Moreover, we look for a third order method being L-stable, with the previous property. We will obtain a method with the above properties later.

## 11 Three-stage methods of rational type.

Now, following our notations in section 8, we can describe any three-stage method of rational type by

$$y_{n+1} = y_n + hk_1 \tilde{G}_4(s_2, \tilde{s}_3), \quad (61)$$

with  $\tilde{s}_3 = s_3 - s_2$ , and where as usually  $s_2$  and  $s_3$  are given through (17) in terms of the stages  $k_i$  ( $1 \leq i \leq 3$ )

$$k_1 = f(y_n), \quad k_2 = f(y_n + hk_1 G_2), \quad k_3 = f(y_n + hk_1 G_3(s_2)),$$

Now the rational functions are given by

$$\begin{aligned} G_2 &= c_2, \\ G_3(s_2) &= c_3 \left( \frac{1 + \sum_{i=1}^{n_3^*} n_i s_2^i}{1 + \sum_{i=1}^{d_3^*} d_i s_2^i} \right), \\ \tilde{G}_4(s_2, \tilde{s}_3) &= c_4 \left( \frac{1 + \sum_{i+2j=1}^{\tilde{n}_4^*} n_{ij} s_2^i \tilde{s}_3^j}{1 + \sum_{i+2j=1}^{\tilde{d}_4^*} d_{ij} s_2^i \tilde{s}_3^j} \right), \end{aligned} \quad (62)$$

and it is easily seen that when looking for fifth order formulas, it is enough to consider the parameters in (62) with  $n_3^* = d_3^* = 3$  and  $\tilde{n}_4^* = \tilde{d}_4^* = 4$ , that is, the order conditions (for order five) are completely given in terms of the parameters  $c_i$  ( $2 \leq i \leq 4$ ),  $n_i$  and  $d_i$  ( $1 \leq i \leq 3$ ), and  $n_{ij}$  and  $d_{ij}$  ( $1 \leq i + 2j \leq 4$ ). In fact, we can obtain the order conditions for order five, by comparing (as in the two stage case) the expansions of the functions in (62) with those of the functions in (29) with the parameters given by (42). The  $c_i$  are given as in (42). The  $n_i$  are given in terms of  $a_1$  in (42) and the free parameters  $a_2, a_3, d_1, d_2$  and  $d_3$ , through relations

$$n_i = \sum_{j=0}^i a_j d_{i-j}, \quad 1 \leq i \leq 3, \quad (63)$$

where obviously we take  $a_0 = d_0 = 1$ . The terms  $n_{ij}$  are given by

$$n_{i_1 i_2} = \sum_{j_k=0}^{i_k} a_{j_1 j_2} d_{i_1 - j_1 i_2 - j_2}, \quad 1 \leq i_1 + 2i_2 \leq 4, \quad (64)$$

with  $a_{00} = d_{00} = 1$ , and where the parameters  $a_{ij}$  are those of (29). The  $d_{ij}$  can be arbitrarily chosen.

In order to minimize the principal part of the local truncation error, we take  $n_3^* = d_3^* = 4$  and  $\tilde{n}_4^* = \tilde{d}_4^* = 5$  in (62), obtaining the additional conditions

$$n_4 = \sum_{j=0}^4 a_j d_{4-j}, \quad n_{i_1 i_2} = \sum_{j_k=0}^{i_k} a_{j_1 j_2} d_{i_1 - j_1 i_2 - j_2}, \quad i_1 + 2i_2 = 5, \quad (65)$$

where  $a_i$  ( $1 \leq i \leq 2$ ) and  $a_{ij}$  ( $1 \leq i+2j \leq 5$ ) are given as in (42), (50) and (51), and the other parameters are free (as before, we take  $a_0 = d_0 = a_{00} = d_{00} = 1$ ).

## 12 Linear stability properties of the methods.

Now we are going to study the linear stability properties of the methods. When we apply a  $p$ -stage method (7) of our family to the scalar test equation

$$y' = \lambda y, \quad \lambda \in \mathbf{C},$$

we get

$$y_{n+1} = R(z) y_n,$$

where  $R(z)$  is the associated stability function, with  $z = h\lambda$ . Moreover, from (6–9) we obtain recursively

$$\begin{aligned} k_1 &= \lambda y_n \\ s_2 &= z G_2 \\ s_3 &= z G_3(z G_2) \\ &\vdots \\ s_p &= z G_p(z G_2, z G_3(z G_2), \dots, z G_{p-1}(z G_2, \dots, z G_{p-2}(\dots))), \end{aligned} \quad (66)$$

from which we obtain the following expression for the stability function

$$R(z) = 1 + z G_{p+1}(s_2, s_3, \dots, s_p),$$

where the  $s_i$  are given in terms of  $z$  by the relations (66).

With our change of notation for the  $s_i$  in (17), it is not difficult to see that in place of (66) we have that

$$\begin{aligned} k_1 &= \lambda y_n \\ s_2 &= z \\ s_3 &= z G_3(z) \\ &\vdots \\ s_p &= z G_p(z, z G_3(z), \dots, z G_{p-1}(z, \dots, z G_{p-2}(\dots))) \end{aligned} \quad (67)$$

(with  $z = h\lambda$ ), and so the associated stability function is now given by

$$R(z) = 1 + zG_{p+1}(s_2, s_3, \dots, s_p),$$

in terms of  $z$  through relations (67).

### 13 Two and three-stage methods being A-stable.

From the above section it is clear that the stability function of a method of polynomial type is always a polynomial function. Therefore it is not possible to obtain formulas of polynomial type with good linear stability properties such as A-stability or L-stability. However, it is also clear that we can obtain methods of rational type whose associated stability function is a rational function. Moreover, we will show in this section that we can obtain A-stable and L-stable methods of rational type, without losing the highest attainable order for a given number of stages.

As we have commented in section 10, it is possible to obtain a third order method (of rational type) being L-stable, in which the denominator of the associated stability function is given by  $(1 - as_2)^\alpha$  ( $\alpha \in \mathbb{N}$ ) so that only one LU-decomposition per step is needed. All we need is to take  $\alpha = 3$  and let the free parameters in (60) satisfy

$$d_1 = -3a, \quad d_2 = 3a^2, \quad d_3 = -a^3, \quad (68)$$

where  $a$  is the root of the polynomial  $6x^3 - 18x^2 + 9x - 1 = 0$  given by

$$\begin{aligned} a &= 1 + \frac{\sqrt{6}}{2} \sin\left(\frac{1}{3} \arctan\left(\frac{\sqrt{2}}{4}\right)\right) - \frac{\sqrt{2}}{2} \cos\left(\frac{1}{3} \arctan\left(\frac{\sqrt{2}}{4}\right)\right) \\ &\approx 0.435866521508459, \end{aligned} \quad (69)$$

and we take  $n_i = 0$  for  $i \geq 3$ , and  $d_i = 0$  for  $i \geq 4$ .

Now we will give some three-stage methods of order five being A-stable (or L-stable). For all such methods we consider the  $c_i$  given as in (42) (with the upper sign), that is

$$c_2 = \frac{6 - \sqrt{6}}{10}, \quad c_3 = \frac{6 + \sqrt{6}}{10}, \quad c_4 = 1. \quad (70)$$

We will also take  $d_i = 0$  ( $i \geq 1$ ) in all cases, obtaining that  $n_i = a_i$ .

For example, taking

$$\begin{aligned}
d_{10} &= \frac{-3}{5}, & d_{01} &= \frac{3 - 7\sqrt{6}}{30}, & d_{20} &= \frac{77 - 18\sqrt{6}}{100}, & d_{11} &= \frac{153 + 29\sqrt{6}}{360}, \\
d_{30} &= \frac{27 - 73\sqrt{6}}{600}, & d_{21} &= \frac{-44 + 3\sqrt{6}}{120}, & d_{40} &= \frac{-168 + 97\sqrt{6}}{600}, \\
n_{10} &= \frac{-3 + 2\sqrt{6}}{5}, & n_{10} &= \frac{-1}{10}, & n_{01} &= \frac{63 - 37\sqrt{6}}{180}, \\
n_{20} &= \frac{216 - 79\sqrt{6}}{300}, & n_{11} &= \frac{44 - 3\sqrt{6}}{120}, & n_{30} &= \frac{168 - 97\sqrt{6}}{600},
\end{aligned} \tag{71}$$

and the other parameters equal to zero, we obtain from (61–62) a method whose associated stability function is given by the (2, 3)-Padé approximation to the exponential function (see e.g. [9]), and thus being L-stable.

Taking

$$\begin{aligned}
d_{10} &= \frac{-2}{3}, & d_{01} &= \frac{3 - 7\sqrt{6}}{30}, & d_{20} &= \frac{41 - 9\sqrt{6}}{50}, \\
d_{11} &= \frac{431 - 59\sqrt{6}}{600}, & d_{30} &= \frac{1396 - 619\sqrt{6}}{750}, & d_{21} &= \frac{1436 - 709\sqrt{6}}{3600}, \\
d_{40} &= \frac{432353 - 178017\sqrt{6}}{180000}, & d_{31} &= \frac{-20769 + 7966\sqrt{6}}{21600}, \\
d_{50} &= \frac{127698 - 38147\sqrt{6}}{1080000}, & d_{60} &= \frac{-7193669 + 2942716\sqrt{6}}{2160000}, \\
n_{10} &= \frac{-3 + 2\sqrt{6}}{5}, & n_{20} &= \frac{-519 + 226\sqrt{6}}{300}, & n_{10} &= \frac{-1}{6}, \\
n_{01} &= \frac{63 - 37\sqrt{6}}{180}, & n_{20} &= \frac{221 - 79\sqrt{6}}{300}, & n_{11} &= \frac{3474 - 1111\sqrt{6}}{5400}, \\
n_{30} &= \frac{43409 - 18001\sqrt{6}}{18000}, & n_{21} &= \frac{20769 - 7966\sqrt{6}}{21600}, \\
n_{40} &= \frac{1892669 - 781091\sqrt{6}}{540000}, & n_{50} &= \frac{7193669 - 2942716\sqrt{6}}{2160000},
\end{aligned} \tag{72}$$

(the other parameters are zero) we get another L-stable method, being optimal with respect to the local truncation error. The associated stability function is given by the (2, 4)-Padé approximation to the exponential function.

Finally, if we take

$$\begin{aligned}
d_{10} &= \frac{-1}{2}, & d_{01} &= \frac{3 - 7\sqrt{6}}{30}, & d_{20} &= \frac{36 - 9\sqrt{6}}{50}, & d_{11} &= \frac{1323 - 247\sqrt{6}}{1800}, \\
d_{30} &= \frac{5969 - 2566\sqrt{6}}{3000}, & d_{21} &= \frac{1159 - 486\sqrt{6}}{2400}, \\
d_{40} &= \frac{480158 - 199037\sqrt{6}}{180000}, & d_{31} &= \frac{-3729 + 1411\sqrt{6}}{3600},
\end{aligned}$$

$$\begin{aligned}
d_{50} &= \frac{135777 - 46528\sqrt{6}}{720000}, & d_{60} &= \frac{-1282889 + 525021\sqrt{6}}{360000}, \\
n_1 &= \frac{-3 + 2\sqrt{6}}{5}, & n_2 &= \frac{-519 + 226\sqrt{6}}{300}, & n_{01} &= \frac{63 - 37\sqrt{6}}{180}, \\
n_{20} &= \frac{216 - 79\sqrt{6}}{300}, & n_{11} &= \frac{421 - 144\sqrt{6}}{600}, \\
n_{30} &= \frac{45569 - 18791\sqrt{6}}{18000}, & n_{21} &= \frac{3729 - 1411\sqrt{6}}{3600}, \\
n_{40} &= \frac{694953 - 286792\sqrt{6}}{180000}, & n_{50} &= \frac{1282889 - 525021\sqrt{6}}{360000}, \tag{73}
\end{aligned}$$

the resulting method is A-stable, with the property of being optimal with respect to the local truncation error, and with associated stability function given by the (3, 3)-Padé approximation to the exponential function.

## 14 Numerical experiments.

In order to show the behaviour as  $h \rightarrow 0$  for the methods explained in the last section, we will consider the following simple problem (taken from [9], pp. 134)

$$\begin{aligned}
y'(x) &= \frac{y(x)(1 - y(x))}{2y(x) - 1}, \\
y(0) &= \frac{5}{6},
\end{aligned}$$

for which the solution is

$$y(x) = \frac{1}{2} + \sqrt{\frac{1}{4} - \frac{5}{36} e^{-x}}.$$

With fixed step size  $h = 2^{-n}$  for various  $n = 1, 2, 3, \dots, 9$  over  $2^n$  steps, the value of  $y(1)$  was computed using our methods and two Runge-Kutta methods. The magnitude of the error  $E$  for different  $h$  and for each of these methods is shown in Figure 1 in logarithmic scale. The third-order method of the last section is marked *MS3*, and the fifth-order methods of this section with associated stability function given by the (2, 3), (2, 4) and (3, 3)-Padé approximation to the exponential function are marked *M23*, *M24* and *M33* respectively. For comparison purposes we also include in figure 1 a three stage third-order Runge-Kutta method (see e.g. [9], pp. 134) marked *RK3* and a six stage fifth-order Runge-Kutta method (see e.g. [9], pp. 202) marked *RK5*.

On the logarithmic scale used for this figure, the error for each method is represented very closely by a straight line whose slope equals the order of the method. It can be seen that methods marked *M24* and *M33* perform very

similarly for this problem, and the slope for the associated graphs is bigger than 5 (in fact  $\approx 5.5$ ). It is easy to explain this behaviour by noting that both methods share the property of being optimal with respect to the local truncation error. It is also clear that our methods perform better than the Runge-Kutta methods marked *RK3* and *RK5* of the same order, and this with less function evaluations (at the cost of more arithmetical operations).

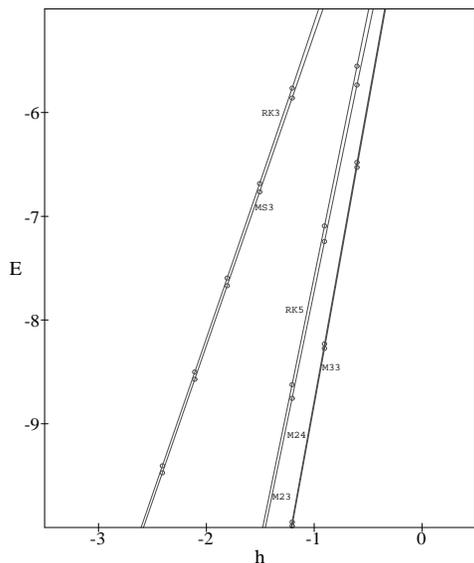


Fig. 1. Error versus step size (double logarithmic scale) for various methods.

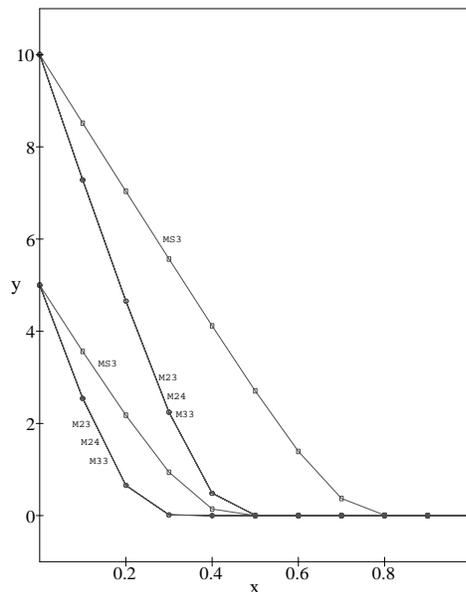


Fig. 2. Some numerical solutions for  $a = 5, 10, b = 10, c = 3000$ , with  $h = 0.1, x \in [0, 1]$ .

To show the good behaviour of the methods from a stability point of view, we will also consider the following problem

$$\begin{aligned} y'(x) &= -b y(x) \sqrt{c^2 + y^2(x)}, \\ y(0) &= a, \end{aligned} \tag{74}$$

depending on the three parameters  $a, b > 0$  and  $c > 0$ , for which the exact solution is given by

$$y(x) = \frac{a c}{c \operatorname{ch}(b c x) + \sqrt{a^2 + c^2} \operatorname{sh}(b c x)}.$$

The derivative with respect to  $y$  of the function  $f(y) = -b y \sqrt{c^2 + y^2}$  in (74) is

$$f_y = -b \frac{c^2 + 2y^2}{\sqrt{c^2 + y^2}}, \tag{75}$$

and so function  $f$  is one-sided Lipschitz continuous ( $f_y < 0$ ) with one-sided Lipschitz constant 0. Therefore, the true solutions of this non-linear problem

show a contractive behaviour. In fact, the solutions after a transient are virtually identical to the steady-state solution  $y(x) \equiv 0$  (the solution of (74) when  $a = 0$ ). We also have from (75) and from our previous comment that  $f_y \approx -bc$  along the integration (at least after the transient).

To illustrate the behaviour of the methods marked *MS3*, *M23*, *M24* and *M33* when applied to a non-linear stiff problem, we take  $b = 10$  and  $c = 3000$  in (74), and integrate this problem over  $x \in [0, 1]$  with initial conditions  $a = 5, 10$  and fixed step size  $h = 0.1$ . Figure 2 shows the good qualitative behaviour of the numerical solutions we get in this manner. The three fifth-order methods perform very similarly for this problem. Note that for the range of values of the initial condition  $a$  we are considering, we have  $f_y \approx -30000$  along the integration, and therefore the two explicit Runge-Kutta methods marked *RK3* and *RK5* give numerical overflow when applied to this problem with fixed step size  $h \geq 0.0001$ .

The numerical solutions we get from our methods when applied to problem (74) for a wide range of values of the parameters  $a$ ,  $b$  and  $c$  ( $bc \gg 1$  in order to retain the stiffness of the problem), and fixed step size  $h = 0.1$ , show that the qualitative behaviour is not always as good. For example, when  $f_{yy}$  is too big (that is,  $f_y$  is a rapidly varying function) the numerical solutions tend to the steady-state slowly. However, decreasing the step size the situation is easily solved. In fact, the step size to be used seems to depend much more on the nonlinear character of  $f$  than on the stiffness of the problem.

## 15 Two-stage methods for some special systems.

In what follows, we will obtain a generalization of our methods in order to integrate some special systems of ODEs. We begin considering autonomous systems given by

$$\begin{aligned} y'_{(1)} &= f_{11}(y_{(1)}) + f_{12}(y_{(2)}) + \dots + f_{1m}(y_{(m)}), \\ y'_{(2)} &= f_{21}(y_{(1)}) + f_{22}(y_{(2)}) + \dots + f_{2m}(y_{(m)}), \\ &\vdots \\ y'_{(m)} &= f_{m1}(y_{(1)}) + f_{m2}(y_{(2)}) + \dots + f_{mm}(y_{(m)}), \end{aligned} \tag{76}$$

that is, systems in which  $f : \mathbb{R}^m \rightarrow \mathbb{R}^m$  can be written in the form  $f(y) = F(y)\mathbb{1}$  with  $y = (y_{(1)}, y_{(2)}, \dots, y_{(m)})$ ,  $\mathbb{1} = (1, 1, \dots, 1)^T$  and where  $F$  is the

matrix

$$F(y) = \begin{pmatrix} f_{11}(y(1)) & f_{12}(y(2)) & \cdots & f_{1m}(y(m)) \\ f_{21}(y(1)) & f_{22}(y(2)) & \cdots & f_{2m}(y(m)) \\ \vdots & \vdots & & \vdots \\ f_{m1}(y(1)) & f_{m2}(y(2)) & \cdots & f_{mm}(y(m)) \end{pmatrix},$$

with  $f_{ij} : \mathbb{R} \rightarrow \mathbb{R}$ . Now we will briefly describe how we can obtain from any consistent two-stage method for scalar autonomous problems, a two-stage method being of the same order for problem (76) and retaining the linear stability properties.

From the preceding sections, it is clear that any two-stage method (of order  $p \geq 1$ ) for the scalar autonomous problem can be written in the form

$$y_{n+1} = y_n + hk_1 G_3(s_2), \quad (77)$$

with

$$k_1 = f(y_n), \quad k_2 = f(y_n + hc_2 k_1), \quad s_2 = \frac{k_2 - k_1}{c_2 k_1}.$$

The associated method for problem (76) (of the same order) is given by

$$y_{n+1} = y_n + hG_3(S_2)k_1, \quad (78)$$

where

$$k_1 = f(y_n) = F(y_n) \mathbb{1}, \quad k_2 = f(y_n + hc_2 k_1) = F(y_n + hc_2 k_1) \mathbb{1},$$

and  $S_2$  is the matrix given by

$$\begin{pmatrix} \frac{f_{11}(y_{n(1)} + hc_2 k_{1(1)}) - f_{11}(y_{n(1)})}{c_2 k_{1(1)}} & \cdots & \frac{f_{1m}(y_{n(m)} + hc_2 k_{1(m)}) - f_{1m}(y_{n(m)})}{c_2 k_{1(m)}} \\ \vdots & & \vdots \\ \frac{f_{m1}(y_{n(1)} + hc_2 k_{1(1)}) - f_{m1}(y_{n(1)})}{c_2 k_{1(1)}} & \cdots & \frac{f_{mm}(y_{n(m)} + hc_2 k_{1(m)}) - f_{mm}(y_{n(m)})}{c_2 k_{1(m)}} \end{pmatrix}$$

where  $y_n = (y_{n(1)}, y_{n(2)}, \dots, y_{n(m)})$  and  $k_1 = (k_{1(1)}, k_{1(2)}, \dots, k_{1(m)})$ .

Note that if  $G_3$  in (77) is a rational function then method (78) is linearly implicit.

It is also possible to integrate with our methods some non-autonomous scalar ODEs and systems. In fact, systems given by

$$\begin{aligned}
y'_{(1)}(x) &= f_{11}(y_{(1)}(x)) + f_{12}(y_{(2)}(x)) + \dots + f_{1m}(y_{(m)}(x)) + f_{1\ m+1}(x), \\
y'_{(2)}(x) &= f_{21}(y_{(1)}(x)) + f_{22}(y_{(2)}(x)) + \dots + f_{2m}(y_{(m)}(x)) + f_{2\ m+1}(x), \\
&\vdots \\
&\vdots \\
y'_{(m)}(x) &= f_{m1}(y_{(1)}(x)) + f_{m2}(y_{(2)}(x)) + \dots + f_{mm}(y_{(m)}(x)) + f_{m\ m+1}(x),
\end{aligned} \tag{79}$$

with  $f_{ij} : \mathbb{R} \rightarrow \mathbb{R}$ , can be rewritten in autonomous form in the usual way (that is, adding the trivial equation  $x' = 1$ ) and then integrated following our previous comments.

When we rewrite the non-autonomous problem (80) in autonomous form, the resulting system is one dimension higher. However, noting that the last row in  $S_2$  has all entries equal to zero, it can be seen that our methods can be implemented without increasing this dimension.

## 16 Numerical experiments.

Next we study the following family of problems, taken from [18], pp. 34 (see also [26], pp. 233 and [16]).

$$\begin{aligned}
y'_1 &= -(b + an)y_1 + by_2^n & y_1(0) &= c^n \\
y'_2 &= y_1 - ay_2 - y_2^n & y_2(0) &= c
\end{aligned} \tag{80}$$

for which the solution is

$$y_1(x) = c^n e^{-anx}, \quad y_2(x) = ce^{-ax}.$$

We chose  $a = 0.1$ ,  $c = 1$ ,  $n = 4$  and  $b = 1, 100, 10000, 1000000$  in the interval  $0 \leq x \leq 10$ . For increasing  $b$  the system becomes stiffer (in fact, for the eigenvalues along the true solution it holds  $\lambda_1 \approx -b$  and  $\lambda_2 \approx -a$ ) and for increasing  $n$  more nonlinear (see e.g. [18] for more details).

With fixed step size  $h = 2^{-k}$  for various  $k = 0, 1, 2, \dots, 8$  over  $10 \cdot 2^k$  steps, the value of  $y$  was computed using our two-stage method. The magnitude of the global error  $E$  ( $L_2$ -norm) for different  $h$  and for each of these problems is shown in Figure 3 in double logarithmic scale. We observe from Figure 3 that the order of convergence for problem (80) is smaller than three (in fact two) when we take  $b = 1000000$  (crosses) and  $b = 10000$  (circles) for the range of step sizes we are considering. When we take  $b = 100$  (diamonds) the order increases from two to three when we decrease the step size. This order reduction phenomenon is related to the concept of B-convergence and many implicit methods show this behaviour when applied to some stiff differential

equations. For the value  $b = 1$  (boxes) the problem is not stiff and our method shows order three as expected.

We will also consider two inhomogeneous stiff problems. Problem A is taken from [25], pp. 27, and is given by

$$y' = -10^6 y + \cos x + 10^6 \sin x, \quad y(0) = 1,$$

with solution

$$y(x) = \sin x + e^{-10^6 x}.$$

Note that this problem is a particular case of the family of scalar equations proposed by Prothero and Robinson in [23]. Problem B is given by the following inhomogeneous system

$$\begin{aligned} y_1' &= -2y_1 + y_2 + 2\sin x & y_1(0) &= 2 \\ y_2' &= 998y_1 - 999y_2 + 999(\cos x - \sin x) & y_2(0) &= 3 \end{aligned}$$

taken from [21], pp. 213, for which the exact solution is given by

$$y_1(x) = 2e^{-x} + \sin x, \quad y_2(x) = 2e^{-x} + \cos x.$$

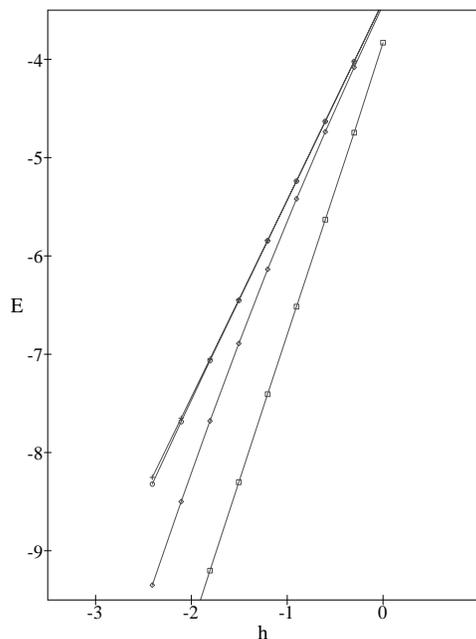


Fig. 3. Error versus step size, in double logarithmic scale, for our method (autonomous problem).

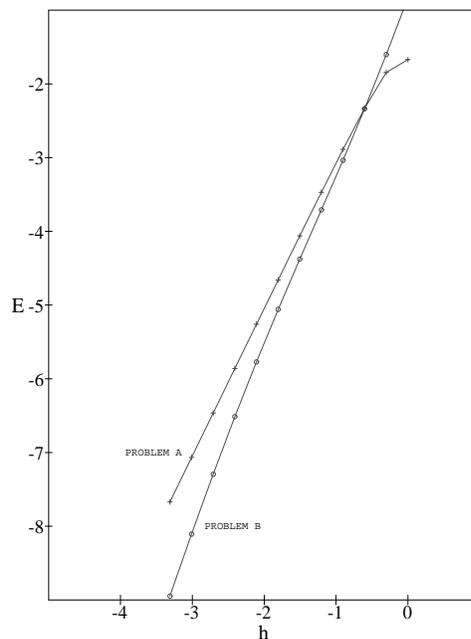


Fig. 4. Error versus step size, in double logarithmic scale, for our method (non-autonomous problems).

Note that both problems are stiff (in fact, the eigenvalues are  $\lambda = -1000000$  for problem A and  $\lambda_1 = -1$ ,  $\lambda_2 = -1000$  for problem B). We apply our two-stage method to these problems with fixed step size  $h = 2^{-k}$  for various

$k = 0, 1, 2, \dots, 11$  over  $10 \cdot 2^k$  steps. Figure 4 shows the magnitude of the global error  $E$  ( $L_2$ -norm) versus the step size in double logarithmic scale. As in a previous example, we can see that the effective order for problem A appears to be 2 for the range of values of the step size we are considering in Figure 4. Taking small enough step sizes the order changes to 3 as expected. For problem B, the effective order changes progressively from 2 to 3 when we decrease the step size, as can be seen in Figure 4. As before, this can be explained in terms of the concept of B-convergence.

## 17 Conclusions.

Our new methods are primarily useful when applied to some stiff problems for which no accurate evaluation of a Jacobian is available or the evaluation of the Jacobian is too expensive. For some non stiff problems for which function evaluations are expensive, our explicit methods might be more efficient than the usual explicit Runge-Kutta methods. This follows from the fact that for a given number of stages (or function evaluations per step), our explicit methods attain bigger order than the Runge-Kutta ones, as we have pointed out before.

Though several practical questions remain to be solved, for example, how to extend our methods in order to integrate a wider class of problems, these new methods seem quite promising, for instance in the context of solving some nonlinear parabolic equations (by the method of lines).

## References

- [1] R. Alexander, Diagonally implicit Runge-Kutta methods for stiff O.D.E.'s, *SIAM J. Numer. Anal.* **14** (1977) 1006–1021.
- [2] J. Álvarez, Obtaining New Explicit Two-Stage Methods for the Scalar Autonomous IVP with Prefixed Stability Functions, *Intl. Journal of Applied Sc. & Computations* **6** (1999) 39–44.
- [3] J. Álvarez and J. Rojo, New A-stable explicit two-stage methods of order three for the scalar autonomous IVP, in: P. de Oliveira, F. Oliveira, F. Patrício, J.A. Ferreira, A. Araújo, eds., *Proc. of the 2nd. Meeting on Numerical Methods for Differential Equations, NMDE'98* (Coimbra, Portugal, 1998) 57–66.
- [4] J. Álvarez and J. Rojo, A New Family of Explicit Two-Stage Methods of order Three for the Scalar Autonomous IVP, *Intl. Journal of Applied Sc. & Computations* **5** (1999) 246–251.
- [5] K. Burrage, The dichotomy of stiffness: Pragmatism versus theory, *Appl. Math. Comput.* **31** Spec. Issue (1989) 92–111.

- [6] J.C. Butcher, Implicit Runge-Kutta processes, *Math. Comp.* **18** (1964) 50–64.
- [7] J.C. Butcher, On the convergence of numerical solutions to ordinary differential equations, *Math. Comp.* **20** (1966) 1–10.
- [8] J.C. Butcher, An algebraic theory of integration methods, *Math. Comput.* **26** (1972) 79–106.
- [9] J.C. Butcher, *The Numerical Analysis of Ordinary Differential Equations: Runge-Kutta and General Linear Methods* (Wiley, Chichester, 1987).
- [10] R. Caira; C. Costabile; F. Costabile, A Class of pseudo Runge-Kutta methods, *BIT* **30** (1990) 642–649.
- [11] G. Dahlquist, A special stability problem for linear multistep methods, *BIT* **3** (1963) 27–43.
- [12] J.D. Day and D.N.P. Murthy, Two Classes of Internally  $S$ -Stable Generalized Runge-Kutta Processes Which Remain Consistent With an Inaccurate Jacobian, *Math. Comp.* **39** (1982) 491–509.
- [13] E. Hairer, S.P. Nørset and G. Wanner, *Solving Ordinary Differential Equations I. Nonstiff Problems* (Springer-Verlag, Berlin, 1993).
- [14] E. Hairer; G. Wanner, On the Butcher group and general multi-value methods, *Computing* **13** (1974) 1–15.
- [15] E. Hairer and G. Wanner, *Solving Ordinary Differential Equations II. Stiff and Differential-Algebraic Problems* (Springer-Verlag, Berlin, 1996).
- [16] P. Kaps, Rosenbrock-type methods, in: G. Dahlquist and R. Jeltsch, eds., *Numerical methods for solving stiff initial value problems, Proceedings, Oberwolfach 28/6–4/7 1981*, Bericht Nr. 9, Inst. für Geometrie und Praktische Mathematik der RWTH Aachen (Aachen, Germany, 1981) 5 pp.
- [17] P. Kaps and A. Ostermann, Rosenbrock Methods using few LU-Decompositions, *IMA J. Numer. Anal.* **9** (1989) 15–27.
- [18] P. Kaps, S. W. H. Poon and T. D. Bui, Rosenbrock Methods for Stiff ODEs: A Comparison of Richardson Extrapolation and Embedding Technique, *Computing* **34** (1985) 17–40.
- [19] P. Kaps and P. Rentrop, Generalized Runge-Kutta Methods of Order Four with Step-size Control for Stiff Ordinary Differential Equations, *Numer. Math.* **33** (1979) 55–68.
- [20] P. Kaps and G. Wanner, A Study of Rosenbrock-Type Methods of High Order, *Numer. Math.* **38** (1981) 279–298.
- [21] J.D. Lambert, *Numerical Methods for Ordinary Differential Systems. The Initial Value Problem* (Wiley, Chichester, 1991).
- [22] S.P. Nørsett and A. Wolfbrandt, Order Conditions for Rosenbrock Type Methods, *Numer. Math.* **32** (1979) 1–15.

- [23] A. Prothero; A. Robinson, On the stability and accuracy of one-step methods for solving stiff systems of ordinary differential equations, *Math. Comput.* **28** (1974) 145–162.
- [24] H.H. Rosenbrock, Some general implicit processes for the numerical solution of differential equations, *Comput. J.* **5** (1963) 329–330.
- [25] J.M. Sanz-Serna and M.P. Calvo, *Numerical Hamiltonian Problems* (Chapman-Hall, London, 1993).
- [26] S. Scholz, Order Barriers for the B-Convergence of ROW Methods, *Computing* **41** (1989) 219–235.
- [27] T.E. Simos, Some modified Runge-Kutta methods for the numerical solution of initial-value problems with oscillatory solutions, *J. Sci. Comput.* **13** (1998) 51–63.
- [28] T. Steihaug and A. Wolfbrandt, An Attempt to Avoid Exact Jacobian and Nonlinear Equations in the Numerical Solution of Stiff Differential Equations, *Math. Comp.* **33** (1979) 521–534.
- [29] A.M. Urbani, Metodi Row modificati per problemi stiff, *Calcolo* **27** (1990) 89–102.
- [30] J. G. Verwer, *S*-Stability Properties for Generalized Runge-Kutta Methods, *Numer. Math.* **27** (1977) 359–370.
- [31] J. G. Verwer and S. Scholz, Rosenbrock methods and time-lagged Jacobian matrices, *Beitr. Numer. Math.* **11** (1983) 173–183.
- [32] J. G. Verwer, S. Scholz, J. G. Blom and M. Louter-Nool, A Class of Runge-Kutta-Rosenbrock Methods for Solving Stiff Differential Equations, *Z. Angew. Math. Mech.* **63** (1983) 13–20.
- [33] H. Zedan, Avoiding the exactness of the Jacobian matrix in Rosenbrock formulae, *Comput. Math. Appl.* **19** (1990) 83–89.