



A Parameter Study of Explicit Runge-Kutta Pairs of Orders 6(5)

CH. TSITOURAS

Department of Mathematics, National Technical University
 Zografou Campus, GR 15780, Athens, Greece
 tsitoura@math.gsd.ntua.gr

(Received March 1997; accepted April 1997)

Communicated by P. Kaps

Abstract—Improvements over a Runge-Kutta pair of orders six and five are presented in this paper. Methods with minimised truncation errors, phase-lag errors, dissipation errors, and with extended stability regions are given and tested in the standard test problems within each category.

Keywords—Ordinary differential equations, Initial value problems, Runge-Kutta pairs, Phase-lag, Dissipation error.

1. INTRODUCTION

Explicit Runge-Kutta (RK) pairs are widely used for the numerical solution of the initial value problem $y' = f(x, y)$, $y(x_0) = y_0 \in \mathbb{R}^m$, $x \in [x_0, x_e]$, where $f : \mathbb{R} \times \mathbb{R}^m \mapsto \mathbb{R}^m$. These pairs are characterised by the extended Butcher tableau

$$\begin{array}{c|c} c & A \\ \hline & b \\ & \hat{b} \end{array}$$

with $b^\top, \hat{b}^\top, c \in \mathbb{R}^s$, and $A \in \mathbb{R}^{s \times s}$ is strictly lower triangular. The procedure that advances the solution from (x_n, y_n) to $x_{n+1} = x_n + h_n$ computes at each step two approximations y_{n+1}, \hat{y}_{n+1} to $y(x_{n+1})$ of orders p and $p - 1$, respectively, given by $y_{n+1} = y_n + h_n \sum_{i=1}^s b_i f_{ni}$ and $\hat{y}_{n+1} = y_n + h_n \sum_{i=1}^s \hat{b}_i f_{ni}$ with $f_{ni} = f(x_n + c_i h_n, y_n + h_n \sum_{j=1}^{i-1} a_{ij} f_{nj})$, $i = 1, 2, \dots, s$. From this embedded form we can obtain an estimate $E_{n+1} = y_{n+1} - \hat{y}_{n+1}$ of the local truncation error of the $p - 1$ order formula. So the step-size control algorithm $h_{n+1} = 0.9 \cdot h_n \cdot (\text{TOL}/E_{n+1})^{1/p}$ is in common use, with TOL being the requested tolerance. The above formula is used even if TOL is exceeded by E_{n+1} . Then h_{n+1} is used as current step-size.

2. PAIRS OF ORDERS 6(5)

In order to construct a 6(5) pair, 37 equations for the 6th order formula and 17 equations for the 5th order formula have to be solved. These nonlinear equations involve b, A, c for the higher order and \hat{b}, A, c for the lower order formula, and can be found easily in the bibliography,

The author would like to thank S. N. Papakostas for helpful suggestions and discussions. The algorithm and the methods presented in this paper can be requested from the author by e-mail.

i.e., [1]. Usually in parallel, 48 equations of 7th order are to be minimised in order to reduce the truncation error of the 6th order method which is used to advance the solution. Many authors in the last few years have dealt with RK pairs of orders 6 and 5. See [2–5]. Some new families of solutions for these sets of equations have been discovered and especially [2,3] belong to the same one, while [4] is a special case of the family studied recently in [5]. All these families use the FSAL device (First Stage As Last), so even if $s = 9$, only eight stages are used effectively every step. According now to the size of the truncation error for each individual pair suggested in [2–5] and exhaustive numerical tests between CMR6(5) (see [4]), DLMP6(5) [3], VE6(5)a [2], and PTP6(5) [5], performed in [5], the later pair is the one recommended. This is due to the one extra free parameter that this family offers in order to satisfy the RK design criteria.

Unfortunately, no explicit algorithm furnishing the coefficients of that pair can be derived. This is a drawback in the construction of pairs with various properties, such as minimized truncation errors or extended stability regions. Thereafter, using the compact theory which appeared in [6], we can give an explicit (symbolic and numerical) algorithm of the Verner-Dormand, Lockyer, McCorrigan and Prince family, which demands only the solution of linear equations.

ALGORITHM. The free parameters are $c_2, c_4, c_5, c_6, c_7, \hat{b}_9$. It is known that for this family $c_8 = c_9 = 1$, $b_2 = b_3 = b_9 = \hat{b}_2 = \hat{b}_3 = 0$, and $a_{i2} = 0$, $i = 4, 5, 6, 7, 8$.

1. Solve¹ $be = 1$, $bc = 1/2$, $bc^2 = 1/3$, $bc^3 = 1/4$, $bc^4 = 1/5$, $bc^5 = 1/6$, for b_1, b_4, b_5, b_6, b_7 , and b_8 .
2. Put $c_3 = 2/3c_4$, $a_{43} = c_4^2/(2c_3)$, $a_{32} = c_3^2/(2c_2)$.
3. Solve² $(Ac)_5 = c_5^2/2$ and $(Ac^2)_5 = c_5^3/3$ for a_{53}, a_{54} .
4. Substitute³ a_{87} from $(b(A + C - I))_7 = 0$.
5. Since $(b(C - I)A)_3 = 0$, evaluate a_{76} from

$$b(C - I)A(C - c_4I)(C - c_5I)c - (b(C - I)A)_3 = \int_0^1 (x - 1) \int_0^x (y - c_4)(y - c_5)y \, dy \, dx.$$

6. a_{86} is given from $(b(A + C - I))_6 = 0$.
7. Solve simultaneously for $\hat{b}_1, \hat{b}_4, \hat{b}_5, \hat{b}_6, \hat{b}_7, \hat{b}_8, a_{63}, a_{73}$, and a_{83} the equations:

$$\begin{aligned} \hat{b}e = 1, \quad \hat{b}c = \frac{1}{2}, \quad \hat{b}c^2 = \frac{1}{3}, \quad \hat{b}c^3 = \frac{1}{4}, \\ \hat{b}c^4 = \frac{1}{5}, \quad (b(A + C - I))_3 = 0, \quad (\hat{b}A)_3 = 0, \end{aligned}$$

$$\begin{aligned} b(C - I)A(C - c_4I)(C - c_5I)c &= \int_0^1 (x - 1) \int_0^x (y - c_5)(y - c_4)y \, dy \, dx, \\ \hat{b}A(C - c_5I)(C - c_4I)c &= \int_0^1 \int_0^x (y - c_5)(y - c_4)y \, dy \, dx. \end{aligned}$$

8. From $(Ac)_6 = c_6^2/2$, $(Ac^2)_6 = c_6^3/3$, evaluate a_{64} and a_{65} .
9. From $(Ac)_7 = c_7^2/2$, $(Ac^2)_7 = c_7^3/3$, evaluate a_{74} and a_{75} .
10. From $(Ac)_8 = c_8^2/2$, $(Ac^2)_8 = c_8^3/3$, evaluate a_{84} and a_{85} .
11. Finally from $Ae = c$, evaluate $a_{21}, a_{31}, \dots, a_{81}$.

Using a symbolic manipulation package [7], we can derive expressions for all coefficients that depend only to the free parameters.

¹ c^i is the vector with the components of c raised in i^{th} power and $e = [1, 1, \dots, 1]^T \in \mathbb{R}^s$.

² $(Ac)_5$ is the 5th component of Ac . See [5] for more details.

³ $C = \text{diag}(c)$ and I is the identity matrix of proper dimension.

3. MINIMISATION OF LOCAL TRUNCATION ERROR

Using the explicit expressions of b, \hat{b}, A, c , then the 7th order (or principal) local truncation error $\|T^{(7)}\|_2$ can be found explicitly. The expression is a little lengthy and independent of \hat{b}_9 , but it can be used easily with a minimisation package in order to find an optimal value for $\|T^{(7)}\|_2$. If we use the `constr` routine of constrained minimisation of Matlab [8], and require that the value $D_\infty = \max(\max_{i,j=1}^s |a_{ij}|, \|b\|_\infty, \|\hat{b}\|_\infty, \|c\|_\infty)$ is small enough, then the selection $c_2 = 17/183, c_4 = 18/83, c_5 = 71/125, c_6 = 42/59, c_7 = 199/200, \hat{b}_9 = 1/20$, given in [6], leads to the method P6(5). Its principal truncation error is of the same level with the error of the method PTP6(5), but it is almost four times smaller than the principal truncation error of other methods in this family [2,3]. The rest of the characteristics of the method can be found in Table 1. This method does not give any clear advantage over PTP6(5) in any of the numerical tests we have tried.

Table 1. The main characteristics of the RK pairs discussed in this paper.

Method	$\ T^{(7)}\ _2$	Stability Region	D_∞
PTP6(5)S	$4.32 \cdot 10^{-5}$	(-6.6, 0)	35.9
PTP6(5)	$1.25 \cdot 10^{-5}$	(-4.4, 0)	33.1
DLMP6(5)	$4.37 \cdot 10^{-5}$	(-4.2, 0)	12.5
CMR6(5)	$6.00 \cdot 10^{-5}$	(-4.4, 0)	16.8
VE6(5)a	$4.93 \cdot 10^{-5}$	(-4.2, 0)	29.6
NEW6(5)	$2.87 \cdot 10^{-6}$	(-4.9, 0)	208.2
P6(5)	$1.23 \cdot 10^{-5}$	(-4.4, 0)	18.4
NEW6(5)S	$5.06 \cdot 10^{-4}$	(-8.1, 0)	50.7
NEW6(5)P6A11	$1.44 \cdot 10^{-4}$	(-4.9, 0)	48.9
NEW6(5)P8A9	$4.9 \cdot 10^{-4}$	(-4.3, 0)	4.5
PTP6(5)P10A7	$1.55 \cdot 10^{-4}$	(-4.5, 0)	9.3

If we admit a little greater coefficients, say no greater than 200, then the selection $c_2 = 1/11, c_4 = 20/139, c_5 = 88/177, c_6 = 35/36, c_7 = 544/545, \hat{b}_9 = 1/20$ leads to a method with principal truncation error which is about five times smaller than all the methods known until now. We have applied the PTP6(5) and the new method to the DETEST set of test problems [9] for tolerances $10^{-10}, 10^{-12}, 10^{-14}, 10^{-16}, 10^{-18}$. According to the tests developed in [5], we notify the percentage difference in the number of function evaluations required for achieving a given maximum global error over the range of integration. This percentage is called efficiency gain, and it is recorded for each problem and accuracy in Table 2 in units of 10%. In that table, positive numbers mean that the second of the two methods is superior. The final row gives the mean value of efficiency gain for each problem. The final row's first number is the average efficiency gain for all problems. The empty places are due to unavailability of data for the respective errors.

Table 2. The efficiency gains of PTP6(5) relative to NEW6(5) for the 25 problems of DETEST and for tolerances $10^{-10}, 10^{-12}, \dots, 10^{-18}$.

	A_1	A_2	A_3	A_4	A_5	B_1	B_2	B_3	B_4	B_5	C_1	C_2	C_3	C_4	C_5	D_1	D_2	D_3	D_4	D_5	E_1	E_2	E_3	E_4	E_5	
-10						-1													2	1	5					
-12			-5	3		0	2		4	-4	3	2			0	1	2	2	1	4	2	-2	3			
-14	5	3	-5	3	5	0	3	3	5	-2	3	3	3	3	0	0	2	1	1	3	2	-1	4	4	1	
-16	4	3	-5	1	6	1	3	4	6	0	3	4	3	3	0	-1	1	1	1		2	0	5	4	3	
-18	3	3	-4	1	7		3	4			4	4	3	3	0						3			3		
19.9%	4	3	-5	2	5	0	3	4	5	-2	3	3	3	3	0	0	2	2	1	4	2	-1	4	4	2	

We observe that an almost 20% reduction of the cost has been gained by the new method. This difference is remarkable for methods of the same algebraic order.

4. METHODS WITH EXTENDED STABILITY REGIONS

For mildly stiff problems an extended stability region is needed. Again, the explicit algorithm can help us to evaluate bA^5c and bA^6c in terms of c_2, c_4, c_5, c_6, c_7 . So the stability polynomial $p(x) = 1 + x + x^2/2! + x^3/3! + x^4/4! + x^5/5! + x^6/6! + x^7bA^5c + x^8bA^6c$ depends only on the five nodes, and $|p(x)| < 1$ must hold for as much as possible values of $x \in C^-$. Choosing in advance $bA^5c = 1/6331$ and $bA^6c = 1/128550$, we ensure an area of 43.95 for the stability region in the left complex plane. We also observe that for $z \in [-8.1, 0] \subset \mathbb{R}^-$, the required inequality holds for all $p(z)$. This interval is called real stability interval and is the longest from all known methods. Then the two equations are solved directly for

$$c_5 = \frac{303888c_4}{(151944 - 3085200c_4 + 27128335c_4^2 - 54256670c_4^3)}$$

and

$$\begin{aligned} c_6 = & 24 (438490994730048 - 14165385440558976c_4 + 324642117388786080c_4^2 \\ & - 4945954966284568320c_4^3 + 44911241045914807915c_4^4 - 220996382998067560650c_4^5 \\ & + 538926866232165410350c_4^6 - 508883495861727676300c_4^7) / \\ & (27128335c_4 (455832 - 561456c_4 - 18011695c_4^2 + 48178910c_4^3) \\ & (151944 - 3996864c_4 + 34818175c_4^2 - 108513340c_4^3 + 108513340c_4^4)). \end{aligned}$$

Again using the routine `constr` of Matlab, we conclude that the selection $c_2 = -25/204$, $c_4 = 620/261$, $c_7 = 452/385$, and $\hat{b}_9 = 1/20$ leads to the minimum $\|T^{(7)}\|_2$ for the new method NEW6(5)S as we see in Table 1. Testing PTP6(5) and NEW6(5)S in the linear DETEST problems A₁, B₂, C₁, C₂, C₃, and C₄, where it is expected such methods to perform better than the conventional ones, we get the corresponding efficiency gains in the left-hand of Table 3.

Table 3. On the left, we present the efficiency gains of PTP6 (5) relative to NEW6(5)S for the six problems of linear DETEST for tolerances $10^{-2}, \dots, 10^{-7}$, while on the right, we give PTP6(5)P10A7 vs. NEW6(5)P8A9 efficiency gains table.

	A ₁	B ₂	C ₁	C ₂	C ₃	C ₄
-3	0	3	3	2	1	1
-4	0	2	2	1	4	4
-5	-1	3	0	2	3	3
-6	-1	2	-1	2	3	3
-7	-1	1	-2	2	4	3
15.7%	-1	2	0	2	3	3

	G ₁	G ₂	G ₃	G ₄	G ₅
-3		2	1		
-4	1	1	1	1	-1
-5	1	1	1	1	1
-6	2	2	1	1	
-7	2	2	1	2	0
-8	2	3	1	2	0
-9				0	
10.4%	1	2	1	1	0

It is obvious that the extended stability region helps NEW6(5)S to give better results. In [5] PTP6(5)S, a method with extended stability interval was also suggested, but it was only 5.7% better than PTP6(5). The area in the left complex plane for that method was about 32.5 and it was the largest among the known pairs 6(5).

5. METHODS FOR PERIODIC PROBLEMS

For these problems, it is instructive to examine the performance of a RK pair to the test problem $y' = iwy$, $y(0) = c_0$, and $w, c_0 \in \mathbb{R}$ with exact solution $y(x) = c_0 \exp(iwx)$. The application of a RK pair to this problem leads to a numerical scheme of the form $y_n = p(v_n)y_{n-1}$, with $v_n = wh_n$ and $p(v_n) = q(v_n) + ir(v_n)$. For a generic $v = v_n$, the quantities $\delta(v) = v - \arg(p(v))$

and $a(v) = 1 - |p(v)|$ are called *phase-lag* and *amplification* or *dissipation error*, respectively. If $\delta(v) = 0(v^{p+1})$, we have phase-lag order p , while $a(v) = 0(v^{t+1})$ implies dissipation order t . Recently the method PTP6(5)P10A7 derived from the family discussed here was presented in [10], and it was of the highest possible phase-lag order 10, but only of 7th amplification order. PTP6(5) and P6(5) share 7th dissipation order and 6th phase-lag order only. Here we present the method NEW6(5)P8A9 of 8th phase-lag order and 9th amplification error order which outperforms the previous method in a periodic set of test problems [10], including Model, Inhomogeneous, Wave, Bessel, and Duffing equations.

The theory appearing in [10,11] informs us that two equations have to be solved for achieving the NEW6(5)P8A9: $bA^5c = 1/7!$ and $bA^6c = 1/8!$. These equations are solved giving $c_5 = -2c_4/(-1 + 8c_4 - 56c_4^2 + 112c_4^3)$ and

$$c_6 = -\frac{-3 + 60c_4 - 792c_4^2 + 6616c_4^3 - 34032c_4^4 + 109760c_4^5 - 211456c_4^6 + 188160c_4^7}{56(1 - 4c_4)c_4(1 - 10c_4 + 42c_4^2 - 56c_4^3)(3 - 16c_4 + 4c_4^2 + 72c_4^3)}.$$

The minimisation of $\|T^{(7)}\|_2$ results the coefficients $c_2 = 3/31, c_4 = 8/29, c_7 = 2/17$, with $\hat{b}_9 = 1/20$ as always. As it was done in [10], we compare for tolerances $10^{-3}, 10^{-4}, \dots, 10^{-9}$, the method PTP6(5)P10A7 appeared there with the new method and we obtain an average 10% superiority of the later one. More details are given in the right-hand of Table 3. The new formula is also about 25% better than the conventional method PTP6(5) of [5].

Another method NEW6(5)P6A11 with 11th order dissipation error was also found sharing the coefficients $c_2 = -5/18, c_4 = 2/15, c_5 = 180/389, c_6 = 163679/242280, c_7 = 125/126, \hat{b}_9 = 1/20$, that satisfy $bA^5c = 11/57600$ and $bA^6c = 1/57600$. The results of this method were inferior than the other methods with special properties for periodic problems. So it is obvious that the method balancing the increment in the order of the two types of “periodic” errors is the best choice.

REFERENCES

1. E. Hairer, S.P. Nørsett and G. Wanner, *Solving Ordinary Differential Equations I*, Second edition, Springer, Berlin, (1993).
2. J.H. Verner, Some Runge-Kutta formula pairs, *SIAM J. Numer. Anal.* **28**, 496–511 (1991).
3. J.R. Dormand, M.R. Lockyer, N.E. McCarrigan and P.J. Prince, Global error estimation with Runge Kutta triples, *Computers Math. Applic.* **18** (9), 835–846 (1989).
4. M. Calvo, J.I. Montijano and L. Randez, A new embedded pair of Runge-Kutta formulas of order 5 and 6, *Computers Math. Applic.* **20** (1), 15–24 (1990).
5. S.N. Papakostas, Ch. Tsitouras and G. Papageorgiou, A general family of explicit Runge-Kutta pairs of orders 6(5), *SIAM J. Numer. Anal.* **33**, 917–936 (1996).
6. S.N. Papakostas, Ph.D. Dissertation, National Technical University Athens, Athens, (1996).
7. S. Wolfram, *Mathematica. A System for Doing Mathematics by Computer*, Second edition, Addison-Wesley, Redwood City, CA, (1991).
8. *Matlab, User's Guide*, The MathWorks Inc., Natick, MA, (1991).
9. W.H. Enright and J.D. Pryce, Two FORTRAN packages for assessing initial value methods, *ACM Trans. Math. Software* **13**, 1–27 (1987).
10. G. Papageorgiou, Ch. Tsitouras and S.N. Papakostas, Runge-Kutta pairs for periodic initial value problems, *Computing* **51**, 151–163 (1993).
11. P.J. van der Houwen and B.P. Sommeijer, Explicit Runge-Kutta (-Nystrom) methods with reduced phase-errors for computing oscillating solutions, *SIAM J. Numer. Anal.* **24**, 407–442 (1987).