

ONE- AND TWO-STAGE IMPLICIT INTERVAL METHODS OF RUNGE-KUTTA TYPE

ANDRZEJ MARCINIAK¹, BARBARA SZYSZKA²

¹ *Institute of Computing Science,* ² *Institute of Mathematics*
Poznań Univeristy of Technology
Piotrowo 3a, 60-965 Poznań, Poland

Abstract: The paper presents one- and two-stage implicit interval methods of Runge-Kutta type. It is shown that the exact solution of the initial value problem belongs to interval-solutions obtained by both kinds of these methods. Moreover, some approximations of the widths of interval-solutions are given.

1. INTRODUCTION

Interval methods for solving the initial value problem are interesting due to interval-solutions obtained by such methods which contain their errors. Computer implementations of interval methods in floating-point interval arithmetic together with the representation of initial data in the form of minimal machine intervals, i. e. by intervals which ends are equal or neighboring machine numbers, yield interval solutions which contain all possible numerical errors.

Explicit interval methods of Runge-Kutta type have been considered and analysed by Šokin [3,7]. In this paper we try to extend his approach for implicit methods. A reason to do this follows from a well-known fact concerning conventional implicit Runge-Kutta methods - higher orders of accuracy can be obtained than for explicit methods.

This paper is dealt with one- and two-stage implicit interval methods of Runge-Kutta type, which are presented in sections 3 and 4. We prove that the exact solution of the initial value problem belongs to interval-solutions obtained by both kinds of these methods (section 5). In section 6 some approximations of the widths of interval-solution are given.

2. THE INITIAL VALUE PROBLEM AND CONVENTIONAL RUNGE-KUTTA METHODS

As is well-known (see e. g. [4]), the initial value problem consists in finding the function $y = y(x)$, such that

$$y' = f(t, y(t)), \quad y(0) = y_0, \quad (1)$$

where $t \in [0, T]$, $y \in \mathbb{R}^N$ and $f: [0, T] \times \mathbb{R}^N \rightarrow \mathbb{R}^N$. We will assume that the solution of (1) exists and is unique. From the theory of ordinary differential equations it is known that these conditions are fulfilled if the function f is determined and continuous in the set $\{(t, y): 0 \leq t \leq T, y \in \mathbb{R}^N\}$ and there exists a constant $L > 0$ such that for each $t \in [0, T]$ and all $y_1, y_2 \in \mathbb{R}^N$ we have

$$\|f(t, y_1) - f(t, y_2)\| \leq L \|y_1 - y_2\|.$$

To carry out a single step by a conventional, m -stage Runge-Kutta method we apply the formula (see e. g. [1])

$$y_{k+1} = y_k + h \sum_{i=1}^m w_i \kappa_i, \quad k = 0, 1, \dots, \quad (2)$$

where

$$\kappa_i = f(t_k + c_i h, y_k + h \sum_{j=1}^s a_{ij} \kappa_j), \quad i = 1, 2, \dots, m, \quad (3)$$

$$\sum_{i=1}^m w_i = 1, \quad c_i = \sum_{j=1}^s a_{ij}, \quad (4)$$

and where $s = i - 1$ for an explicit method, and $s = m$ for an implicit one. The set of numbers w_i, c_i, a_{ij} are constants which characterize a particular method.

The local truncation error of step $k + 1$ for a Runge-Kutta method (explicit and implicit) of order p can be written in the form (see e.g. [1] or [4])

$$\begin{aligned} r_{k+1}(h) &= \Psi(t_k, y(t_k)) h^{p+1} + O(h^{p+2}) \\ &= r_{k+1}^{(p+1)}(0) \frac{h^{p+1}}{(p+1)!} + r_{k+1}^{(p+1)}(\theta h) \frac{h^{p+2}}{(p+2)!}, \end{aligned} \quad (5)$$

where

$$0 < \theta < 1, \quad \left| \frac{r_{k+1}^{(p+2)}(\theta h)}{(p+2)!} \right| \leq M. \quad (6)$$

This error is equal to the difference between the exact value $y(t_k + h)$ and its approximation evaluated on the basis of the exact value $y(t_k)$. The function $\Psi(t, y(t))$ depends on coefficients w_i, c_i, a_{ij} , and on partial derivatives of the function $f(t, y)$ occurring in (1). The form of $\Psi(t, y(t))$ is very complicated and cannot be written in a general form for an arbitrary p (see e.g. [1], [4] or [5]).

3. ONE-STAGE IMPLICIT INTERVAL METHODS

Let us denote:

Δ_t and Δ_y - sets in which the function $f(t, y)$ is defined, i. e.

$$\begin{aligned} \Delta_t &= \{t \in \mathbf{R}: 0 \leq t \leq a\}, \\ \Delta_y &= \{y = (y_1, y_2, \dots, y_N)^T \in \mathbf{R}^N: \underline{b}_i \leq y_i \leq \bar{b}_i, \quad i = 1, 2, \dots, N\}, \end{aligned}$$

$F(T, Y)$ - an interval extension of $f(t, y)$
 $\Psi(T, Y)$ - an interval extension of $\Psi(t, y)$ (see (5)),

and let us assume that:

- the function $F(T, Y)$ is defined and continuous for all $T \subset \Delta_t$ and $Y \subset \Delta_y$,
- the function $F(T, Y)$ is monotonic with respect to inclusion, i. e.

$$T_1 \subset T_2 \wedge Y_1 \subset Y_2 \Rightarrow F(T_1, Y_1) \subset F(T_2, Y_2),$$

- for each $T \subset \Delta_t$ and for each $Y \subset \Delta_y$ there exists a constant $L > 0$ such that

$$d(F(T, Y)) \leq L(d(T) + d(Y)),$$

where $d(A)$ denotes the width of A (if $A = (A_1, A_2, \dots, A_N)^T$, then the number $d(A)$ is defined by $d(A) = \max_{i=1, 2, \dots, N} d(A_i)$),

- the function $\Psi(T, Y)$ is defined for all $T \subset \Delta_t$ and $Y \subset \Delta_y$,
- the function $\Psi(T, Y)$ is monotonic with respect to inclusion.

From the theory of conventional Runge-Kutta methods it is known that a one-stage implicit method is of order 2 if and only if $c_1 - a_{11} = 1/2$ (for other values we get methods of the first order, and this case is not interesting because of the Euler method, which is explicit one). Moreover, from (4) it follows that in the case of one-stage method we have $w_1 = 1$. For these values of coefficients, $t_0 = 0$ and $y_0 \in Y_0$, implicit one-stage interval method of Runge-Kutta type we define as follows:

$$\begin{aligned} Y_n(t_0) &= Y_n(0) = Y_0, \\ Y_n(t_{k+1}) &= Y_n(t_k) + h K_{1,k}(h) + (\Psi(T_k, Y_n(t_k)) + [-\alpha, \alpha]) h^3, \\ k &= 0, 1, \dots, n-1, \end{aligned} \tag{7}$$

where

$$\begin{aligned} K_{1,k}(h) &= F\left(T_k + \frac{h}{2}, Y_n(t_k) + \frac{h}{2} K_{1,k}(h)\right), \\ \alpha &= Mh_0. \end{aligned} \tag{8}$$

¹⁾ An interval extension of the function

$$f: \mathbf{R} \times \mathbf{R}^N \supset \Delta_t \times \Delta_y \rightarrow \mathbf{R}^N$$

we call a function

$$F: I(\mathbf{R}) \times I(\mathbf{R}^N) \supset I(\Delta_t) \times I(\Delta_y) \rightarrow I(\mathbf{R}^N)$$

such that

$$(t, y) \in (T, Y) \Rightarrow f(t, y) \in F(T, Y).$$

$I(\mathbf{R})$ and $I(\mathbf{R}^N)$ denote the space of real intervals and the space of N -dimensional real interval vectors, respectively.

The step-size h of the method (7), where $0 < h \leq h_0$, h_0 denotes a given number, is calculated from the formula

$$h = \frac{\eta_1^*}{n}, \quad (9)$$

where

$$\eta_1^* = \min \{ \eta_0, \eta_1 \}, \quad (10)$$

The number $\eta_1 > 0$ is evaluated in such a way that

$$Y_0 + \frac{\eta_1}{2} F(\Delta_t, \Delta_y) \subset \Delta_y, \quad (11)$$

and the number η_0 — from the relation

$$Y_0 + \eta_0 F(\Delta_t, \Delta_y) + (\Psi(\Delta_t, \Delta_y) + [-\alpha, \alpha])h_0^2 \subset \Delta_y, \quad (12)$$

for $Y_0 \subset \Delta_y$ and $y_0 \in Y_0$.

The interval $[0, \eta_1^*]$ we divide into n parts by the points $t_k = kh$ ($k = 0, 1, \dots, n$), and the intervals T_k , which appear in the method (7) - (8), we choose in such a way that

$$t_k = kh \in T_k \subset [0, \eta_1^*].$$

From (8) it follows that in each step of the method we have to solve a (vector) interval equation of the form

$$X = G(T, X),$$

where

$$T \in \Delta_t \subset I(\mathbf{R}), X = (X_1, X_2, \dots, X_N)^T \in I(\Delta_y) \subset I(\mathbf{R}^N), \\ G: I(\Delta_t) \times I(\Delta_y) \rightarrow I(\mathbf{R}^N).$$

If we assume that the function G is a contraction mapping, then the well-known fixed-point theorem implies that the iteration process

$$X^{(l+1)} = G(T, X^{(l)}), \quad l = 0, 1, \dots, \quad (13)$$

is convergent to X^* , i. e. $\lim_{l \rightarrow \infty} X^{(l)} = X^*$, for an arbitrary choice of $X^{(0)} \in I(\Delta_y)$. Let us

remind that G is called a contraction mapping if

$$\rho(G(T, X_{(1)}), G(T, X_{(2)})) \leq \alpha \rho(X_{(1)}, X_{(2)}),$$

where ρ is a metric¹⁾, and $\alpha < 1$ denotes a constant.

For the equation (8) the process (13) is of the form

$$K_{1,k}^{(l+1)}(h) = F\left(T_k + \frac{h}{2}, Y_n(t_k) + \frac{h}{2} K_{1,k}^{(l)}\right),$$

$$k = 0, 1, \dots, n-1, \quad l = 0, 1, \dots,$$

where we choose

$$K_{1,k}^{(0)} = F\left(T_k + \frac{h}{2}, Y_n(t_k)\right).$$

4. TWO-STAGE IMPLICIT INTERVAL METHODS

From the theory of conventional Runge-Kutta methods it follows that two-stage implicit methods, which are characterized by the set of coefficients w_i, c_i, a_{ij} ($i, j = 1, 2$), can have order up to four, and the maximum order condition ($p = 4$) is fulfilled if

$$w_1 = w_2 = \frac{1}{2}, \quad c_1 = \frac{1}{2} \pm \frac{\sqrt{3}}{6}, \quad c_2 = \frac{1}{2} \mp \frac{\sqrt{3}}{6},$$

$$a_{11} = a_{22} = \frac{1}{4}, \quad a_{12} = \frac{1}{4} \pm \frac{\sqrt{3}}{6}, \quad a_{21} = \frac{1}{4} \mp \frac{\sqrt{3}}{6}.$$

In general, for $t_0 = 0$ and $y_0 \in$ conventional ones we determine by the following formulas:

$$Y_n(t_0) = Y_n(0) = Y_0,$$

$$Y_n(t_{k+1}) = Y_n(t_k) + h\left(w_1 K_{1,k}(h) + w_2 K_{2,k}(h)\right) + \left(\Psi(T_k, Y_n(t_k)) + [-\alpha, \alpha]\right) h^{p+1}, \quad (14)$$

$$k = 0, 1, \dots, n-1,$$

¹⁾ In the space $I(\mathbf{R})$ the distance between intervals $A = [\underline{a}, \bar{a}]$ and $B = [\underline{b}, \bar{b}]$ is determined by

$$\rho(A, B) = \max \left\{ \left| \underline{a} - \underline{b} \right|, \left| \bar{a} - \bar{b} \right| \right\},$$

where $\rho : I(\mathbf{R}) \times I(\mathbf{R}) \rightarrow \mathbf{R}$ defines a metric. The space $I(\mathbf{R})$ with the metric ρ is the complete metric space. If A and B are interval vectors, i.e.

$$A = (A_1, A_2, \dots, A_N)^T \in I(\mathbf{R}^N) \quad \text{and} \quad B = (B_1, B_2, \dots, B_N)^T \in I(\mathbf{R}^N),$$

then the distance between them is defined by the formula

$$\rho(A, B) = \max_{i=1, 2, \dots, N} \rho(A_i, B_i).$$

where $p \leq 4$, and

$$\begin{aligned} K_{1,k}(h) &= F\left(T_k + c_1 h, Y_n(t_k) + h(a_{11} K_{1,k}(h) + a_{12} K_{2,k}(h))\right), \\ K_{2,k}(h) &= F\left(T_k + c_2 h, Y_n(t_k) + h(a_{21} K_{1,k}(h) + a_{22} K_{2,k}(h))\right), \\ i &= 1, 2, \dots, m, \\ \alpha &= M h_0. \end{aligned} \quad (15)$$

The step-size h , $0 < h \leq h_0$, where h_0 is given, can be found from the formula

$$h = \frac{\eta_2^*}{n}, \quad (16)$$

where

$$\eta_2^* = \min \{\eta_0, \eta_1, \eta_2\}. \quad (17)$$

The numbers $\eta_1 > 0$ and $\eta_2 > 0$ should fulfill the conditions

$$Y_0 + \eta_i c_i F(\Delta_t, \Delta_y) \subset \Delta_y, \quad i = 1, 2, \quad (18)$$

and the number η_0 should be chosen in such a way that

$$Y_0 + \eta_0 \left(w_1 F(\Delta_t, \Delta_y) + w_2 F(\Delta_t, \Delta_y) \right) + \left(\Psi(\Delta_t, \Delta_y) + [-\alpha, \alpha] \right) h_0^p \subset \Delta_y, \quad (19)$$

for $Y_0 \subset \Delta_y$ and $y_0 \in Y_0$.

As for the one-stage method described in the previous section, the interval $[0, \eta_2^*]$ is divided into n parts by the points $t_k = kh$ ($k = 0, 1, \dots, n$), and the intervals T_k in the method (14) - (15) should be such that

$$t_k = kh \in T_k \subset [0, \eta_2^*].$$

As previously, at each step of the method (14)-(15) one should apply the process (13).

5. THE EXACT SOLUTION VS. INTERVAL SOLUTIONS

For the methods (7)-(8) and (14)-(15) we can prove that the exact solution of the initial value problem (1) belongs to the intervals obtained by these methods. Let us note that in the proof of the theorem below there are no restrictions to one- or two-stage implicit interval methods of Runge-Kutta type, and in the same way we can prove this theorem for any arbitrary number of stages.

Theorem 1. For the exact solution $y(t)$ of the initial value problem (1) we have

$$y(t_k) \in Y_n(t_k) \quad (k = 0, 1, \dots, n), \quad \text{where } Y_n(t_k) \text{ are obtained from the method (7)-(8) or (14)-(15).}$$

Proof (induction with respect to k). For $k = 0$ we have

$$y(t_0) = y_0 \in Y_0 = Y_n(0) = Y_n(t_0).$$

Let us assume that $y(t_k) \in Y_n(t_k)$. For $y(t_{k+1})$ we have

$$\begin{aligned} y(t_{k+1}) &= y(t_k + h) = y(t_k) + \int_0^h y'(t_k + \tau) d\tau \\ &= y(t_k) + h \int_0^1 f(t_k + ch, y(t_k + ch)) dc. \end{aligned}$$

If we substitute an interpolation polynomial for the integrand above and then use the quadrature formula analogous to, we get

$$y(t_{k+1}) = y(t_k) + h \sum_{i=1}^m w_i \kappa_{ik}(h) + \rho_k(h),$$

where

$$\kappa_{ik}(h) = f(t_k + c_i h, y(t_k)) + h \sum_{j=1}^m a_{ij} \kappa_{jk}(h), \quad i = 1, \dots, m,$$

and

$$\rho_k(h) = \left[\Psi(t_k, y(t_k)) + \frac{r_k^{(p+2)}(\theta h) h}{(p+2)!} \right] h^{p+1} \quad (20)$$

is a summarized error of interpolation and integration. But $\Psi(t_k, y(t_k)) \in \Psi(T_k, Y_n(t_k))$, and from the assumption (about our methods) it follows that

$$\left| \frac{r_k^{(p+2)}(\theta h)}{(p+2)!} \cdot h \right| \leq M h \leq M h_0 = \alpha.$$

This implies that $\frac{r_k^{(p+2)}(\theta h)}{(p+2)!} \cdot h \in [-\alpha, \alpha]$. Hence, taking into account (20), we have

$$\rho_k(h) \in (\Psi(T_k, Y_n(t_k)) + [-\alpha, \alpha]) h^{p+1}.$$

Moreover, $f(t, y) \in F(T, Y)$ for each $t \in \Delta_t$ and $y \in \Delta_y$, and from the induction assumption we have $y(t_k) \in Y_n(t_k)$. Thus, we get

$$y(t_{k+1}) \in Y_n(t_k) + h \sum_{i=1}^m w_i K_{i,k}(h) + (\Psi(T_k, Y_n(t_k)) + [-\alpha, \alpha]) h^{p+1}.$$

But on the basis of (7) (in the case of one-stage method, i.e. with $m = 1$) or (14) (for two-stages methods, i.e. with $m = 2$) the interval on the right-hand side of membership operator is equal to $Y_n(t_{k+1})$.

6. WIDTHS OF INTERVAL SOLUTIONS

Before we estimate the widths of interval solutions obtained by the methods (7) and (14), let us consider the widths of intervals $K_{i,k}(h)$ given by (8) and (15). From these formulas and properties of the function F it follows that

$$d(K_{i,k}(h)) \leq L[d(T_k) + d(Y_n(t_k))] + hL \sum_{j=1}^m |a_{ij}| d(K_{j,k}(h)), \quad (21)$$

where $i = m = 1$ for the method (7)-(8), and $m = 2, i = 1, 2$ for the method (14)-(15). The inequalities (21) are of the form

$$x_i \leq \beta + \sum_{j=1}^m \alpha_{ij} x_j, \quad i = 1, 2, \dots, m,$$

and can be also written as

$$(1 - \alpha_{ii}) x_i - \sum_{\substack{j=1 \\ j \neq i}}^m \alpha_{ij} x_j \leq \beta, \quad i = 1, 2, \dots, m. \quad (22)$$

For $m = 1$ and 2 we get the following solutions of (22):

- $m = 1$

$$x \leq \frac{\beta}{1 - \alpha}, \quad 1 - \alpha > 0 \quad (23)$$

- $m = 2$

$$x_1 \leq \frac{(1 + \alpha_{12} - \alpha_{22}) \beta}{(1 - \alpha_{11})(1 - \alpha_{22}) - \alpha_{12} \alpha_{21}}, \quad (24)$$

$$x_2 \leq \frac{(1 + \alpha_{21} - \alpha_{11}) \beta}{(1 - \alpha_{11})(1 - \alpha_{22}) - \alpha_{12} \alpha_{21}},$$

where

$$1 - \alpha_{11} > 0, \quad 1 - \alpha_{22} > 0, \quad (1 - \alpha_{11})(1 - \alpha_{22}) - \alpha_{12} \alpha_{21} > 0$$

On the basis of (23) for the method (7) from (21) we have

$$d(K_{1,k}(h)) \leq \frac{L[d(T_k) + d(Y_n(t_k))]}{1 - \frac{hL}{2}}, \quad (25)$$

if

$$1 - \frac{hL}{2} > 0. \quad (26)$$

Using the inequality (25) we can prove

Theorem 2. *If $Y_n(t_k)$ ($k = 0, 1, \dots, n$) are obtained from (17)-(18), then for $h_0 < \frac{2}{L}$ we have*

$$d(Y_n(t_k)) \leq Qh^2 + Rd(Y_0) + S \cdot \max_{l=1,2,\dots,n} d(T_l), \quad (27)$$

where Q , R and S are some nonnegative constants.

Proof. From (7) we get

$$d(Y_n(t_{k+1})) \leq d(Y_n(t_k)) + hd(K_{1,k}(h)) + [d(\Psi(\Delta_t, \Delta_y)) + 2\alpha] h^3. \quad (28)$$

Since $h \leq h_0$, then from the assumption that $h_0 < \frac{2}{L}$ it follows the inequality (31), and also (30). For $h \leq h_0$ from (30) we have

$$d(K_{1,k}(h)) \leq \frac{2L}{2 - h_0L} [d(T_k) + d(Y_n(t_k))].$$

Insertion of this estimate into (33) yields

$$d(Y_n(t_{k+1})) \leq d(Y_n(t_k)) + \frac{2hL}{2 - h_0L} [d(T_k) + d(Y_n(t_k))] + [d(\Psi(\Delta_t, \Delta_y)) + 2\alpha] h^3.$$

Denoting

$$v_1 = \frac{2L}{2 - h_0L},$$

we can write the last inequality in the form

$$d(Y_n(t_{k+1})) \leq d(Y_n(t_k))(1 + v_1hL) + v_1hLd(T_k) + [d(\Psi(\Delta_t, \Delta_y)) + 2\alpha] h^3, \quad (29)$$

$$k = 0, 1, \dots, n - 1.$$

From (29) it follows that

$$d(Y_n(t_1)) \leq d(Y_n(t_0))(1 + v_1hL) + v_1hLd(T_0) + [d(\Psi(\Delta_t, \Delta_y)) + 2\alpha] h^3,$$

$$\begin{aligned}
d(Y_n(t_2)) &\leq d(Y_n(t_1))(1 + v_1 hL) + v_1 hLd(T_1) + [d(\Psi(\Delta_t, \Delta_y)) + 2\alpha] h^3 \\
&\leq \left(d(Y_n(t_0))(1 + v_1 hL) + v_1 hLd(T_0) + [d(\Psi(\Delta_t, \Delta_y)) + 2\alpha] h^3 \right) (1 + v_1 hL) \\
&\quad + v_1 hLd(T_1) + [d(\Psi(\Delta_t, \Delta_y)) + 2\alpha] h^3 \\
&= d(Y_n(t_0))(1 + v_1 hL)^2 \\
&\quad + \left(v_1 hL \cdot \max_{l=0,1} d(T_l) + [d(\Psi(\Delta_t, \Delta_y)) + 2\alpha] h^3 \right) [1 + (1 + v_1 hL)],
\end{aligned}$$

$$\begin{aligned}
d(Y_n(t_3)) &\leq d(Y_n(t_2))(1 + v_1 hL) + v_1 hLd(T_2) + [d(\Psi(\Delta_t, \Delta_y)) + 2\alpha] h^3 \\
&\leq \left[d(Y_n(t_0))(1 + v_1 hL)^2 \right. \\
&\quad \left. + \left(v_1 hL \cdot \max_{l=0,1} d(T_l) + [d(\Psi(\Delta_t, \Delta_y)) + 2\alpha] h^3 \right) (1 + (1 + v_1 hL)) \right] (1 + v_1 hL) \\
&\quad + v_1 hLd(T_2) + [d(\Psi(\Delta_t, \Delta_y)) + 2\alpha] h^3 \\
&\leq d(Y_n(t_0))(1 + v_1 hL)^3 \\
&\quad + \left(v_1 hL \cdot \max_{l=0,1,2} d(T_l) + [d(\Psi(\Delta_t, \Delta_y)) + 2\alpha] h^3 \right) \\
&\quad \cdot \left(1 + (1 + v_1 hL) + (1 + v_1 hL)^2 \right)
\end{aligned}$$

.....

Thus, for each $k=1,2,\dots, n$ we have

$$\begin{aligned}
d(Y_n(t_k)) &\leq d(Y_n(t_0))(1 + v_1 hL)^k \\
&\quad + \left(v_1 hL \cdot \max_{l=0,1,\dots,k-1} d(T_l) + [d(\Psi(\Delta_t, \Delta_y)) + 2\alpha] h^3 \right) \sum_{i=0}^{k-1} (1 + v_1 hL)^i.
\end{aligned}$$

But

$$\begin{aligned}
\sum_{i=0}^{k-1} (1 + v_1 hL)^i &= \frac{(1 + v_1 hL)^k - 1}{v_1 hL} \\
&\leq \frac{\exp(v_1 k hL) - 1}{v_1 hL} \leq \frac{\exp(v_1 n hL) - 1}{v_1 hL} = \frac{\exp(v_1 \eta_1^* L) - 1}{v_1 hL},
\end{aligned}$$

where, according to (9) and (10),

$$\eta_1^* = \min \{ \eta_0, \eta_1 \}.$$

Hence

$$d(Y_n(t_k)) \leq R d(Y_n(t_0)) + S \cdot \max_{l=0,1,\dots,k} d(T_l) + Q h^2, \quad (30)$$

where

$$R = \exp(v_1 \eta_1^* L), \quad S = R - 1, \quad Q = \frac{\exp(v_1 \eta_1^* L) - 1}{v_1 L} [d(\Psi(\Delta_t, \Delta_y)) + 2\alpha].$$

Taking into account that $T_0 = [0, 0]$, i. e. $d(T_0) = 0$, the inequality (27) follows immediately from (30). •

For the two-stage implicit interval method of the Runge-Kutta type we can prove

Theorem 3. *If $Y_n(t_k)$ ($k = 0, 1, \dots, n$) are obtained on the basis of the method (14)- (15), then for h_0 such that*

$$h_0 < \min \left\{ 1, \frac{1}{L|a_{11}|}, \frac{1}{L|a_{22}|}, \frac{1}{L(|a_{11}| + |a_{22}|) + L^2|a_{12}||a_{21}|} \right\} \quad (31)$$

we have

$$d(Y_n(t_k)) \leq Q h^p + R d(Y_0) + S \cdot \max_{l=1,2,\dots,n} d(T_l), \quad (32)$$

where Q , R and S denote some nonnegative constants.

Proof. The formulas (15) yield

$$d(Y_n(t_{k+1})) \leq d(Y_n(t_k)) + h \left(|w_1| d(K_{1k}(h)) + |w_2| d(K_{2k}(h)) \right) + [d(\Psi(\Delta_t, \Delta_y)) + 2\alpha] h^{p+1}, \quad (33)$$

where $p \leq 4$, $k = 0, 1, \dots, n - 1$. On the basis of (15) we have

$$\begin{aligned} d(K_{1k}(h)) &\leq L[d(T_k) + d(Y_n(t_k))] \\ &\quad + hL|a_{11}|d(K_{1k}(h)) + hL|a_{12}|d(K_{2k}(h)), \\ d(K_{2k}(h)) &\leq L[d(T_k) + d(Y_n(t_k))] \\ &\quad + hL|a_{21}|d(K_{1k}(h)) + hL|a_{22}|d(K_{2k}(h)). \end{aligned} \quad (34)$$

From (24) it follows that the solution of the inequalities (34) is of the form

$$d(K_{1,k}(h)) \leq \frac{1 - hL|a_{22}| + hL|a_{12}|}{(1 - hL|a_{11}|)(1 - hL|a_{22}|) - h^2 L^2 |a_{12}||a_{21}|} L[d(T_k) + d(Y_n(t_k))], \quad (35)$$

$$d(K_{2,k}(h)) \leq \frac{1 - hL|a_{11}| + hL|a_{21}|}{(1 - hL|a_{11}|)(1 - hL|a_{22}|) - h^2L^2|a_{12}||a_{21}|} L[d(T_k) + d(Y_n(t_k))], \quad (35 \text{ cont.})$$

if

$$h < \frac{1}{L|a_{11}|}, \quad h < \frac{1}{L|a_{22}|},$$

$$1 - hL(|a_{11}| + |a_{22}|) + h^2L^2(|a_{11}||a_{22}| - |a_{12}||a_{21}|) > 0.$$

The first two inequalities are fulfilled from the assumption (31) and because of $h \leq h_0$. The third inequality also follows from (31), because for $h \leq h_0$ we have

$$h < \frac{1}{L(|a_{11}| + |a_{22}|) + L^2|a_{12}||a_{21}|},$$

i. e.

$$1 - hL(|a_{11}| + |a_{22}|) - hL^2|a_{12}||a_{21}| > 0. \quad (36)$$

Since $h < 1$ (as a consequence of $h_0 < 1$), then $h^2 < h$. Thus, from (36) it also follows that

$$1 - hL(|a_{11}| + |a_{22}|) - h^2L^2|a_{12}||a_{21}| > 0,$$

and hence, obviously

$$1 - hL(|a_{11}| + |a_{22}|) - h^2L^2|a_{12}||a_{21}| + h^2L^2|a_{11}||a_{22}| > 0.$$

Taking into account that $h \leq h_0$, from (35) we get

$$d(K_{1,k}(h)) \leq \frac{1 + h_0L|a_{21}|}{(1 - h_0L|a_{11}|)(1 - h_0L|a_{22}|) - h_0^2L^2|a_{12}||a_{21}|} L[d(T_k) + d(Y_n(t_k))],$$

$$d(K_{2,k}(h)) \leq \frac{1 + h_0L|a_{21}|}{(1 - h_0L|a_{11}|)(1 - h_0L|a_{22}|) - h_0^2L^2|a_{12}||a_{21}|} L[d(T_k) + d(Y_n(t_k))].$$

Using these estimate, from the inequality (33) we obtain

$$d(Y_n(t_{k+1})) \leq d(Y_n(t_k))(1 + v_2hL) + v_2hLd(T_k) + [d(\Psi(\Delta_t, \Delta_y)) + 2\alpha]h^{p+1}, \quad (37)$$

$$k = 0, 1, \dots, n-1,$$

where

$$v_2 = \frac{|w_1|(1 + h_0L|a_{12}|) + |w_2|(1 + h_0L|a_{21}|)}{(1 - h_0L|a_{11}|)(1 - h_0L|a_{22}|) - h_0^2L^2|a_{12}||a_{21}|}. \quad (38)$$

Proceeding further as in the proof of the previous theorem we get

$$d(Y_n(t_k)) \leq R d(Y_n(t_0)) + S \cdot \max_{l=0,1,\dots,k} d(T_l) + Q h^p, \quad (39)$$

where

$$R = \exp(v_2 \eta_2^* L), \quad S = R - 1, \quad Q = \frac{\exp(v_2 \eta_2^* L) - 1}{v_2 L} [d(\Psi(\Delta_t, \Delta_y)) + 2\alpha],$$

$$\eta_2^* = \min \{ \eta_0, \eta_1, \eta_2 \}.$$

Since $d(\theta) = 0$, the inequality (31) is an obvious consequence of (39). •

7. REMARKS

Theoretical justifications presented in this paper must be accompanied by a practical realization of the methods on the computer. An appropriate object-oriented system, called OOIRK (Object-oriented interval Runge-Kutta methods), is just developed by the authors [5]. Currently this system is fully functional for a number of explicit interval methods of Runge-Kutta type, and makes possible to provide calculations in standard floating-point arithmetic (sometimes called naive arithmetic) and in interval floating-point arithmetic together with interval representations of data in the form of machine intervals.

We plan to add to this system not only one- and two-stages implicit interval methods presented in this paper, but also three- and four-stage methods, including symplectic ones. Some theoretical results for such methods already have been obtained [6], but other still wait for considerations. In our opinion, one of the main problems which should be solved concerns an iteration process used in the implicit methods. Such a process cannot be too complicated and should be possible to apply to a wide range of interval functions. The assumption about such functions in the process (13) (to make them contraction mappings) seems to be too strong.

References

- [1] J. C. Butcher, *The Numerical Analysis of Ordinary Differential Equations. Runge-Kutta and General Linear Methods*, J. Wiley & Sons, Chichester 1987.
- [2] E. Hairer, S. P. Nørsett, G. Wanner, *Solving Ordinary Differential Equations I. Nonstiff Problems*, Springer-Verlag, Berlin, Heidelberg 1987.
- [3] S. A. Kalmykov, Ju. I. Šokin, Z. H. Juldašev, *Methods of Interval Analysis* [in Russian], Nauka, Novosibirsk 1986.
- [4] A. Krupowicz, *Numerical Methods of Initial Value Problems of Ordinary Differential Equations* [in Polish], PWN, Warsaw 1986.
- [5] A. Marciniak, K. Gajda, A. Marlewski, B. Szyszka, *The Concept of an Object-Oriented System for Solving the Initial Value Problem by Interval Methods of Runge-Kutta Type* [in Polish], *Pro Dialog* 8 (1999), 39-82.
- [6] A. Marciniak, A. Marlewski, *Interval Representations of Non-Machine Numbers in Object Pascal* [in Polish], *Pro Dialog* 7 (1998), 75-100.
- [6] A. Marciniak, *Interval Methods of Runge-Kutta Type in Floating-Point Interval Arithmetics* [in Polish], Technical Report RB-027/99, Poznań University of Technology, Institute of Computing Science, Poznań 1999.
- [7] Ju. I. Šokin, *Interval Analysis* [in Russian], Nauka, Novosibirsk 1981.