

Multistep Interval Methods of Nyström and Milne-Simpson Types

Andrzej Marciniak

*Poznań University of Technology, Institute of Computing Science
Piotrowo 3a, 60-965 Poznań, Poland*

*Adam Mickiewicz University, Faculty of Mathematics and Computer Science
Umultowska 87, 61-614 Poznań, Poland*

e-mail: anmar@sol.put.poznan.pl

(Rec. 24 January 2006)

Abstract: The paper is dealt with two kinds of multistep intervals methods which can be used to solve the initial value problem in the form of intervals containing all possible numerical errors. The interval methods of Nyström type are explicit, while the methods of Milne-Simpson are implicit. It appears that we can get two families of interval methods of the second kind. For both kinds of interval methods numerical examples are presented and compared with other interval multistep method considered in previous papers of the author.

Key words: initial value problem, interval methods, floating-point interval arithmetic

I. INTRODUCTION

Interval methods for solving the initial value problem in floating-point interval arithmetic give solutions in the form of intervals which contain all possible numerical errors, i.e. representation errors, rounding errors, and errors of methods.

In a number of our previous papers we have presented interval methods of Runge-Kutta type [2, 3, 13, 15, 16, 20], and interval multistep methods of Adams type (explicit of Adams-Bashforth type [5, 18] and implicit of Adams-Moulton type [4, 6, 18]). This paper is dealt with other interval multistep methods based on the well-known conventional methods of Nyström (Sec. III) and Milne-Simpson (Sec. IV). We show that in the case of implicit interval methods of Milne-Simpson type there exist two families of such methods and one of them is better (a similar case has been shown for interval methods of Adams-Moulton type). For both kinds of the interval methods introduced, i.e. of Nyström type and of Milne-Simpson type) we prove theorems that the exact solution of the initial value problem belongs to the intervals obtained (see Theorem 1 in Sec. III and Theorem 3 in Sec. IV). We also estimate the widths of interval solutions obtained by the considered methods (see Theorem 2 in Sec. III and Theorem 5 in Sec. IV). All the proofs of these theorems are original.

On the basis of a number of numerical tests concerning one- and multidimensional problems we can conclude that for the same number of steps explicit interval methods of Nyström type are somewhat better (i.e. give the interval solution with a smaller width) than the methods of Adams-Bashforth type, and implicit interval methods of Milne-Simpson type give somewhat better results than the methods of Adams-Moulton type. Two examples that confirm this conclusion are presented in Sec. V.

II. THE INITIAL VALUE PROBLEM AND CONVENTIONAL METHODS OF NYSTRÖM AND MILNE-SIMPSON

The initial value problem is of the form

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

where $y = y(t) \in \mathbf{R}^N$. We assume that the solution of (1) exists and is unique.

As is well-known, in order to construct the Nyström explicit methods (see e.g. [1]) we start from the fact that for $t \in [t_{n-2}, t_n]$ the equation $y' = f(t, y)$ is equivalent to the equation

$$y(t_n) = y(t_{n-2}) + \int_{t_{n-2}}^{t_n} f(x, y(x)) dx. \quad (2)$$

Replacing $f(x, y(x))$ with $W(x) + r(x)$, where $W(x)$ is the interpolation polynomial of degree $k-1$ and $r(x)$ denotes the interpolation error, and then integrating the equation obtained, finally from (2) we get

$$y(t_n) = y(t_{n-2}) + h \sum_{j=0}^{k-1} v_j \nabla^j f(t_{n-1}) + h^{k+1} \left[v_k^* y^{(k+1)}(\eta_n^*) + v_k^{**} y^{(k+1)}(\eta_n^{**}) \right], \quad (3)$$

where ∇ is the backward difference operator,

$$f(t_{n-1}) = f(t_{n-1}, y(t_{n-1})),$$

h denotes a step size, $\eta_n^*, \eta_n^{**} \in (t_0, t_n)$, and

$$v_0 = 2, \quad v_j = \frac{1}{j!} \int_{-1}^1 t(t+1) \dots (t+j-1) dt,$$

$$j = 1, 2, \dots, k-1,$$

$$v_k^* = \frac{1}{k!} \int_{-1}^0 t(t+1) \dots (t+k-1) dt,$$

$$v_k^{**} = \frac{1}{k!} \int_0^1 t(t+1) \dots (t+k-1) dt.$$

The coefficients v_k^* and v_k^{**} are very important in the interval methods considered (see Sec. III).

From (3) the conventional k -step method of Nyström follows immediately. We have

$$y_n = y_{n-2} + h \sum_{j=0}^{k-1} v_j \nabla^j f_{n-1},$$

where y_n is an approximation to $y(t_n)$ and

$$f_{n-1} = f(t_{n-1}, y_{n-1}).$$

In order to obtain the Milne-Simpson implicit methods (see e.g. [1]) we can start from the exact relation containing either backward differences, i.e.

$$y(t_n) = y(t_{n-2}) + h \sum_{j=0}^k \bar{v}_j \nabla^j f(t_n) + h^{k+2} \left[\bar{v}_{k+1}^* y^{(k+2)}(\theta_n^*) + \bar{v}_{k+1}^{**} y^{(k+2)}(\theta_n^{**}) \right], \quad (4)$$

or only the values of the function, i.e.

$$y(t_n) = y(t_{n-2}) + h \sum_{j=0}^k \bar{\delta}_{kj} f(t_{n-1}) + h^{k+2} \left[\bar{v}_{k+1}^* y^{(k+2)}(\theta_n^*) + \bar{v}_{k+1}^{**} y^{(k+2)}(\theta_n^{**}) \right], \quad (5)$$

where θ_n^* and θ_n^{**} are some points in (t_0, t_n) ,

$$\bar{v}_0 = 2,$$

$$\bar{v}_j = \frac{1}{j!} \int_{-2}^0 t(t+1) \dots (t+j-1) dt, \quad j = 1, 2, \dots, k,$$

$$\bar{v}_{k+1}^* = \frac{1}{(k+1)!} \int_{-2}^{-1} t(t+1) \dots (t+k) dt,$$

$$\bar{v}_{k+1}^{**} = \frac{1}{(k+1)!} \int_{-1}^0 t(t+1) \dots (t+k) dt,$$

$$\bar{\delta}_{kj} = (-1)^j \sum_{l=j}^k \binom{l}{j} \bar{v}_l, \quad j = 0, 1, \dots, k.$$

From both of these formulas, i.e. from (4) or (5), we can get the conventional k -step method of Milne-Simpson:

$$y_n = y_{n-2} + h \sum_{j=0}^k \bar{v}_j \nabla^j f_n.$$

In interval case the formulas (4) and (5) give quite different multistep methods (see Sec. IV for details).

III. EXPLICIT INTERVAL METHODS OF NYSTRÖM TYPE

Let us denote:

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

$$\Delta_t = \{t \in \mathbf{R} : 0 \leq t \leq a\},$$

$$\Delta_y = \{y = (y_1, y_2, \dots, y_N)^T \in \mathbf{R}^N : \underline{b}_i \leq y_i \leq \bar{b}_i, \quad (6)$$

$$i = 1, 2, \dots, N\},$$

$F(T, Y)$ and $\Psi(T, Y)$ – interval extensions of $f(t, y)$ and $\psi(t, y) = f^{(k)}(t, y) \equiv y^{(k+1)}(t)$, respectively.

Let us assume that:

- $F(T, Y)$ is determined and continuous for all $T \subset \Delta_t$, and $Y \subset \Delta_y$,
- $F(T, Y)$ 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 $Y \subset \Delta_y$ there exists a constant $L > 0$ such that

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

where $d(A)$ denotes the diameter of interval A (if $A = (A_1, A_2, \dots, A_N)^T$, then $d(A)$ is defined as the maximum of $d(A_i)$, $i = 1, 2, \dots, N$),

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

The explicit interval methods of Nyström type we define as follows:

$$Y_n = Y_{n-2} + h \sum_{j=0}^{k-1} v_j \nabla^j F_{n-1} + h^{k+1} (v_k^* \Psi_k + v_k^{**} \Psi_k), \quad (7)$$

$$n = k, k+1, \dots, m,$$

where $F_{n-1} = F(T_{n-1}, Y_{n-1})$, and

$$\Psi_k = \Psi(T_{n-1} + [-(k-1)h, h],$$

$$Y_{n-1} + [-(k-1)h, h]F(\Delta_t, \Delta_y)).$$

In (7) it is assumed that for the integration interval $[0, \xi]$ the intervals Y_i such that $y(t_i) \in Y_i$ for $i = 0, 1, \dots, k-1$ are known, and that

$$h = \frac{\xi}{m}, \quad t_i = ih \in T_i, \quad i = 0, 1, \dots, m.$$

Let us note that we cannot write $(v_k^* + v_k^{**})\Psi_k$ instead of $v_k^*\Psi_k + v_k^{**}\Psi_k$, because in general $|v_k^* + v_k^{**}|$ may be different from $|v_k^*| + |v_k^{**}|$. Moreover, the formula (7) can be written in more convenient form:

$$Y_n = Y_{n-2} + h \sum_{j=1}^k \delta_{kj} F_{n-j} + h^{k+1} (v_k^* \Psi_k + v_k^{**} \Psi_k), \quad (8)$$

$$n = k, k+1, \dots, m,$$

where

$$\delta_{kj} = (-1)^{j-1} \sum_{l=j-1}^{k-1} \binom{l}{j-1} v_l,$$

$$j = 1, 2, \dots, k.$$

In particular, for a given k from (7) and (8) we have the following methods:

- $k = 1$

$$Y_n = Y_{n-2} + 2hF_{n-1} + \frac{h^2}{2}(\Psi_1 - \Psi_1),$$

where

$$\Psi_1 = \Psi(T_{n-1} + [0, h], Y_{n-1} + [0, h]F(\Delta_t, \Delta_y)),$$

- $k = 2$ (in the conventional case we have the same method as for $k = 1$)

$$Y_n = Y_{n-2} + 2hF_{n-1} + \frac{h^3}{12}(5\Psi_2 - \Psi_2),$$

where

$$\Psi_2 = \Psi(T_{n-1} + [-h, h], Y_{n-1} + [-h, h]F(\Delta_t, \Delta_y)),$$

- $k = 3$

$$Y_n = Y_{n-2} + \frac{h}{3}(7F_{n-1} - 2F_{n-2} + F_{n-3}) + \frac{h^4}{24}(9\Psi_3 - \Psi_3),$$

where

$$\Psi_3 = \Psi(T_{n-1} + [-2h, h], Y_{n-1} + [-2h, h]F(\Delta_t, \Delta_y)).$$

For explicit interval methods of Nyström type we can prove that the exact solution of the initial value problem belongs to the intervals obtained with these methods. We have

Theorem 1. *If $y(0) \in Y_0$ and $y(t_i) \in Y_i$ for $i = 1, 2, \dots, k-1$, then for the exact solution $y(t)$ of the initial value problem (1) we have $y(t_n) \in Y_n$ for $n = k, k+1, \dots, m$, where $Y_n = Y(t_n)$ are obtained from (7).*

Proof. First, let us prove that if $(t_i, y(t_i)) \in (T_i, Y_i)$ for $i = n-k, n-k+1, \dots, n-1$, where $Y_i = Y(t_i)$, then for any $j = 0, 1, \dots, k-1$ we have

$$\nabla^j f(t_{n-1}) \in \nabla^j F_{n-1}, \quad (9)$$

where $f(t_{n-1}) = f(t_{n-1}, y(t_{n-1}))$. Since $F(T, Y)$ is an interval extension of $f(t, y)$ and $(t_i, y(t_i)) \in (T_i, Y_i)$ for $i = n-k, n-k+1, \dots, n-1$, we can write

$$f(t_{n-1-m}, y(t_{n-1-m})) \in F(T_{n-1-m}, Y_{n-1-m}) \equiv F_{n-1-m},$$

$$m = 0, 1, \dots, j.$$

From this relation it follows that

$$\sum_{m=0}^j (-1)^m \binom{j}{m} f(t_{n-1-m}, y(t_{n-1-m})) \in$$

$$\in \sum_{m=0}^j (-1)^m \binom{j}{m} F_{n-1-m}. \quad (10)$$

But

$$\sum_{m=0}^j (-1)^m \binom{j}{m} f(t_{n-1-m}, y(t_{n-1-m})) = \nabla^j f(t_{n-1}), \quad (11)$$

$$\sum_{m=0}^j (-1)^m \binom{j}{m} F_{n-1-m} = \nabla^j F_{n-1}.$$

From (10) and (11) the inclusion (9) follows immediately.

Let us consider the formula (3) for $n = k$, i.e.

$$y(t_k) = y(t_{k-2}) + h \sum_{j=0}^{k-1} v_j \nabla^j f(t_{k-1}) + h^{k+1} \left[v_k^* y^{(k+1)}(\eta_k^*) + v_k^{**} y^{(k+1)}(\eta_k^{**}) \right], \quad (12)$$

where

$\eta_k^*, \eta_k^{**} \in [t_0, t_k]$, and where $y^{(k+1)}(\eta_k) \equiv \psi(\eta_k, y(\eta_k))$ for $\eta_k = \eta_k^*$ and $\eta_k = \eta_k^{**}$. From the assumption we have $y(t_{k-2}) \in Y_{k-2}$, and from (9) and the fact that all coefficients v_j are nonnegative it follows that

$$h \sum_{j=0}^{k-1} v_j \nabla^j f(t_{k-1}) \in h \sum_{j=0}^{k-1} v_j \nabla^j F_{k-1}. \quad (13)$$

Applying Taylor's formula we have

$$y(\eta_k) = y(t_{k-1}) + (\eta_k - t_{k-1}) y'(t_{k-1} + \vartheta(\eta_k - t_{k-1})), \quad (14)$$

where $\vartheta \in [0, 1]$. Because $\eta_k \in [t_0, t_k]$ and $t_i = ih$ for $i = 0, 1, \dots, m$, we get

$$\eta_k - t_{k-1} \in [-(k-1)h, h]. \quad (15)$$

Moreover, since $y'(t) = f(t, y(t))$ and

$$f(t_{k-1} + \vartheta(\eta_k - t_{k-1}), y(t_{k-1} + \vartheta(\eta_k - t_{k-1}))) \in F(\Delta_t, \Delta_y)$$

then

$$y'(t_{k-1} + \vartheta(\eta_k - t_{k-1})) \in F(\Delta_t, \Delta_y). \quad (16)$$

From (14)-(16) it follows that

$$y(\eta_k) \in Y_{k-1} + [-(k-1)h, h] F(\Delta_t, \Delta_y). \quad (17)$$

From the assumption the function Ψ is an interval extension of ψ . Thus, on the basis of (15) and (17) we have

$$v_k \Psi(\eta_k, y(\eta_k)) \in v_k \Psi(T_{k-1} + [-(k-1)h, h], Y_{k-1} + [-(k-1)h, h] F(\Delta_t, \Delta_y)),$$

where (v_k, η_k) equals either (v_k^*, η_k^*) or (v_k^{**}, η_k^{**}) . It means that

$$h^{k+1} \left[v_k^* y^{(k+1)}(\eta_k^*) + v_k^{**} y^{(k+1)}(\eta_k^{**}) \right] \in h^{k+1} (v_k^* \Psi_k + v_k^{**} \Psi_k).$$

Thus, we have shown that $y(t_k)$ belongs to the interval

$$Y_{k-2} + h \sum_{j=0}^{k-1} v_j \nabla^j F_{k-1} + h^{k+1} (v_k^* \Psi_k + v_k^{**} \Psi_k),$$

but, according to the formula (7), this is the interval Y_k . This conclusion ends the proof for $n = k$. In the same way we can show that $y(t_n) \in Y_n$ for $n = k+1, k+2, \dots, m$. ■

We can also prove the following

Theorem 2. *If the intervals Y_n are known for $n = 0, 1, \dots, k-1$, $t_i = ih \in T_i$ for $i = 0, 1, \dots, m$, $h = \xi/m$, and the intervals Y_n are obtained from (7), then*

$$d(Y_n) \leq A \max_{q=0,1,\dots,k-1} d(Y_q) + B \max_{j=1,2,\dots,m-1} d(T_j) + Ch^k, \quad (18)$$

where the constants A , B and C are independent of h .

Proof. From (8) we obtain

$$d(Y_n) \leq d(Y_{n-2}) + h \sum_{j=1}^k |\delta_{kj}| d(F_{n-j}) + h^{k+1} \left(|v_k^*| + |v_k^{**}| \right) d(\Psi_k). \quad (19)$$

Since the function Ψ is monotonic with respect to inclusion, then if the step size h is chosen in such a way that

$$T_{n-1} + [-(k-1)h, 0] \subset \Delta_t,$$

$$Y_{n-1} + [-(k-1)h, 0] \subset F(\Delta_t, \Delta_y),$$

we have $\Psi_k \subset \Psi(\Delta_t, \Delta_y)$, and hence $d(\Psi_k) \leq d(\Psi(\Delta_t, \Delta_y))$. Moreover, on the basis of the assumptions about the function F it follows that there exists a constant L such that

$$d(F_{n-j}) \leq L(d(T_{n-j}) + d(Y_{n-j})).$$

Thus, from (19) we have

$$d(Y_n) \leq d(Y_{n-2}) + hL\delta_k \sum_{j=1}^k (d(T_{n-j}) + d(Y_{n-j})) + h^{k+1} \left(|v_k^*| + |v_k^{**}| \right) d(\Psi(\Delta_t, \Delta_y)), \quad (20)$$

where $\delta_k = \max_{j=1,2,\dots,k} |\delta_{kj}|$.

For further simplicity, let us denote:

$$\delta = hL\delta_k, \quad \alpha = 1 + \delta, \quad \nu = h^{k+1} \left(|v_k^*| + |v_k^{**}| \right). \quad (21)$$

From (20) we get

$$\begin{aligned}
 d(Y_n) &\leq d(Y_{n-2}) + \delta \sum_{j=1}^k d(Y_{n-j}) + \\
 &+ \delta \sum_{j=1}^k d(T_{n-j}) + \nu d(\Psi(\Delta_t, \Delta_y)) \leq \\
 &\leq \alpha \sum_{j=1}^k d(Y_{n-j}) + \delta \sum_{j=1}^k d(T_{n-j}) + \nu d(\Psi(\Delta_t, \Delta_y)).
 \end{aligned} \tag{22}$$

Hence, for $n = k$ we have

$$\begin{aligned}
 d(Y_k) &\leq \alpha \sum_{j=1}^k d(Y_{k-j}) + \\
 &+ \delta \sum_{j=1}^k d(T_{k-j}) + \nu d(\Psi(\Delta_t, \Delta_y)),
 \end{aligned}$$

and for $n = k + 1$ we get

$$\begin{aligned}
 d(Y_{k+1}) &\leq \alpha d(Y_k) + \alpha \sum_{j=1}^{k-1} d(Y_{k-j}) + \\
 &+ \delta \sum_{j=1}^k d(T_{k+1-j}) + \nu d(\Psi(\Delta_t, \Delta_y)).
 \end{aligned}$$

If in the above inequality we take into considerations the estimation (23), then we obtain

$$\begin{aligned}
 d(Y_{k+1}) &\leq (\alpha^2 + \alpha) \sum_{j=1}^k d(Y_{k-j}) + \\
 &+ \delta \left(\alpha \sum_{j=1}^k d(T_{k-j}) + \sum_{j=1}^k d(T_{k+1-j}) \right) + \\
 &+ \nu(\alpha + 1) d(\Psi(\Delta_t, \Delta_y)).
 \end{aligned} \tag{24}$$

For $n = k + 2$ from (22) we have

$$\begin{aligned}
 d(Y_{k+2}) &\leq \alpha d(Y_{k+1}) + \alpha d(Y_k) + \\
 &+ \alpha \sum_{j=1}^{k-2} d(Y_{k-j}) + \delta \sum_{j=1}^k d(T_{k+2-j}) + \\
 &+ \nu d(\Psi(\Delta_t, \Delta_y))
 \end{aligned}$$

and taking into account (23) and (24) we get

$$\begin{aligned}
 d(Y_{k+2}) &\leq (\alpha^2 + 2\alpha + \alpha) \sum_{j=1}^k d(Y_{k-j}) + \\
 &+ \delta \left((\alpha^2 + \alpha) \sum_{j=1}^k d(T_{k-j}) + \alpha \sum_{j=1}^k d(T_{k+1-j}) + \sum_{j=1}^k d(T_{k+2-j}) \right) + \\
 &+ \nu(\alpha^2 + 2\alpha + 1) d(\Psi(\Delta_t, \Delta_y)).
 \end{aligned} \tag{25}$$

From (23), (24) and (25) it follows that in general (for each $i = 1, 2, \dots, m - 1$) we have

$$\begin{aligned}
 d(Y_{k+i}) &\leq \left(\sum_{l=0}^i \binom{i}{l} \alpha^{l+i} \right) \left(\sum_{j=1}^k d(Y_{k-j}) \right) + \\
 &+ \delta \sum_{p=0}^i \left(\sum_{l=0}^{p-1} \binom{p-1}{l} \alpha^{l+1} \right) \left(\sum_{j=1}^k d(T_{k+i-p-j}) \right) + \\
 &+ \left(\nu \sum_{l=0}^i \binom{i}{l} \alpha^l \right) d(\Psi(\Delta_t, \Delta_y)),
 \end{aligned}$$

(23) i.e., taking into account the notations (21),

$$\begin{aligned}
 d(Y_{k+i}) &\leq k \sum_{l=0}^i \binom{i}{l} (1 + hL\delta_k)^{l+1} \max_{q=0,1,\dots,k-1} d(Y_q) + \\
 &+ hL\delta_k k \sum_{p=0}^i \sum_{l=0}^{p-1} \binom{p-1}{l} (1 + hL\delta_k)^{l+1} \max_{j=0,1,\dots,k+i-1} d(T_j) + \\
 &+ h^{k+1} (|\nu_k^*| + |\nu_k^{**}|) \sum_{l=0}^i \binom{i}{l} (1 + hL\delta_k)^l d(\Psi(\Delta_t, \Delta_y)).
 \end{aligned} \tag{26}$$

Since

$$\begin{aligned}
 \binom{i}{j} &\leq i! \leq (m-k)!, \quad l = 0, 1, \dots, i, \\
 \binom{p-1}{l} &\leq p! \leq i! \leq (m-k)!, \quad l = 0, 1, \dots, p-1,
 \end{aligned} \tag{27}$$

$$(1 + hL\delta_k)^{l+1} \leq \exp(mhL\delta_k) = \exp(\xi L\delta_k),$$

$$\sum_{l=0}^{p-1} (1 + hL\delta_k)^{l+1} \leq \frac{\exp(mhL\delta_k) - 1}{hL\delta_k} = \frac{\exp(\xi L\delta_k) - 1}{hL\delta_k},$$

we have

$$k \sum_{l=0}^i \binom{i}{l} (1 + hL\delta_k)^{l+1} \leq m(m-k+1)(m-k)! \exp(\xi L\delta_k),$$

$$k \sum_{p=0}^i \sum_{l=0}^{p-1} \binom{p-1}{l} (1 + hL\delta_k)^{l+1} \leq$$

$$\leq m(m-k+1)(m-k)! \frac{\exp(\xi L\delta_k) - 1}{hL\delta_k}$$

$$\sum_{l=0}^i \binom{i}{l} (1 + hL\delta_k)^l \leq (m-k+1)(m-k)! \frac{\exp(\xi L\delta_k) - 1}{hL\delta_k}.$$

Thus, from (26) we finally get

$$d(Y_{k+i}) \leq A \max_{q=0,1,\dots,k-1} d(Y_q) + B \max_{j=0,1,\dots,m-1} d(T_j) + Ch^k,$$

where

$$A = m(m-k+1)(m-k)! \exp(\xi L \delta_k),$$

$$B = m(m-k+1)(m-k)! (\exp(\xi L \delta_k) - 1),$$

$$C = \frac{(|v_k^*| + |v_k^{**}|) (m-k+1)(m-k)!}{L \delta_k} (\exp(\xi L \delta_k) - 1).$$

Taking into account that $T_0 = [0, 0]$ i.e. $d(T_0) = 0$, from (28) the inequality (18) follows immediately. ■

IV. IMPLICIT INTERVAL METHODS OF MILNE-SIMPSON TYPE

Let $\bar{\Psi}(T, Y)$ denotes the interval extension of

$$\bar{\psi}(t, y) = f^{(k+1)}(t, y) \equiv y^{(k+2)}(t),$$

and let us assume that $\bar{\Psi}(T, Y)$ is monotonic with respect to inclusion and determined for all $T \subset \Delta_t$ and $Y \subset \Delta_y$. Using (4) we get the following implicit interval methods:

$$Y_n = Y_{n-2} + h \sum_{j=0}^k \bar{v}_j \nabla^j F_n + h^{k+2} (\bar{v}_{k+1}^* \bar{\Psi}_k + \bar{v}_{k+1}^{**} \bar{\Psi}_k), \quad (29)$$

$$n = k, k+1, \dots, m,$$

where $F_n = F(T_n, Y_n)$, and where

$$\bar{\Psi}_k = \bar{\Psi}(T_n + [-kh, 0], Y_n + [-kh, 0]F(\Delta_t, \Delta_y)).$$

The second kind of interval methods of Milne-Simpson type, based on (5), are as follows:

$$Y_n = Y_{n-2} + h \sum_{j=0}^k \bar{\delta}_{kj} F_{n-j} + h^{k+2} (\bar{v}_{k+1}^* \bar{\Psi}_k + \bar{v}_{k+1}^{**} \bar{\Psi}_k), \quad (30)$$

$$n = k, k+1, \dots, m.$$

In particular, for a given k from (29) we get the following methods of the first kind:

- $k = 1$

$$Y_n = Y_{n-2} + 2h(F_n - F_n + F_{n-1}) + \frac{h^3}{12}(5\bar{\Psi}_1 - \bar{\Psi}_1),$$

where $\bar{\Psi}_1 = \bar{\Psi}(T_n + [-h, 0], Y_n + [-h, 0]F(\Delta_t, \Delta_y))$,

- $k = 2$

$$Y_n = Y_{n-2} + \frac{h}{3}(7F_n - 6F_n + 6F_{n-1} - 2F_{n-1} + F_{n-2}) + \frac{h^4}{24}(\bar{\Psi}_2 - \bar{\Psi}_2),$$

where $\bar{\Psi}_2 = \bar{\Psi}(T_n + [-2h, 0], Y_n + [-2h, 0]F(\Delta_t, \Delta_y))$,

- $k = 3$ (in the conventional case we have the same method as for $k = 2$)

$$Y_n = Y_{n-2} + \frac{h}{3}(7F_n - 6F_n + 6F_{n-1} - 2F_{n-1} + F_{n-2}) + \frac{h^5}{720}(11\bar{\Psi}_3 - 19\bar{\Psi}_3),$$

where $\bar{\Psi}_3 = \bar{\Psi}(T_n + [-3h, 0], Y_n + [-3h, 0]F(\Delta_t, \Delta_y))$.

Below there are examples of the methods of the second kind (obtained from (30)).

- $k = 1$

$$Y_n = Y_{n-2} + 2hF_{n-1} + \frac{h^3}{12}(5\bar{\Psi}_1 - \bar{\Psi}_1),$$

- $k = 2$

$$Y_n = Y_{n-2} + \frac{h}{3}(F_n + 4F_{n-1} + F_{n-2}) + \frac{h^4}{24}(\bar{\Psi}_2 - \bar{\Psi}_2),$$

- $k = 3$

$$Y_n = Y_{n-2} + \frac{h}{3}(F_n + 4F_{n-1} + F_{n-2}) + \frac{h^5}{720}(11\bar{\Psi}_3 - 19\bar{\Psi}_3).$$

If we denote:

Y_n^1 – the interval solution obtained from (29), i.e. from the formula with backward interval differences,

Y_n^2 – the interval solution obtained from (30), i.e. from the formula without backward interval differences,

then we can prove

Theorem 3. $Y_n^2 \subseteq Y_n^1$,

which means that the second kind of implicit interval formulas gives the interval solution with a smaller diameter (width).

Proof. The proof of the theorem 3 follows immediately from the inclusion

$$\sum_{j=0}^k \bar{\delta}_{kj} F_{n-j} \subseteq \sum_{j=0}^k \bar{v}_j \nabla^j F_n. \quad \blacksquare$$

Let us note that we can get only one kind of explicit interval methods of Nyström type. It follows from the fact that for these methods we have

$$\sum_{j=1}^k \delta_{kj} F_{n-j} = \sum_{j=0}^{k-1} \nu_j \nabla^j F_{n-1}.$$

(In explicit methods all coefficients ν_j are non-negative!)

In each step of the interval methods of Milne-Simpson type (of both kinds) we have to solve a system of nonlinear interval equations of the form

$$Y = G(T, Y)$$

where

$$T \in I(\Delta_t) \subset I(\mathbf{R}), \quad Y \in I(\Delta_y) \subset I(\mathbf{R}^N),$$

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

and where $I(\mathbf{R})$ and $I(\mathbf{R}^N)$ denote the sets of all real intervals and all n -dimensional real interval vectors, respectively. An iteration process follows immediately from the well-known fixed-point theorem. If we assume that G is the contracting mapping, i.e.

$$\rho(G(T, Y), G(T, \bar{Y})) \leq \alpha \rho(Y, \bar{Y}),$$

$$T \subset \Delta_t, \quad Y, \bar{Y} \subset \Delta_y,$$

where ρ denotes a metric¹, and $\alpha < 1$ is a constant, then the sequence

$$Y^{(l+1)} = G(T, Y^{(l)}), \quad l = 1, 2, \dots, \quad (31)$$

is convergent to a unique element Y^* with an arbitrary choice of $Y^0 \in I(\Delta_y)$.

For the second kind of interval methods of Milne-Simpson type the process (31) is as follows:

$$Y_n^{(l+1)} = Y_n + h \bar{\delta}_{k0} F(T_n, Y_n^{(l)}) + h \sum_{j=1}^k \bar{\delta}_{kj} F_{n-j} + h^{k+2} \left(\bar{\nu}_{k+1}^* \bar{\Psi}^{(l)} + \bar{\nu}_{k+1}^{**} \bar{\Psi}^{(l)} \right),$$

$$l = 0, 1, \dots,$$

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

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

where $\rho : I(\mathbf{R}) \times I(\mathbf{R}) \rightarrow \mathbf{R}$ is a metric. If A and B denote interval vectors, i.e.

$$A = (A_1, A_2, \dots, A_N)^T \in I(\mathbf{R}^N),$$

$$B = (B_1, B_2, \dots, B_N)^T \in I(\mathbf{R}^N),$$

then the distance between A and B is defined by

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

where

$$\bar{\Psi}^{(l)} = \bar{\Psi}(T_n + [-kh, 0], Y_n^{(l)} + [-kh, 0]F(\Delta_t, \Delta_y)),$$

and we can choose $Y_n^{(0)} = Y_{n-1}$.

As for interval methods of Nyström type, for implicit interval methods of Milne-Simpson type we can prove that the exact solution of the initial value problem (1) belongs to the intervals obtained. We have

Theorem 4. *If $y(0) \in Y_0$ and $y(t_i) \in Y_i$ for $i = 1, 2, \dots, k-1$, then for the exact solution $y(t)$ of the initial value problem (1) we have $y(t_n) \in Y_n$ for $n = k, k+1, \dots, m$, where $Y_n = Y(t_n)$ are obtained from (29) or (30).*

The proof is similar to that of Theorem 1.

Moreover, we can estimate the widths of intervals obtained.

Theorem 5. *If the intervals Y_n for $n = 0, 1, \dots, k-1$ are known, $t_i = ih \in T_i$, $i = 0, 1, \dots, m$, $h = \xi/m$, $0 < h \leq h_0$, where*

$$h_0 < \frac{1}{L \bar{\delta}_k}, \quad \bar{\delta}_k = \max_{j=0, 1, \dots, k} |\bar{\delta}_{kj}|,$$

and Y_n for $n = k, k+1, \dots, m$ are obtained from (29) or (30), then

$$d(Y_n) \leq \bar{A} \max_{q=0, 1, \dots, k-1} d(Y_q) + \bar{B} \max_{j=1, 2, \dots, m} d(T_j) + \bar{C} h^{k+1}, \quad (32)$$

where the constants \bar{A} , \bar{B} and \bar{C} are independent of h .

Proof. We prove the estimation of $d(Y_n)$ for Y_n given by (30) (for Y_n given by (29) the proof is similar). From (30) it follows that

$$d(Y_n) \leq d(Y_{n-2}) + h \sum_{j=0}^k |\bar{\delta}_{kj}| d(F_{n-j}) + h^{k+2} \left(\left| \bar{\nu}_{k+1}^* \right| + \left| \bar{\nu}_{k+1}^{**} \right| \right) d(\bar{\Psi}_k). \quad (33)$$

We have assumed that $\bar{\Psi}$ is monotonic with respect to inclusion. Moreover, if the step size h is such that satisfies the condition

$$T_n + [-kh, 0] \subset \Delta_t,$$

$$Y_n + [-kh, 0]F(\Delta_t, \Delta_y) \subset \Delta_y,$$

then $\bar{\Psi}_k \subset \bar{\Psi}(\Delta_t, \Delta_y)$, and hence $d(\bar{\Psi}_k) \leq d(\bar{\Psi}(\Delta_t, \Delta_y))$. We have also assumed that for the function F there exists

a constant $L > 0$ such that $d(F_{n-j}) \leq L(d(T_{n-j}) + d(Y_{n-j}))$.
Therefore, from (33) we get

$$\begin{aligned} d(Y_n) &\leq d(Y_{n-2}) + \\ &+ hL\bar{\delta}_k \sum_{j=0}^k (d(T_{n-j}) + d(Y_{n-j})) + \\ &+ h^{k+2} \left(\left| \bar{v}_{k+1}^* \right| + \left| \bar{v}_{k+1}^{**} \right| \right) d(\bar{\Psi}(\Delta_t, \Delta_y)), \end{aligned} \quad (34)$$

where

$$\bar{\delta}_k = \max_{j=0,1,\dots,k} |\bar{\delta}_{kj}|.$$

If we denote

$$\begin{aligned} \bar{\delta} &= hL\bar{\delta}_k, \quad \alpha = 1 + \bar{\delta}, \\ \bar{v} &= h^{k+2} \left(\left| \bar{v}_{k+1}^* \right| + \left| \bar{v}_{k+1}^{**} \right| \right), \end{aligned} \quad (35)$$

then we can write (34) in the form

$$\begin{aligned} d(Y_n) &\leq d(Y_{n-2}) + \bar{\delta}d(Y_n) + \\ &+ \bar{\delta} \sum_{j=1}^k d(Y_{n-j}) + \bar{\delta} \sum_{j=0}^k d(T_{n-j}) + \bar{v}d(\Psi(\Delta_t, \Delta_y)), \end{aligned}$$

from which it follows that

$$\begin{aligned} (1 - \bar{\delta})d(Y_n) &\leq \bar{\alpha} \sum_{j=1}^k d(Y_{n-j}) + \\ &+ \bar{\delta} \sum_{j=0}^k d(T_{n-j}) + \bar{v}d(\Psi(\Delta_t, \Delta_y)). \end{aligned} \quad (36)$$

Let us note that

$$1 - \bar{\delta} = 1 - hL\bar{\delta}_k > 0. \quad (37)$$

The condition (37) is satisfied, since from the assumptions $0 < h < h_0$ and $h_0 = 1/L\bar{\delta}_k$. Thus, the inequality (36) can be written in the form

$$\begin{aligned} d(Y_n) &\leq \beta\bar{\alpha} \sum_{j=1}^k d(Y_{n-j}) + \\ &+ \beta\bar{\delta} \sum_{j=0}^k d(T_{n-j}) + \beta\bar{v}d(\Psi(\Delta_t, \Delta_y)), \end{aligned}$$

where

$$\beta = \frac{1}{1 - h_0 L \bar{\delta}_k}.$$

From (38) for $n = k$ we have

$$\begin{aligned} d(Y_k) &\leq \beta\bar{\alpha} \sum_{j=1}^k d(Y_{k-j}) + \\ &+ \beta\bar{\delta} \sum_{j=0}^k d(T_{k-j}) + \beta\bar{v}d(\Psi(\Delta_t, \Delta_y)), \end{aligned} \quad (39)$$

while for $n = k + 1$ we get

$$\begin{aligned} d(Y_{k+1}) &\leq \beta\bar{\alpha}d(Y_k) + \beta\bar{\alpha} \sum_{j=1}^{k-1} d(Y_{k-j}) + \\ &+ \beta\bar{\delta} \sum_{j=0}^k d(T_{k+1-j}) + \beta\bar{v}d(\Psi(\Delta_t, \Delta_y)). \end{aligned}$$

Applying (39) to the above inequality we obtain

$$\begin{aligned} d(Y_{k+1}) &\leq \left((\beta\bar{\alpha})^2 + \beta\bar{\alpha} \right) \sum_{j=1}^k d(Y_{k-j}) + \\ &+ \beta\bar{\delta} \left(\beta\bar{\alpha} \sum_{j=0}^k d(T_{k-j}) + \sum_{j=0}^k d(T_{k+1-j}) \right) + \\ &+ \beta\bar{v}(\beta\bar{\alpha} + 1)d(\Psi(\Delta_t, \Delta_y)). \end{aligned} \quad (40)$$

From (38) for $n = k + 2$ we get

$$\begin{aligned} d(Y_{k+2}) &\leq \\ &\leq \beta\bar{\alpha}d(Y_{k+1}) + \beta\bar{\alpha}d(Y_k) + \beta\bar{\alpha} \sum_{j=1}^{k-2} d(Y_{k-j}) + \\ &+ \beta\bar{\delta} \sum_{j=0}^k d(T_{k+2-j}) + \beta\bar{v}d(\Psi(\Delta_t, \Delta_y)). \end{aligned}$$

Insertion of (39) and (40) into the above inequality yields

$$\begin{aligned} d(Y_{k+2}) &\leq \left((\beta\bar{\alpha})^3 + 2(\beta\bar{\alpha}) + \beta\bar{\alpha} \right) \sum_{j=1}^k d(Y_{k-j}) + \\ &+ \beta\bar{\delta} \left(\left((\beta\bar{\alpha})^2 + \beta\bar{\alpha} \right) \sum_{j=0}^k d(T_{k-j}) + \beta\bar{\alpha} \sum_{j=0}^k d(T_{k+1-j}) + \right. \\ &\left. + \sum_{j=0}^k d(T_{k+2-j}) \right) + \beta\bar{v} \left((\beta\bar{\alpha})^2 + 2\beta\bar{\alpha} + 1 \right) d(\Psi(\Delta_t, \Delta_y)). \end{aligned} \quad (41)$$

From (39), (40) and (41) it follows that in general for each $i = 0, 1, \dots, m - k$ we have

$$\begin{aligned}
d(Y_{k+i}) &\leq \left(\sum_{l=0}^i \binom{i}{l} (\beta \bar{\alpha})^{l+1} \right) \left(\sum_{j=1}^k d(Y_{k-j}) \right) + \\
&+ \beta \bar{\delta} \sum_{p=0}^i \left(\sum_{l=0}^{p-1} \binom{p-1}{l} (\beta \bar{\alpha})^{l+1} \right) \left(\sum_{j=0}^k d(T_{k+i-p-j}) \right) + \\
&+ \beta \bar{v} \sum_{l=0}^i \binom{i}{l} (\beta \bar{\alpha})^l d(\Psi(\Delta_t, \Delta_y)).
\end{aligned}$$

Applying the notations (35) we obtain

$$\begin{aligned}
d(Y_{k+i}) &\leq k \sum_{l=0}^i \binom{i}{l} (\beta(1+hL\bar{\delta}_k))^{l+1} \max_{q=0,1,\dots,k-1} d(Y_q) + \\
&+ \beta h L \bar{\delta}_k (k+1) \times \\
&\times \sum_{p=0}^i \sum_{l=0}^{p-1} \binom{p-1}{l} (\beta(1+hL\bar{\delta}_k))^{l+1} \max_{j=0,1,\dots,k+i} d(T_j) + \quad (42) \\
&+ \beta h^{k+2} \left(\left| v_{k+1}^* \right| + \left| v_{k+1}^{**} \right| \right) \times \\
&\times \sum_{l=0}^i \binom{i}{l} (\beta(1+hL\bar{\delta}_k))^l d(\Psi(\Delta_t, \Delta_y)).
\end{aligned}$$

Using (27) and the following estimation:

$$\sum_{l=0}^{p-1} \beta^{l+1} (1+hL\bar{\delta}_k) \leq \frac{\beta^m \exp(\xi L \bar{\delta}_k) - 1}{\beta h L \bar{\delta}_k},$$

it is easy to show that

$$\begin{aligned}
k \sum_{l=0}^i \binom{i}{l} (\beta(1+hL\bar{\delta}_k))^{l+1} &\leq m(m-k)! \beta \exp(\xi L \bar{\delta}_k) \frac{1-\beta^m}{1-\beta}, \\
(k+1) \sum_{p=0}^i \sum_{l=0}^{p-1} \binom{p-1}{l} \beta^{l+1} (1+hL\bar{\delta}_k)^{l+1} &\leq \\
&\leq (m+1)(m-k+1)(m-k)! \frac{\beta^m \exp(\xi L \bar{\delta}_k) - 1}{\beta h L \bar{\delta}_k}, \\
\sum_{l=0}^i \binom{i}{l} (\beta(1+hL\bar{\delta}_k))^l &\leq (m-k)! \frac{\beta^m \exp(\xi L \bar{\delta}_k) - 1}{\beta h L \bar{\delta}_k}.
\end{aligned}$$

Thus, from (42) we finally have

$$\begin{aligned}
d(Y_{k+i}) &\leq \bar{A} \max_{q=0,1,\dots,k-1} d(Y_q) + \\
&+ \bar{B} \max_{j=0,1,\dots,m} d(T_j) + \bar{C} h^{k+1}, \quad (43)
\end{aligned}$$

for each $i = 0, 1, \dots, m-k$, where

$$\bar{A} = m(m-k)! \beta \exp(\xi L \bar{\delta}_k) \frac{1-\beta^m}{1-\beta},$$

$$\bar{B} = (m+1)(m-k+1)(m-k)! (\beta^m \exp(\xi L \bar{\delta}_k) - 1),$$

$$\bar{C} = \frac{\left| v_{k+1}^* \right| + \left| v_{k+1}^{**} \right|}{L \bar{\delta}_k} (m-k)! \times$$

$$\times (\beta^m \exp(\xi L \bar{\delta}_k) - 1) d(\Psi(\Delta_t, \Delta_y)).$$

Since $d(T_0) = 0$, the inequality (32) follows immediately from (43). ■

V. NUMERICAL RESULTS

All calculations in the examples presented below have been performed using the *IntervalArithmetic* unit written by the author in the Delphi Pascal language². This unit takes advantage of the floating-point *Extended* type. The *IntervalArithmetic* unit makes it possible to:

- represent any input dat in the form of machine interval (the ends of this interval are equal or are two subsequent machine numbers),
- perform all calculations in floating-point interval arithmetic,
- use some standard interval functions,
- give results in the form of proper intervals (if the ends of an interval are not the same machine numbers, you see the difference in the output).

Example 1

As the first example let us consider the commonly used test problem, i.e. the following initial value problem:

$$y' = \lambda y, \quad y(0) = 1,$$

with the exact solution

$$y = \exp(\lambda t).$$

For $\lambda = 0.5$ we have (17 digits after decimal point)

$$y(0.5) \approx 1.28402541668774149,$$

$$y(1.0) \approx 1.64872127070012815.$$

In all methods we have assumed

$$\Delta_t = [0, 1], \quad \Delta_y = [1, 1.65], \quad h = 5 \cdot 10^{-4},$$

$$2000 \text{ steps}, \quad T_0 = [0, 0], \quad \text{and } Y_0 = [1, 1].$$

² The *IntervalArithmetic* has been presented and discussed in [12] and [14], and can be obtain from the author by e-mail.

Table 1. Interval solutions of the test problem obtained by explicit interval methods

k	Method type	Results
2	Nyström	T(1000) = [4.999999999999999E-0001, 5.0000000000000001E-0001]
		Y(1000) = [1.2840254166859112E+0000, 1.2840254166895718E+0000] width = 3.66E-0012
		T(2000) = [9.999999999999994E-0001, 1.0000000000000001E+0000]
		Y(2000) = [1.6487212706959476E+0000, 1.6487212707043086E+0000] width = 8.36E-0012
	Adams-Bashforth	T(1000) = [4.999999999999999E-0001, 5.0000000000000001E-0001]
		Y(1000) = [1.2840254166844718E+0000, 1.2840254166914287E+0000] width = 6.96E-0012
4	Nyström	T(1000) = [4.999999999999999E-0001, 5.0000000000000001E-0001]
		Y(1000) = [1.2840254166877413E+0000, 1.2840254166877417E+0000] width = 2.93E-0016
		T(2000) = [9.999999999999994E-0001, 1.0000000000000001E+0000]
		Y(2000) = [1.6487212707001277E+0000, 1.6487212707001285E+0000] width = 7.01E-0016
	Adams-Bashforth	T(1000) = [4.999999999999999E-0001, 5.0000000000000001E-0001]
		Y(1000) = [1.2840254166877410E+0000, 1.2840254166877419E+0000] width = 8.01E-0016
		T(2000) = [9.999999999999994E-0001, 1.0000000000000001E+0000]
		Y(2000) = [1.6487212707001259E+0000, 1.6487212707001305E+0000] width = 4.51E-0015

In Table 1 we present the results obtained by the interval methods of Nyström type for $k = 2$ and $k = 4$, and compare them to the results obtained by the interval methods of Adams-Bashforth type, presented e.g. in [5], for the same values of k . One can observe that interval methods of Nyström type give somewhat better solutions, because the diameters of the intervals obtained are smaller.

In Table 2 similar results for implicit interval methods of Milne-Simpson type and Adams-Moulton type are presented. Implicit interval methods of Adams-Moulton type have been presented in [4], [6] and [17]. It is easy to observe that for the same number of steps implicit interval methods of Milne-Simpson type, presented in this paper, give somewhat better results than the methods of Adams-Moulton type. Let us note that the number of iterations in both of these implicit methods is the same.

Figure 1 presents widths of interval solutions as functions of the number of methods steps. It appears that

for each of the method considered there exists an optimal number of method steps. From the second figure (Fig. 2) it

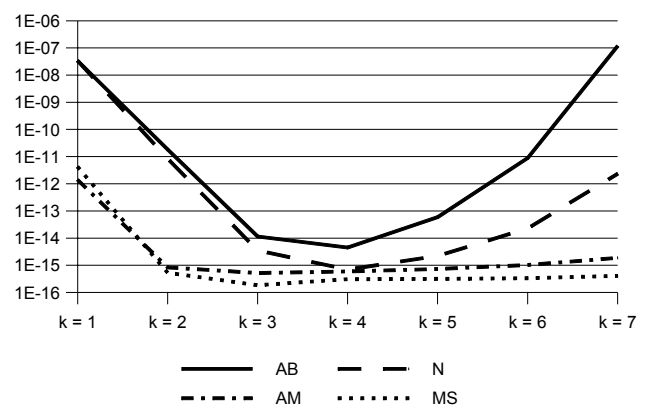


Fig. 1. Widths of interval solutions vs. the number k of method steps at $t = 1$, $h = 0.0005$ (AB – Adams-Bashforth, N – Nyström, AM – Adams-Moulton, MS – Milne-Simpson)

Table 2. Interval solutions of the test problem obtained by implicit interval methods

k	Method type	Results
2	Milne-Simpson	T(1000) = [4.999999999999999E-0001, 5.0000000000000001E-0001] Y(1000) = [1.2840254166877413E+0000 , 1.2840254166877417E+0000] width = 2.34E-0016, number of iterations: 5
		T(2000) = [9.999999999999994E-0001, 1.0000000000000001E+0000] Y(2000) = [1.6487212707001279E+0000 , 1.6487212707001285E+0000] width = 5.32E-0016, number of iterations: 5
		T(1000) = [4.999999999999999E-0001, 5.0000000000000001E-0001] Y(1000) = [1.2840254166877413E+0000 , 1.2840254166877418E+0000] width = 4.15E-0016, number of iterations: 5
	Adams-Moulton	T(2000) = [9.999999999999994E-0001, 1.0000000000000001E+0000] Y(2000) = [1.6487212707001278E+0000 , 1.6487212707001287E+0000] width = 8.37E-0016
		T(1000) = [4.999999999999999E-0001, 5.0000000000000001E-0001] Y(1000) = [1.2840254166877414E+0000 , 1.2840254166877416E+0000] width = 8.12E-0017, number of iterations: 5
		T(2000) = [9.999999999999994E-0001, 1.0000000000000001E+0000] Y(2000) = [1.6487212707001280E+0000 , 1.6487212707001283E+0000] width = 1.85E-0016, number of iterations: 5
3	Milne-Simpson	T(1000) = [4.999999999999999E-0001, 5.0000000000000001E-0001] Y(1000) = [1.2840254166877413E+0000 , 1.2840254166877417E+0000] width = 2.73E-0016, number of iterations: 5
		T(2000) = [9.999999999999994E-0001, 1.0000000000000001E+0000] Y(2000) = [1.6487212707001279E+0000 , 1.6487212707001285E+0000] width = 5.20E-0016, number of iterations: 5
		T(1000) = [4.999999999999999E-0001, 5.0000000000000001E-0001] Y(1000) = [1.2840254166877413E+0000 , 1.2840254166877417E+0000] width = 2.73E-0016, number of iterations: 5
	Adams-Moulton	T(2000) = [9.999999999999994E-0001, 1.0000000000000001E+0000] Y(2000) = [1.6487212707001279E+0000 , 1.6487212707001285E+0000] width = 5.20E-0016, number of iterations: 5
		T(1000) = [4.999999999999999E-0001, 5.0000000000000001E-0001] Y(1000) = [1.2840254166877413E+0000 , 1.2840254166877417E+0000] width = 2.73E-0016, number of iterations: 5
		T(2000) = [9.999999999999994E-0001, 1.0000000000000001E+0000] Y(2000) = [1.6487212707001279E+0000 , 1.6487212707001285E+0000] width = 5.20E-0016, number of iterations: 5

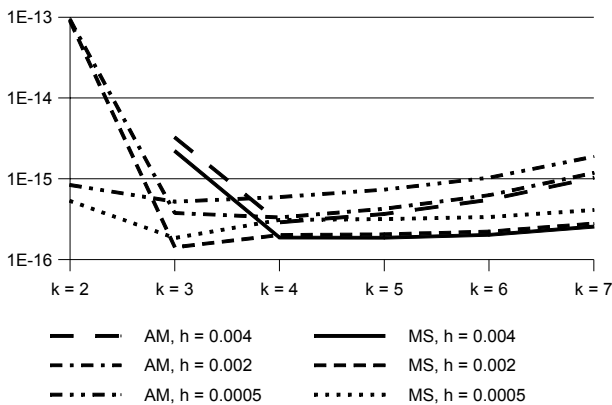


Fig. 2. Widths of interval solutions vs. the number k of method steps for different step sizes at $t = 1$ (AM – Adams-Moulton, MS – Milne-Simpson)

follows that for a given number of method steps there exists the best step size.

Example 2

As an example of two-dimensional problem let us consider a plane motion of two bodies. Let us assume that the bodies are material points with the masses m_1 and m_2 , and are in a constant mutual distance α . From the well-known equations of motion of N bodies in inertial frame of references, i.e. from the following equations:

$$\ddot{x}_{li} = -G \sum_{\substack{j=1 \\ j \neq i}}^N m_j \frac{x_{li} - x_{lj}}{r_{ij}^3}, \quad r_{ij} = \sqrt{\sum_{l=1}^3 (x_{li} - x_{lj})^2},$$

$$l = 1, 2, 3; \quad i = 1, 2, \dots, N,$$

where G denotes the gravitational constant, it follows that the considered motion is described by the following differential equations of the first order:

$$\begin{aligned} \dot{x}_{li} &= v_{li}, \quad \dot{v}_{l1} = -\frac{Gm_2}{\alpha^3}(x_{l1} - x_{l2}), \\ \dot{v}_{l2} &= -\frac{Gm_1}{\alpha^3}(x_{l2} - x_{l1}), \quad l=1, 2; \quad i=1, 2. \end{aligned} \quad (44)$$

Assuming that

$$\begin{aligned} x_{11}^0 &= \alpha, \quad v_{21}^0 = 2\pi\alpha, \\ x_{li}^0 &= v_{li}^0 = 0 \quad \text{for other } l \text{ and } i \text{ at } t_0 = 0, \end{aligned}$$

from (44) we can get the exact solution in the form (see [10])

$$\begin{aligned} x_{11} &= \frac{\alpha}{m_1 + m_2}(m_1 + m_2 \cos 2\pi t), \\ x_{21} &= \frac{\alpha}{m_1 + m_2}(2\pi t m_1 + m_2 \sin 2\pi t), \end{aligned}$$

$$x_{12} = \frac{\alpha m_1}{m_1 + m_2}(1 - \cos 2\pi t),$$

$$x_{22} = \frac{\alpha m_1}{m_1 + m_2}(2\pi t - \sin 2\pi t).$$

$$G = \frac{4\pi^2 \alpha^3}{m_1 + m_2}.$$

Taking the astronomical unit (149597.9×10^6 m), mass of the Earth with the Moon ($1/328900.1$ of the Sun mass) and astronomical year at the moment 1950.0 (365.24219572 days) as units, we have $G \approx 1.20021974563227948 \times 10^{-4}$. Hence, for $m_1 = 1$ and $m_2 = 328900.1$ from (46) it follows that $\alpha \approx 0.999974178082659804$.

In order to apply explicit interval methods of Nyström type and both kinds of implicit interval methods of Milne-Simpson type for $k = 1, 2, 3$ with $h = 0.0001$, on the basis of (45) we have taken the following initial intervals:

$n = 0:$	$x_{11} = [9.9997417808265980E-0001,$	$9.9997417808265981E-0001],$
	$x_{21} = [0.0000000000000000E+0000,$	$0.0000000000000000E+0000],$
	$x_{12} = [0.0000000000000000E+0000,$	$0.0000000000000000E+0000],$
	$x_{22} = [0.0000000000000000E+0000,$	$0.0000000000000000E+0000],$
	$v_{11} = [0.0000000000000000E+0000,$	$0.0000000000000000E+0000],$
	$v_{21} = [6.2830230632879513E+0000,$	$6.2830230632879514E+0000],$
	$v_{12} = [0.0000000000000000E+0000,$	$0.0000000000000000E+0000],$
	$v_{22} = [0.0000000000000000E+0000,$	$0.0000000000000000E+0000],$
$n = 1:$	$x_{11} = [9.9997398069627545E-0001,$	$9.9997398069627546E-0001],$
	$x_{21} = [6.2830226498828692E-0004,$	$6.2830226498828693E-0004],$
	$x_{12} = [6.0014084624780502E-0013,$	$6.0014084624780503E-0013],$
	$x_{22} = [1.2569320656687768E-0016,$	$1.2569320656687769E-0016],$
	$v_{11} = [3.9477275570247397E-0003,$	$3.9477275570247396E-0003],$
	$v_{21} = [6.2830218230727214E+0000,$	$6.2830218230727215E+0000],$
	$v_{12} = [1.2002816530079314E-0008,$	$1.2002816530079315E-0008],$
	$v_{22} = [3.7707961473825318E-0012,$	$3.7707961473825319E-0012],$
$n = 2:$	$x_{11} = [9.9997338853720034E-0001,$	$9.9997338853720035E-0001],$
	$x_{21} = [1.2566042819335441E-0003,$	$1.2566042819335443E-0003],$
	$x_{12} = [2.4005631480651783E-0012,$	$2.4005631480651784E-0012],$
	$x_{22} = [1.0055455929885264E-0015,$	$1.0055455929885265E-0015],$
	$v_{11} = [7.8954535555491599E-0003,$	$7.8954535555491598E-0003],$
	$v_{21} = [6.2830181024275211E+0000,$	$6.2830181024275212E+0000],$
	$v_{12} = [2.4005628321636751E-0008,$	$2.4005628321636752E-0008],$
	$v_{22} = [1.5083183100879902E-0011,$	$1.5083183100879903E-0011].$

Table 3. Interval solutions of the two-body problem obtained by interval methods of Nyström type

k	Step	Results	
1	2000	$x[1, 1] = [3.0901054026171765E-0001,$ width = 9.94E-0007	3.0901153413474893E-0001]
		$x[2, 1] = [9.5103241587730308E-0001,$ width = 9.94E-0007	9.5103340975033433E-0001]
		$v[1, 1] = [-5.9754951387795067E+0000,$ width = 6.24E-0006	5.9754888988231943E+0000]
		$v[2, 1] = [1.9415704887021080E+0000,$ width = 6.24E-0006	1.9415767286584202E+0000]
	10000	$x[1, 1] = [9.9986838859246625E-0001,$ width = 2.12E-0004	1.0000799675726825E+0000]
		$x[2, 1] = [-8.6273007868630888E-0005,$ width = 2.12E-0004	1.2530597234547504E-0004]
		$v[1, 1] = [-6.6728861987543796E-0004,$ width = 1.32E-0003	6.6209361737553662E-0004]
		$v[2, 1] = [6.2823583721687961E+0000,$ width = 1.32E-0003	6.2836877544060329E+0000]
2	2000	$x[1, 1] = [3.0901111563729169E-0001,$ width = 3.12E-0010	3.0901111594921569E-0001]
		$x[2, 1] = [9.5103288710799053E-0001,$ width = 3.12E-0010	9.5103288741991450E-0001]
		$v[1, 1] = [-5.9754918592474882E+0000,-$ width = 1.96E-0009	-5.9754918572861251E+0000]
		$v[2, 1] = [1.9415741015266600E+0000,$ width = 1.96E-0009	1.9415741034880229E+0000]
	10000	$x[1, 1] = [9.9997414490329201E-0001,$ width = 6.64E-0008	9.9997421126202760E-0001]
		$x[2, 1] = [1.9069897700167120E-0005,$ width = 6.64E-0008	1.9136256432965659E-0005]
		$v[1, 1] = [-2.0847327440342510E-0007,$ width = 4.17E-0007	2.0847337791437236E-0007]
		$v[2, 1] = [6.2830228548146321E+0000,$ width = 4.17E-0007	6.2830232717612706E+0000]
3	2000	$x[1, 1] = [3.0901111579295821E-0001,$ width = 5.55E-0013	3.0901111579351298E-0001]
		$x[2, 1] = [9.5103288726368103E-0001,$ width = 5.55E-0013	9.5103288726423578E-0001]
		$v[1, 1] = [-5.9754918582685861E+0000,$ width = 3.48E-0012	-5.9754918582651012E+0000]
		$v[2, 1] = [1.9415741025054845E+0000,$ width = 3.49E-0012	1.9415741025089710E+0000]
	10000	$x[1, 1] = [9.9997417671168836E-0001,$ width = 2.74E-0009	9.9997417945363130E-0001]
		$x[2, 1] = [1.9101705914996750E-0005,$ width = 2.74E-0009	1.9104448408562482E-0005]
		$v[1, 1] = [-8.6146142851223850E-0009,$ width = 1.72E-0008	8.6135215736874046E-0009]
		$v[2, 1] = [6.2830230546720339E+0000,$ width = 1.72E-0008	6.2830230719038694E+0000]

Table 4. Interval solutions of the two-body problem obtained by interval methods of Milne-Simpson type of the first kind (with backward differences)

k	Step	Results
1	2000	$x[1, 1] = [3.09011111492869487E-0001, 3.09011111668988840E-0001]$ width = 1.76E-0009
		$x[2, 1] = [9.5103288643274413E-0001, 9.5103288819393640E-0001]$ width = 1.76E-0009
		$v[1, 1] = [-5.9754918641104098E+0000, -5.9754918530438282E+0000]$ width = 1.10E-0008
		$v[2, 1] = [1.9415740970748240E+0000, 1.9415741081413977E+0000]$ width = 1.10E-0008 number of iterations : 4
	10000	$x[1, 1] = [9.9672257365831662E-0001, 1.0032257830265033E+0000]$ width = 6.50E-0003
		$x[2, 1] = [-3.2324994014764165E-0003, 3.2707055557727574E-0003]$ width = 6.50E-0003
		$v[1, 1] = [-2.0430434776658366E-0002, 2.0430434775736423E-0002]$ width = 4.09E-0002
		$v[2, 1] = [6.2625926440011807E+0000, 6.3034534858388382E+0000]$ width = 4.09E-0002 number of iterations : 7
2	2000	$x[1, 1] = [3.0901111579289397E-0001, 3.0901111579360825E-0001]$ width = 7.14E-0013
		$x[2, 1] = [9.5103288726359762E-0001, 9.5103288726430908E-0001]$ width = 7.11E-0013
		$v[1, 1] = [-5.9754918582690561E+0000, -5.9754918582645677E+0000]$ width = 4.49E-0012
		$v[2, 1] = [1.9415741025050898E+0000, 1.9415741025095606E+0000]$ width = 4.47E-0012 number of iterations : 6
	10000	$x[1, 1] = [9.9993689042939006E-0001, 1.0000114657359296E+0000]$ width = 7.46E-0005
		$x[2, 1] = [-1.8040700405988558E-0005, 5.6246854566328955E-0005]$ width = 7.43E-0005
		$v[1, 1] = [-2.3428523519754089E-0004, 2.3428523513016957E-0004]$ width = 4.69E-0004
		$v[2, 1] = [6.2827896820509971E+0000, 6.2832564445249057E+0000]$ width = 4.67E-0004 number of iterations : 8
3	2000	$x[1, 1] = [3.0901111579324826E-0001, 3.0901111579325603E-0001]$ width = 7.76E-0015
		$x[2, 1] = [9.5103288726395081E-0001, 9.5103288726395523E-0001]$ width = 4.40E-0015
		$v[1, 1] = [-5.9754918582668345E+0000, -5.9754918582667852E+0000]$ width = 4.92E-0014
		$v[2, 1] = [1.9415741025073175E+0000, 1.9415741025073459E+0000]$ width = 2.82E-0014 number of iterations : 6
	10000	$x[1, 1] = [9.9997377360945398E-0001, 9.9997458255586562E-0001]$ width = 8.09E-0007
		$x[2, 1] = [1.8869511256453819E-0005, 1.9336642893008227E-0005]$ width = 4.67E-0007
		$v[1, 1] = [-2.5413801037084521E-0006, 2.5413801047086599E-0006]$ width = 5.08E-0006
		$v[2, 1] = [6.2830215957506335E+0000, 6.2830245308252693E+0000]$ width = 2.94E-0006 number of iterations : 8

Table 5. Interval solutions of the two-body problem obtained by interval methods of Milne-Simpson type of the second kind (without backward differences)

k	Step	Results
1	2000	$x[1, 1] = [3.09011111565332963E-0001, 3.09011111596525363E-0001]$ width = 3.12E-0010
		$x[2, 1] = [9.5103288715737828E-0001, 9.5103288746930225E-0001]$ width = 3.12E-0010
		$v[1, 1] = [-5.9754918595578006E+0000, -5.9754918575964375E+0000]$ width = 1.96E-0009
		$v[2, 1] = [1.9415741016274294E+0000, 1.9415741035887923E+0000]$ width = 1.96E-0009 number of iterations : 3
10000	10000	$x[1, 1] = [9.9997414516304228E-0001, 9.9997421152177755E-0001]$ width = 6.64E-0008
		$x[2, 1] = [1.9069897781524307E-0005, 1.9136256514813561E-0005]$ width = 6.64E-0008
		$v[1, 1] = [-2.0847378713724874E-0007, 2.0847286520361374E-0007]$ width = 4.17E-0007
		$v[2, 1] = [6.2830228564466912E+0000, 6.2830232733933277E+0000]$ width = 4.17E-0007 number of iterations : 4
2	2000	$x[1, 1] = [3.0901111579321830E-0001, 3.0901111579328392E-0001]$ width = 6.56E-0014
		$x[2, 1] = [9.5103288726392055E-0001, 9.5103288726398616E-0001]$ width = 6.56E-0014
		$v[1, 1] = [-5.9754918582670182E+0000, -5.9754918582666057E+0000]$ width = 4.12E-0013
		$v[2, 1] = [1.9415741025071190E+0000, 1.9415741025075314E+0000]$ width = 4.12E-0013 number of iterations : 6
10000	10000	$x[1, 1] = [9.9997417807568387E-0001, 9.9997417808963573E-0001]$ width = 1.39E-0011
		$x[2, 1] = [1.9103070104057486E-0005, 1.9103084056286399E-0005]$ width = 1.39E-0011
		$v[1, 1] = [-4.3864799169172155E-0011, 4.3797457369853193E-0011]$ width = 8.77E-0011
		$v[2, 1] = [6.2830230632441190E+0000, 6.2830230633317838E+0000]$ width = 8.77E-0011 number of iterations : 7
3	2000	$x[1, 1] = [3.0901111579325198E-0001, 3.0901111579325231E-0001]$ width = 3.16E-0016
		$x[2, 1] = [9.5103288726395287E-0001, 9.5103288726395316E-0001]$ width = 2.81E-0016
		$v[1, 1] = [-5.9754918582668110E+0000, -5.9754918582668086E+0000]$ width = 2.35E-0015
		$v[2, 1] = [1.9415741025073306E+0000, 1.9415741025073329E+0000]$ width = 2.19E-0015 number of iterations : 6
10000	10000	$x[1, 1] = [9.9997417808262542E-0001, 9.9997417808269419E-0001]$ width = 6.88E-0014
		$x[2, 1] = [1.9103077041422683E-0005, 1.9103077108037477E-0005]$ width = 6.66E-0014
		$v[1, 1] = [-2.1575457406226130E-0013, 2.1675752618362943E-0013]$ width = 4.33E-0013
		$v[2, 1] = [6.2830230632877418E+0000, 6.2830230632881610E+0000]$ width = 4.19E-0013 number of iterations : 7

Moreover, we have assumed that $m = 10000$ steps, $T_0 = [0, 0]$, $\Delta_t = [0, 1]$ and

$$\Delta_y = \begin{pmatrix} \Delta_{x_{11}} \\ \Delta_{x_{21}} \\ \Delta_{x_{12}} \\ \Delta_{x_{22}} \\ \Delta_{v_{11}} \\ \Delta_{v_{21}} \\ \Delta_{v_{12}} \\ \Delta_{v_{22}} \end{pmatrix} = \begin{pmatrix} [-1, 1] \\ [-1, 1] \\ [-2 \cdot 10^{-5}, 2 \cdot 10^{-5}] \\ [-2 \cdot 10^{-5}, 2 \cdot 10^{-5}] \\ [-6.3, 6.3] \\ [-6.3, 6.3] \\ [-4 \cdot 10^{-5}, 4 \cdot 10^{-5}] \\ [-4 \cdot 10^{-5}, 4 \cdot 10^{-5}] \end{pmatrix}$$

All the results presented in Tables 3, 4 and 5 have been obtained for the above data.

It is easy to observe that for the same number of steps the interval solutions obtained by the implicit methods of the second kind are the best (the diameters of intervals are the smallest), and that the methods (30) give better solutions than the methods (29). One can also check that in each case and at each moment the exact solution given by (45) belongs to the intervals obtained. Moreover, the implicit interval methods of the first kind are good enough at the beginning of integration interval, but the diameters of intervals grow rapidly at the end of this interval.

VI. CONCLUSIONS

Interval methods for solving the initial value problem in floating-point interval arithmetic give solutions in the form of interval which contain all possible numerical errors, i.e. representation errors, rounding errors, and errors of methods. For the method considered in this paper it follows that for the same number of steps explicit interval methods of Nyström type are somewhat better (i.e. give the interval solution with a smaller width) than the methods of Adams-Bashforth type [5, 9, 18, 21], and implicit interval methods of Milne-Simpson type give somewhat better results than the methods of Adams-Moulton type [4, 6, 17, 18]. Another conclusion concerning the interval methods of Milne-Simpson type is that the methods based on backward interval differences give somewhat worse results than the methods based only on the combinations of interval function values at different points (see Theorem 4 and Table 3). Moreover, for each particular problem one should choose the appropriate step size and the number of method steps to obtain the interval solution with the smallest width (diameter). It appears that for a given step size there exists

the optimal number of method steps, and for a given number of method steps there exists the best step size.

References

- [1] J. C. Butcher, *The Numerical Analysis of Ordinary Differential Equations. Runge-Kutta and General Linear Methods*, J. Wiley & Sons, Chichester 1987.
- [2] K. Gajda, A. Marciniak, A. Marlewski, B. Szyszka, *A Layout of an Object-Oriented System for Solving the Initial Value Problem by Interval Methods of Runge-Kutta Type* [in Polish], *Pro Dialog* **8**, 39-62 (1999).
- [3] K. Gajda, A. Marciniak, B. Szyszka, *Three- and Four-Stage Implicit Interval Methods of Runge-Kutta Type*, *Computational Methods in Science and Technology* **6**, 41-59 (2000).
- [4] M. Jankowska, A. Marciniak, *Implicit Interval Multistep Methods for Solving the Initial Value Problem*, *Computational Methods in Science and Technology* **8(1)**, 17-30 (2002).
- [5] M. Jankowska, A. Marciniak, *On Explicit Interval Methods of Adams-Bashforth Type*, *Computational Methods in Science and Technology* **8(2)**, 46-57 (2002).
- [6] M. Jankowska, A. Marciniak, *On Two Families of Implicit Interval Methods of Adams-Moulton Type*, *Computational Methods in Science and Technology* **12(2)**, 109-113 (2006).
- [7] S. A. Kalmykov, Ju. I. Šokin, E. Ch. Juldašev, *On an Interval-Analytical Method of Second Order for Solving Ordinary Differential Equations* [in Russian], *Izv. AN UzSSR, Ser. Fiz.-Mat. Nauk* **3** (1976).
- [8] S. A. Kalmykov, Ju. I. Šokin, E. Ch. Juldašev, *Some Interval Methods for Solving Ordinary Differential Equations* [in Russian], *Èislennoje Metody Mehaniki Splošnoj Sredy* **7(6)**, (1976).
- [9] S. A. Kalmykov, Ju. I. Šokin, E. Ch. Juldašev, *Solving Ordinary Differential Equations by Interval Methods* [in Russian], *Doklady AN SSSR*, **230(6)** (1976).
- [10] F. Krückeberg, *Ordinary Differential Equations*, in: E. Hansen (Ed.), *Topics in Interval Analysis*, Oxford University Press, 91-97 (1969).
- [11] A. Marciniak, *Numerical Solutions of the N-body Problem*, D. Reidel Publishing Co., Dordrecht 1985.
- [12] A. Marciniak, *0,1*, *Pro Dialog* **5**, 55-82 (1997).
- [13] A. Marciniak, B. Szyszka, *One- and Two-Stage Implicit Interval Methods of Runge-Kutta Type*, *Computational Methods in Science and Technology* **5**, 53-65 (1999).
- [14] A. Marciniak, *Borland Delphi 5 Professional. Object Pascal* [in Polish], NAKOM Publishers, Poznań 2000.
- [15] A. Marciniak, *Finding the Integration Interval for Interval Methods of Runge-Kutta Type in Floating-Point Interval Arithmetic*, *Pro Dialog* **10**, 35-45 (2000).
- [16] A. Marciniak, B. Szyszka, *On Representations of Coefficients in Implicit Interval Methods of Runge-Kutta Type*, *Computational Methods in Science and Technology* **10(1)**, 57-71 (2004).
- [17] A. Marciniak, *Implicit Interval Methods for Solving the Initial Value Problem*, *Numerical Algorithms* **37**, 241-251 (2004).

- [18] A. Marciniak, *On Multistep Interval Methods for Solving the Initial Value Problem*, Journal of Computational and Applied Mathematics (2006) (in press).
- [19] R. E. Moore, *Interval Analysis*, Prentice Hall, 1966.
- [20] B. Szyszka, *Implicit Interval Methods of Runge-Kutta Type* [in Polish], Ph. D. Thesis, Poznan University of Technology, 2003.
- [21] Ju. I. Šokin, *Interval Analysis* [in Russian], Izdatel'stvo "Nauka", Novosibirsk, 1981.



ANDRZEJ MARCINIAK, Associate Professor (PhD 1981, Dr. habil. 1993), in years 1977-1986 in the Institute of Mathematics, Adam Mickiewicz University of Poznan, in years 1986-1992 in the Institute of Mathematics, Poznan University of Technology, since 1992 in the Institute of Computing Science, and since 2000 also in the Faculty of Mathematics and Computer Science, Adam Mickiewicz University of Poznan. In years 1982-1983 visiting professor in the University of Florida in Gainesville, and visiting professor in the University of Texas at Arlington in 1990. Specialist in computer programming and numerical analysis, especially in numerical methods for solving differential equations (including interval methods) and their application in celestial mechanics. Author or co-author of 3 monographs, 30 textbooks, and over 40 papers in professional journals and conference proceedings. Editor-in-Chief of the Microcomputer User's Library of NAKOM Publisher, and Pro Dialog journal – an official scientific journal of the Polish Information Processing Society. Vice-president (in 1999-2005) and President (since 2005) of the Polish Information Processing Society. Member of the State Accreditation Committee (since 2005).