# Simulation of Dynamic Systems Using Global Calculation Error

Yaroslav Shmygelsky
Department of radiophysics and computer technology
Ivan Franko National University of Lviv,
107 Tarnavsky Str, 79017 Lviv, Ukraine
shmygelsky@ukr.net

Volodymyr Brygilevych
The Bronisław Markiewicz State Higher School of Technology
and Economics in Jarosław (PWSTE in Jarosław)
Jarosław, Poland

*Abstract*—**A numerical value of the estimate of global discretization error is used in the simulation process of transient and steady-state modes in dynamic systems. This allowed the calculation accuracy to by increased significantly with low computational expanses.**

*Index Terms*—**dynamic systems; transient and steady-state modes; numerical methods; global discretization error**

## I. INTRODUCTION

One of the most important characteristics that determine the reliability of the process of numerical simulation of dynamic systems is the global calculation error, that is, the difference between accurate and calculated solutions over the whole range of integration.

The significance of such a characteristic for the researcher is undeniable. In general, it is believed that a good matrix for solving the Cauchy problem for ordinary differential equations should, at the request of the user, provide information about the global calculation error at each starting point [1].

Naturally, the exact definition of global error in modeling programs is not real. We can only talk about some of its approximation. In many cases, the investigator can only satisfy the order of the magnitude of the error and its sign. But if the global error is determined by the corresponding algorithms sufficiently reliable, then its calculated value can be tried to use to correct the results of the calculation conducted with low accuracy.

## II. PRELIMINARY REZULTS

Consider the mathematical model of a dynamical system in the form implicit differential-algebraic equations:

$$F(x,\dot{x},t) = 0, \tag{1}$$

where $x \in R^n$ - is the vector of independent variables, $F : R^{2n+1} \to R^n$ - is the vector-function, continuous by $t$ and continuously differentiated with respect to $x$ and $\dot{x}$.

Let $x(t) = \varphi(t,t_0,x_0)$, $t \in [t_0,t_{fin}]$ - the solution of the system (1) satisfying the initial condition $x(t_0) = x_0$, and

$x_m \approx x(t_m)$ - the numerical approximation of this solution on the grid $\Omega : \{t_{m+1} = t_m + h_m, \ h_m > 0, \ m = 0,1,...\}$.

The global discretization error on a grid $\Omega$ is the value [3]:

$$\delta(t_m) = x_m - x(t_m), \quad t_m \in \Omega. \tag{2}$$

The value of the global discretization error depends on the type of method used for numerical integration of the system (1) and the value of its local truncation error.

If for an algebraization of the derivative in (1) on the grid $\Omega$ an implicit k-step method of the backward differentiation formulas (the BDF method [4]) is used

$$\dot{x}_{m+1} = -\frac{1}{h_{m+1}} \sum_{i=0}^{k} \alpha_i x_{m+1-i} \tag{3}$$

then, as shown in [2], the numerical approximation $\delta_m$ of the global error $\delta(t_m)$ on the grid $\Omega$ satisfies the equation

$$A_{m+1}\delta_{m+1} = B_{m+1}\frac{1}{h_{m+1}}\left[\sum_{i=0}^{k}\alpha_i\delta_{m+1-i} + e_{m+1}\right], \tag{4}$$

where

$$A_{m+1} = \left[\frac{\partial F(x,\dot{x},t)}{\partial x} + \frac{\alpha_0}{h_{m+1}}\frac{\partial F(x,\dot{x},t)}{\partial \dot{x}}\right]_{x=x_{m+1},\dot{x}=\dot{x}_{m+1},t=t_{m+1}},$$

$$B_{m+1} = \left[\frac{\partial F(x,\dot{x},t)}{\partial \dot{x}}\right]_{x=x_{m+1},\dot{x}=\dot{x}_{m+1},t=t_{m+1}}, \tag{5}$$

$e_{m+1}$ - is the local truncation error of the BDF method at the time moment $t_{m+1}$, calculated by the Brayton formula [2]:

$$e_{m+1} = \frac{h_{m+1}}{t_{m+1}-t_{m-k}}\left(x_{m+1} - x_{m+1}^0\right), \tag{6}$$

where $x_{m+1}^0$ - is the predicted value for $x_{m+1}$.

The use of global error in modeling programs is justified in the case where the procedure for determining it is cheap enough. If the discrete system model (1)

$$F(x_{m+1}, \dot{x}_{m+1}, t_{m+1}) = 0, \tag{7}$$

where $\dot{x}_{m+1}$ is determined by the expression (3), is solved at each time step by the Newton method, then the matrices $A$ and $B$ are calculated in the process of numerical integration of the system (1), and the matrix $A$ is usually represented by its LU-expansion. Therefore, the computational cost for finding $\delta_{m+1}$ from equation (4) does not exceed the cost of one Newtonian iteration for equation (7).

Equation (4) is obtained by linearization of the corresponding nonlinear equation, which correlates the true value of the global error $\delta(t_{m+1})$ with the true value of the local error $e(t_{m+1})$ of the numerical method [3]. Therefore, in the case of the linear system (1), equation (4) is valid for any of the values of the global error and the accuracy of the determination $\delta_{m+1}$ is determined only by the error of estimation $e_{m+1}$. In the case of a nonlinear system (1), equation (4) is valid only for small values $\delta_{m+1}$ at the points of the grid $\Omega$, and the domain of its applicability, generally speaking, can be determined only experimentally.

### III. CALCULATION OF TRANSIENT REGIMES

The numerical experiments have shown that the global error for nonlinear systems by using BDF methods is determined by (4) very reliably, even when its value reached 10-15% of the calculated variable value. Therefore, there are every reason to use a numerical approximation of global error to correct the calculation results in order to increase their accuracy.

If $\delta_m$ is a good approximation for $\delta(t_m)$, then, at is follows from (2),

$$\hat{x}_m = x_m - \delta_m \tag{8}$$

will be a better approximation for $x(t_m)$, than $x_m$ at the grid $\Omega$.

**Example 1**. Consider the calculation of the transient process in the sequential linear oscillatory *RLC*-circuit, which is acted upon by the harmonic e.m.f. $E(t) = A\sin(\frac{2\pi}{T}t)$. Mathematical model of *RLC*-circuit:

$$\begin{aligned}
C\dot{U}_C - I_L &= 0, \\
L\dot{I}_L + RI_L + U_C - A\sin\left(\frac{2\pi}{T}\right) &= 0,
\end{aligned} \tag{9}$$

where $U_C$ and $I_L$ - voltage on the capacitor and current through inductance

The circuit parameters are: $R$=10($\Omega$), $L$=10$^{-3}$ (H), $C$=10$^{-7}$ (F); The e.m.f. parameters are: $A$=10(V), $T$=62,5·10$^{-6}$ (s).

The proximity of the external force frequency ($\omega$ = 1.0053 10$^5$ s$^{-1}$) to the resonance circuit frequency ($\omega_0$ = 10$^5$ s$^{-1}$) increases the effect of numerical errors on the calculation results. The simulation was carried out using the BDF method with various fixed orders of formulas (3) from zero initial conditions. In all experiments, the limit of local error at the step was assumed to be 10$^{-4}$. Together with the solution, its global error was determined in accordance with (4). Table 1 present the values of the circuit state variables $U_C(t_M)$ and $I_L(t_M)$ at the instant of time $t_M = 4T$ : exact (obtained from the analytical expression for the transient process), calculated (obtained by numerical integration of the circuit equations) and improved (taking into account the global error in accordance with the formula (8)).

TABLE I. CALCULATION OF TRANSIENT REGIM IN AN OSCILLATORY CIRCUIT

| The order of the BDF method | Value $U_C$ (V) exact: -71,284 | | Value $I_L$ (mA) exact: -28,820 | |
|---|---|---|---|---|
| | calculation | improvement | calculation | improvement |
| 1 | -66,287 | -71,210 | -23,245 | -28,600 |
| 2 | -71,021 | -71,339 | -38,548 | -28,953 |
| 3 | -71,636 | -71,302 | -30,121 | -28,552 |
| 4 | -71,394 | -71,261 | -27,274 | -28,621 |

It follows from the above results, in this example, the global error is determined using equation (4) very confidently. Moreover, the calculated values of the global error make not less than 80% of the true values, and their consideration can significantly improve the accuracy of the simulation.

### IV. CALCULATION OF STEADY-STATE REGIMES

Let the system (1) be time periodic with the period T> 0, that is:

$$F(x, \dot{x}, t+T) = F(x, \dot{x}, t). \tag{10}$$

The search for the periodic regime of the system (1) can be reduced to the solution of the equation [5,6]

$$y - P(y) = 0 , \qquad (11)$$

where $y = x(t_0)$ - vector of initial conditions, $P(y)$ - mapping of a point along the trajectory $x(t) = \varphi(t, t_0, y)$ of system (1) for period $T$.

In the simulation of system (1) we have to deal with not exact mapping $P(y)$, but with its numerical approximation $\tilde{P}(y)$, which is calculated by means of some formulas of numerical integration on a grid $\Omega$, and instead of (11) solve the equation

$$y - \tilde{P}(y) = 0 , \qquad (12)$$

Assume that in the discrete system (7) there is a periodic mode, that is, the mapping $\tilde{P}$ has a fixed point. Let $y^*$ - a fixed point of the exact mapping $P$, and $\tilde{y}^*$ - a fixed point of the calculated mapping $\tilde{P}$. For these fixed points, the relation:

$$y^* = P(y^*), \quad \tilde{y}^* = \tilde{P}(\tilde{y}^*) + d , \qquad (13)$$

where $d$ - some vector satisfying the condition $\|d\| \le \varepsilon$, $\varepsilon > 0$ - the accuracy of solving the equation (12).

We introduce the error of the calculation of the mapping $P$ and the fixed point:

$$\delta = \tilde{P}(y) - P(y) , \qquad (14)$$

$$\eta = \tilde{y}^* - y^* , \qquad (15)$$

The value $\delta$ is simply defined in (2) the global error of a numerical method at a time moment $t_N = t_0 + T$ when integrating the system (1) from the initial conditions $x(t_0) = y$ on the interval $[t_0, t_0 + T]$. From (15), taking into account (13), we have:

$$\eta = \tilde{y}^* - y^* = \tilde{P}(\tilde{y}^*) - P(y^*) + d = \left(\tilde{P}(\tilde{y}^*) - \tilde{P}(y^*)\right) + \left(\tilde{P}(y^*) - P(y^*)\right) + d , \qquad (16)$$

If the BDF method (3) is used for the system (1) discretization, then, under the assumptions made about the $F(x, \dot{x}, t)$ smoothness, it can be shown [5] that the calculated mapping $\tilde{P}$ will be differentiable whit the derivative continuous in the sense of Lipschitz, that is|:

$$\tilde{P}(y) - \tilde{P}(z) = \tilde{P}'(y)(y - z) + O\left(\|y - z\|^2\right), \quad \forall y, z \in R^n , \quad (17)$$

under the condition that $\tilde{P}(y)$ i $\tilde{P}(z)$ are calculated on the same grid $\Omega$ and in the same order of $k$ formulas of the numerical method (3) in the corresponding steps, and the matrix $A_{m+1}$ defined in (5) is regular for all $m = 0, 1, \ldots N - 1$.

Then we obtain from (16) taking into account (14) and (17)

$$\eta \approx \left(I - \tilde{P}'(y^*)\right)^{-1} (\delta + d) , \qquad (18)$$

where $I$ - is a unit matrix, and the sign $\approx$ means "with precision to members $O\left(\|\eta\|^2\right)$".

If we do not take into account the rounding errors, then in the linear case it is possible to put $d = 0$, and in the nonlinear case, with the iterative process of the equation (12) is convergent (for example, Newton's method), the value $\|d\|$ can be very low.

Consequently, the error $\eta$ for calculating a fixed point $y^*$ is mainly defined by the value $\left\|\left(I - \tilde{P}'(y^*)\right)^{-1}\right\|$ and global error $\delta$ of the mapping calculation.

For weakly damped systems $\left\|\left(I - \tilde{P}'(y^*)\right)^{-1}\right\| \gg 1$. So, in the case of linear system (1), it may be shown [6], that $\left\|\left(I - \tilde{P}'(y^*)\right)^{-1}\right\| \sim Q$, where $Q$ - is the maximum figure of merit of the system (1). Therefore, when looking for a periodic regime of weakly damped systems to obtain reliable results, it is necessary to calculate with very high accuracy, which naturally leads to significant computational costs.

For this to be eliminated, in the search process of the steady-state regime the improved value of the calculated mapping $\tilde{P}(y)$ may be used as

$$\hat{P}(y) = \tilde{P}(y) - \delta_N , \qquad (19)$$

where $\delta_N$ - is the numerical approximation of global error $\delta(t_N)$.

**Example** 2. Consider the problem of periodic states definition of the Duffing equation:

$$\begin{aligned} \dot{x}_1 &= x_2 \\ \dot{x}_2 &= -x_1^3 - 0,2x_2 + 0,3\cos t \end{aligned} \qquad (20)$$

As is known, this equation has three periodic states, two of which are stable and one unstable. In [6] the periodic states of equation (20) are calculated, with high accuracy. Whit the initial conditions $(x_1(0) = 0,027; x_2(0) = 1,10)$ the stable periodic state $(x_1^*(0) = 0,6263873; x_2^*(0) = 1,03347995)$ is obtained, and with $(x_1(0) = -0,027; x_2(0) = 0,729)$ the unstable periodic state an unstable periodic state

$\left( x_1^*(0) = -0,71598261; \ x_2^*(0) = 0,74740203 \right)$. Let's take these values as benchmarks.

In our case, the mapping was calculated by numerical integration of the equation (20) on the interval $[0.2\pi]$ by the BDF method of the first order (implicit Euler method) with a constant step. The derivative of the reflection for the period $\tilde{P}'(y)$ was determined by solving the equation in the variations for (20) simultaneously with $\tilde{P}(y)$. The equation (12) was solved by Newton's method with $\varepsilon = 10^{-4}$ accuracy. The accuracy of the calculation of the mapping $\tilde{P}(y)$ from experiment to experiment was increased by reducing the value of the integration step. Experiment results are given in Table 2

In the first column of Table 2, the number of steps for integrating N in the period of forced fluctuations T in the various experiments is indicated. The last line shows the calculation results for N = 1000, taking into account calculated value of global sampling error. The first line of each part of Table 3 gives the initial conditions from which the search for periodicity began. The third column of each part of the table specifies the processor time for solving the task (along with the output of the results).

TABLE II.       CALCULATION OF PERIODIC STATES OF DUFFING EQUATION

| N=T/h | $x_1(0) = 0{,}027; \ x_2(0) = 1{,}1$ | | | $x_1(0) = -0{,}027; \ x_2(0) = 0{,}729$ | | |
|---|---|---|---|---|---|---|
| | $x_1^*(0)$ | $x_2^*(0)$ | $t$ s | $x_1^*(0)$ | $x_2^*(0)$ | $t$ s |
| 1000 | 0,575 | 1,039 | 30 | -0,689 | 0,767 | 26 |
| 5000 | 0,620 | 1,034 | 129 | -0,711 | 0,750 | 100 |
| 10000 | 0,629 | 1,032 | 254 | -0,715 | 0,748 | 198 |
| 1000+gl.error | 0,627 | 1,033 | 32 | -0,717 | 0,747 | 28 |

Comparing the obtained values $x_1^*(0)$, $x_2^*(0)$ of periodic states with reference ones, it is not difficult to make sure that the use of numerical approximation of global sampling error in the process of solution allowed, at small computational costs, to increase the accuracy of determining the periodic state of the system (20) and significantly accelerate the process of its search.

REFERENCES

[1] H. J. Stetter, "Global Error Estimation in Ordinary Initial Value Problems," Lect. Notes Math., 1982, vol. 968, pp. 269–279.

[2] "Modern Numerical Methods for Ordinary Differential Equation," Edited by G. Hall and J. M. Watt, Clarendon Press, Oxford, 1976 (Russian translation, 1979).

[3] R. K. Brayton, F. G. Gustavson, G. D. Hechtel, "A New Efficient Algorithm for Solving Differential-Algebraic Systems Using Implicit Backward Differentiation Formula," Proc. IEEE, vol. 60, No. 1, January 1972, pp. 98-108.

[4] Ya. Shmygelsky, "On Definition of Global Error in Numerical Calculation of Electronic Circuit Equations," Teor. Electrotekhn., 1892, vyp. 34, pp.124–129 (in Russian).

[5] S. Scalboe, "Time Stationary Analisis of Nonlinear Electrical Systems," Proc. IEEE, vol. 70, No. 10, October 1982, pp. 89-111.

[6] T. J. Aprill and T. N. Trick, "Steady State Analisis of Nonlinear Circuits with Periodic Inputs," Proc. IEEE, vol. 60, No. 1, January 1972, pp. 108-114.