

NEW COLLOCATION INTEGRATOR FOR SOLVING DYNAMIC PROBLEMS. I. THEORETICAL BACKGROUND

V. A. Avdyushev

UDC 519.62:531

A new collocation integrator with Lobatto spacings is proposed for numerical solving mixed systems of first- and second-order differential equations for dynamic problems. The general theory of collocation integrators is outlined from which the basic formulas of the new integrator are derived.

Keywords: numerical integrators, collocation methods, ordinary differential equations, dynamical systems.

INTRODUCTION

Collocation integrators [1–4] are simple, graceful, convenient, and at the same time powerful tools for solving differential equations of dynamics. The idea of the collocation integrators is rather transparent, and knowledge of other integrators [3], for example, Runge–Kutta, Gragg–Bulirsch–Stoer, or Adams integrators are not required to understand it. It is suffice to know only what are differential equations, integration, and interpolation. Therefore, it is well justified to consider the collocation integrators as a whole as an independent class, the most known representatives of which are the implicit Runge–Kutta collocation integrators.

A remarkable feature of the collocation integrators is that their theoretical basis, as well as software implementation, is universal for an arbitrary order [5]. Practically, the order is defined by spacing within a step, namely, by the number and specificity of distribution of nodal values in terms of which all other integrator constants are expressed. In addition, with Gaussian Legendre or Lobatto quadrature spacings, the collocation integrators become geometrical [4]: symmetric and orbitally stable,¹ and with Legendre spacings, also symplectic. It should also be noted that, unlike others, the collocation integrators allow one to design easily within each step an approximate analytical solution, which is convenient to use for obtaining results on a dense time grid.

In the present work, a new collocation integrator Lobbie is proposed with Lobatto spacings. Actually, its prototype is the Everhart integrator widely used in dynamical astronomy [6, 7]. More precisely, it is the result of cardinal revision of the predecessor, although the theory, algorithmization, and program code of the new integrator can only conceptually remind of the author's version of the glorified Everhart integrator.

The general theory of collocation integrators with application to solving the differential equations of the first and second orders is briefly outlined. Particular examples, including the Everhart integrator, are given followed by the derivation of the main formulas for the integrator Lobbie and the special features of their software implementation. The procedure Lobbie is described.

¹ This definition with reference to the Runge–Kutta method was introduced in [5].

1. COLLOCATION METHODS

1.1. First-order differential equations

Let a dynamic state \mathbf{x} be described as a function of time t by the first-order vector differential equation

$$\mathbf{x}' = \mathbf{f}(t, \mathbf{x}) \quad (1)$$

with the known initial dynamic state at time t_0 :

$$\mathbf{x}_0 = \mathbf{x}(t_0) . \quad (2)$$

Here the prime designates the full time derivative and \mathbf{f} is the known vector function of time and dynamic state. It is required to determine the dynamic state of the system at time $t_0 + h$:

$$\mathbf{x}(t_0 + h) , \quad (3)$$

where h is a small parameter (size of the integration step).

We represent the approximate solution of Eq. (1) in the form of a polynomial

$$\mathbf{u}(t) \approx \mathbf{x}(t) , \quad (4)$$

which must precisely satisfy to the equation at some intermediate moments of time $t_i \in [t_0, t_0 + h]$ ($i = 1, \dots, s$) (at collocation points)

$$\mathbf{u}'(t_i) = \mathbf{f}(t_i, \mathbf{u}(t_i)) \quad (i = 1, \dots, s) , \quad (5)$$

and to initial condition (2)

$$\mathbf{x}_0 = \mathbf{u}(t_0) . \quad (6)$$

Then solution (3), with allowance for (4), is approximately determined as

$$\mathbf{x}_1 = \mathbf{u}(t_0 + h) \approx \mathbf{x}(t_0 + h) . \quad (7)$$

The geometric sense of (5) is that the tangents to the polynomial at the collocation points must coincide (collocate) with the directions of the vector field generated by the function of the differential equation \mathbf{f} (Fig. 1). Though the values of the polynomial can differ considerably from the exact solution.

Collocation conditions (5) can be considered as the Lagrange conditions imposed on the derivative of polynomial (4):

$$\mathbf{u}'(t_0 + h\tau) \equiv \mathbf{p}(\tau) , \quad (8)$$

which in this case plays the role of the polynomial interpolant of the function \mathbf{f} with respect to the dimensionless variable τ . According to (5), the Lagrange conditions for the nodal values c_1, \dots, c_s of the dimensionless variable τ can be represented in the form

$$\mathbf{p}(c_i) = \mathbf{f}_i \equiv \mathbf{f}(t_i, \mathbf{u}_i), \quad \mathbf{u}_i \equiv \mathbf{u}(t_i), \quad t_i = t_0 + hc_i \quad (i = 1, \dots, s) . \quad (9)$$

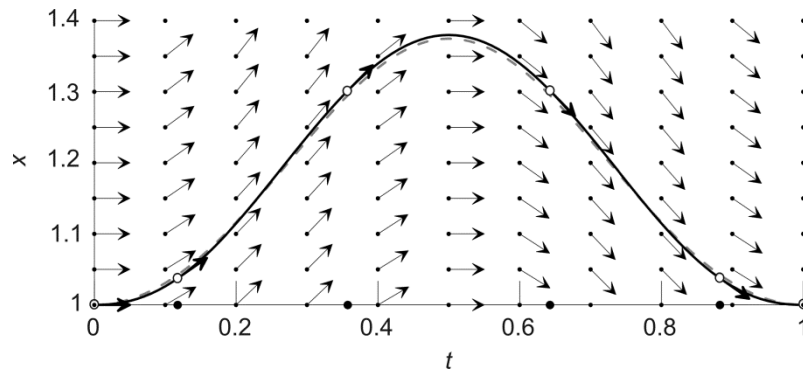


Fig. 1. Collocation conditions for the differential equation $x' = f(t, x) = \sin(2\pi t)x$ with the initial condition $x_0 = x(0) = 1$. The arrows indicate directions of the vector field $(\cos \alpha, \sin \alpha)$ at points (t, x) , where $\tan \alpha = f(t, x)$. The black solid curve shows the collocation polynomial u (the approximated solution) with the Lobatto spacings (filled circles) for $s = 6$; the dashed curve shows the exact solution $x = e^{(1-\cos(2\pi t))/2\pi}$; the collocation points for planes (t, x) are indicated by open circles.

Forming the polynomial interpolant \mathbf{p} from conditions (9), then integrating relation (8) over τ on the segment $[0, 1]$, and taking into account (6) and (7):

$$\int_0^1 \mathbf{u}'(t_0 + h\tau) d\tau = \frac{\mathbf{u}(t_0 + h\tau)|_0^1}{h} = \frac{\mathbf{x}_1 - \mathbf{x}_0}{h} = \int_0^1 \mathbf{p}(\tau) d\tau,$$

we obtain the approximate solution

$$\mathbf{x}_1 = \mathbf{x}_0 + h \int_0^1 \mathbf{p}(\tau) d\tau.$$

The interpolant \mathbf{p} (8) is constructed from the intermediate approximate solutions $\mathbf{u}_1, \dots, \mathbf{u}_s$ (9), that is, $\mathbf{p} = \mathbf{p}(\tau, \mathbf{u}_1, \dots, \mathbf{u}_s)$. Each i th solution is also determined by integrating (8) over τ , but on the segment $[0, c_i]$. Thus, the collocation method of solving differential equation (1) can be represented by the set of formulas

$$\mathbf{x}_1 = \mathbf{x}_0 + h \int_0^1 \mathbf{p}(\tau, \mathbf{u}_1, \dots, \mathbf{u}_s) d\tau, \quad \mathbf{u}_i = \mathbf{x}_0 + h \int_0^{c_i} \mathbf{p}(\tau, \mathbf{u}_1, \dots, \mathbf{u}_s) d\tau \quad (i = 1, \dots, s). \quad (10)$$

The intermediate solutions in (10) are implicitly expressed through the nonlinear equations. For this reason, all collocation methods are implicit.² As a rule, the equations are solved by the fixed-point iteration in the Seidel modifications, that is, refining one by one the intermediate solutions at each iteration step. In fact, the iterative solution of nonlinear equations is the kernel of any collocation integrator, and its efficiency depends in many respects on how successfully organized is the iterative process.

The order p of the method is determined by the number s of collocation points and the specificity of their distribution. The collocation principle allows one to obtain practically any order. So, for any arbitrary distribution

² Except the explicit Euler method (of the first order) which is the collocation method with the Radau I spacings for $s = 1$.

c_1, \dots, c_s (for example, uniform) at least $p = s$ [3]. However, using the nodal values of the Gaussian Legendre, Radau, or Lobatto quadratures, the order can be increased to $p = 2s$, $p = 2s - 1$, and $p = 2s - 2$, respectively [3, 4, 8, 9]. These nodal values are solutions of the algebraic equations

$$\begin{aligned} \frac{d^s}{d\tau^s}[\tau^s(\tau-1)^s] &= 0 \quad (\text{Legendre}), \\ \frac{d^{s-1}}{d\tau^{s-1}}[\tau^s(\tau-1)^{s-1}] &= 0 \quad (\text{Radau I}); \quad \frac{d^{s-1}}{d\tau^{s-1}}[\tau^{s-1}(\tau-1)^s] = 0 \quad (\text{Radau II}), \\ \frac{d^{s-2}}{d\tau^{s-2}}[\tau^{s-1}(\tau-1)^{s-1}] &= 0 \quad (\text{Lobatto}). \end{aligned} \tag{11}$$

The collocation methods are remarkable in that, in fact, they give the analytical solution within the step in the form of the collocation polynomial

$$\mathbf{u}(\tau) = \mathbf{x}_0 + h \int_0^\tau \mathbf{p}(\tau) d\tau$$

(not only the solution \mathbf{x}_1 at the end of the step and several intermediate solutions $\mathbf{u}_1, \dots, \mathbf{u}_s$), which is very convenient for derivation of approximate solutions for a dense time grid. However, it should be borne in mind that the order of accuracy within the step is reduced to $p = s$ [3]. In addition, the presence of the interpolant allows one to obtain a sufficiently good initial approximation of the function \mathbf{f} at the next step by extrapolation:

$$\mathbf{f}_i = \mathbf{p}(1 + c_i) \quad (i = 1, \dots, s)$$

for the subsequent iterative determination of intermediate solutions (10).

It should especially be noted that the collocation methods with the Legendre and Lobatto spacings become geometrical [4]: symmetric and orbitally stable, and the Legendre spacing, also symplectic. It seems that the Legendre spacing is more preferable, because with the same number of collocation points, its order is higher by 2. However, in spite of the fact that the order of the method with the Lobatto spacings is lower, it works a little faster. Indeed, at each step with ni iterations for determining intermediate solutions, the number of calculations of the equation function is equal to $ncf = ni \cdot (s - 1)$, whereas with the Legendre spacings, $ncf = ni \cdot s$. In addition, the interpolant of the equation function is constructed for the entire integration segment $[t_0, t_0 + h]$ with $c_1 = 0$ and $c_s = 1$, unlike the Legendre spacings for which all nodal points lie inside the segment. Hence, the predictor with the Lobatto spacings is better, which is very important for its software implementation.

As simple examples of the collocation methods, consider the Runge–Kutta methods well-known for a long time and including the explicit and implicit Euler methods (Radau I and II: $s = 1$, $p = 1$), the midpoint method (Legendre: $s = 1$, $p = 2$), the trapezoidal rule (Lobatto: $s = 2$, $p = 2$), and the Simpson method (Lobatto: $s = 3$, $p = 4$). If we take the Lagrange polynomial

$$\mathbf{p}(\tau) = \sum_{j=1}^s \mathbf{f}_j \prod_{k \neq j} \frac{\tau - c_k}{c_j - c_k}$$

as an interpolant of equation function (1), collocation method (10) acquires the classical form of the Runge–Kutta method [1–4]:

$$\mathbf{x}_1 = \mathbf{x}_0 + h \sum_{j=1}^s b_j \mathbf{f}_j, \quad \mathbf{u}_i = \mathbf{x}_0 + h \sum_{j=1}^s a_{ij} \mathbf{f}_j \quad (i=1, \dots, s), \quad (12)$$

where the constants are expressed through the integrals of the basic Lagrange functions:

$$a_{ij} = \int_0^{c_i} \prod_{k \neq j} \frac{\tau - c_k}{c_j - c_k} d\tau, \quad b_j = \int_0^1 \prod_{k \neq j} \frac{\tau - c_k}{c_j - c_k} d\tau.$$

Soon after the publication of the first works on the Runge–Kutta collocation methods [1, 2], E. Everhart [6, 7] proposed to use the polynomial in the canonical form³

$$\mathbf{p}(\tau) = \sum_{j=1}^s \mathbf{a}_j \tau^{j-1}$$

as an interpolant. Here $\mathbf{a}_j = \mathbf{a}_j(\mathbf{f}_1, \dots, \mathbf{f}_s)$ ($j=1, \dots, s$). The simple interpolation yields a sufficiently simple form of the approximate solution

$$\mathbf{x}_1 = \mathbf{x}_0 + h \sum_{j=1}^s \frac{\mathbf{a}_j}{j}, \quad \mathbf{u}_i = \mathbf{x}_0 + h \sum_{j=1}^s \frac{\mathbf{a}_j}{j} c_i^j \quad (i=1, \dots, s). \quad (13)$$

However, to express the coefficients of the polynomial \mathbf{a} through the collocation values of the function \mathbf{f} , the author expressed the divided differences of the Newton polynomial \mathbf{a} directly through the collocation values. Meanwhile, the coefficients of the canonical polynomial are expressed in terms of the divided differences by means of the linear relations

$$\mathbf{a}_j = \sum_{i=j}^s c_{ij} \mathbf{a}_i \quad (j=1, \dots, s),$$

where the constants are calculated from the nodal values of spacing c_1, \dots, c_s using the recurrent formulas

$$c_{ii} = 1 \quad (i > 0), \quad c_{ij} = c_{i-1, j-1} - c_{i-1} c_{i-1, j} \quad (i > j > 0).$$

1.2. Second-order differential equations

Suppose now that the dynamical system is described by the second-order vector differential equation

$$\mathbf{x}'' = \mathbf{f}(t, \mathbf{x}, \mathbf{x}') \quad (14)$$

with the known initial dynamic state at time t_0 :

$$\mathbf{x}_0 = \mathbf{x}(t_0), \quad \mathbf{x}'_0 = \mathbf{x}'(t_0). \quad (15)$$

It is required to determine the dynamic state of the system at time $t_0 + h$:

³ Though the author himself did not present his method as a collocation one.

$$\mathbf{x}(t_0 + h), \quad \mathbf{x}'(t_0 + h). \quad (16)$$

Here \mathbf{x} and \mathbf{x}' are the coordinate and velocity vectors, respectively.

We represent the approximate solutions for the coordinates and velocity in the form of polynomials

$$\mathbf{u}(t) \approx \mathbf{x}(t), \quad \mathbf{v}(t) = \mathbf{u}'(t) \approx \mathbf{x}'(t), \quad (17)$$

which should satisfy to Eq. (14) at collocation times $t_i \in [t_0, t_0 + h]$ ($i = 1, \dots, s$):

$$\mathbf{u}''(t_i) = \mathbf{v}'(t_i) = \mathbf{f}(t_i, \mathbf{u}(t_i), \mathbf{v}(t_i)) \quad (i = 1, \dots, s), \quad (18)$$

and to the initial condition (15):

$$\mathbf{x}_0 = \mathbf{u}(t_0), \quad \mathbf{x}'_0 = \mathbf{v}(t_0). \quad (19)$$

Then solutions (16) with allowance for (17) are approximately determined as

$$\mathbf{x}_1 = \mathbf{u}(t_0 + h) \approx \mathbf{x}(t_0 + h), \quad \mathbf{x}'_1 = \mathbf{v}(t_0 + h) \approx \mathbf{x}'(t_0 + h). \quad (20)$$

According to (18), the Lagrange conditions for the interpolant \mathbf{p} of the function \mathbf{f}

$$\mathbf{u}''(t_0 + h\tau) = \mathbf{v}'(t_0 + h\tau) \equiv \mathbf{p}(\tau) \quad (21)$$

can be represented in the form

$$\mathbf{p}(c_i) = \mathbf{f}_i \equiv \mathbf{f}(t_i, \mathbf{u}_i, \mathbf{v}_i), \quad \mathbf{u}_i \equiv \mathbf{u}(t_i), \quad \mathbf{v}_i \equiv \mathbf{v}(t_i), \quad t_i = t_0 + hc_i \quad (i = 1, \dots, s). \quad (22)$$

Forming the polynomial interpolant \mathbf{p} from conditions (22) and then integrating (21) over τ on the segments $[0, 1]$ and $[0, c_i]$ ($i = 1, \dots, s$), we obtain a set of formulas of the collocation method for differential equation (14):

$$\mathbf{x}_1 = \mathbf{x}_0 + \mathbf{x}'_0 h + h^2 \int_0^1 \int_0^\tau \mathbf{p}(\tau) d\tau^2, \quad \mathbf{x}'_1 = \mathbf{x}'_0 + h \int_0^1 \mathbf{p}(\tau) d\tau, \quad (23)$$

$$\mathbf{u}_i = \mathbf{x}_0 + \mathbf{x}'_0 hc_i + h^2 \int_0^{c_i} \int_0^\tau \mathbf{p}(\tau) d\tau^2, \quad \mathbf{v}_i = \mathbf{x}'_0 + h \int_0^{c_i} \mathbf{p}(\tau) d\tau \quad (i = 1, \dots, s),$$

where the subintegral polynomial \mathbf{p} depends also on the intermediate solutions $\mathbf{u}_1, \dots, \mathbf{u}_s$ and $\mathbf{v}_1, \dots, \mathbf{v}_s$.

1.3. Mixed systems of the first- and second-order differential equations

The problem described by (14) and (15) can be represented in the form of (1) and (2):

$$\mathbf{x}'' = \dot{\mathbf{x}}, \quad \dot{\mathbf{x}}' = \mathbf{f}(t, \mathbf{x}, \dot{\mathbf{x}}), \quad \mathbf{x}_0 = \mathbf{x}(t_0), \quad \dot{\mathbf{x}}_0 = \dot{\mathbf{x}}(t_0) = \mathbf{x}'_0. \quad (24)$$

In spite of the fact that both problems describe the same dynamical system, according to the collocation principle, their numerical solutions will differ significantly. Indeed, according to (10), for the alternative problem we have the approximate solutions

$$(\mathbf{x}_1, \dot{\mathbf{x}}_1) = (\mathbf{x}_0, \dot{\mathbf{x}}_0) + h \int_0^1 (\mathbf{q}(\tau), \mathbf{p}(\tau)) d\tau, \quad (\mathbf{u}_i, \mathbf{v}_i) = (\mathbf{x}_0, \dot{\mathbf{x}}_0) + h \int_0^{c_i} (\mathbf{q}(\tau), \mathbf{p}(\tau)) d\tau \quad (i=1, \dots, s), \quad (25)$$

where the polynomial \mathbf{q} interpolates the right-hand side of the equation for the coordinate vector. The accuracies of solutions (25) (for coordinate and velocity vectors) are of the same order p . However, the solution for the coordinates in (23) has the order $p+1$ owing to double integration; in other words, the coordinates in (23) are determined more precisely than in (25).

Nevertheless, sometimes the user is compelled to use the representation of the dynamical system in the form of Eqs. (24), thereby using the integrator for first-order equations (25) understanding that this inevitably reduces the accuracy of the approximate solution. This is necessary when Eq. (14) describing the dynamical system is supplemented with the first-order equations for some auxiliary dynamic quantities. For example, the mixed systems are used to investigate dynamical chaos [10] as well as for linearization, regularization, and stabilization of the dynamic equations [11–14]. In this case, Eqs. (24) are used instead of Eq. (14) in order that to reduce all equations of the mixed system to the same order. The alternative variant is to differentiate the additional equations so that the system as a whole was described by Eq. (14), though for complex dynamical systems, taking derivatives of functions of the additional equations is often practically impossible.

The natural and effective approach to solving mixed system of equations is the application of the hybrid integrator. If it is required to solve the first-order vector equation for auxiliary dynamic quantities \mathbf{z} together with Eq. (14):

$$\mathbf{z}' = \mathbf{g}(t, \mathbf{x}, \mathbf{x}', \mathbf{z}) \quad (26)$$

with the initial condition $\mathbf{z}_0 = \mathbf{z}(t_0)$, the set of formulas (23) should be supplemented with the formulas

$$\mathbf{z}_1 = \mathbf{z}_0 + h \int_0^1 \mathbf{r}(\tau) d\tau, \quad \mathbf{w}_i = \mathbf{z}_0 + h \int_0^{c_i} \mathbf{r}(\tau) d\tau \quad (i=1, \dots, s). \quad (27)$$

Here the interpolant \mathbf{r} of the function \mathbf{g} is constructed by analogy with the interpolant \mathbf{p} from the intermediate solutions $\mathbf{u}_1, \dots, \mathbf{u}_s$, $\mathbf{v}_1, \dots, \mathbf{v}_s$, and $\mathbf{w}_1, \dots, \mathbf{w}_s$. Thus, the hybrid collocation method for Eqs. (14) and (26) has the form of (23) and (27).

2. INTEGRATOR LOBBIE

For the mixed system of differential equations (14) and (26)

$$\mathbf{x}'' = \mathbf{f}(t, \mathbf{x}, \mathbf{x}', \mathbf{z}), \quad \mathbf{z}' = \mathbf{g}(t, \mathbf{x}, \mathbf{x}', \mathbf{z}), \quad \mathbf{x}_0 = \mathbf{x}(t_0), \quad \mathbf{x}'_0 = \mathbf{x}'(t_0), \quad \mathbf{z}_0 = \mathbf{z}(t_0), \quad (28)$$

we take as interpolants the Newton polynomials

$$\mathbf{p}(\tau) = \sum_{j=1}^s \boldsymbol{\alpha}_j \prod_{k=1}^{j-1} (\tau - c_k) \quad \text{and} \quad \mathbf{r}(\tau) = \sum_{j=1}^s \boldsymbol{\beta}_j \prod_{k=1}^{j-1} (\tau - c_k) \quad (29)$$

with the Lobatto spacings ($c_1 = 0, c_s = 1$). Here $\prod_{k=1}^0 = 1$. The divided differences in (29) are determined in terms of the nodal values of the functions \mathbf{f} and \mathbf{g} using the recurrent formulas

$$\begin{aligned} \boldsymbol{\alpha}_j = \mathbf{f}_j, \quad \boldsymbol{\beta}_j = \mathbf{g}_j, \quad \boldsymbol{\alpha}_j := (\boldsymbol{\alpha}_j - \boldsymbol{\alpha}_k) / (c_j - c_k), \quad \boldsymbol{\beta}_j := (\boldsymbol{\beta}_j - \boldsymbol{\beta}_k) / (c_j - c_k) \\ (j = 1, \dots, s; k = 1, \dots, j-1). \end{aligned} \quad (30)$$

Substituting interpolants (29) into (23) and (27) as well as considering that $c_1 = 0$ and $c_s = 1$, we obtain the collocation method for system (28) in the form

$$\begin{aligned} \mathbf{u}_1 &= \mathbf{x}_0, \quad \mathbf{v}_1 = \mathbf{x}'_0, \quad \mathbf{w}_1 = \mathbf{z}_0, \\ \mathbf{u}_i &= \mathbf{x}_0 + \mathbf{x}'_0 h c_i + h^2 \sum_{j=1}^s a_{ij} \mathbf{a}_j, \quad \mathbf{v}_i = \mathbf{x}'_0 + h \sum_{j=1}^s b_{ij} \mathbf{a}_j, \quad \mathbf{w}_i = \mathbf{z}_0 + h \sum_{j=1}^s b_{ij} \mathbf{\beta}_j \quad (i = 2, \dots, s), \\ \mathbf{x}_1 &= \mathbf{u}_s, \quad \mathbf{x}'_1 = \mathbf{v}_s, \quad \mathbf{z}_1 = \mathbf{w}_s, \end{aligned} \quad (31)$$

where

$$a_{ij} = \int_0^{c_i} \int_0^\tau \prod_{k=1}^{j-1} (\tau - c_k) d\tau^2, \quad b_{ij} = \int_0^{c_i} \prod_{k=1}^{j-1} (\tau - c_k) d\tau \quad (i, j = 1, \dots, s). \quad (32)$$

We designate the k -fold integral of the j th Newton basic function by

$$\gamma_{jk}(\tau) = \int_0^\tau \dots \int_0^\tau (\tau - c_1) \dots (\tau - c_{j-1}) d\tau^k. \quad (33)$$

If one considers $\gamma_{jk}(\tau)$ ($j, k = 1, \dots, s+1$) as elements of some matrix Γ of size $(s+1) \times (s+1)$, the constants of the collocation method – integrals of the Newton basic functions (32) – will form its first two columns:

$$a_{ij} = \gamma_{j2}(c_i), \quad b_{ij} = \gamma_{j1}(c_i) \quad (i, j = 1, \dots, s).$$

Meanwhile, the elements of the matrix Γ for any arbitrary value τ are calculated row-by-row using the recurrent relations:

$$\gamma_{1k} = \tau^k / k! \quad (k = 1, \dots, s+1); \quad \gamma_{jk} = (\tau - c_{j-1}) \gamma_{j-1,k} - k \gamma_{j-1,k+1} \quad (j = 2, \dots, s; k = 1, \dots, s-j+2). \quad (34)$$

Nonlinear equations (31) for $\mathbf{u}_1, \dots, \mathbf{u}_s$, $\mathbf{v}_1, \dots, \mathbf{v}_s$, and $\mathbf{w}_1, \dots, \mathbf{w}_s$ are solved at each step by the fixed-point iteration in the Seidel modification. Before iterations, the solutions \mathbf{u}_1 , \mathbf{v}_1 , and \mathbf{w}_1 are known together with the divided differences

$$\mathbf{a}_1 = \mathbf{f}_1 = \mathbf{f}(t_0, \mathbf{u}_1, \mathbf{v}_1, \mathbf{w}_1) \quad \text{and} \quad \mathbf{\beta}_1 = \mathbf{g}_1 = \mathbf{g}(t_0, \mathbf{u}_1, \mathbf{v}_1, \mathbf{w}_1)$$

at the first collocation point $c_1 = 0$. In the beginning of the iterative process at the second collocation point c_2 , the group of solutions \mathbf{u}_2 , \mathbf{v}_2 , and \mathbf{w}_2 is determined, and from them \mathbf{f}_2 and \mathbf{g}_2 , by which the divided differences \mathbf{a}_2 and $\mathbf{\beta}_2$ are refined using recurrent formulas (30). Then the divided differences \mathbf{a}_3 and $\mathbf{\beta}_3$ at the third collocation point c_3 are refined in the same way. After successive refinement of all divided differences at the step, the iteration is repeated. The iterative process is continued until the inequality

$$\|\mathbf{u}_s - \mathbf{u}_s^*\| < \varepsilon \|\mathbf{u}_s\| \quad (35)$$

is satisfied. Here \mathbf{u}_s^* is the solution \mathbf{u}_s at the preceding iteration and ε is a small parameter that determines the accuracy of convergence. By the way, the number of iterations at the step can be preset disregarding condition (35). In

this case, it should be borne in mind that the obtained solution \mathbf{u}_s will not correspond to the preset order, and the method will lose the geometrical properties.

As initial values of the divided differences \mathbf{a} and \mathbf{b} , one accepts their estimates obtained from recurrent formulas (30) using extrapolated values of the functions \mathbf{f} and \mathbf{g} :

$$\mathbf{f}_i = \mathbf{p}(1 + c_i), \quad \mathbf{g}_i = \mathbf{r}(1 + c_i) \quad (i = 1, \dots, s). \quad (36)$$

At the first step, in the absence of such estimates, the iterative process begins with zero values of the divided differences, and continues until condition (35) is satisfied.

For multiple derivation of solutions on a dense time grid, it is convenient and expedient to use an appropriate collocation polynomial covering the times of the grid at a certain step rather than to carry out step-by-step integration at each of these times. The collocation polynomials for system (28) at any arbitrary time $t_0 + h\tau$ can be represented in the form

$$\mathbf{u}(\tau) = \mathbf{x}_0 + \mathbf{x}'_0 h\tau + h^2 \sum_{j=1}^s \gamma_{j2}(\tau) \mathbf{a}_j, \quad \mathbf{v}(\tau) = \mathbf{x}'_0 + h \sum_{j=1}^s \gamma_{j1}(\tau) \mathbf{a}_j, \quad \mathbf{w}(\tau) = \mathbf{z}_0 + h \sum_{j=1}^s \gamma_{j1}(\tau) \mathbf{b}_j, \quad (37)$$

where $\gamma_{j1}(\tau)$ and $\gamma_{j2}(\tau)$ ($j = 1, \dots, s$) are calculated by formulas (34).

An alternative method of obtaining the time series of approximate solutions is polynomial interpolation on intermediate solutions, for example, of type (29):

$$(\mathbf{u}, \mathbf{v}, \mathbf{w})(\tau) = \sum_{j=1}^s (\mathbf{u}_j, \mathbf{v}_j, \mathbf{w}_j) \prod_{k=1}^{j-1} (\tau - c_k). \quad (38)$$

However, the accuracy of interpolation polynomials (38) between the collocation points appears much lower than the accuracy of collocation polynomials (37) obtained by direct integration of polynomials (29).

The step size h as an integrator parameter is assigned by the user. However, the automatic choice of the step size is possible during step-by-step integration. The h value is chosen so that the approximate estimate of the s -order term of the Taylor series for the velocity vector was retained [5]:

$$\|\mathbf{e}\|_{\text{cal}} = \frac{h}{s} \|\mathbf{a}_s\| \approx \frac{h^s}{s!} \|\mathbf{x}^{(s)}\|, \quad (39)$$

that is, if the term of the series is considered as the principal error term, the step size should be determined as for a method of order $p = s - 1$, though the order with the Gaussian quadrature spacings is significantly higher (not less than twice). Unfortunately, it is practically impossible to obtain the estimate of higher order.

Suppose that estimate (39) should be equal to a constant $\|\mathbf{e}\|_{\text{tol}}$ assigned by the user. Since

$$\|\mathbf{e}\|_{\text{cal}} \approx \frac{h^s}{s!} \|\mathbf{x}^{(s)}\| \quad \text{and} \quad \|\mathbf{e}\|_{\text{tol}} \approx \frac{\tilde{h}^s}{s!} \|\mathbf{x}^{(s)}\|,$$

the step size \tilde{h} needed to provide $\|\mathbf{e}\|_{\text{tol}}$ can be estimated as

$$\tilde{h} = h \left(\frac{\|\mathbf{e}\|_{\text{tol}}}{\|\mathbf{e}\|_{\text{cal}}} \right)^{1/s}. \quad (40)$$

After obtaining estimate (40) at the current step, integration is not repeated with the new step size \tilde{h} , but at the next step, it is used according to the formula

$$h_{n+1} = rh_n, \quad r = \left(\frac{s \|e\|_{\text{tol}}}{h_n \|\alpha_s\|} \right)^{1/s}, \quad (41)$$

where n and $n+1$ are the numbers of the current and subsequent steps, respectively. We note that the automatic choice of the step size implies the modification of predictor (36), namely

$$f_i = p(1 + rc_i), \quad g_i = r(1 + rc_i) \quad (i = 1, \dots, s). \quad (42)$$

Algorithm (41) is effective if h_n and h_{n+1} for any n differ insignificantly. Noticeable changes in the sequence of estimates (41) are observed when the function f behaves irregularly, for example, during integration in the vicinity of its singularities. To avoid large differences between h_n and h_{n+1} , damping of the ratio $r = h_{n+1}/h_n$ is required, that is, the restrictions

$$\sigma^{-1/s} < r < \sigma^{1/s} : \quad r < \sigma^{-1/s} \Rightarrow r = \sigma^{-1/s} \quad \text{or} \quad r > \sigma^{1/s} \Rightarrow r = \sigma^{1/s} \quad (43)$$

should be imposed on it. Here σ defines the range of admissible variations of $\|e\|_{\text{cal}}$. In order that $\|e\|_{\text{cal}}$ changed within one order of magnitude, we should set $\sigma = \sqrt{10} \approx 3.16$. If the left side of inequality (43) is not fulfilled, the step is repeated.

The initial step size h_1 is determined from estimate (41) for $s = 2$ [5]:

$$h_1 = \sqrt{2\eta \frac{\|e\|_{\text{tol}}}{\|f_2 - f_1\|}}, \quad f_1 = f(t_0, x_0, x'_0, z_0), \quad g_1 = g(t_0, x_0, x'_0, z_0), \quad f_2 = f(t_1, x_1, x'_1, z_1), \quad (44)$$

$$x_1 = x_0 + x'_0\eta + \frac{1}{2}f_1\eta^2, \quad x'_1 = x'_0 + f_1\eta, \quad z_1 = z_0 + g_1\eta, \quad t_1 = t_0 + \eta.$$

Here η is a small value. If η is so small that in computer arithmetic $f_2 = f_1$, it is increased by an order of magnitude, and estimate (44) is repeated. Estimate (44) of the starting step size is adequate only for the Euler method of the first order. Therefore, for the method of any other order, the starting step size will be much less than that expected. Nevertheless, the step-by-step integration begins with estimate (44), but then with the use of algorithm (41) with allowance for restrictions (43), the step size will gradually reach the proper value of the regular operating mode of the integrator with automatic choice of the step size.

The last step is revealed from the condition

$$(t_0 + \Delta t - t_n)/h_{n+1} < 1, \quad (45)$$

where Δt is the length of the entire interval of integration and t_n is the time of the n th step. Then to obtain the time $t_0 + \Delta t$, the last step size is chosen compulsorily as $h_{n+1} = t_0 + \Delta t - t_n$, and its ratio to the current step size h_n is redefined: $r = h_{n+1}/h_n$.

3. PROCEDURE LOBBIE

The integrator was implemented in Fortran up to order 32 in computer arithmetic with double and quadruple precision. The program procedure of the integrator Lobbie is called by the command

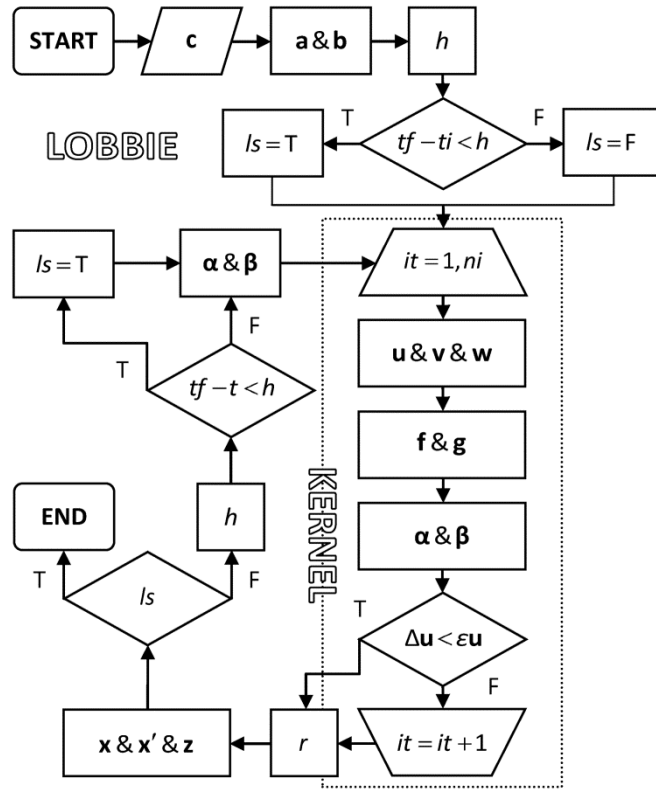


Fig. 2. Flow chart of the procedure Lobbie.

call lobbie (x, y, z, ts, tf, step, etol, nxy, nz, ns, ni, nst, ncf, fun).

Here \mathbf{x} , \mathbf{y} , \mathbf{z} are the arrays of integral variables $\mathbf{x}, \mathbf{x}', \mathbf{z}$, respectively: at the input, their values at an initial moment of time t_0 (\mathbf{ts}); at the output, their values at a final moment of time $t_0 + \Delta t$ (\mathbf{tf}); **step** is an initial step size h_1 : for automatic choice (41) its output value is h_n (the size at the penultimate step), while if **step** is zero, h_1 is assigned by the integrator according to estimate (44); **etol** is $\|\mathbf{e}\|_{\text{tol}}$: at zero value, the mode of a constant integration step; **nxy** and **nz** are dimensions of the arrays \mathbf{x} , \mathbf{y} ($\dim \mathbf{x} = \dim \mathbf{x}'$), and \mathbf{z} ($\dim \mathbf{z}$), respectively; **ns** is the number of nodal values s ; **ni** is the maximal number of iterations at the step for solving nonlinear equations (31) for $\mathbf{u}_1, \dots, \mathbf{u}_s$, $\mathbf{v}_1, \dots, \mathbf{v}_s$, and $\mathbf{w}_1, \dots, \mathbf{w}_s$; **nst** and **ncf** are the numbers of the steps and calls of the procedure **fun** for evaluation of functions \mathbf{f} and \mathbf{g} over the entire interval of integration. The procedure **fun** is defined as

subroutine fun(t, x, y, z, f),

where \mathbf{t} is a current time t ; \mathbf{x} , \mathbf{y} , \mathbf{z} are the arrays of integrated variables with values at time t ; \mathbf{f} is the output array of values of functions \mathbf{f} and \mathbf{g} with dimensions $\dim \mathbf{f} + \dim \mathbf{g}$.

The computational process in the procedure Lobbie is carried out according to the flow chart presented in Fig. 2. We now briefly describe the main stages of step-by-step integration in the case of a variable step with an automatic choice of its starting value, while time increases ($\Delta t > 0$).

I. From the data block attached to the procedure, the array of the nodal values c_1, \dots, c_s is read, and the integrator constants (32) are recurrently calculated by using (34). The starting step size h_1 is estimated by (44). If the condition $h_1 > \Delta t$ is fulfilled, the starting step is considered to be the last, and then the size is set $h_1 = \Delta t$.

II. The intermediate solutions u_1, \dots, u_s , v_1, \dots, v_s , and w_1, \dots, w_s (31) are iteratively determined together with the functions f_1, \dots, f_s and g_1, \dots, g_s (28) as well as the divided differences $\alpha_1, \dots, \alpha_s$ and β_1, \dots, β_s . Iterations ends when condition (35) is satisfied or when the number of iterations reaches its maximum value ni .

III. When the iterative process ends, the solutions x_1, x'_1, z_1 (31) are formed taking into account the refined divided differences α_s and β_s at the last iteration. The scaling multiplier r (41) is estimated with allowance for damping conditions (43).

IV. If the step is the last one, the procedure ends. Otherwise, the size of the successive step (41) is defined. The fulfillment of condition (45) establishes the last integration step, and then its size is redefined so that to reach the final time $t_0 + \Delta t$.

V. The values of the divided differences (30) are extrapolated using the values of the functions of differential equations (42) and the computing process is repeated from item II, where the obtained solutions x_1, x'_1, z_1 are accepted to be the initial ones x_0, x'_0, z_0 .

CONCLUSIONS

Thus, the theoretical background of the new collocation integrator intended for solving the mixed systems of first- and second-order differential equations of dynamics have been outlined in the paper. The practical implementation of the integrator was considered, and the software procedure Loblie was described. In future we plan to present results of testing of the new integrator on the example of simple dynamical systems and show its efficiency compared to other integrators widely used in practice.

The research was carried out within the State Assignment of the Ministry of Science and Higher Education of the Russian Federation No. 0721-2020-0049.

REFERENCES

1. A. Guillou and J. L. Soule, Rev. Francaise Informat. Recherche Opérationnelle 3, Ser. R-3, 17–44 (1969).
2. K. Wright, BIT, **10**, 217–227 (1970).
3. E. Hairer, S. P. Nørsett, and G. Wanner, Solving Ordinary Differential Equations I. Nonstiff Problems, Springer (2008).
4. E. Hairer, C. Lubich, and G. Wanner, Geometric Numerical Integration. Structure-Preserving Algorithms for Ordinary Differential Equations, Springer (2006).
5. V. A. Avdyushev, Numerical Simulation of Orbits of Celestial Bodies [in Russian], Publishing House of Tomsk State University, Tomsk (2015).
6. E. Everhart, Celest. Mech., **10**, 35–55 (1974).
7. E. Everhart, in: Proc. IAU Colloq. 83 (Rome 1984), ed. by A. Carusi and G. B. Valsecchi, D. Reidel Publishing Co., Dordrecht (1985), pp. 185–202.
8. J. Kuntzmann, Z. Angew. Math. Mech., **41**, 28–31 (1961).
9. J. C. Butcher, Math. Comput., **18**, 50–64 (1964).
10. P. M. Cincotta, C. M. Giordano, and C. Simó, Physica D, **182**, 151–178 (2003).
11. C. A. Burdet, Z. Angew. Math. Phys., **19**, 345–368 (1968).
12. P. Kustaanheimo and E. Stiefel, J. Reine Angew. Math., **218**, 204–219 (1965).
13. V. A. Shefer, Astron. Zh., **68**, 197–205 (1991).
14. J. Baumgarte, Comp. Math. Appl. Mech. Eng., **1**, 1–16 (1972).