



ELSEVIER

Available online at www.sciencedirect.com

SCIENCE @ DIRECT®

Journal of Computational and Applied Mathematics 185 (2006) 203–211

JOURNAL OF
COMPUTATIONAL AND
APPLIED MATHEMATICS

www.elsevier.com/locate/cam

Verifying approximate solutions to differential equations[☆]

W.H. Enright

Department of Computer Science, University of Toronto, Toronto, Ont., Canada M5S 1A4

Received 7 March 2003

Abstract

It is now standard practice in computational science for large-scale simulations to be implemented and investigated in a problem solving environment (PSE) such as MATLAB or MAPLE. In such an environment, a scientist or engineer will formulate a mathematical model, approximate its solution using an appropriate numerical method, visualize the approximate solution and verify (or validate) the quality of the approximate solution. Traditionally, we have been most concerned with the development of effective numerical software for generating the approximate solution and several efficient and reliable numerical libraries are now available for use within the most widely used PSEs. On the other hand, the visualization and verification tasks have received little attention, even though each often requires as much computational effort as is involved in generating the approximate solution.

In this paper, we will investigate the effectiveness of a suite of tools that we have recently introduced in the MATLAB PSE to verify approximate solutions of ordinary differential equations. We will use the notion of ‘effectivity index’, widely used by researchers in the adaptive mesh PDE community, to quantify the credibility of our verification tools. Numerical examples will be presented to illustrate the effectiveness of these tools when applied to a standard numerical method on two model test problems.

© 2005 Elsevier B.V. All rights reserved.

MSC: 65L05; 65Y20

Keywords: Verification; ODEs; Defect; Numerical solutions; Reliability

[☆]This research was supported by the Natural Science and Engineering Research Council of Canada.

E-mail address: enright@cs.utoronto.ca.

1. Introduction

In the last decade there has been a significant change in the way scientific computing is carried out. Those of us developing numerical methods must be aware of this change if our software is ever going to be widely used and have an impact. The expectations and measures we use to evaluate software must evolve to make the results more relevant and applicable to modern scientific computing.

Two environments where numerical software is now widely used are:

- **Large-Scale Scientific Computation:**
Characterized by simulations generating massive amounts of data and involving expensive runs on High Performance distributed systems (often located at remote sites). Applications require data compression, interactive data viewing (or data mining), remote visualization and the distributed collaboration of experts.
- **Problem Solving Environments (PSEs):**
Often where students are first exposed to scientific computing and where scientists and engineers formulate their mathematical models and investigate approximate solutions of these problems.

In this investigation we focus on the latter environment but the tools and approaches we consider and the observations we make are also applicable to tools that can be developed for the former environment. We will use the numerical solution of ordinary differential equations (ODEs), in particular initial value problems (IVPs), to illustrate the tools and evaluation criteria that have been developed and implemented. The approach can be adopted in a similar fashion for other problem classes.

In the next section we will discuss the concept of a PSE and identify the characteristics or features that make a PSE most effective. We will consider the importance of verification tools (as an essential component of a PSE) and describe four such tools that we have implemented for use with ODE solvers. In Section 3, we describe how one can quantify the effectiveness of these tools by introducing the notion of an ‘effectivity index’ (a concept which has been developed and applied in the PDE community). Finally, we illustrate the credibility of these tools on two model problems and the method `ode45` of the MATLAB PSE.

2. The components of a PSE

In a PSE the focus is on visualizing approximate solutions of mathematical models. This generally involves graphical representation of a ‘subset’ of the results. (e.g. Phase plane or standard plots.) It can also involve the use of colour, texture, lighting, animation and sound. Stringent accuracy is not usually necessary, but off-mesh accuracy is often required.

Ease of use and robustness is very important (as important as efficiency) in a PSE. Information should be hidden unless necessary to define the problem. Options and additional parameters (if required) should be specified in the same way for all methods that perform similar tasks.

Similarly, there should be a standard representation (adopted by all methods of the PSE performing similar tasks) of what is meant by an approximate solution to a particular problem. For example, an approximate solution of an ODE can be represented in a PSE by a vector of piecewise polynomials. Such a representation makes it easy to compare, visualize, verify and manipulate the results of simulations without having to pay attention to or understand how the approximations were generated. It must be

acknowledged that, without an assumption of this type, it would be difficult to define and/or interpret verification tools suitable for use with ODE software.

A PSE should provide tools to:

- (1) Formulate and alter a mathematical model of the problem.
- (2) Approximate the solution of the mathematical model.
- (3) Visualize the *approximate* solution.
- (4) Verify that an approximate solution is consistent with the mathematical model and the requested accuracy.
- (5) Verify that the mathematical model is well-posed. (That is, quantify the underlying inherent conditioning.)

Generic tools for the first three of these tasks have been widely investigated and implemented as the key components of PSEs.

We have recently [3] developed a class of tools that can be used for the latter two tasks when the underlying mathematical model is a system of ODEs and the numerical solution is a continuous approximation to the true solution. In this investigation we will illustrate the effectiveness of these tools.

The key assumption we make of any method which generates an approximate solution is that, when applied to the IVP,

$$y' = f(x, y), \quad y(a) = y_0, \quad \text{on } [a, b],$$

with a specified accuracy, TOL , the numerical method, M , generates a piecewise polynomial, $z_{TOL}(x)$, defined for $x \in [a, b]$ and satisfying,

$$\|z_{TOL}(x) - y(x)\| \leq K_M TOL.$$

This assumption is satisfied by most state-of-the-art ODE software and certainly by the built-in numerical methods of MATLAB [4]. Note that K_M will always depend on the underlying mathematical conditioning of the problem, but it can also be very sensitive to the particular method. K_M can be interpreted as the ‘numerical condition number’ associated with method M .

The restrictions we imposed on the verification tools we developed were:

- (1) The tool must be easy to use—If possible the calling sequence should be the same as that used to generate the approximate solution with no additional information required.
- (2) The ‘cost’ associated with applying the tool should not be much greater than that involved in generating the approximate solution (or visualizing it). (Note that this requirement rules out techniques that provide strict bounds on the global error as such techniques are much more expensive to apply.)
- (3) The calling sequence should be the same for any method it is applied with. The tool may require some knowledge of the method but this information should be transparent to the user.

The four verification tools we have implemented for ODEs in MATLAB (see [3] for details) are:

Check 1: A consistency Check based on solving the problem with a more stringent value for TOL and returning the difference in the two approximations as a piecewise polynomial. That is we compute

$$E_{TOL}(x) = z_{TOL}(x) - z_{TOL1}(x),$$

where $TOL1 < TOL$. The ‘cost’ of applying this tool is about double the cost of generating $z_{TOL}(x)$. (Note that, in our implementation, we have used the default value $TOL1 = TOL/10$.)

Check 2: A consistency Check based on solving the problem with an alternative method, \tilde{M} . The piecewise polynomial returned by this routine is

$$\tilde{E}_{TOL}(x) = z_{TOL}(x) - \tilde{z}_{TOL}(x),$$

where $\tilde{z}_{TOL}(x)$ is the approximate solution associated with \tilde{M} . (Note that this Check will be most suitable if the method \tilde{M} is very different from M and the value of $K_{\tilde{M}}$ is not very sensitive to the method. For example, in our implementation of this tool, we use a variable order Adams method when assessing a Runge–Kutta method from MATLAB.)

Check 3: A Check which computes the defect,

$$D_{TOL}(x) = z'_{TOL}(x) - f(x, z_{TOL}(x)).$$

Note that $D_{TOL}(x)$ is not a piecewise polynomial, but it is a function that can be evaluated at any point in $[a, b]$. It is an inexpensive Check most suitable for methods designed to control the magnitude of the defect.

Check 4: A consistency Check based on an idea of Zadunaisky [5] where we determine a piecewise polynomial

$$EST_{TOL}(x) = w_{TOL}(x) - z_{TOL}(x),$$

where $w_{TOL}(x)$ is the approximate solution, produced by method M applied to the perturbed IVP, $z' = f(x, z) + D_{TOL}(x)$ with known exact solution $z_{TOL}(x)$ on the same mesh that was used to determine $z_{TOL}(x)$. That is, we have $z_{TOL}(x)' = f(x, z_{TOL}(x)) + D_{TOL}(x)$. (Note that this Check is most appropriate when the defect is small in magnitude as it will likely be for methods that directly estimate and control the size of the defect.)

For each of the vector-valued functions returned by Check 1, Check 2 and Check 4 we can interpret its max value for $x \in [a, b]$ to be an estimate of the maximum global error. Therefore, a crude indicator of the numerical condition number is the ratio of this max value to TOL . Similarly for Check 3, the maximum magnitude of the components of $D_{TOL}(x)/TOL$ can be considered an indication of the contribution of the numerical method to the overall numerical conditioning of the problem (it does not reflect the underlying mathematical conditioning of the problem).

Each of these verification tools are scalable with respect to the dimension of the problem and the length of the integration (in the sense that the cost associated with applying the tool remains a small multiple of the cost associated with generating the approximate solution). Each is also straightforward to apply in a parallel computing environment where the underlying numerical method and the verification tool can be run simultaneously on independent processors with very little communication required between the processors.

In the numerical solution of partial differential equations (PDEs) there has been a similar (although more advanced) effort in developing efficient numerical methods and corresponding verification tools for special classes of problems. In the last decade there has been considerable interest in developing verification tools based on a posteriori error estimates for assessing the validity of a particular numerical solution. Such estimates are intended to be computed as a post processing task, and have been implemented and packaged with some of the widely used PDE PSEs. In particular, the finite element

and adaptive mesh refinement research communities have adopted this point of view (see for example [1,2]). They have introduced the concept of an ‘effectivity index’ to quantify the credibility of their tools. We will adopt this approach in quantifying the credibility of the ODE verification tools we have introduced.

3. Effectivity indices for ODE verification tools

For the ODE verification tool, Check 3, the only error that arises is due to round-off and therefore this verification Check will be reliable at all reasonable error tolerances. Because it does involve the subtraction of near equals (the pointwise evaluation of the defect) it will become unreliable at stringent tolerances (values of TOL near the unit round-off).

Each of the other verification tools (Check 1, Check 2 and Check 4) can be interpreted as an estimate of the global error, $EST \approx TRUERR$. In this case the natural definition for the ‘Effectivity Index’ is,

$$EI = \|EST\| / \|TRUERR\|.$$

Clearly, the closer this value is to 1, the more accurate and useful it will be. Too large a value indicates an over-estimate while too small a value indicates an under-estimate and the possibility that an inaccurate approximate solution will not be recognized. We will compute the value of EI for Check 1, Check 2 and Check 4 for two test problems, a range of values of TOL , and monitor its behaviour over the entire range of integration ($x \in [a, b]$).

In our numerical results we report,

$$CEST \equiv \frac{\max_{x \in [x_0, x_F]} \|EST(x)\|}{TOL},$$

$$MnEI \equiv \min_{x \in [x_0, x_F]} EI(x),$$

$$MxEI \equiv \max_{x \in [x_0, x_F]} EI(x).$$

Our main objective, in developing these verification tools, is to produce an inexpensive but reasonably reliable tool for assessing the credibility of a particular approximate solution. For Check 1 and Check 2 the ‘reliability’ of these tools can be improved by setting the accuracy parameter associated with the secondary solution to be more severe. We have been somewhat arbitrary in setting the default for these values ($TOL/10$ for Check 1 and TOL for Check 2) but the overall credibility of these tools is quite good and not very sensitive to this choice. For Check 4, the reliability of the tool is sensitive to the definition of the piecewise polynomial associated with the approximate solution and we have assumed that this is part of the definition of the method. For ode45 there are several possible choices for how this piecewise polynomial can be defined and the MATLAB implementation, while appropriate for most purposes, does not necessarily result in a small magnitude defect (especially at relaxed tolerances). Alternative definitions for this piecewise polynomial (for the same underlying discrete Runge–Kutta formula) are available which, at a small increase in cost per step, will result in a smaller magnitude defect and hence a more robust and credible verification tool.

4. Numerical results and discussion

The first test problem corresponds to a recently discovered stable orbit that arises in the simulation of the restricted three-body problem (where the orbits are planar). The bodies have equal mass (in our case we assume $m_1 = m_2 = m_3 = 1.0$) and, with the appropriate starting conditions, will follow the same figure-eight orbit as a periodic steady-state solution.

The two spatial coordinates of the j th body are y_{1j}, y_{2j} for $j = 1, 2, 3$. Each of the six coordinates satisfy a second-order differential equation,

$$y''_{ij} = \sum_{k=1, k \neq j}^3 m_k \left(\frac{y_{ik} - y_{ij}}{d_{jk}^3} \right),$$

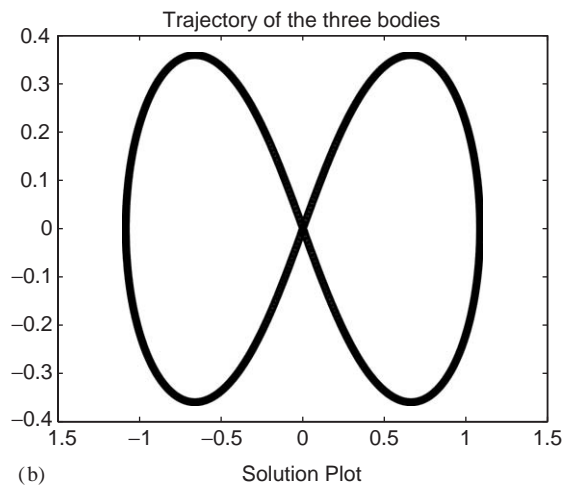
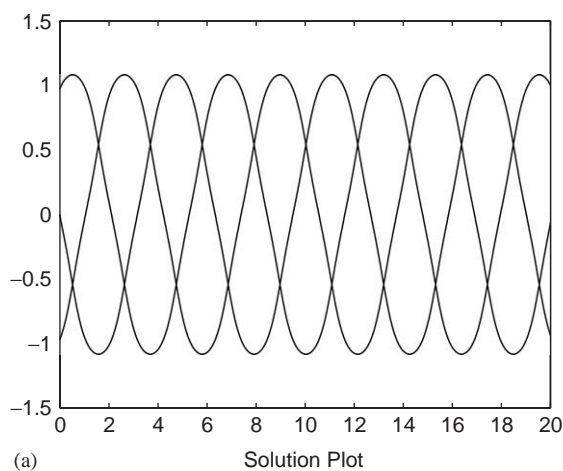


Fig. 1. The approximate solution produced by ode45 for the Three Body problem: (a) solution plot; (b) phase portrait.

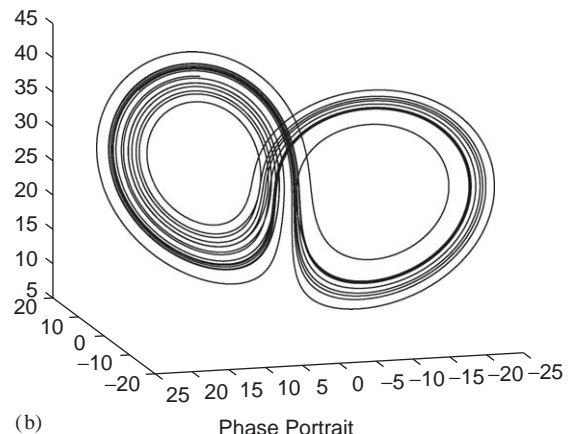
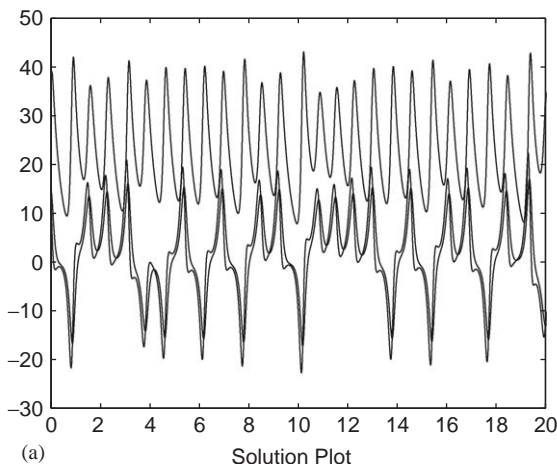


Fig. 2. The approximate solution produced by ode45 for the Lorentz problem: (a) solution plot; (b) phase portrait.

Table 1
Effectiveness indices for the Three Body problem

TOL	CEST	MnEI	MxEI	CEST	MnEI	MxEI	CEST	MnEI	MxEI
	Check 1			Check 2			Check 4		
10^{-3}	801	0.67	1.01	778	0.56	15.5	1078	0.42	3.59
10^{-4}	1370	0.98	1.01	1450	0.81	5.43	1411	0.39	1.53
10^{-6}	120	0.99	1.02	188	0.90	11.9	76	0.62	1.22
10^{-8}	190	0.98	1.03	247	0.63	5.80	223	0.21	2.75
10^{-9}	220	0.99	1.02	220	0.99	1.02	258	0.28	2.09
10^{-10}	230	0.13	1.00	218	0.19	1.92	273	0.05	1.32

where

$$d_{kj}^2 = \sum_{i=1}^2 (y_{ik} - y_{ij})^2, \quad k, j = 1, 2, 3.$$

When this system is re-written, as a first-order system, the dimension of the problem is 12 and the initial conditions, at $x = 0$, are given by

$$\begin{aligned} y_{11} &= -0.97000436, & y'_{11} &= 0.466203685, \\ y_{21} &= 0.24308753, & y'_{21} &= 0.43236573, \\ y_{12} &= 0.0, & y'_{12} &= -0.93240737, \\ y_{22} &= 0.0, & y'_{22} &= -0.86473146, \\ y_{13} &= 0.97000436, & y'_{13} &= 0.466203685, \\ y_{23} &= -0.24308753, & y'_{23} &= 0.43236573. \end{aligned}$$

We solve the problem for $x \in [0, 20]$.

The second problem is the well-known Lorentz problem which arises in the study of dynamical systems and is known to have solutions which are potentially poorly conditioned,

$$\begin{aligned} y'_1 &= 10(y_2 - y_1), \\ y'_2 &= y_1(28 - y_3) - y_2, \\ y'_3 &= y_1 y_2 - \frac{8}{3} y_3, \end{aligned}$$

with

$$y_1(0) = 15, \quad y_2(0) = 15, \quad y_3(0) = 36, \quad \text{and } x \in [0, 20].$$

Using ode45 to investigate solutions of these problems allows us to illustrate the ability of our tools to identify potential conditioning difficulties. Figs. 1 and 2 present standard visualizations of numerical solutions to these problems determined by ode45.

Tables 1 and 2 and Fig. 3 report statistics on the behaviour of the effectivity indices quantifying the credibility of these verification tools on these two test problems. As these tables indicate the tools are

Table 2
Effectiveness indices for the Lorentz problem

TOL	CEST	MnEI	MxEI	CEST	MnEI	MxEI	CEST	MnEI	MxEI
	Check 1			Check 2			Check 4		
10^{-3}	3.8×10^4	0.20	2.22	4.0×10^4	0.06	6.21	3.3×10^4	0.18	6.21
10^{-4}	4.0×10^5	0.26	18.1	3.7×10^5	0.10	6.83	3.6×10^5	0.09	8.44
10^{-6}	3.6×10^7	0.92	1.12	3.1×10^7	0.30	14.3	3.6×10^7	0.28	3.77
10^{-8}	5.1×10^8	0.99	1.00	2.8×10^9	0.70	9.85	9.1×10^8	0.39	2.05
10^{-9}	4.4×10^8	0.98	1.00	5.7×10^8	0.41	4.33	4.4×10^8	0.69	1.24
10^{-10}	4.3×10^8	0.24	1.08	9.4×10^8	0.60	10.4	4.5×10^8	0.18	1.19

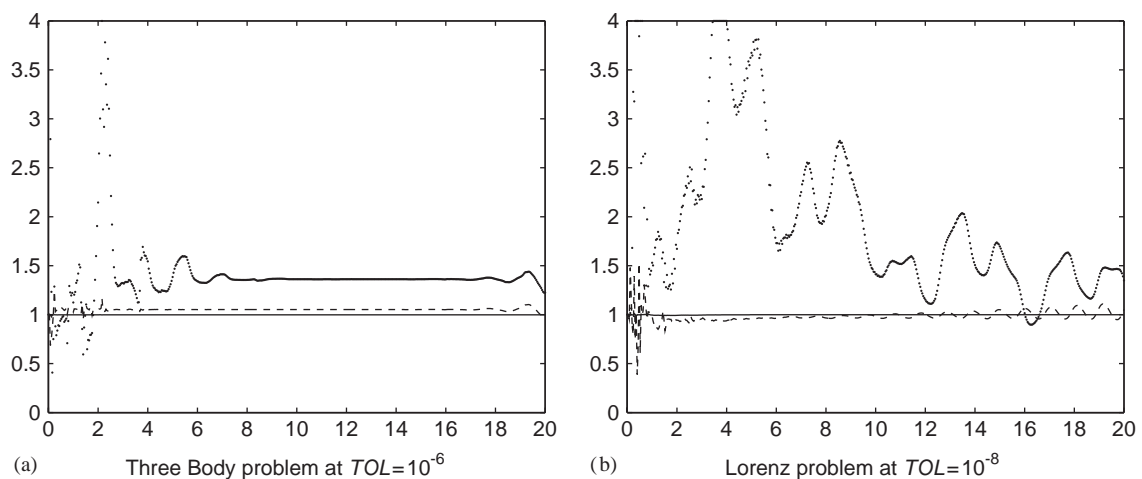


Fig. 3. Graphs of the effectiveness indices over the range of integration for the verification tools. The solid line corresponds to Check 1, dashed line to Check 2 and dotted line to Check 4: (a) Three Body problem at $TOL = 10^{-6}$; (b) Lorentz problem at $TOL = 10^{-8}$.

able to provide a consistent and reliable estimate of the conditioning of the mathematical model over a range of accuracy requests for a method that does not attempt to control the size of the defect. Note that even on poorly conditioned problems, such as the Lorentz problem (where the condition number is exponential in the length of the integration interval), all three tools reflect this. While the effectivity index values associated with Check 1 are particularly good (all values very close to 1), the values are also quite reasonable for the other tools as well (the respective estimates of the norm of the maximum error is rarely out by a factor of 5).

From the reported results it is also clear that, when the errors are potentially underestimated (small values of *MnEI*), the reported values are usually large in magnitude (the corresponding relative error is estimated to be much greater than 100%) and this is a reliable signal that the approximate solution cannot be trusted. The results summarized in Fig. 3 illustrate that the tools provide reliable and credible information about the approximate solution over the entire interval of integration.

These results indicate that verification tools, such as those we have developed, can provide inexpensive confirmation that a particular approximate solution can be trusted. It is also clear that, when the

reported error is small enough for the approximation to be accurate to a few significant digits, then the associated estimate of the condition number is generally a good indication of the underlying numerical conditioning.

References

- [1] M. Ainsworth, J. Tinsley Oden, *A Posteriori Error Estimation in Finite Element Analysis*, Wiley-Interscience, New York, 2000.
- [2] I. Babuska, T. Strouboulis, *The Finite Element Method and its Reliability*, Oxford University Press, Oxford, 2001.
- [3] W.H. Enright, *Tools for the Verification of Approximate Solutions to Differential Equations*, Department of Computer Science Technical Report, University of Toronto, January 2002.
- [4] L.F. Shampine, M.W. Reichelt, The MATLAB ODE suite, *SIAM J. Sci. Comput.* 18 (1997) 1–22.
- [5] P.E. Zadunaisky, On the estimation of errors propagated in the numerical integration of ordinary differential equations, *Numer. Math.* 27 (1976) 21–39.