

## OPTIMAL CONVERGENCE OF THE SIMPLIFIED NEWTON ITERATIONS FOR NON STIFF PROBLEMS

Ronald Tshelametse\*

### ABSTRACT:

The paper deals with the numerical solution of IVP's for systems of stiff ODE's with particular emphasis on implicit linear multistep methods (LMM), particularly the backward differentiation formulae (BDF). In this paper we intend to find the value of the optimal convergence rate factor that will minimize the computational costs in the context of terminating simplified Newton iterations for non-stiff problems. We refer to this convergence rate factor,  $\eta_{ref}$  as the optimal convergence rate. We conclude that for non-stiff problems the optimal convergence rate is such that  $0.3 \leq \eta_{ref} \leq 0.5$ . This compares favourably with the optimal convergence rate for non-stiff problems using small stepsizes as obtained in Gustaffsson and Soderlind [1, Fig2, p27]. We further conduct experiments using several values of  $\eta_{ref}$  in solving the van de Pol problem which uses small stepsizes when solved using ode15s with ,  $\sigma = 10$  for  $t \in [0,30]$ . We monitor the number of function calls for various convergence rate. The results obtained confirm our lower bound on  $\eta_{ref}$ .

\* Department of Mathematics, University of Botswana, Private Bag 0022, Gaborone, Botswana.

## 1 Introduction

The paper is concerned with the numerical integration of

$$\frac{du}{dt} = f(t, u), \quad 0 < t \leq T, \quad u(0) = u_0, \quad f: R \times R^m \rightarrow R^m \quad (1)$$

In the literature some initial value problems (1) are referred to as stiff. A prominent feature for these problems is that they are extremely difficult to solve by standard explicit methods. The time integration of stiff systems is usually achieved using implicit methods, and for many codes by linear multistep methods. A linear multistep method also called linear k-step methods [2], [3], [4], [5], [6] in standard constant stepsize form a linear multistep or k-step method is defined thus:

$$\sum_{i=0}^k \alpha_i u_{n-i} = h \sum_{i=0}^k \beta_i f_{n-i}, \quad (2)$$

where  $\alpha_i$  and  $\beta_i$  are constants and  $\alpha_0 = 1$ .  $f_{n-i}$  denotes  $f(t_{n-i}, u_{n-i})$ ,  $t_{n-i} = t_n - ih$ ,  $i = 0, 1, \dots, k$  and  $h$  is the stepsize. The condition that  $\alpha_0 = 1$  removes the arbitrariness that arises from the fact that both sides of the IVP could be multiplied by the same constant without altering the method. The linear multistep method (2) is said to be explicit if  $\beta_0 = 0$  and implicit if  $\beta_0 \neq 0$ . This leads to

$$F(u_n) \equiv u_n + \varphi_n - h_n \beta_0 f(u_n) = 0, \quad (4)$$

where  $\varphi_n = \sum_{i=1}^k \alpha_i u_{n-i}$  is a known value. [7], [8], [9], [10], [11], [12]. At each integration step  $t_n$  we must solve the nonlinear equation (4). To solve for  $u_n$  most codes use the Newton iterative method and its variants in the following form

$$W_n^{(l)} \varepsilon_n^{(l)} = -F(u_n^{(l)}), \quad u_n^{(l+1)} = u_n^{(l)} + \varepsilon_n^{(l)} \quad l = 0, 1, 2, \dots \quad (5)$$

with the starting value  $u_n^{(0)}$  known and “fairly” accurate. For the full Newton method

$$W_n^{(l)} = F'(u_n^{(l)}) = I - h_n \beta_0 f'(u_n^{(l)}) \quad (6)$$

The use of the Newton method is due to the stiffness phenomenon. For large problems evaluating the Jacobian,  $J = f'(u_n^{(l)})$  (and hence the Newton iteration matrix  $W_n^{(l)}$ ) and solving the linear algebraic system are by far the most computationally expensive operations in the integration.

There are various strategies used in practice to try and minimise the cost of computing the Jacobian and the Newton matrix. These measures are mainly centred on administering the iteration matrix in (6) leading to simplified Newton methods amongst others [13], [14], [15]. In this paper we focus on optimally terminating the simplified Newton iterations.

## 2 Theory

We intend to find the value of  $\eta$  (set to 0.9 in ode15s) that will minimize the computational costs in the context of simplified Newton. We refer to this convergence rate factor,  $\eta_{ref}$  as the optimal convergence rate. For functional iteration this task is considered in Gustafsson and Soderlind [1], ideal for nonstiff problems, where  $\eta_{ref} \approx 0.4$ . To simulate the nonstiff implementation in the context of simplified Newton we consider the case when the stepsize  $h_n$  is small. We emphasise here that the convergence rate of an iterative solver usually depends on the stepsize. Thus the step size selection strategy is important.

It can be shown that with  $e_l = u_n^{(l)} - u_n$ ,

$$e_{l+1} = W^{-1}(h_n \beta_0 \tilde{J} - h_0 \beta_0 J) e_l, \quad (7)$$

where  $W = I - h_0 \beta_0 J$  and  $\tilde{J}$  is the mean value Jacobian at  $t_n$  such that

$$\tilde{J} e_l = f(u_n^{(l)}) - f(u_n). \quad (8)$$

If the iteration is convergent, there follows

$$\|e_{l+1}\| \leq \eta \|e_l\|, \quad (9)$$

where  $\|W^{-1}(h_n \beta_0 \tilde{J} - h_0 \beta_0 J) e_l\| \leq \eta$ . Now since we assume that  $\eta < 1$ ,

$$\|e_l\| \leq \|u_n^{(l)} - u_n^{(l-1)}\| + \|e_{l+1}\| \leq \eta \|e_l\| + \|u_n^{(l)} - u_n^{(l-1)}\|, \quad (10)$$

yielding in (9)

$$\|e_{l+1}\| \leq \frac{\eta}{1 - \eta} \|u_n^{(l)} - u_n^{(l-1)}\|, \quad (11)$$

In practice  $\eta$  is replaced with the estimate of the convergence rate factor to get the practical convergence test of the form

$$\|e_{l+1}\| \sim \leq \frac{\eta}{1-\eta} \|u_n^{(l)} - u_n^{(l-1)}\| \leq \tau, \quad (12)$$

for some appropriate  $\tau$ . We now follow Gustafsson and Soderlind [1]. Let  $l$  be the number of iterations needed to satisfy (12). Then

$$\tau \geq \frac{\eta}{1-\eta} \|u_n^{(l)} - u_n^{(l-1)}\| \geq \sim \|e_l\| \approx \eta^l \|e_0\|, \quad (13)$$

where  $e_0$  is the error in the initial  $u_n^{(0)}$ . Thus

$$l \sim \geq \frac{\log \tau - \log \|e_0\|}{\log \eta}, \quad (\eta < 1) \quad (14)$$

implying that, the number of function evaluations per unit step (i.e. per unit time of integration) is proportional to  $\bar{l}$ , where approximately

$$\bar{l} = \frac{l}{h_n} = \frac{\log \tau - \log \|e_0\|}{h_n \log \eta}, \quad (15)$$

For modelling the relationship between  $\eta$  and  $h_n$  we have from (9) the approximation

$$\eta \approx \|W^{-1}(h_n \beta_0 \tilde{J} - h_0 \beta_0 J)\|. \quad (16)$$

Hence if  $h_n = h_0$

$$\eta \approx h_n \beta_0 \|(I - h_n \beta_0 J)^{-1}(\tilde{J} - J)\|. \quad (17)$$

For small  $h_n$ ,

$$(I - h_n \beta_0 J)^{-1} \approx I. \quad (18)$$

We may therefore use the model

$$\eta \approx Kh_n \quad (19)$$

for some constant  $K$ . In ode15s the iteration matrix is updated and refactored every time the step size,  $h$  is changed, see Shampine and Reichelt [16], so that for every successful time step the condition  $h_n = h_0$  holds.

This model is also used by Gustafsson and Soderlind [1] who analyse the case for functional iteration, namely,  $J = 0$ . Therefore

$$\bar{l} \approx \frac{\log \tau - \log \|e_0\|}{\eta \log \eta}, \quad (20)$$

The starting value  $u_n^0$  is usually constructed using an interpolation formula, which makes the initial error  $e_0$  depend on  $h_n$  (and hence  $\eta$ ). Assuming that the dependence is usually not strong,  $\log \tau - \log \|e_0\|$  is assumed to be a constant, so that we have the proportionality relationship

$$\bar{l} \approx -\frac{1}{\eta \log \eta}, \quad (21)$$

where the minus sign is a result of assuming  $e_0 > \tau$ . The expression (21) has a minimum at  $\eta = e^{-1} \approx 0.36788$  which indicates the appropriate selection of the convergence rate as, say  $\eta_{ref} \approx 0.4$ . The function (21) is fairly flat around the minimum, any value  $0.3 \leq \eta_{ref} \leq 0.5$  would probably be acceptable. See Figure 1 below.

It should be noted that this model, for both the functional iteration [1] and our analysis for small  $h_n$ , relies on the assumption that the number of iterations  $l$ , is not limited and  $l$  is significantly large. From [1, Fig2, p.27] for functional iteration for a non stiff problem van de Pol problem

with  $\sigma = 10$ , if the maximum allowed iterations is 10 then any value  $0.2 \leq \eta_{ref} \leq 0.95$  is acceptable. Their theoretical expectations are approximately met when  $l = 100$ . We conduct experiments to investigate the range of our  $\eta_{ref}$  in solving the same problem.

### 3 Numerical Experiments

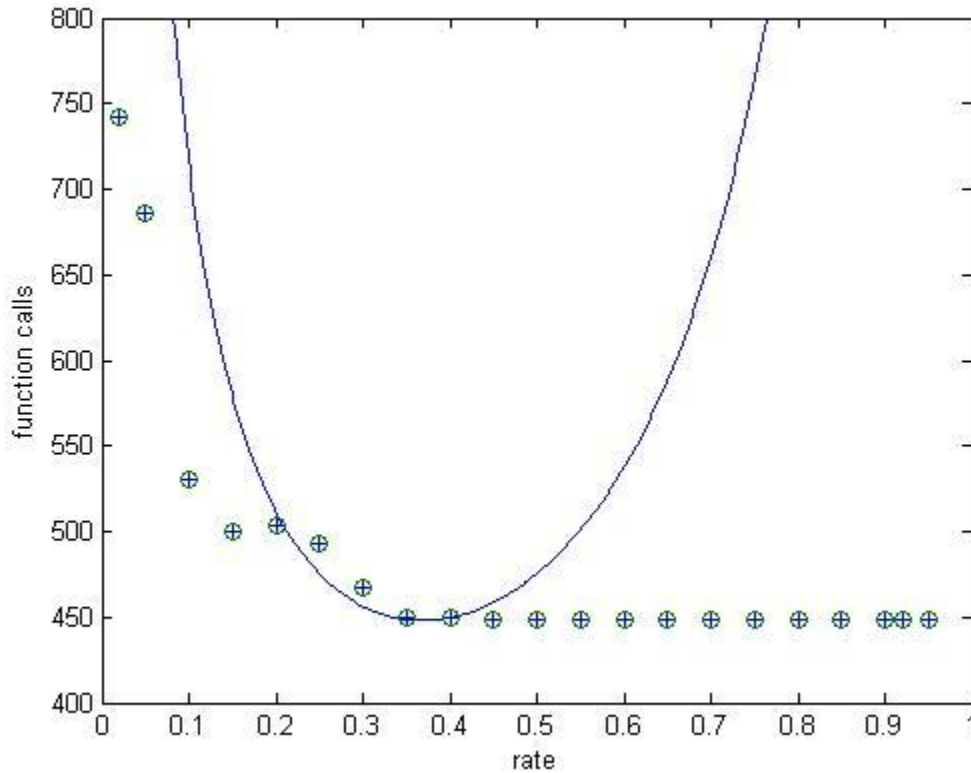
In this experiment we aim to obtain a plot similar to the plots obtained by Gustafsson and Soderlind [1, Fig.2,p.27] for the simplified Newton iteration for small  $h_n$ . That is, we intend to obtain a plot of the total work (number of function calls) as a function of the convergence rate when integrating the nonstiff van de Pol problem with  $\sigma = 10$ , for  $t \in [0, 30]$  with `ode15s`. The code `ode15s` uses small stepsizes (maximum stepsize is about 0.9) throughout the integration when solving this problem. This will then simulate the above small  $h_n$  theory for the simplified Newton method. We have a nonstiff problem solved by the stiff methods of `ode15s` verifying that the theory for small  $h_n$  corresponds to using functional iteration in the case of Gustafsson and Soderlind [1].

In `ode15s` the acceptable convergence rate is  $\eta_{ref} \leq 0.9$  and the maximum allowed number of iterations is 4. We use the weighted infinity norm and experiment with various convergence rates  $\eta_{ref}$  in the range  $[0.01, 0.97]$ . For the default `ode15s` settings discussed above the plot of function calls against convergence rate is plotted as rings (o) in Figure 1. We further increase the maximum number of allowed iterations in `ode15s` to 100 and investigate the behaviour of function evaluations versus convergence rate. The behaviour is plotted as crosses (+) in Figure 1. To determine whether `ode15s` conforms to the theory we also include the plot of the function (21) for the optimal convergence rate as a continuous line in Figure 1. For the plots to fit on the same scale the expression (11) was normalised with a constant  $k$  that makes its minimum equal to the minimum number of function calls experienced in practice, that is,

$$k = \frac{448}{\left[ \frac{1}{-0.36788 \log(0.36788)} \right]} = 164.81$$

We obtain **Figure 1**.

#### 4 Results and Conclusions



**Figure1:** Total work (number of function calls) as a function of the convergence rate set-point when integrating the van de Pol problem with  $\sigma = 10$ . The equation solver includes a limit on the number of allowed iterations. The crosses and the rings correspond to setting this limit as 100 and 4 respectively. The solid line is the graph of the function (21). We use default tolerances.

The expression (21) has a minimum at  $\eta = e^{-1} \approx 0.36788$  which indicates the appropriate selection of the convergence rate as, say  $\eta_{ref} \approx 0.4$ . The function (21) is fairly flat around the minimum, any value  $0.3 \leq \eta_{ref} \leq 0.5$  would probably be acceptable. See Figure 1 below. This is in agreement with  $\eta_{ref} \approx 0.4$  as obtained by Gustafsson and Soderlind [1].

Regarding our experiments, From Figure 1, for the simplified Newton iteration with small  $h_n$  it is clear that for the nonstiff problem, van de Pol problem if the maximum number of allowed iterations is 4 or 100 then any value  $0.2 \leq \eta_{ref} \leq 0.95$  is acceptable. The largest observed



number of iterations performed throughout the entire integration was 3. As seen in (13), the key requirement is that  $\|e_l\| \approx \eta^l \|e_0\|$ , which is more meaningful if  $l$  is large. Moreover the result (21) is only valid when the dependence of the initial error  $e_0$  on the step size is not strong. We note that in ode15s,  $u_n^{(0)}$  is constructed using an interpolation polynomial, so that the initial error  $e_0$  depends on  $h_n$ . This in part explains the difference between our plots and the plots in Gustafsson and Soderlind [1] for the van de Pol problem using a functional iteration based code. The convergence rate  $\eta_{ref} \leq 0.9$  will (for the non-stiff van de Pol problem) yield the minimum number of function evaluations. The same conclusions can be drawn even when the maximum number of allowed iterations is increased to 100, see Figure 1.

Meanwhile, we conclude that the optimal convergence rate for non-stiff problems in the solution of odes using BDFs in the simplified Newton context is such that:

$$0.3 \leq \eta_{ref} \leq 0.5$$

To be specific  $\eta_{ref} = e^{-1} = 0.36788$ .



## References

- [1] Kjell Gustafsson, and Gustaf Soderlind. Control strategies for the iterative solution of nonlinear equations in ODE solver. *SIAM J. Sci. Stat. Comput.*, 18, No. 1:23–40, January 1997.
- [2] John C. Butcher. Numerical methods for ordinary differential equations. *John Wiley*, 2003
- [3] W. H. Enright, T. E. Hull, and B. Linberg. Comparing numerical methods for stiff systems of ODEs. *BIT*, 15:10–48, 1975.
- [4] E. Hairer. Backward error analysis for linear multistep methods. *Numer. Math.*, 84:2:199–232, 1999.
- [5] E. Hairer. Conjugate-symplecticity of linear multistep methods. *J. Computational Mathematics*, 26:5:657–659, 2008.
- [6] E. Hairer and C. Lubich. Symmetric multistep methods over long times. *Numer. Math.*, 97:4:699–723, 2004.
- [7] J. D. Lambert. Numerical Methods for Ordinary Differential Systems. *John Wiley and Sons*, 1991.
- [8] Peter N. Brown, George D. Byrne, and Alan C. Hindmarsh. VODE: a variable-coefficient ODE solver. *SIAM J. Sci. Stat. Comput.*, 10, No. 5:1038–1051, September 1989.
- [9] C. T. Kelley. Iterative methods for linear and nonlinear equations. *Society for Industrial and Applied Mathematics, Philadelphia, PA, USA*, 1995.
- [10] Kenneth R. Jackson. The numerical solution of stiff IVPs for ODEs. *J. Applied Numerical Mathematics*, 1995.
- [11] L. F. Shampine and P. Bogacki. The effect of changing the step size in the linear multistep codes. *SIAM J. Sci. Stat. Comput.*, 10:1010–1023, September 1989.
- [12] R. Tshelametse. Terminating simplified newton iterations: A modified strategy. *International Journal of Management, IT and Engineering*, 4, Issue 4: pp. 246–266, 2014.
- [13] G. D. Byrne and A. C. Hindmarsh. A polyalgorithm for the numerical solution of ordinary differential equation. *Comm. ACM*, 1, No. 1:71–96, March 1975.
- [14] G. Gheri and P. Marzulli. Parallel shooting with error estimate for increasing the accuracy. *J. Comput. and Appl. Math.*, 115, Issues 1-2:213–227, March 2000.
- [15] Lawrence F. Shampine. Numerical solution of ordinary differential equations. *Chapman and Hall*, 1994.
- [16] Lawrence F. Shampine and Mark W. Reichelt. The MATLAB ODE suite. *SIAM J. Sci. Stat. Comput.*, 18, No 1:103–118, September 1989.