

ON A CLASS OF NON-LINEAR METHODS FOR ORDINARY DIFFERENTIAL EQUATIONS*

SUN GENG (孙 耿)

(Institute of Mathematics, Academia Sinica, Beijing, China)

MAO ZU-FAN (毛祖范)

(Computing Centre, Academia Sinica, Beijing, China)

Abstract

In this paper, a class of non-linear methods proposed in [1] is discussed. A new derivation of the methods is given. The analysis based on the new derivation shows that this class of methods is not suitable for stiff problems. The numerical tests support our argument.

§ 1. Introduction

In [1], Qin Zeng-fu put forward a class of non-linear methods for numerical integration of ordinary differential equations. The derivation of the formulas is based on the Frenet frame and the regular representation of curves.

Let the initial value problem be in the form

$$\begin{cases} \frac{dy_i}{dx} = f_i(x, y_1, \dots, y_m), \\ y_i(x_0) = y_{i0}, \quad i = 1, 2, \dots, m. \end{cases} \quad (1.1)$$

By introducing

$$y_0 = x,$$

$$f_0(y_0, y_1, \dots, y_m) = 1$$

and writing

$$Y = (y_0, y_1, \dots, y_m)^T,$$

$$F = (f_0, f_1, \dots, f_m)^T,$$

the initial value problem (1.1) can be rewritten in the form

$$\begin{cases} \frac{dY}{dx} = F(Y), \\ Y(x_0) = Y_0. \end{cases} \quad (1.2)$$

The solution of (1.2) is a curve in the space R^{m+1} . With the aid of the Frenet frame and the regular representation of the solution curve a class of non-linear formulas can be constructed. The derivation is rather complicated, for detail see [1]. Two of the non-linear formulas are as follows:

$$(I) \quad Y_{n+1} = Y_n + \frac{h}{2} \left(\frac{1}{l_n} F_n + \frac{1}{l_n^*} F_n^* \right),$$

where

* Received October 16, 1985.

$$l_n = \|F_n\|_2,$$

$$F_n^* = F\left(Y_n + h \frac{1}{l_n} F_n\right),$$

$$l_n^* = \|F_n^*\|_2,$$

and

$$(II) \quad Y_{n+1} = Y_n + \frac{h}{l_n} F_n + \frac{h^2}{6} \left(\left(\frac{1}{l_n} U_n - \frac{q}{l_n^4} F_n \right) + 2 \left(\frac{1}{l_n^{*2}} U_n^* - \frac{q}{l_n^{*4}} F_n^* \right) \right),$$

where

$$l_n = \|F_n\|_2,$$

$$F_n^* = F\left(Y_n + \frac{h}{2l_n} F_n + \frac{h^2}{8} \left(\frac{1}{l_n^2} U_n - \frac{q}{l_n^4} F_n \right)\right),$$

$$U_n = \frac{\partial F(Y_n)}{\partial Y} F_n,$$

$$q_n = F_n^T \cdot U_n,$$

$$l_n^* = \|F_n^*\|_2,$$

$$U_n^* = U\left(Y_n + \frac{h}{2l_n} F_n + \frac{h^2}{8} \left(\frac{1}{l_n^2} U_n - \frac{q}{l_n^4} F_n \right)\right),$$

$$q_n^* = U_n^{*T} \cdot F_n^*.$$

It should be pointed out that the h in the formulas is a step-size in arc length along the solution curve rather than a common step-size in the independent variable. It has been shown that formula (I) is a two-stage method of order 2 and formula (II) a two-stage method of order 4. Formula (II) is recommended by Qin Zeng-fu for solving stiff ordinary differential equations. In addition the step-size criterion derived from (I) or (II) (for the standard test equation $y = \lambda y$) is

$$h < \frac{2l^3}{\alpha(l^2 + 1)}, \quad \alpha = \operatorname{Re} \lambda, \quad (1.3)$$

$$h_{\text{perm}} = \frac{4(l^2 - 1)}{\kappa l^2(l^2 + 1)} \quad (1.4)$$

where κ is the curvature.

In this paper, we give another derivation of the non-linear methods by which a wider class of non-linear formulas can be easily constructed. However it can be seen from the new derivation that the non-linear formulas of [1] are essentially the results obtained by applying certain "explicit linear methods" to the ordinary differential equations which have been transformed with an independent variable transformation. One can expect that such kind of methods will have the same restriction on step-size as a general explicit method. The analysis in § 3 verifies this expectation and the numerical tests in § 4 is identical with the analysis.

§ 2. A New Derivation of the Formulas

Consider the initial value problem

$$\begin{cases} \frac{dY}{dx} = F(Y), \\ Y(x_0) = Y_0. \end{cases} \quad (2.1)$$

Introducing an independent transformation

$$x = x(s),$$

where s is a parameter of arc length along the solution curve, we have

$$\begin{aligned} Y(x) &= Y(x(s)) = Y(s), \\ \frac{dY(s)}{ds} &= \frac{dY(x(s))}{dx} \cdot \frac{dx}{ds} = \frac{1}{l} F(Y(s)), \end{aligned} \quad (2.2)$$

where

$$l = \left(\sum_{i=0}^m F_i(Y(s))^2 \right)^{1/2} = \|F\|_2. \quad (2.3)$$

Therefore the initial value problem becomes

$$\begin{cases} \frac{dY}{ds} = \tilde{F}(Y), \\ Y(s_0) = Y_0, \end{cases} \quad (2.4)$$

where

$$\tilde{F}(Y) = \frac{1}{l} F(Y(s)). \quad (2.5)$$

It is worthwhile to note that formally (2.4) is identical with a general initial value problem. By applying various types of "linear methods" to (2.4) we can get various non-linear formulas with respect to the original problem (2.1). The following are some examples.

1. Runge-Kutta type formulas

(1) Euler's formula

$$Y_{n+1} = Y_n + h\tilde{F}_n = Y_n + \frac{h}{l_n} F_n,$$

where $l_n = \|F_n\|_2$. This is formula (8) of [1].

(2) The second order formulas

$$(i) \quad Y_{n+1} = Y_n + \frac{h}{2} (K_1 + K_2),$$

$$K_1 = \tilde{F}_n = \frac{1}{l_1} F_n,$$

$$K_2 = \tilde{F}(Y_n + hK_1) = \frac{1}{l_2} F\left(Y_n + \frac{h}{l_1} F_n\right),$$

$$l_i = \|K_i\|_2, \quad \text{for } i=1, 2.$$

This is the two-stage method (10) of [1] or (I) in § 1.

(ii)

$$Y_{n+1} = Y_n + hK_2,$$

$$K_1 = \tilde{F}_n = \frac{1}{l_1} F(Y_n),$$

$$K_2 = \tilde{F}\left(Y_n + \frac{h}{2} K_1\right) = \frac{1}{l_2} F\left(Y_n + \frac{h}{2l_1} F_n\right),$$

$$l_i = \|K_i\|_2, \quad i=1, 2.$$

(3) The third order formulas

$$(i) \quad Y_{n+1} = Y_n + \frac{h}{4} (K_1 + K_3),$$

$$K_1 = \tilde{F}(Y_n) = \frac{1}{l_1} F(Y_n),$$

$$K_2 = \tilde{F}\left(Y_n + \frac{h}{3} K_1\right) = \frac{1}{l_2} F\left(Y_n + \frac{h}{3l_1} F_n\right),$$

$$K_3 = \tilde{F}\left(Y_n + \frac{2}{3} h K_2\right) = \frac{1}{l_3} F\left(Y_n + \frac{2}{3} h K_2\right),$$

$$l_i = \|K_i\|_2, \quad i=1, 2, 3.$$

$$(ii) \quad Y_{n+1} = Y_n + \frac{h}{6} (K_1 + 4K_2 + K_3),$$

$$K_1 = \tilde{F}(Y_n) = \frac{1}{l_1} F(Y_n),$$

$$K_2 = \tilde{F}\left(Y_n + \frac{h}{2} K_1\right) = \frac{1}{l_2} F\left(Y_n + \frac{h}{2} K_1\right),$$

$$K_3 = \tilde{F}(Y_n - hK_1 + 2hK_2) = \frac{1}{l_3} F(Y_n - hK_1 + 2hK_2),$$

$$l_i = \|K_i\|_2, \quad i=1, 2, 3.$$

(4) The fourth order formulas

$$Y_{n+1} = Y_n + \frac{h}{6} (K_1 + 2K_2 + 2K_3 + K_4),$$

$$K_1 = \tilde{F}(Y_n) = \frac{1}{l_1} F(Y_n),$$

$$K_2 = \tilde{F}\left(Y_n + \frac{h}{2} K_1\right) = \frac{1}{l_2} F\left(Y_n + \frac{h}{2} K_1\right),$$

$$K_3 = \tilde{F}\left(Y_n + \frac{h}{2} K_2\right) = \frac{1}{l_3} F\left(Y_n + \frac{h}{2} K_2\right),$$

$$K_4 = \tilde{F}(Y_n + hK_3) = \frac{1}{l_4} F(Y_n + hK_3),$$

$$l_i = \|K_i\|_2, \quad i=1, 2, 3, 4.$$

It is clear that we can construct other fourth order or higher order Runge-Kutta type formulas, and implicit Runge-Kutta formulas as well.

2. The Runge-Kutta type formulas with second derivatives

Set

$$G(Y(s)) = \tilde{F}'(Y(s)).$$

Then we have the following formulas:

$$(1) \quad Y_{n+1} = Y_n + h\tilde{F}_n + \frac{1}{2} h^2 G_n = Y_n + \frac{1}{l_1} F_n + \frac{1}{2} h^2 \left(\frac{1}{l_1^2} U_n - \frac{q_n}{l_1^4} F_n \right),$$

where

$$l_1 = \|F(Y_n)\|_2,$$

$$U_n = \frac{\partial F(Y_n)}{\partial Y} \cdot F_n,$$

and

$$q_n = F_n^T \cdot U_n,$$

which is the second order formula (9) with the second derivatives.

$$(2) \quad Y_{n+1} = Y_n + h\tilde{F}_n + \frac{h^2}{6} (2G_n + \bar{G}_n),$$

where

$$\bar{G}_n = G(Y_n + h\tilde{F}_n).$$

Since

$$\bar{G}_n = G(Y_n + h\tilde{F}_n) = G(Y_n) + \frac{\partial G_n}{\partial Y} \cdot h\tilde{F}_n + O(h^2)$$

we get

$$\begin{aligned} Y_{n+1} &= Y_n + h\tilde{F}_n + \frac{h^2}{2} G_n + \frac{h^3}{6} \frac{\partial G_n}{\partial Y} \tilde{F}_n + O(h^4) \\ &= Y_n + h\tilde{F}_n + \frac{h^2}{2} \dot{\tilde{F}}_n + \frac{h^3}{6} \ddot{\tilde{F}}_n + O(h^4). \end{aligned}$$

Clearly, the third order formula (2) is formula (11) of [1].

(3) The two-stage fourth order formula.

Let

$$\begin{aligned} \bar{G}_n &= G(Y_n + ah\tilde{F}_n + bh^2\dot{\tilde{F}}_n) \\ &= G(Y_n) + \frac{\partial G(Y_n)}{\partial Y} (ah\tilde{F}_n + bh^2\dot{\tilde{F}}_n) + \frac{1}{2} \frac{\partial^2 G(Y_n)}{\partial Y^2} (ah\tilde{F}_n + bh^2\dot{\tilde{F}}_n)^2 + O(h^3), \end{aligned}$$

where a and b are undetermined parameters. Then we consider a linear combination of G_n and \bar{G}_n , $c_1G_n + c_2\bar{G}_n$. We have

$$\begin{aligned} c_1G_n + c_2\bar{G}_n &= (c_1 + c_2)G_n + ac_1h \frac{\partial G_n}{\partial Y} \tilde{F}_n \\ &\quad + h^2 \left(bc_2 \frac{\partial G_n}{\partial Y} \dot{\tilde{F}}_n + \frac{1}{2} a^2c_2 \frac{\partial^2 G_n}{\partial Y^2} \cdot \tilde{F}_n^2 \right) + O(h^3). \end{aligned}$$

In order to obtain a fourth order formula, the following equations should be satisfied

$$\begin{cases} c_1 + c_2 = \frac{1}{2}, \\ ac_2 = \frac{1}{6}, \\ b = \frac{1}{2} a^2, \\ bc_2 = \frac{1}{24}, \end{cases}$$

which have a solution:

$$a = \frac{1}{2}, \quad b = \frac{1}{8}, \quad c_1 = \frac{1}{3}, \quad c_2 = \frac{1}{6}.$$

Now we have a two-stage formula

$$Y_{n+1} = Y_n + h\tilde{F}_n + \frac{h^2}{6} (G_n + \bar{G}_n),$$

where

$$\bar{G}_n = G\left(Y_n + \frac{1}{2} h\tilde{F}_n + \frac{1}{8} h^2\dot{\tilde{F}}_n\right).$$

This is a fourth order formula. It is formula (15) of [1] or (II) in § 1.

It can be seen from the derivation above that all the formulas given are only

the results obtained by applying some explicit linear methods to the system of ordinary differential equations for which an independent transformation has been performed. Using arc length as a parameter, we can set up in an explicit way the relation between (2.4) and (2.1). In the next section we will show that this kind of methods has the same restriction on step-size as a general explicit method.

§ 3. The Step-Size Criterion

The step-size criterion (1.4) has been derived for the standard test equation

$$\begin{cases} \frac{dy_0}{dx} = 1, \\ \frac{dy_1}{dx} = -\alpha y_1 - \beta y_2, \\ \frac{dy_2}{dx} = \beta y_1 - \alpha y_2. \end{cases} \quad \alpha > 0, \beta \text{ real number}, \quad (3.1)$$

Note that (3.1) is derived from the scalar test equation $y' = \lambda y$, $\lambda = -\alpha + i\beta$. First we point out that the criterion is applicable to a system of linear equations $Y' = AY$ which can be reduced to a diagonal system under unitary transformations. That means that there exists a unitary matrix U such that

$$Z' = \Lambda Z$$

where $\Lambda = U^{-1}AU = \text{diag}(\lambda_i)$, $Z = U^{-1}Y$. Since the 2-norm is invariant under unitary transformation, we have

$$l = \|\Lambda Z\|_2 = \|U^{-1}AUU^{-1}Y\|_2 = \|U^{-1}AY\|_2 = \|AY\|_2 = l. \quad (3.2)$$

Therefore it can be shown that the following communications hold:

$$\begin{array}{ccc} Y' = AY & \xrightarrow{U} & Z' = \Lambda Z \\ \downarrow \text{transformation} & & \downarrow \text{transformation} \\ \frac{dY}{ds} = \frac{1}{l} AY & & \frac{dZ}{ds} = \frac{1}{l} \Lambda Z \\ \downarrow \text{linear method} & & \downarrow \text{linear method} \\ Y_1 & \xrightarrow{U} & Z_1 \end{array}$$

Thus we now can analyze the step-size criterion for a diagonal system

$$\frac{dZ}{ds} = \Lambda Z, \quad \Lambda = \text{diag}(\lambda_i). \quad (3.3)$$

For formula (I) in § 1 the step-size h must satisfy the relation

$$h < \frac{2l^3}{\alpha(1+l^2)} \quad (3.4)$$

where

$$\alpha = \max_i |\text{Re } \lambda_i|,$$

$$l = \left(\sum_{i=0}^m |\lambda_i Z_i|^2 \right)^{1/2}.$$

In the transient phase the solution components are relatively large. Without loss of generality, let $Z_i = O(1)$ and $l \doteq \max |\lambda_i|$. For a stiff problem $\alpha \doteq \max |\lambda_i|$ so that from (3.4) the restriction on step-size h is roughly as follows:

$$h < 2. \quad (3.5)$$

This means that in the transient phase we can take a rather big step-size. However, when the values of Z_i are almost negligible in the stationary phase, $l \rightarrow 1$ and the right-hand side of (3.4) tends to $1/\alpha$. In this case the restriction on step-size is approximately

$$\alpha h < 1 \quad (3.6)$$

or

$$\max |\lambda_i| h < 1,$$

which is the same as the restriction on step-size to a general explicit method.

At first glance it is strange that the step-size h can take a big value in transient phase while it must be very small in stationary phase. It will become clear if we return to the original independent variable x . Notice that we have

$$\frac{dx}{ds} = \frac{1}{l}. \quad (3.7)$$

It can be approximately expressed as

$$\frac{\Delta x}{\Delta s} \doteq \frac{1}{l} \quad \text{or} \quad \Delta s \doteq l \Delta x. \quad (3.7)$$

Now, it can be seen that although Δs can be quite big from the viewpoint that a length Δx is only $\Delta s/l$, it is still very small. It turns out from (3.5) that $\Delta x < 2/l$. Therefore, from the viewpoint of the independent variable x the step-size is always subjected to the restriction which has been imposed on general explicit methods during the whole computation. In our opinion, this class of methods is not suitable for solving stiff ordinary differential equations.

§ 4. Numerical Tests

The computations have been carried out for the following two examples, using the method (II).

Example 1. The diagonal system of equations

$$\begin{cases} y_1' = -y_1, & y_1(0) = 1, \\ y_2' = -\lambda y_2, & y_2(0) = 1. \end{cases}$$

Example 2. The non-linear equations

$$\begin{cases} y_1' = 0.01 - (0.01 + y_1 + y_2)(1 + (y_1 + 1000)(y_1 + 1)), & y_1(0) = 0, \\ y_2' = 0.01 - (0.01 + y_1 + y_2)(1 + y_2^2), & y_2(0) = 0. \end{cases}$$

In the computations the step-size Δs is determined from the step-size criterion (1.4). At each step, we also compute the solution using the classical fourth order Runge-Kutta method with Δx produced by (II), in order to compare the two solutions. From the results of computations we obtain the following observations:

(1) In both examples

$$\Delta s \doteq l \Delta x$$

always holds.

(2) The step-size Δs determined by (1.4) is not so big as expected in (3.5). In transient phase Δs is often too small. For example 1, when we take $\lambda=100$ the value of Δs is between 0.026 and 0.029 for $0 < x < 1$, and the value of Δx slowly increases from 0.000285 to 0.02 as l decreases from 100 to 1. When $\lambda=1000$ it is even more typical. In this case, Δs is between 0.00283 and 0.0025 and Δx between 0.00000283 and 0.000004. We did not complete the computation, because it took too many steps. It is worthwhile to note that the value of Δs is something like the restriction for the Runge-Kutta method.

For Example 2, Δs is also smaller than the step-size which is necessary in transient phase. At $x=0.005$, $\Delta s = \Delta x = 0.00007324$, which is much smaller than the step-size for the Runge-Kutta method as shown in [1].

(3) When Δs determined by (1.4) is big, the accuracy of results is poor. In example 2 when $0.02 < \Delta s < 0.05$ the accuracy is about 10^{-2} . When $\Delta s > 0.05$ we get wrong results. [1] has suggested using

$$h_{\text{working}} = \min(h_{\text{max}}, h_{\text{perm}}),$$

but the selection of h_{max} remains a problem.

References

- [1] Qin Zeng-fu, A class of non-linear methods for ordinary differential equations, JCM, 3: 4 (1985), 320—327.
- [2] E. Hairer, Unconditionally stable explicit methods for parabolic equations, *Numer. Math.*, V. 135, pp. 57—68.