

107

SURVEY OF DIGITAL SOLUTION
OF DIFFERENTIAL EQUATIONS

by

914

ING-WEN HWANG

B.S.E.E., National Taiwan University, 1963

A MASTER'S REPORT

submitted in partial fulfillment of the

requirement for the degree

MASTER OF SCIENCE

Department of Electrical Engineering

KANSAS STATE UNIVERSITY
Manhattan, Kansas

1966

Approved by:

Charles A. Halijak

Major Professor

LD
2668
R4
1966
H991

TABLE OF CONTENT

I. INTRODUCTION 1

II. SURVEY OF METHODS FOR DIGITAL SOLUTION OF DIFFERENTIAL EQUATION . 3

III. CLASSICAL NUMERICAL TECHNIQUES..... 3

 1). One-Step Methods 3

 A. Taylor Series Method 4

 B. Euler Method 6

 C. Runge-Kutta Methods 10

 2). Multiple Step Methods 30

 A. Adams-Bashforth Method 32

 B. Adams-Moulton Method 36

 C. Milne Method 40

 D. Hamming Method 47

 3). Comparison between One-Step Method and Multiple Step-Method. 54

IV. NUMERICAL TRANSFORM TECHNIQUES 54

 A. Tustin Program 56

 B. Single-Integrator Program 58

 C. Multiple-Integrator Substitution Program 60

V. SUMMARY 62

APPENDIX 1 64

REFERENCES 66

ACKNOWLEDGMENT 68

INTRODUCTION

Developments of modern science have confronted the scientist and the engineer with a variety of problems which cannot be solved formally. Hence the widening interest in numerical analysis, a branch of mathematics which leads to approximate solutions by repeated application of the four basic operations of algebra. Numerical analysis has been applied to scientific and technological problems from the very beginning of applied science, but has been given new impetus by development of the electronic digital computer. Many classical numerical analysis techniques are available for engineering applications to solve differential equations. There is another group of procedures unique to the engineering literature which generates approximating recurrence relations from approximate Z-transforms. No matter how different all these techniques are in terms of their mathematical approaches and the algebra involved, they all have two things in common; the calculations are performed with discrete values and on a step-by-step basis. Consequently the time interval between steps is an important factor which affects the accuracy as well as the speed of the computation.

Up to a point, the smaller the step size used, the longer is the solution time required and the more accurate is the solution obtained. Since significant digits are limited in computation, the increased solution accuracy produced by the reduced step size is lost. This effect is termed, "round-off error" Therefore there exists an optional step size to be used in producing the numerical solution of optimal accuracy. In general no single technique is best in all cases. In fact, the effectiveness of approximate methods hinges on the type of function in question and the goodness of each

method is measured by a number of factors, namely, the program set-up time, the solution speed, and solution accuracy and stability. In other words, it is a consideration of economy and accuracy and these two quantities are contradictory in nature. Therefore a brief examination of strong points and weaknesses of various types of digital techniques and comparisons among them would provide a guidepost for selection of an optimal technique.

The objective of this report is to review previous methods for digital solution of differential equations, and to illustrate their applications for approximating the solutions of ordinary differential equations.

II. Survey of Methods for Digital Solution of Differential Equations.

Methods for approximating solution of ordinary differential equations are based on the principle of discretization. These methods have the common feature that no attempt is made to approximate the exact solution $y(t)$ over a continuous range of the independent variable. Approximating values are sought only on a set of discrete points $t_0, t_1, t_2, \dots, t_n$. Generally speaking, a discrete variable method for solving a differential equation consists of an algorithm which furnishes a number y_n corresponding to each lattice point t_n . Which is to be regarded as an approximation to the value $y(t_n)$ of the exact solution at the point t_n .

Discrete variable methods fall into two classes: classical numerical analysis techniques; and numerical transform techniques. Classical numerical analysis techniques yield one-step methods and multiple step methods. In a one-step method the value of y_n can be found if only one initial value is known. In a multiple step method the calculation of y_{n+1} requires explicit knowledge of more than one starting value. Oftentimes, a one-step method is used to start a multiple step method.

The numerical transform techniques were first introduced by Tustin and were further expanded by Madwed, Boxer-Thaler, and Halijak. These techniques generate approximating recurrence relations from approximate z-transform of $1/s^n$ and \bar{f}/s^n .

III. Classical Numerical Analysis Techniques.

1). One-Step Methods

Let a typical first-order differential equation be given by

$$\frac{dy}{dt} = f(t, y)$$

$$y = y_0 \quad \text{at } t = t_0 \quad (3.1)$$

and let it satisfy Lipschitz conditions in some closed region D. There exists a single-valued function $y(t)$ continuous in D such that it satisfies Eq. (3.1), then it can be solved approximately by these methods.

One-step methods are usually divided into two classes; The first class includes Taylor series method and Picard's method. In these methods, y in Eq. (3.1) is approximated by a truncation series, the individual terms of which are functions of the independent variable t . The second class is represented by the methods of Euler and Runge-Kutta methods.

A. Taylor Series Method.

Consider the differential equation (3.1) with the initial condition $y=y_0$ at $t=t_0$. Let the required solution be

$$y = y(t) \quad (3.2)$$

If $t=t_0$ is not a singular point of the function, $y(t)$ can be expanded in Taylor series about this point. Thus with

$$y_0^{(m)} = \left[\frac{d^m y}{dt^m} \right]_{t=t_0} \quad (3.3)$$

$$y = y_0 + (t-t_0)y_0^{(1)} + \frac{1}{2!}(t-t_0)^2 y_0^{(2)} + \frac{1}{3!}(t-t_0)^3 y_0^{(3)} + \dots \quad (3.4)$$

a power series in t that converges over some range $t_0 \leq t \leq t_b$.

The value of $y^{(1)}$ is evaluated by making use of the differential equation (3.1) which can be written as

$$y^{(1)} = f(t, y) = g_1(t, y), \quad (3.5)$$

differentiating Eq. (3.5) yields

$$y^{(2)} = \frac{\partial g_1}{\partial t} + \frac{\partial g_1}{\partial y} y^{(1)} = g_2(t, y, y^{(1)}), \quad (3.6)$$

and by repeated differentiation

$$y^{(m)} = g_m(t, y, y^{(1)}, y^{(2)}, \dots, y^{(m-1)}) \quad m = 1, 2, \dots \quad (3.7)$$

where

$$g_m = \frac{\partial g_{m-1}}{\partial t} + \sum_{i=0}^{m-2} \frac{\partial g_{m-1}}{\partial y^{(i)}} y^{(i+1)}, \quad (3.8)$$

Setting $t=t_0$ and $y=y_0$, Eq. (3.7) becomes

$$y_0^{(m)} = g_m(t_0, y_0, y_0^{(1)}, \dots, y_0^{(m-1)}), \quad m = 1, 2, \dots \quad (3.9)$$

The truncation error after the m th term is given by

$$R_m = \frac{f^{(m)}(\xi)}{m!} (t - t_0)^m, \quad \text{where} \quad t_0 \leq \xi \leq t \quad (3.10)$$

Replacing $f(\xi)$ with an upper bound $M^{(m)}$ to its value in the interval (t, t_0) . The truncation error is then bounded by

$$R_m \leq \left| \frac{M^{(m)}}{m!} (t - t_0)^m \right|, \quad (3.11)$$

In particular, for a convergent Taylor series with alternating sign, the truncation error after m th term cannot exceed the $(m+1)$ th term Eq.(3.11) becomes

$$|R_m| \leq \left| \frac{f^{(m)}(t_0)}{m} (t - t_0)^m \right|, \quad (3.12)$$

B. 1. Euler's Method.

This method is of very little practical importance, but it illustrates in simple form the basic idea of those numerical methods which seek to determine the change of Δy in y corresponding to a small increment of the argument.

Consider Eq.(3.1), the left-hand side of the first part is, by definition

$$\frac{dy}{dt} = \lim_{\Delta t \rightarrow 0} \frac{\Delta y}{\Delta t} \quad (3.13)$$

therefore, for small value of Δt

$$\Delta y \approx \frac{dy}{dt} \Delta t$$

Thus, the increase in $y(t)$ when t increases to $t + \Delta t$ is approximately

$$y_{m+1} - y_m = f(t_m, y_m)T \quad (3.14)$$

where

$$T = \Delta t$$

However, the exact expression for y at $t = t_{m+1}$, using Taylor series formula is

$$z_{m+1} = z_m + Tf(t, z) + \frac{1}{2} T^2 y^{(2)}(\xi)_{m+1} \quad (3.15)$$

where

$$t_m < \xi_{m+1} < t_{m+1}.$$

Subtracting Eq. (3.14) from Eq. (3.15) results in

$$z_{m+1} - y_{m+1} = (z_m - y_m) + T(f(t_m, z_m) - f(t_m, y_m)) \\ + \frac{1}{2} T^2 y^{(2)}(\xi_{m+1}). \quad (3.16)$$

The left-hand side of Eq.(3.16) is by definition the truncation error of y at $t=t_{m+1}$,

$$\text{Let } \epsilon_{m+1} = z_{m+1} - y_{m+1} \quad (3.17)$$

$$\text{and } k_m = \frac{f(t_m, z_m) - f(t_m, y_m)}{z_m - y_m} \approx \frac{\partial f(t_m, z_m)}{\partial y} \quad (3.18)$$

Eq.(3.16) becomes

$$\epsilon_{m+1} = \epsilon_m + Tk_m \epsilon_m + \frac{1}{2} T^2 y^{(2)}(\xi_{m+1}). \quad (3.19)$$

For $m = 0$

$$\epsilon_0 = 0 \quad (3.20)$$

$m = 1$

$$\epsilon_1 = \frac{1}{2} R^2 y^{(2)}(\xi_1), \quad t_0 < \xi_1 < t_1 \quad (3.21)$$

$$m = 2$$

$$\epsilon_2 = (t + TK_1) \frac{1}{2} T^2 y^{(2)}(\xi_1) + \frac{1}{2} T^2 y^{(2)}(\xi_2) \quad (3.22)$$

where

$$t_1 < \xi_2 < t_2 ;$$

proceeding in this manner, the truncation error at n th step is

$$\begin{aligned} \epsilon_n = T^2 \left[\frac{1}{2} y^{(2)}(\xi_1) \prod_{j=1}^{n-1} (1 + TK_j) + \frac{1}{2} y^{(2)}(\xi_2) \prod_{j=2}^{n-1} (1 + TK_j) \right. \\ \left. + \dots + \frac{1}{2} y^{(2)}(\xi_n) \right] \quad (3.23) \end{aligned}$$

$$\text{or} \quad \epsilon_n = T^2 \sum_{i=0}^{n-1} A_i a_i \quad (3.24)$$

where

$$\begin{cases} A_i = \prod_{j=1}^{n-1} (1 + TK_j) & i \neq n \\ A_i = 1 & i = n \end{cases}$$

$$\text{and} \quad a_i = \frac{1}{2} y^{(2)}(\xi_i) \quad (3.25)$$

Thus errors $\frac{1}{2} y^n (\xi_1)$ introduced at each step because of the inaccuracy of formula (3.14) are to be multiplied by amplification factors A_1 before being summed.

C. Runge-Kutta Method.

This method is an algorithm designed to approximate the Taylor's series solution.

Consider the Taylor's Series

$$y_{n+1} = y_n + T y'_n + \frac{T^2}{2} y''_n + \dots, \quad (3.26)$$

where

$$y_{n+1} = y((n+1)T) \quad \text{and} \quad y_n = y(nT).$$

If the Runge-Kutta formula is derived by retaining terms in the Taylor's series expression up to m th power of T , this formula is called the Runge-Kutta method of m th order accuracy or the Runge-Kutta m th order method.

C. 1. Second Order Runge-Kutta Formula.

Consider increments in y defined by the equations

$$\Delta y = R_1 \Delta'y + R_2 \Delta''y \quad (3.27)$$

where

$$\Delta'y = f(t_n, y_n) T,$$

$$\Delta y = f(t_n + aT, y_n + bT) T,$$

and R_1, R_2, a and b are constants.

Expanding Eq. (3.27) with respect to Taylor's series of two variables, results in

$$\Delta y = T(R_1 + R_2) f_n + (a(f_t)_n + b(f_y)_n) T^2 R_2 + \dots \quad (3.28)$$

where

$$f_n = f(t_n, y_n) = y_n, \quad (f_t)_n = \left. \frac{\partial f}{\partial t} \right|_{t=t_n, y=y_n},$$

and

$$(f_y)_n = \left. \frac{\partial f}{\partial y} \right|_{t=t_n, y=y_n}$$

Eq. (3.28) will be equal to Eq. (3.26) up to second order of T if its parameters take the values

$$R_1 + R_2 = 1 \quad (3.29)$$

$$R_2 a = 1/2$$

$$R_2 b = 1/2$$

There are four constants to be determined and only three equations. Choose R_2 as parameter, the constants $a, b,$ and R_1 can be determined in terms of R_2 . Thus

$$a = b = \frac{1}{2R_2} \quad (3.30)$$

$$R_1 = 1 - R_2$$

Consequently, an infinite number of forms of Eq.(3.28) can be established. Formulas of higher order can be obtained in the similar manner although no formula beyond fifth order has ever been developed.

C. II. Third Order Runge-Kutta Formula.

Again take the differential equation

$$\frac{dy}{dt} = f(t, y). \quad (3.31)$$

and expand it with respect to (t_n, y_n) by Taylor series for two variables to obtain

$$\begin{aligned} \Delta y = f_n T + \frac{1}{2} (f_t + f_y f)_n T^2 + \frac{1}{6} \left\{ f_{tt} + 2f_{ty} f + f_{yy} f^2 \right. \\ \left. + (f_t + f_y f) f_y \right\}_n T^3 + \dots \end{aligned} \quad (3.32)$$

Define the increment of y by

$$\Delta y = R_1 \Delta'y + R_2 \Delta''y + R_3 \Delta'''y \quad (3.33)$$

where

$$\Delta'y = f(t_n, y_n) T$$

$$\Delta''y = f(t_n + mT, y_n + m\Delta'y) T$$

$$\Delta''y = f(t_n + \lambda T, y_n + p\Delta'y + (\lambda - p)\Delta'y) T$$

The symbols m , λ , and p are constants. Expanding Eq. (3.33) with respect to Taylor's series at (t_n, y_n) , results in

$$\Delta'y = Tf_n$$

$$\Delta''y = T \left\{ f_n + mT(f_t)_n + m\Delta'y(f_y)_n + \frac{1}{2!} \left[(mT)^2(f_{tt})_n + 2m^2T\Delta'y(f_{ty})_n + (m\Delta'y)^2(f_{yy})_n \right] + \dots \right\}$$

$$\begin{aligned} \Delta''y = T \left\{ f_n + \lambda T(f_t)_n + \left[\rho\Delta''y + (\lambda - \rho)\Delta'y \right] (f_y)_n \right. \\ \left. + \frac{1}{2!} \left\{ (\lambda T)^2(f_{tt})_n + 2\lambda T \left[\rho\Delta''y + (\lambda - \rho)\Delta'y \right] (f_{ty})_n \right. \right. \\ \left. \left. + \left[\rho\Delta''y + (\lambda - \rho)\Delta'y \right]^2 (f_{yy})_n \right\} + \dots \right\} \quad (3.34) \end{aligned}$$

where

$$(f_{yy})_n = \left(\frac{\partial^2 f}{\partial y^2} \right)_{t=t_n, y=y_n} \quad (f_y)_n = \left(\frac{\partial f}{\partial y} \right)_{t=t_n, y=y_n}$$

$$(f_t)_n = \left(\frac{\partial f}{\partial t} \right)_{t=t_n, y=y_n} \quad f_n = f(t_n, y_n)$$

Eliminating $\Delta'y$ and $\Delta''y$ from the right hand side of the second and third equations of Eq.(3.34) and substiting into Eq.(3.33), the result coincides with Taylor's series expansion up to the third power of T provided the following relations hold among the constants

$$R_1 + R_2 + R_3 = 1 \quad (3.35)$$

$$R_2 m + R_3 \lambda = \frac{1}{2}$$

$$R_2 m^2 + R_3 \lambda^2 = \frac{1}{3}$$

$$R_3 \rho m = \frac{1}{6}$$

Since there are six constants to be determined, namely, R_1 , R_2 , R_3 , m , λ , and ρ , and only four equations, the constants m and λ can be chosen as parameters. The constants R_1 , R_2 , R_3 and ρ as calculated from Eq.(3.35) are

$$R_1 = \frac{6m\lambda - 3(m + \lambda) + 2}{6m\lambda} \quad (3.36)$$

$$R_2 = \frac{2 - 3\lambda}{6m(m - \lambda)}$$

$$R_3 = \frac{2 - 3m}{6\lambda(\lambda - m)}$$

$$\rho = \frac{\lambda(\lambda - m)}{m(2 - m)}$$

From Eq.(3.35), it is clear that R_3 cannot be zero and R_1 and R_2 cannot both zero. If $R_2=0$, then by Eq.(3.35)

$$\lambda = \frac{2}{3} \quad R_3 = \frac{3}{4} \quad R_1 = \frac{1}{4} \quad \text{and} \quad \rho = \frac{2}{3} \quad (3.37)$$

Moreover, assume ρ to be $1/3$, then m is equal to $2/3$. Therefore

$$\Delta y = \frac{1}{4} (\Delta'y + 3 \Delta''y) \quad (3.38)$$

$$\Delta'y = f(t_n, y_n) T$$

$$\Delta''y = f\left(t_n + \frac{2}{3} T, y_n + \frac{2}{3} \Delta'y\right) T$$

$$\Delta'''y = f\left[t_n + \frac{2}{3} T, y_n + \frac{2}{3} (\Delta'y + \Delta''y)\right] T$$

On the other hand, if $\lambda=0$, then

$$m = \frac{2}{3} \quad R_2 = \frac{3}{4} \quad R_3 = \frac{1}{4\rho} \quad \text{and} \quad R_1 = \frac{\rho - 1}{4\rho} \quad (3.39)$$

If it is assumed that $\rho=1$, then another form can be constructed from Eq.(3.39) and Eq.(3.33); it is

$$\Delta y = \frac{1}{4} (3 \Delta''y + \Delta'''y) \quad (3.40)$$

where

$$\Delta'y = f(t_n, y_n) T$$

$$\Delta''y = f\left(t_n + \frac{2}{3}T, y_n + \frac{2}{3}\Delta'y\right) T$$

$$\Delta''''y = f\left(t_n, y_n + \Delta''y - \Delta'y\right) T$$

C. IV. Fourth Order Runge-Kutta Formula.

The Runge-Kutta fourth order formula has the form

$$\Delta y = R_1 \Delta'y + R_2 \Delta''y + R_3 \Delta''''y + R_4 \Delta''''''y \quad (3.41)$$

where

$$\Delta'y = f(t_n, y_n) T$$

$$\Delta''y = f\left(t_n + mT, y_n + m\Delta'y\right) T$$

$$\Delta''''y = f\left(t_n + \lambda T, y_n + \rho\Delta''y_n + (\lambda - \rho)\Delta'y\right) T$$

$$\Delta''''''y = f\left(t_n + \mu T, y_n + \sigma\Delta''''y + \tau\Delta''y + (\mu - \sigma - \tau)\Delta'y\right) T.$$

Expand second, third, fourth and fifth equations of Eq. (3.41) to obtain

$$\Delta'y = f(t_n, y_n) T \quad (3.42)$$

$$\begin{aligned}
\Delta''y &= T \left\{ f(t_n, y_n) + mT(f_t)_n + m\Delta'y(f_y)_n \right. \\
&+ \frac{1}{2!} \left[(mT)^2 (f_{tt})_n + 2(mT)(m\Delta'y)(f_{ty})_n + (m\Delta'y)^2 (f_{yy})_n \right] \\
&+ \frac{1}{3!} \left[(mT)^3 (f_{ttt})_n + 3(mT)^2(m\Delta'y)(f_{tty})_n \right. \\
&\left. + 3(mT)(m\Delta'y)^2 (f_{tyy})_n + (m\Delta'y)^3 (f_{yyy})_n \right] + \dots \left. \right\}
\end{aligned}$$

$$\begin{aligned}
\Delta''y &= T \left\{ f(t_n, y_n) + \lambda T(f_t)_n + (\rho\Delta''y + (\lambda - \rho)\Delta'y)(f_y)_n \right. \\
&+ \frac{1}{2!} \left[(\lambda T)^2 (f_{tt})_n + 2(\lambda T)(\rho\Delta''y + (\lambda - \rho)\Delta'y)(f_{ty})_n \right. \\
&\left. + (\rho\Delta''y + (\lambda - \rho)\Delta'y)^2 (f_{yy})_n \right] \\
&+ \frac{1}{3!} \left[(\lambda T)^3 (f_{ttt})_n + 3(\lambda T)^2(\rho\Delta''y + (\lambda - \rho)\Delta'y)(f_{tty})_n \right. \\
&+ 3(\lambda T)(\rho\Delta''y + (\lambda - \rho)\Delta'y)^2 (f_{tyy})_n + (\rho\Delta''y + (\lambda - \rho)\Delta'y)^3 \\
&\left. + (\rho\Delta''y + (\lambda - \rho)\Delta'y)^3 (f_{yyy})_n \right] + \dots \left. \right\}
\end{aligned}$$

$$\begin{aligned}
\Delta^{(4)}y &= T \left\{ f_n + \mu T(f_t)_n + (\sigma\Delta^{(3)}y + \tau\Delta^{(2)}y + (\mu - \sigma - \tau)\Delta'y)(f_y)_n \right. \\
&+ \frac{1}{2!} \left[(\mu T)^2 (f_{tt})_n + 2(\mu T)(\sigma\Delta^{(3)}y + \tau\Delta^{(2)}y \right. \\
&+ (\mu - \sigma - \tau)\Delta'y)(f_{ty})_n + ((\mu - \sigma - \tau)\Delta'y + \sigma\Delta^{(3)}y \\
&\left. + \tau\Delta^{(2)}y)^2 (f_{yy})_n \right] + \dots \left. \right\}
\end{aligned}$$

Eliminating $\Delta'y$, $\Delta''y$ and $\Delta'''y$ from the right-hand side of Eq. (3.42) by successive substitution, putting these results into Eq. (3.41) and comparing with Taylor's series yield

$$R_1 + R_2 + R_3 + R_4 = 1 \quad (3.43)$$

$$R_2 m + R_3 \lambda + R_4 \mu = \frac{1}{2}$$

$$R_2 m^2 + R_3 \lambda^2 + R_4 \mu^2 = \frac{1}{3}$$

$$R_3 m \rho + R_4 (m \sigma + \lambda \tau) = \frac{1}{6}$$

$$R_2 m^3 + R_3 \lambda^3 + R_4 \mu^3 = \frac{1}{4}$$

$$R_3 m \lambda \rho + R_4 \mu (m \sigma + \lambda \tau) = \frac{1}{8}$$

$$R_3 m^3 \rho + R_4 (m^2 \sigma + \lambda^2 \tau) = \frac{1}{12}$$

$$R_4 m \lambda \tau = \frac{1}{24}$$

There is a one-parameter family of solutions derived by Kutta as follows:

$$R_1 = \frac{1}{6}, \quad R_2 = \frac{2-t}{3}, \quad R_3 = \frac{t}{3}, \quad R_4 = \frac{1}{6} \quad (3.44)$$

$$m = 1/2, \quad \lambda = 1/2, \quad \rho = 1/2t, \quad \mu = 1, \quad \sigma = 1-t, \quad \tau = t$$

If $t=1$, a Runge-Kutta fourth order formula can be constructed as follows:

$$\Delta y = \frac{1}{6} (\Delta^1 y + 2\Delta^2 y + 2\Delta^3 y + \Delta^4 y) \quad (3.45)$$

where

$$\Delta y = y_{n+1} - y_n,$$

$$\Delta^1 y = T f(t_n, y_n),$$

$$\Delta^{(2)} y = T f\left(t_n + \frac{1}{2} T, y_n + \frac{1}{2} \Delta^1 y\right)$$

$$\Delta^{(3)} y = T f\left(t_n + \frac{1}{2} T, y_n + \frac{1}{2} \Delta^2 y\right)$$

$$\Delta^{(4)} y = T f(t_n + T, y_n + \Delta^3 y)$$

4). Fifth Order Runge-Kutta Formula

The 5th order Runge-Kutta formulas were derived recently by H. A. Luther and H. P. Konen. The derivations are as follows:

$$\text{Let } dy/dt = f(t, y) \text{ together with } y(t_0) = y_0 \quad (3.46)$$

The fifth order Runge-Kutta formulas are phrased as

$$\Delta y = + \sum_1^6 R_i \Delta^{(i)} y$$

$$\Delta^s y = T f(t_n, y_n) \quad (3.47)$$

$$\Delta^{(s)} y = T f(t_n + a_s T, y + \sum_{j=1}^{s-1} b_{sj} \Delta^{(j)} y).$$

where $2 \leq s \leq 6$. and the coefficients R_i are the constants to be determined.

The usual procedure yields 44 equations involving the various parameters.

Assume that

$$a_s = \sum_{j=1}^{s-1} b_{sj}, \quad 2 \leq s \leq 6 \quad (3.48)$$

and eliminate easily identifiable combinations to obtain the following 16 relations

$$\left\{ \begin{array}{l} R_1 + R_2 + R_3 + R_4 + R_5 + R_6 = 1, \\ a_2 R_2 + a_3 R_3 + a_4 R_4 + a_5 R_5 + a_6 R_6 = \frac{1}{2}, \\ a_2^2 R_2 + a_3^2 R_3 + a_4^2 R_4 + a_5^2 R_5 + a_6^2 R_6 = \frac{1}{3}, \\ a_2^3 R_2 + a_3^3 R_3 + a_4^3 R_4 + a_5^3 R_5 + a_6^3 R_6 = \frac{1}{4}, \\ a_2^4 R_2 + a_3^4 R_3 + a_4^4 R_4 + a_5^4 R_5 + a_6^4 R_6 = \frac{1}{5} \end{array} \right. \quad (3.49a)$$

$$c_1 R_3 + c_2 R_4 + c_3 R_5 + c_4 R_6 = \frac{1}{6}$$

$$a_3 c_1 R_3 + a_4 c_2 R_4 + a_5 c_3 R_5 + a_6 c_4 R_6 = 1/8,$$

$$a_3^2 c_1 R_3 + a_4^2 c_2 R_4 + a_5^2 c_3 R_5 + a_6^2 c_4 R_6 = 1/10,$$

$$d_1 R_3 + d_2 R_4 + d_3 R_5 + d_4 R_6 = 1/12,$$

(3.49a)

$$a_3 d_1 R_3 + a_4 d_2 R_4 + a_5 d_3 R_5 + a_6 d_4 R_6 = 1/15,$$

$$c_1^2 R_3 + c_2^2 R_4 + c_3^2 R_5 + c_4^2 R_6 = 1/20,$$

$$a_2^3 b_{32} R_3 + (a_2^3 b_{42} + a_3^3 b_{43}) R_4 + (a_2^3 b_{52} + a_3^3 b_{53} + a_4^3 b_{54}) R_5 \\ + (a_2^3 b_{62} + a_3^3 b_{63} + a_4^3 b_{64} + a_5^3 b_{65}) R_6 = 1/20$$

$$c_1 b_{32} R_3 + (c_1 b_{53} + c_2 b_{54}) R_5 + (c_1 b_{63} + c_2 b_{64} + c_3 b_{65}) R_6 = 1/24,$$

$$d_1 b_{43} R_4 + (d_1 b_{53} + d_2 b_{54}) R_5 + (d_1 b_{63} + d_2 b_{64} + d_3 b_{65}) R_6 = 1/60,$$

$$(a_3 + a_4) c_1 b_{43} R_4 + [(a_3 + a_5) c_1 b_{53} + (a_4 + a_5) c_2 b_{54}] R_5$$

$$+ [(a_3 + a_6) c_1 b_{63} + (a_4 + a_6) c_2 b_{64} + (a_5 + a_6) c_3 b_{65}] R_6 = 7/120$$

$$c_1 b_{43} b_{54} R_5 + [c_1 b_{43} b_{64} + (c_1 b_{53} + c_2 b_{54}) b_{65}] R_6 = 1/120,$$

(3.49a)

where

$$c_i = \sum_{j=2}^{i+1} a_j b_{i+2, j}, \quad d_i = \sum_{j=2}^{i+1} a_j^2 b_{i+2, j}. \quad (3.49b)$$

simplification requires that

$$R_2 = 0 \quad (3.50)$$

$$\text{and } \frac{a_i^2}{2} = c_{i-2}, \quad 3 \leq i \leq 6 \quad (3.51)$$

then Eq. (3.49) can be simplified considerably. After eliminating duplicates and combining some equations solve, in addition to (3.51),

$$R_1 + R_3 + R_4 + R_5 + R_6 = 1, \quad (3.52a)$$

and

$$a_3 R_3 + a_4 R_4 + a_5 R_5 + a_6 R_6 = 1/2,$$

$$\frac{a_3^2}{3} R_3 + \frac{a_4^2}{4} R_4 + \frac{a_5^2}{5} R_5 + \frac{a_6^2}{6} R_6 = 1/3,$$

$$\frac{a_3^3}{3} R_3 + \frac{a_4^3}{4} R_4 + \frac{a_5^3}{5} R_5 + \frac{a_6^3}{6} R_6 = 1/4, \quad (3.57b)$$

$$\frac{a_3^4}{3} R_3 + \frac{a_4^4}{4} R_4 + \frac{a_5^4}{5} R_5 + \frac{a_6^4}{6} R_6 = 1/5,$$

and

$$a_3 b_{43} R_4 + (a_3 b_{53} + a_4 b_{54}) R_5 + (a_3 b_{63} + a_4 b_{64} + a_5 b_{65}) R_6 = 1/6$$

$$a_3^2 b_{43} R_4 + (a_3^2 b_{53} + a_4^2 b_{54}) R_5 + (a_3^2 b_{63} + a_4^2 b_{64} + a_5^2 b_{65}) R_6 = 1/12$$

$$a_3^3 b_{43} R_4 + (a_3^3 b_{53} + a_4^3 b_{54}) R_5 + (a_3^3 b_{63} + a_4^3 b_{64} + a_5^3 b_{65}) R_6 = 1/20,$$

$$a_4 a_3 b_{43} R_4 + a_5 (a_3 b_{53} a_4 b_{54}) R_5 + a_6 (a_3 b_{63} + a_4 b_{64} + a_5 b_{65}) R_6 \\ = 1/8,$$

$$a_4^2 a_3 b_{43} R_4 + a_5 (a_4^2 b_{53} + a_4^2 b_{54}) R_5 + a_6 (a_3^2 b_{63} + a_4^2 b_{64} + a_5^2 b_{65}) R_6 \\ = 1/15,$$

(3.52c)

and

$$a_3 b_{43} b_{54} R_5 + [a_3 b_{43} b_{64} + (a_3 b_{53} + a_4 b_{54}) b_{65}] R_6 = 1/24,$$

$$a_3^2 b_{43} b_{54} R_5 + [a_3^2 b_{43} b_{64} + (a_3^2 b_{53} + a_4^2 b_{54}) b_{65}] R_6 = 1/60,$$

(3.52d)

The situation is now as follows. Equations (3.52b), (3.52c), and (3.52d) are to be solved independently. Then (3.52a) yields R_1 . Equation (3.51) are used to find b_{32} , b_{42} , b_{52} , b_{62} . Then equations (3.48) determine b_{21} , b_{31} , b_{41} , b_{51} and b_{61} . This, with $R_2=0$, completes the solution. The family of solutions due to Kutta may now be found by taking $b_{65}=0$. From (3.52d), $a_3=2/5$. It then develops that the third equation in (3.52c) has the left member equal to $-a_3 a_4$ times the left member of the first of

this group, plus $a_3 + a_4$ times the left member of the second of this group. This forces a_4 to be 1. Equations (3.52c) may now be solved for $b_{64}R_6$ and $b_{54}R_5$. When the results are substituted in $a_3b_{43} (b_{54}R_5 + b_{64}R_6) = 1/24$ (the consequence of (3.52d) and $b_{65}=0$) it is found that $b_{43} = 15/4$. In summary, there results the following description of a family of solutions:

$$R_2 = 0, \quad b_{65} = 0, \quad b_{43} = 15/4, \quad a_3 = 2/5, \quad a_4 = 1, \quad (3.53)$$

with

$$\begin{aligned} R_3 &= 125 \left[10a_5a_6 - 5(a_5 + a_6) + 3 \right] / \left[72(2 - 5a_5)(2 - 5a_6) \right], \\ R_4 &= \left[8a_5a_6 - 7(a_5 + a_6) + 6 \right] / \left[36(1 - a_5)(1 - a_6) \right], \\ R_5 &= \left[1 - a_6 \right] / \left[12a_5(1 - a_5)(5a_5 - 2)(a_5 - a_6) \right], \\ R_6 &= \left[1 - a_5 \right] / \left[12a_6(1 - a_6)(5a_6 - 2)(a_6 - a_5) \right], \\ R_1 &= 1 - R_3 - R_4 - R_5 - R_6, \end{aligned} \quad (3.53)$$

and

$$\begin{aligned} b_{53} &= 5 \left[7 - 10a_6 - 108R_4(1 - a_6) \right] / \left[144R_5(a_5 - a_6) \right], \\ b_{54} &= \left[1 - a_6 \right] / \left[36R_5(a_5 - a_6) \right], \\ b_{63} &= 5 \left[7 - 10a_5 - 108R_4(1 - a_5) \right] / \left[144R_6(a_6 - a_5) \right], \\ b_{64} &= \left[1 - a_5 \right] / \left[36R_6(a_6 - a_5) \right], \end{aligned} \quad (3.53c)$$

and

$$b_{32} = 2 / | 25a_2 | ,$$

$$b_{42} = 1/a_2 ,$$

$$b_{52} = [5a_5^2 - 4b_{53} - 10b_{54}] / [10a_2] \quad (3.53d)$$

$$b_{62} = [5a_6^2 - 4b_{63} - 10b_{64}] / [10a_2] ,$$

and

$$b_{j1} = a_j - \sum_{i=2}^{j-1} b_{ji} , \quad 2 \leq j \leq 6. \quad (3.53e)$$

It becomes very simple to construct a fifth order Runge-Kutta formula based on the coefficients derived from equation (3.53). The result is

$$y_{n+1} = y_n + \left\{ 4 \Delta^1 y + (16 + \sqrt{6}) \Delta^{(5)} y + (16 - \sqrt{6}) \Delta^{(6)} y \right\} / 36 ,$$

$$\Delta^1 y = Tf(t_n, y_n) ,$$

$$\Delta^2 y = Tf(t_n + 4T/11, y_n + 4\Delta^1 y/11) , \quad (3.54)$$

$$\Delta^3 y = Tf(t_n + 2T/5, y_n + \left\{ 9 \Delta^1 y + 11 \Delta^2 y \right\} / 50) ,$$

$$\Delta^{(4)} y = Tf(t_n + T, y_n + \left\{ -11 \Delta^2 y + 15 \Delta^3 y \right\} / 4) .$$

$$\begin{aligned} \Delta^{(5)} y = Tf(t_n + (6 - \sqrt{6})T/10, y_n + \left\{ (81 + 9\sqrt{6}) \Delta^1 y \right. \\ \left. + (255 - 55\sqrt{6}) \Delta^2 y + (24 - 14\sqrt{6}) \Delta^{(4)} y \right\} / 600) , \end{aligned}$$

$$\begin{aligned} \Delta^{(6)} y = Tf(t_n + (6 + \sqrt{6})T/10, y_n + \left\{ (81 - 9\sqrt{6}) \Delta^1 y \right. \\ \left. + (255 + 55\sqrt{6}) \Delta^2 y + (24 + 14\sqrt{6}) \Delta^{(4)} y \right\} / 6000) , \end{aligned}$$

Let $b_{43} = 0$ instead of letting $b_{65} = 0$, then Eqs. (3.52d) yields

$$(a_3 b_{53} + a_4 b_{54}) b_{65} R_6 = 1/24,$$

$$(a_3 b_{53} + a_4 b_{54}) b_{65} R_6 = 1/60,$$

Using these equations in conjunction with the first, second, fourth, and fifth equations of (3.52c), the terms in b_{63} , b_{64} , and b_{65} can be eliminated and there result two equations:

$$R (a_6 - a_5) = (4a_6 - 3) b_{65} R_6,$$

$$R (a_6 - a_5) = (5a_6 - 4) b_{65} R_6.$$

These lead to $a_6 = 1$. Another family of solution occurs with

$$R_2 = 0, \quad b_{43} = 0, \quad a_6 = 1, \quad (3.55a)$$

and

$$R_3 = \left[3 - 5(a_4 + a_5) + 10a_4 a_5 \right] / \left[60a_3(a_4 - a_3)(a_3 - a_5)(a_3 - 1) \right],$$

$$R_4 = \left[3 - 5(a_3 + a_5) + 10a_3 a_5 \right] / \left[60a_4(a_3 - a_4)(a_4 - a_5)(a_4 - 1) \right],$$

$$R_5 = \left[3 - 5(a_3 + a_4) + 10a_3 a_4 \right] / \left[60a_5(a_3 - a_5)(a_5 - a_4)(a_5 - 1) \right],$$

$$R_6 = \left[12 - 15(a_3 + a_4 + a_5) + 20(a_3 a_4 + a_4 a_5 + a_5 a_3) - 30a_3 a_4 a_5 \right] / \left[60(1 - a_3)(1 - a_4)(1 - a_5) \right],$$

$$R_1 = 1 - R_3 - R_4 - R_5 - R_6, \quad (3.55b)$$

and

$$\begin{aligned} b_{53} &= [2 - 5a_4] / [120R_5a_3(1 - a_5)(a_3 - a_4)], \\ b_{54} &= [2 - 5a_3] / [120R_5a_4(1 - a_5)(a_4 - a_3)], \\ b_{63} &= [6 - 2a_3 - 10a_4 - 14a_5 + 5a_3a_4 + 25a_4a_5 + 10a_5^2 \\ &\quad - 20a_5^2a_4] / [120R_6a_3(1 - a_5)(a_3 - a_4)(a_3 - a_5)], \quad (3.55c) \\ b_{64} &= [6 - 2a_4 - 10a_3 - 14a_5 + 5a_3a_4 + 25a_3a_5 + 10a_5^2 \\ &\quad - 20a_5^2a_3] / [120R_6a_4(1 - a_5)(a_4 - a_3)(a_4 - a_5)], \\ b_{65} &= [3 - 5a_3 - 5a_4 + 10a_3a_4] / [60R_6a_5(a_5 - a_3)(a_5 - a_4)], \end{aligned}$$

and

$$\begin{aligned} b_{32} &= a_3^2 / [2a_2], \\ b_{42} &= a_4^2 / [2a_2], \quad (3.55d) \\ b_{52} &= [a_5^2 - 2a_3b_{35} - 2a_4b_{54}] / [2a_2], \\ b_{62} &= [1 - 2a_3b_{63} - 2a_4b_{64} - 2a_5b_{65}] / [2a_2], \end{aligned}$$

and

$$b_{j1} = a_j - \sum_{i=2}^{j-1} b_{ji}, \quad 2 \leq j \leq 6 \quad (3.55e)$$

Thus a formula can be constructed as follows:

$$y_{n+1} = y_n + \left\{ \Delta^0 y + 5 \Delta^1 y + 5 \Delta^2 y + \Delta^3 y \right\} / 12, \quad (3.56)$$

$$k_1 = Tf(t_n, y_n),$$

$$k_2 = Tf(t_n + T/2, y_n + k_1/2),$$

$$k_3 = Tf(t_n + (5 - \sqrt{5})T/10, y_n + \left\{ 2k_1 + (3 - \sqrt{5})k_2 \right\} / 10),$$

$$k_4 = Tf(t_n + T/2, y_n + k_1 + k_2/4),$$

$$k_5 = Tf(t_n + (5 + \sqrt{5})T/10, y_n + \left\{ (1 - \sqrt{5})k_1 - 4k_2(5 + 3\sqrt{5})k_3 + 8k_4 \right\} / 20,$$

$$k_6 = Tf(t_n + T, y_n + (5 - 1)k_1 + (2\sqrt{5} - 2)k_2 + (5 - \sqrt{5})k_3 - 8k_4 + (10 - 2\sqrt{5})k_5 / 4.$$

This fifth order formula is not in Kutta's family. It belongs to Lobatto quadrature formulas which have errors of order T^7 rather than T^6 .

The accumulated truncation error of Runge-Kutta's method can be calculated as follows:

Let

$$z_{n+1} = z_n + T \bar{\Delta} z(t_n, z_n; T) \quad (3.57)$$

be the exact solution of Eq. (3.46),

$$\text{and } y_{n+1} = y_n + T \Delta y(t_n, y_n; T) \quad (3.58)$$

be the calculated value of the solution

Substituting Equation (3.57) from Equation (3.58) yields

$$e_{n+1} = e_n + T [\Delta y(t_n, y_n; T) - \bar{\Delta} z(t_n, z_n; T)] \quad (3.59)$$

by application of the triangle inequality, Equation (3.59) becomes

$$\begin{aligned} \|e_{n+1}\| &\leq \|e_n\| + T \|\Delta y(t_n, y_n; T) - \Delta z(t_n, z_n; T)\| \\ &\quad + T \|\Delta z(t_n, z_n; T) - \bar{\Delta} z(t_n, z_n; T)\| \end{aligned} \quad (3.60)$$

The Lipschitz condition yields

$$\begin{aligned} \|\Delta y(t_n, y_n; T) - \Delta z(t_n, z_n; T)\| &\leq L \|y_n - z_n\| \\ &= L \|e_n\| \end{aligned} \quad (3.61)$$

$$\text{and } \|\Delta z(t_n, z_n; T) - \bar{\Delta} z(t_n, z_n; T)\| \leq N(T^p) \quad (3.62)$$

where

$$N = \frac{1}{(p+1)!} \max \|f^{(p)}(y)\|$$

and p is the order of the Runge-Kutta formula

Hence from Equation (3.60) the required truncation error is

$$\| e_{n+1} \| \leq (1 + LT) \| e_n \| + T^{(p+1)} L \quad (3.63)$$

The value of L can be determined from a Runge-Kutta formula and different formulas correspond to different values of L . All L 's have upper bounds.

The Runge-Kutta method seems to be tedious because values of $f(t,y)$ have to be calculated a number of times per time increment. However, the formulas are systematic and hence can be easily programmed on an automatic machine. No special starting procedure is required and calculations can often be checked by repetition using a different step size. Furthermore, such a method is particularly useful if certain coefficients in the differential equation are empirical formulas for which analytical expressions are not known. The step size can be altered as desired. It should be understood that the derivatives as evaluated by Runge-Kutta process are not the actual derivatives at the various points within the step as commonly assumed.

2). Multiple step Methods.

The one-step methods are necessary for obtaining initial values in the solution of a differential equation. However, they involve too much labor to be used for obtaining a numerical solution over an extended range. This can be offset by using multiple step methods.

Multiple step methods are always expressed in the difference-differential equation form

$$\alpha_k y_{n+k} + \alpha_{k-1} y_{n+k-1} + \dots + \alpha_0 y_n = T \left\{ \beta_k f_{n+k} + \beta_{k-1} f_{n+k-1} + \dots + \beta_0 f_n \right\}$$

$$n = 0, 1, 2, \dots \quad (3.64)$$

where k is a fixed integer, $f_m = f(t_m, y_m)$ ($m = 0, 1, 2, \dots$), and where α_μ and β_μ ($\mu = 0, 1, 2, \dots$) are real constants which do not depend on n . Any α_k is always assumed unequal to zero. Equation (3.64) defines the general linear k -step method. If $\beta_k = 0$, the formula (3.64) is called "closed"; otherwise it is called open.

Unlike one-step methods, multiple step methods are not self-starting; if some values $y_{n+k-1}, y_{n+k-2}, \dots, y_n$ are not known, these methods break down. Such is the case at the beginning of the computation, where the initial condition furnishes only one of the required $k+1$ values, or at places where the step T is changed.

Stability and convergence are two important factors that affect the availability of a multiple step method. To insure its stable and convergence, two rules must be fulfilled;

1. The characteristic polynomial

$$\alpha_k S^k + \alpha_{k-1} S^{k-1} + \dots + \alpha_0 = 0$$

of Equation (3.75) has no root with modulus exceeding 1. and the root of modulus 1 must be simple.

II. The order of the associated difference operator be at least 1.

A. Adams-Bashforth Method

Consider the initial value problem

$$y' = f(t, y) \quad y(t_0) = y_0 \quad (3.65)$$

An exact solution of the differential equation (3.65) by definition satisfies the identity

$$y(t+k) - y(t) = \int_t^{t+k} f(t, y) dt \quad (3.66)$$

for any two values of t in the interval $[a, b]$. Replace $f(t, y)$ in the right-hand side of Equation (3.66) by an interpolating polynomial on a set of points t_n where y_n has already been computed or is just about to be computed. Equate the integral and accept its value as the increment of the approximate values y_n between t and $t+k$. If it is assumed that the interpolating points are $t_p, t_{p-1}, t_{p-2}, \dots, t_{p-q}$, then the polynomial replacing $f(t, y)$ is given by

$$p(t) = \sum_{m=0}^q (-1)^m \binom{-s}{m} \nabla^m f_p \quad (3.67)$$

$$s = \frac{t - t_p}{T}$$

where q is an arbitrary integer. This formula is developed in Appendix.

Equation (3.67) can be substituted into Equation (3.66) to yield

$$y_{p+1} - y_p \approx \int_{t_p}^{t_{p+1}} p(t) dt = T \sum_{m=0}^q r_m \nabla^m f_p \quad (3.68)$$

where

$$r_m = (-1)^m \frac{1}{T} \int_{t_p}^{t_{p+1}} \binom{-s}{m} ds \quad (3.69)$$

and $y_p = y(t_p)$

Construct a generating function $G(k)$ as follows

$$\begin{aligned} G(k) &= \sum_{m=0}^{\infty} r_m k^m = \sum_{m=0}^{\infty} (-1)^m \int_{t_p}^{t_{p+1}} \binom{-s}{m} ds \\ &= \int_{t_p}^{t_{p+1}} \sum_{m=0}^{\infty} (-1)^m \binom{-s}{m} ds = \int_{t_p}^{t_{p+1}} (1-k)^{-s} ds \quad (3.70) \end{aligned}$$

The identity $(1-k)^{-s} = \exp(-s \log(1-k))$ causes Eq. (3.70) to become

$$G(k) = - \frac{k}{(1-k) \log(1-k)}$$

this may be written as

$$- G(k) \frac{\log(1 - k)}{k} = \frac{1}{1 - k}.$$

Since
$$\frac{1}{1 - k} = 1 + k + k^2 + \dots$$

and
$$-\frac{\log(1 - k)}{k} = 1 + \frac{1}{2}k + \frac{1}{3}k^2 + \dots$$

one can conclude that

$$(r_0 + r_1 k + r_2 k^2 + \dots) \left(1 + \frac{1}{2}k + \frac{1}{3}k^2 + \dots\right) = 1 + k + k^2 + \dots$$

By comparing the coefficients of corresponding powers of k , a relation can be found

$$r_m + \frac{1}{2}r_{m-1} + \frac{1}{3}r_{m-2} + \dots + \frac{1}{m+1}r = 1 \quad (3.71)$$

thus it is possible to calculate r_m recursively. Some values of r_m calculated from Eq. (3.71) are

m	0	1	2	3	4	5	6
r_m	1	$\frac{1}{2}$	$\frac{5}{12}$	$\frac{3}{8}$	$\frac{251}{720}$	$\frac{95}{288}$	$\frac{19087}{60480}$

If $y^{(q+2)}(t)$ is continuous in $[a, b]$, then $y'(t)$ can be expressed as

$$y' = \sum_{m=0}^q (-1)^m \binom{-s}{m} {}^m y'(t_p) + (-1)^{(q+1)} \binom{-s}{q+1} T^{q+1} y^{(q+2)}(\xi)$$

where ξ is a point between the largest and the smallest of the values t, t_p , and t_{p-q} . Integrating y' between t_p and t_{p+1} , we get

$$y(t_{p+1}) - y(t_p) = T \sum_{m=0}^q r_m \nabla^m f_p + R_q^{AB} \quad (3.72)$$

where

$$R_q^{AB} = (-1)^{(q+1)} T^{(q+1)} \int_{t_p}^{t_{p+1}} \binom{-s}{q+1} y^{(q+2)}(\xi) dt$$

since $\binom{-s}{q+1}$ is a constant sign in the interval $t_p < t < t_{p+1}$ and $y^{(q+2)}(\xi)$ is a continuous function of t , apply the second mean value theorem of the integral calculus

$$R_q^{AB} = (-1)^{(q+1)} T^{(q+1)} y^{(q+2)}(\xi') \int_{t_p}^{t_{p+1}} \binom{-s}{q+1} dt$$

where $t_{p-q} < \xi' < t_{p+1}$. By definition of r_{q+1} , this may be written

$$R_q^{AB} = T^{(q+2)} y^{(q+2)} \left(\frac{\tau}{T}\right) r_{q+1} \quad (3.73)$$

This is the desired expression for the remainder of the Adams-Bashforth formula.

B. Adams-Moulton Method.

This method uses the form

$$y_p - y_{p-1} = \int_{t_{p-1}}^{t_p} p(t) dt = T \sum_{m=0}^q r_m^* \nabla^m f \quad (3.74)$$

where

$$r_m^* = (-1)^m \frac{1}{T} \int_{t_{p-1}}^{t_p} \binom{-s}{m} dt = (-1)^m \frac{1}{T} \int_{-1}^0 \binom{-s}{m} ds \quad (3.75)$$

the generating function of the coefficients is determined as follows:

$$G^*(k) = \sum_{m=0}^{\infty} r_m^* k^m = - \frac{k}{\log(1-k)} \quad (3.76)$$

or

$$- \frac{\log(1-k)}{k} G^*(k) = 1 \quad (3.77)$$

Expand the left-hand side of Eq. (3.77) to power series

$$\left(1 + \frac{1}{2}k + \frac{1}{3}k^2 + \dots\right)(r_0^* + r_1^*k + r_2^*k^2 + \dots) = 1 \quad (3.78)$$

It follows that

$$r_m^* + \frac{1}{2}r_{m-1}^* + \frac{1}{3}r_{m-2}^* + \dots + \frac{1}{m+1}r_0^* = \begin{cases} 1, & m=0 \\ 0, & m=1,2,3,\dots \end{cases} \quad (3.79)$$

The numerical values of r_m^* are easily found from this recurrence relation. Some values of r_m^* calculated from Eq. (3.79) are

m	0	1	2	3	4	5	6
r_m^*	1	$-\frac{1}{2}$	$-\frac{1}{12}$	$-\frac{1}{24}$	$-\frac{19}{720}$	$-\frac{3}{160}$	$-\frac{863}{60480}$

Since y_p occurs as an argument in $f_p = f(t_p, y_p)$ in the right hand side term of Eq. (3.74), it will not be possible to solve this equation explicitly. A better approach to this solution is by means of an iterative procedure.

Assuming an approximation of a solution of Eq. (3.74) has been obtained to be $y_p^{(0)}$, calculate $f^{(0)} = f(t_p, y_p^{(0)})$ and form the difference

$$\nabla f_p^{(0)} = f_p^{(0)} - f_{p-1}, \quad \nabla^2 f_p^{(0)} = \nabla f_p^{(0)} - \nabla f_{p-1}, \dots \quad \text{A better approxima-}$$

tion is then obtained from

$$y_p^{(1)} = y_{p-1} + T \sum_{m=0}^q r_m \nabla^m f_p^{(0)} \quad (3.80)$$

Calculating $f_p^{(1)} = f(t_p, y_p^{(1)})$ and re-evaluating the differences, a still better value $y_p^{(2)}$ is

$$y_p^{(2)} = y_{p-1} + T \sum_{m=0}^q r_m^* \nabla^m f_p^{(1)} \quad (3.81)$$

Generally a sequence $y^{(r)}$ ($r = 0, 1, 2, 3, \dots$) of approximations is obtained recursively from the relation

$$y_p^{(r)} = y_{p-1} + T \sum_{m=0}^q r_m^* \nabla^m f_p^{(r-1)} \quad (3.82)$$

Since

$$\begin{aligned} \nabla^m f_p^{(i)} &= \nabla^{m-1} f_p^{(i)} - \nabla^{m-1} f_p^{(i-1)}, & (m = 1, 2, \dots) \\ & & (i = 1, 2, \dots) \end{aligned} \quad (3.83)$$

it follows that

$$\nabla^m f_p^{(r)} - \nabla^m f_p^{(r-1)} = f_p^{(r)} - f_p^{(r-1)} \quad (3.84)$$

Thus

$$\begin{aligned}
 y_p^{(r+1)} - y_p^{(r)} &= T \sum_{m=0}^q r_m^* (\nabla^m f_p^{(r)} - \nabla^m f_p^{(r-1)}) \\
 &= T \sum_{m=0}^q r_m^* (f_p^{(r)} - f_p^{(r-1)})
 \end{aligned} \tag{3.85}$$

The Lipschitz condition yields

$$|f_p^{(r)} - f_p^{(r-1)}| \leq L |y_p^r - y_p^{r-1}| \tag{3.86}$$

Equation (3.85) becomes

$$y_p^{(r+1)} - y_p^{(r)} \leq T \sum_{m=0}^q r_m^* L |y_p^{(r)} - y_p^{(r-1)}| \tag{3.87}$$

or

$$y_p^{(r+1)} - y_p^{(r)} \leq (TL A)^r |y_p^{(1)} - y_p^{(0)}| \tag{3.88}$$

where

$$A = \sum_{m=0}^q r_m^*$$

The solution of y_p of Eq. (3.74) is now obtained by summing terms of the series

$$\begin{aligned}
 y_p^* = & y_p^{(0)} + (y_p^{(1)} - y_p^{(0)}) + (y_p^{(0)} - y_p^{(1)}) + \dots \\
 & + (y_p^{(n)} - y_p^{(n-1)}) + \dots \quad (3.89)
 \end{aligned}$$

The series on the right-hand side of Eq.(3.89) will converge absolutely provided $0 \leq |TLA| < 1$. In this case the solution y_p^* will exist and be unique.

The local truncation error is

$$R_q^{AM} = T^{(q+2)} y^{(q+2)}(\xi) r_{q+1}^* \quad (3.90)$$

where

$$t_{p-q} < \xi < t_p.$$

C. Milne's Method

Milne's method requires predictor and corrector formulas. The predictor formula is of the form

$$y_{n+1} - y_{n-3} = \int_{t_{n-3}}^{t_{n+1}} f(t, y) dt = \frac{4T}{3} (2f_{n-2} - f_{n-1} + 2f_n) \quad (3.91)$$

and the corrector is of the form

$$y_{n+1} - y_{n-1} = \int_{t_{n-1}}^{t_{n+1}} f(t, y) dt \approx \frac{T}{3} (f_{n-1} + 4f_n + f_{n+1}) \quad (3.92)$$

The predictor formula has truncation error

$$\text{Truncation error} = \frac{14}{45} T^5 y^{(5)}(\xi_1) \quad (3.93)$$

where

$$t_n < \xi_1 < t_{n+1}$$

The corrector formula has truncation error

$$\text{Truncation error} = \frac{T^5}{90} y^{(5)}(\xi_2) \quad (3.94)$$

where

$$t_{n-1} < \xi_2 < t_{n+1}$$

The procedure is as follows:

Step 1; Takes T small enough to insure that the remainder term involving is small in the predictor formula, then find out y_{n+1} from this formula.

Step 2; Obtain a first approximation to y_{n+1}^i by substituting the value of y_{n+1} obtained from step 1 into the following equation

$$y^i = f(t, y) \quad y = y_0 \text{ at } t = t_0 .$$

Step 3; Obtain a better approximation of y_{n+1} by means of the corrector formula Eq. (3.94).

After repeating step 2 and step 3, the value of y_{n+1} becomes very accurate. When values of y_{n+1} and f_{n+1} have negligible error, the next pair of values y_{n+2} and f_{n+2} may be obtained by a repetition of the process.

Let $y_{n+1}^{(0)}$ be the value of y_{n+1} obtained from step 1, $y_{n+1}^{(1)}$ be the value of y_{n+1} introduced by the corrector at first time, and $y_{n+1}^{(m)}$ be the value of y_{n+1} after introducing the corrector m times, then

$$f_{n+1}^{(1)} - f_{n+1}^{(0)} = k(y_{n+1}^{(1)} - y_{n+1}^{(0)}) \quad (3.95)$$

where

$$k = \frac{f_{n+1}^{(1)} - f_{n+1}^{(0)}}{y_{n+1}^{(1)} - y_{n+1}^{(0)}} \quad (3.96)$$

$$\text{and } f_{n+1}^{(0)} = f(t_{n+1}, y_{n+1}^{(0)}), \quad f_{n+1}^{(1)} = f(t_{n+1}, y_{n+1}^{(1)}) \quad (3.97)$$

If $f(t, y)$ possesses a continuous first derivative $f_y(t, y)$ with respect to y , then by the mean-value theorem

$$k = f_y(t_{n+1}, \eta_{n+1}) \quad (3.98)$$

where

$y_{n+1}^{(2)}$ lies between $y_{n+1}^{(0)}$ and $y_{n+1}^{(1)}$

Equation (3.92) together with Equation (3.95) yield

$$y_{n+1}^{(2)} - y_{n+1}^{(1)} = \frac{T}{3} [f_{n+1}^{(1)} - f_{n+1}^{(0)}] = [k(y_{n+1}^{(1)} - y_{n+1}^{(0)})] \quad (3.99)$$

By the same method

$$y_{n+1}^{(3)} - y_{n+1}^{(2)} = \left(\frac{T}{3} k\right) \left(\frac{T}{3} k\right) (y_{n+1}^{(1)} - y_{n+1}^{(0)})$$

where the change of k is negligible.

Proceeding in this fashion a sequence

$$y_{n+1}^{(r)} - y_{n+1}^{(r-1)} = \left(\frac{T}{3} k\right)^{(r-1)} (y_{n+1}^{(1)} - y_{n+1}^{(0)})$$

is obtained.

Thus

$$\begin{aligned} y_{n+1}^* &= y_{n+1}^{(0)} + (y_{n+1}^{(1)} - y_{n+1}^{(0)}) + (y_{n+1}^{(2)} - y_{n+1}^{(1)}) + \dots \\ &= y_{n+1}^{(0)} + (y_{n+1}^{(1)} - y_{n+1}^{(0)}) \left[1 + \frac{T}{3} k + \left(\frac{T}{3} k\right)^2 + \dots \right] \quad (3.100) \end{aligned}$$

Provided that T is sufficiently small, the value of $\frac{T}{3} k$ can be chosen so that $0 \leq \left| \frac{T}{3} k \right| < 1$ to insure convergence of the power series on the right-hand side of Eq. (3.100).

If iterations are finite, the value of \tilde{y}_{n+1} obtained will differ from y_{n+1}^* by

$$y_{n+1}^* - \tilde{y}_{n+1} = (y_{n+1}^{(1)} - y_{n+1}^{(0)}) \left(\frac{T}{3} k \right)^{n+1} \left(\frac{3}{3 - Tk} \right) \quad (3.101)$$

where n is the number of steps in the finite step process.

Let the true values of the \tilde{y}_i and \tilde{y}_i^0 be z_i and z_i^0 respectively. Define errors in the \tilde{y}_i and \tilde{y}_i^0 by the equations

$$z_i = \tilde{y}_i + e_i \quad (3.102)$$

$$z_i^0 = \tilde{y}_i^0 + e_i^0$$

From the differential equation (3.64)

$$e_i^0 = z_i^0 - \tilde{y}_i^0 = f(t_i, z_i) - f(t_i, \tilde{y}_i) = k_i (z_i - \tilde{y}_i) \quad (3.103)$$

where $k_i = f(t_i, \eta_i)$ and $y_i < \eta_i < z_i$

The z_i and z_i^0 satisfy the equation

$$z_{n+1} - z_{n-1} = \frac{T}{3} (z_{n-1}^3 + 4z_n^3 + z_{n+1}^3) - T^5 z^{(5)} (\xi_{n+1}) / 90$$

where

$$t_n < \xi_{n+1} < t_{n+1} \quad (3.104)$$

Subtracting Equation (3.94) from Equation (3.104) results in

$$e_{n+1} - e_{n-1} = \frac{T}{3} (e_{n-1}^3 + 4e_n^3 + e_{n+1}^3) - T^5 A_{n+1} \quad (3.105)$$

where

$$z^{(5)} (\xi_{n+1}) / 90 = A_{n+1} \quad t_n < \xi_{n+1} < t_{n+1}$$

On making use of Equation (3.103), this difference equation may be written

$$\left(1 - \frac{1}{3} Tk_{n+1}\right) e_{n+1} = \left(\frac{4}{3} Tk_n\right) e_n + \left(1 + \frac{1}{3} Tk_{n-1}\right) e_{n-1} - T^5 A_{n+1} \quad (3.106)$$

Assume $\begin{cases} k_i = k, \\ A_i = A \end{cases}$ for $i = 1, 2, \dots, n$, over a sufficiently restricted range.

Solving the difference equation yields

$$e_n = C_1 \lambda_1^n + C_2 \lambda_2^n + \frac{T^5 A}{\left(6\frac{1}{3} Tk\right)} \quad (3.107)$$

where

$$\lambda_1 = \frac{\frac{2}{3} Tk + \sqrt{1 + \frac{1}{3} T^2 k^2}}{1 - \frac{1}{3} Tk} \approx 1 + Tk \approx e^{3\theta} \quad (3.108)$$

$$\lambda_2 = \frac{\frac{2}{3} Tk - \sqrt{1 + \frac{1}{3} T^2 k^2}}{1 - \frac{1}{3} Tk} \approx - \left(1 - \frac{Tk}{3}\right) \approx -e^\theta \quad (3.109)$$

$$\theta = \frac{1}{3} Tk .$$

Expressing the constants c_1 and c_2 in terms of e_1 and e_0 and setting e_0 to be zero, Equation (3.107) becomes

$$e_n \approx \frac{e_1}{2} (\lambda_1^n - \lambda_2^n) + \frac{T^5 A}{6\theta} \left(1 - \lambda_1^n - \frac{3}{2} \theta \lambda_2^n\right) \quad (3.110)$$

As an approximation setting $\epsilon_1=0$, Eq.(3.110) can be replaced by

$$e_n \approx \frac{T^5 A}{2k} \left[1 - e^{k(t_n - t_0)} + (-1)^{n+1} \frac{1}{2} T k e^{-\frac{2}{3} k(t_n - t_0)} \right] \quad (3.111)$$

If k is positive, the last term decreases as t_n increases. In this case the error is always of the same sign and is multiplied by 10 each time t_n increases by $2.3/k$. This case corresponds to a differential equation in

which the plots of the solutions for various boundary conditions, in the neighborhood of the desired solution, diverge to the right, and the accumulated relative error may decrease even though the accumulated error increases exponentially. Milne's method is thus applicable.

If k is negative, the second term vanishes, and the error is alternately positive and negative. This case corresponds to a differential equation in which the plots of the solutions for various boundary conditions converge to the right. The accumulated relative error increases rapidly without bound. This shows that Milne's method is unstable, and cannot be used.

Milne's method has two virtues: It supplies a running check that the method and interval size are suitable and that the computation is locally accurate enough to warrant going on. Moreover, only two evaluations of the derivatives per step forward are required.

D. Hamming Method

Milne's method is always unstable, because it has two roots with modulus one in the characteristic equation of its corrector formula. Hamming has removed instability by eliminating one of these roots and thus modified the corrector formula without changing the predictor formula.

D.I. Third Order formula

The third order formula is of the form

$$y_{n+1} = ay_n + by_{n-1} + T(cy'_{n+1} + dy'_n + ey'_{n-1}) \quad (3.112)$$

Expanding both sides with respect to y_n by Taylor series expansion and

requiring exact fit for $1, T, T^2, T^3$, the coefficients a, b, c, d , and e can be determined to be

$$\begin{aligned} a &= -4 + 12c & d &= 4 - 8c \\ b &= 5 - 12c & e &= 2 - 5c \\ c &= c \end{aligned} \quad (3.113)$$

The truncation error has the value

$$\text{Truncation error} = \frac{1 - 3c}{6} T^4 y^{(4)}$$

The characteristic equation is

$$f^2 + (4 - 12c)f - (5 + 12c) = 0 \quad (3.114)$$

Solving this equation yields

$$\begin{aligned} f_1 &= 1 \\ f_2 &= -5 + 12c \end{aligned}$$

Stability requires

$$0 \leq |f_2| < 1$$

or

$$\frac{1}{2} > c > \frac{1}{3}$$

Setting $c = 5/12$, Equation (3.112) becomes

$$y_{n+1} = y_n + \frac{T}{12} (5y'_{n+1} + 8y'_n - y'_{n-1}) \quad (3.115)$$

D. II. Fourth Order Formulas

This method generalizes Milne's corrector formula to the form

$$y_{n+1} = ay_n + by_{n-1} + cy_{n-2} + T [dy'_n + ey'_n + fy'_{n-1}] \quad (3.116)$$

and then stabilizes this formula by choosing suitable values of the coefficients to minimize one of the characteristic roots with modulus one in Milne's formula.

Expanding Equation (3.116) with respect y_n by Taylor series expansion, and requiring exact fit for $1, T, T^2, T^3, T^4$, yield

$$\begin{aligned} a &= \frac{27(1-b)}{24} & d &= \frac{9-b}{24} \\ b &= b & e &= \frac{18+14b}{24} \\ c &= \frac{-3(1-b)}{24} & f &= \frac{-9+17b}{24} \end{aligned} \quad (3.117)$$

In this case the truncation error is found to be

$$kT^5 y^{(5)} = \frac{-9 + 5b}{360} T^5 y^{(5)}. \quad (3.118)$$

The characteristic equation of Equation (3.116), using Equation (3.117), is

$$8\rho^3 - 9(1 - b)\rho^2 - 8b\rho + (1 - b) = 0. \quad (3.119)$$

The root loci of ρ with respect to b are shown in Fig. 2. For stable operation b should lie in the range $-0.6 < b < 1$ to insure that no characteristic root exceeds one and that the root with modulus one is simple. Milne's formula is a special case of Equation (3.116) with $b=1$. It is easily seen from Fig. 2 that in this case there are two characteristic roots with modulus one.

For the case $b=0$, a stable predictor-corrector formula can be constructed as follows:

$$\begin{aligned} \text{predictor} \quad p_{n+1} &= y_{n-3} + \frac{4T}{3} (2y'_n - y'_{n-1} + 2y'_{n-2}) \\ \text{modify} \quad m_{n+1} &= p_{n+1} - \frac{112}{121} (p_n - c_n) \\ \text{corrector} \quad c_{n+1} &= \frac{1}{8} (9y_n - y_{n-2}) + 3T(f'_{n+1} - 2y'_n - y'_{n-1}) \\ \text{final value} \quad y_{n+1} &= c_{n+1} + \frac{9}{121} (p_{n+1} - c_{n+1}). \end{aligned} \quad (3.120)$$

Truncation error of corrector formula in this case is

$$\text{truncation error} = - \frac{1}{40} T^5 y^{(5)} .$$

as compared with that of Milne's corrector formula, there is a 125% increment.

In Hamming methods predictor formulas are always the same as Milne's formula because if the predictor is generalized by the same method used above, it would be unstable.

From the above discussion it is easily seen that to gain stability in a predictor-corrector method one must lose some accuracy. This loss in accuracy can be compensated by shortening the interval of integration.

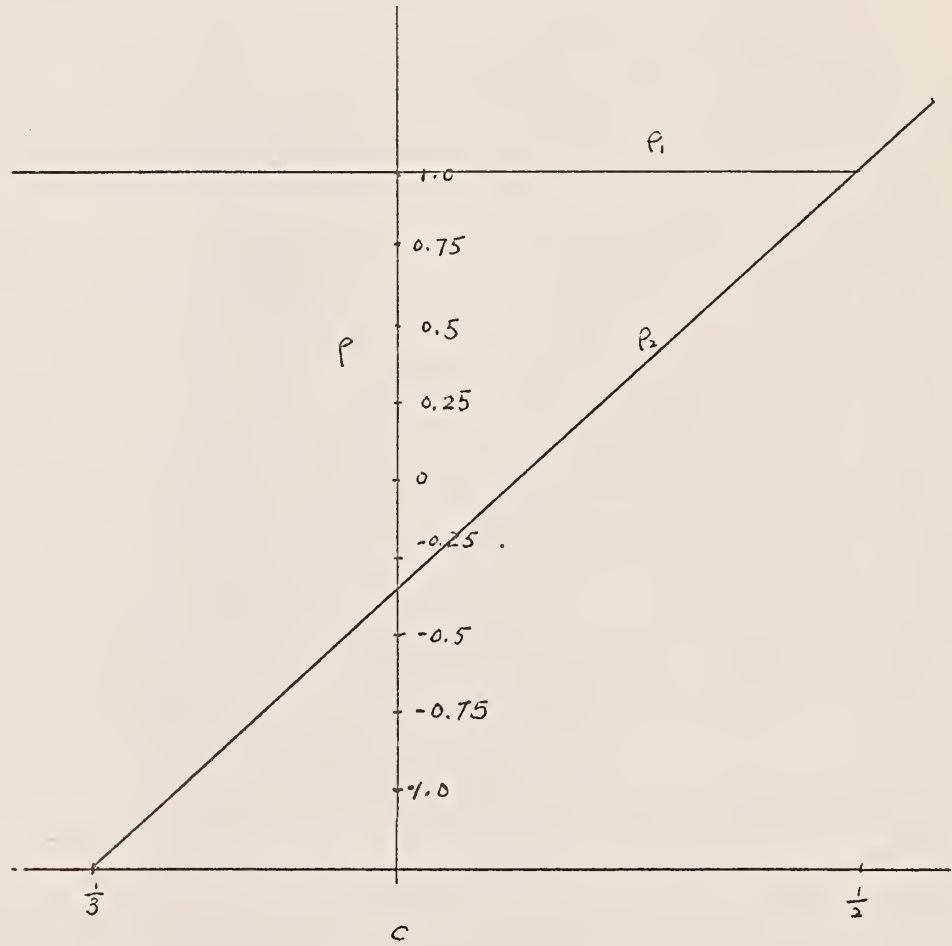


Fig. 1. Root Loci of Eq. (3.114)

$$p^2 + (4 - 12c)p - (5 + 12c) = 0$$

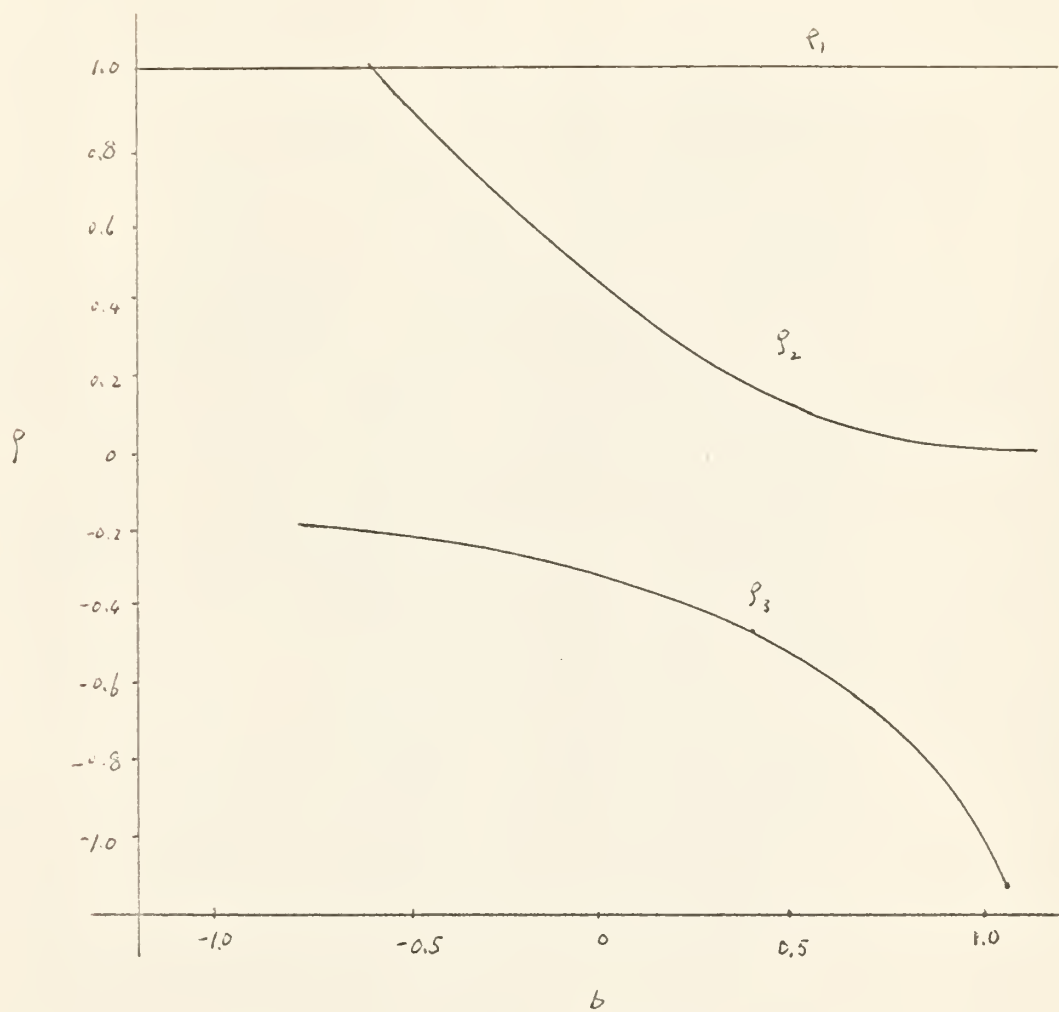


Fig. 2. Root Loci of Eq. (3.119)

$$8p^3 - 9(1 - b)p^2 - 8bp + (1 - b) = 0$$

C. Comparison Between One-Step Methods and Multiple Step Methods

One-step methods are useful for obtaining the first few of the solution of a differential equation, but involve too much labor to be used for obtaining a numerical solution over an extended range. However, multiple step methods require special starting procedures which are furnished by one-step methods.

IV. Numerical Transform Techniques.

Part of the solution of a linear differential equation with constant coefficients has the form of a convolution. The digital computer must approximate convolution on a denumerable set of equally spaced points. The trapezoidal approximation is the commonly used quadrature method.

Trapezoidal approximation proceeds as follows. Let $g(x)$ be a function defined and continuous in a closed interval containing $x=0$, and let its first four derivatives be continuous in the same interval, the Taylor series representation of the function with remainder is

$$g(x) = g(0) + xg^{(1)}(0) + \frac{x^2}{2!} g^{(2)}(0) + \frac{x^3}{3!} g^{(3)}(0) + \int_0^x \frac{(x-\mu)^3}{3!} g^{(4)}(\mu) d\mu \quad (4.1)$$

Integrating the given function results in

$$\frac{1}{h} \int_0^h g(x) dx = \frac{(g_0 + g_1)}{2} - \frac{h^2}{12} \left(\frac{g_0^{(2)}}{2} + \frac{g_1^{(2)}}{2} \right) + R(h) \quad (4.2)$$

The remainder $R(h)$ is given by

$$R(h) = \frac{h^4}{4!} \int_0^1 (\lambda - 2\lambda^3 + \lambda^4) g^{(4)}(\lambda h) d\lambda \quad (4.3)$$

$$\leq \frac{h^4}{120} g^{(4)}(\theta h) \quad 0 < \theta < 1$$

A convolution can be approximated by dividing its interval into a finite number of subintervals each with same width T as follow:

$$\int_0^t f(\tau) g(t - \tau) d\tau = \sum_{k=0}^{n-1} \int_{kT}^{(k+1)T} f(\tau) g(t - \tau) d\tau \quad (4.4)$$

Trapezoidal approximation yields

$$\int_0^t f(\tau) g(t - \tau) d\tau = \frac{T}{2} \left(\sum_{k=0}^{n-1} f_k g_{n-k} + f_{k+1} g_{n-k-1} \right) \quad (4.5)$$

Taking z-transform of both sides of Equation (4.4) results in

$$Z[\bar{f}\bar{g}] = T [Zf][Z\bar{g}] - \frac{T}{2} [f_0 Z\bar{g} + g_0 Z\bar{f}] \quad (4.6)$$

There are many ways of utilizing trapezoidal convolution to solve a given differential equation. Each of these is called a program. Three classes of programs are discernible; The multiple integration substitution program, the Tustin integrator program, and the single integrator program. Differences among them are at the transitions from integration to multiple integration.

A. Tustin Program

This method solves an n-th order differential equation by the state-space approach. A Tustin-like program is demonstrated as follows:

Let an n-th order differential equation of the form

$$D^n y = x(t) \quad (4.7)$$

Define the row vectors

$$w^t = (y, y', \dots, y^{(n-1)})$$

$$v^t = (D, 0, \dots, x(t)) \quad (4.8)$$

Then

$$Dw + \begin{pmatrix} 0 & -I \\ 0 & 0 \end{pmatrix} w = v \quad (4.9)$$

is equivalent to an nth order differential equation. Take the Laplace transform and divided by s to obtain

$$\bar{w} + \frac{1}{s} A\bar{w} = \frac{1}{s} v + \frac{1}{s} w_0 \quad (4.10)$$

Employing the sampling operation and trapezoidal integration yields

$$Z\bar{w} + AZ\bar{w} - Aw_0 = \alpha Z\bar{v} - \beta v_0 + \frac{2}{T} \beta w_0 \quad (4.11)$$

where

$$\alpha = \frac{T}{2} \frac{1+z}{1-z}, \quad \beta = \frac{T}{2(1-z)}, \quad A_{(n \times n)} = \begin{pmatrix} 0 & -I \\ 0 & 0 \end{pmatrix}$$

Since $A^n = 0$, it follows that

$$(I + \alpha A^n) = I \quad \text{or} \quad (I + \alpha A) \left[I - (\alpha A) + (\alpha A)^2 + \dots + (-\alpha A)^{n-1} \right] = I \quad (4.12)$$

Then the desired solution can be written as

$$Z\bar{y} = \alpha^n Z\bar{x} - \beta \alpha^{n-1} x_0 + \frac{1}{1-z} y_0^{(n-1)} + \frac{T^2 z}{(1-z)^2} \sum_{k=2}^n \alpha^{k-2} y_0^{(n-k)} \quad \text{for } n = 1, 2, 3, \dots \quad (4.13)$$

The above equation is the Tustin program, if all initial conditions are zero.

This program's great advantage is that it finds the solution and its

first (n-1) derivatives in one sweep.

B. Single-integrator program.

Halljak has derived a single-integrator program which uses a sequence of ascending order differential equations derived from the given differential equation to solve the same equation. Consider the differential equation.

$$\left[D^n + a_1 D^{n-1} + \dots + a_{n-1} D + a_n \right] y(\tau) = x(\tau), \quad (4.14)$$

First, set up the differential equation

$$(D + a_1) y_1(\tau) = 1 \quad y_1(0) = 0 \quad (4.15)$$

The Laplace transform of this differential equation yields after division by s

$$\left(1 + \frac{a_1}{s}\right) \bar{y}_1(s) = \frac{1}{s^2} \quad (4.16)$$

Taking Z-transform of both sides and employing approximate trapezoidal convolution, yields

$$Z\bar{y}_1 = \frac{2Tz}{[1 - z][(2 + aT) - (2 - aT) z]} = Z\left(\frac{1}{s(s + a)}\right) \quad (4.17)$$

Set up another differential equation of the form;

$$(D^2 + a_1 D + a_2) y_2(t) = 1, \quad y_2(0) = y_2'(0) = 0 \quad (4.18)$$

taking Laplace transform of both sides and dividing by $s(s+a)$ yield

$$\left(1 + \frac{a^2}{s(s+a)}\right) \bar{y}(s) = \frac{1}{s^2(s+a)} \quad (4.19)$$

taking Z-transform yields

$$Z\bar{y}(s) + Ta_2 \left[Z\left(\frac{1}{s(s+a_1)}\right) \right] \left[Z\bar{y}(s) \right] = \left[Z\left(\frac{1}{s}\right) \right] \left[Z\left(\frac{1}{s(s+a)}\right) \right] \quad (4.20)$$

Substituting equation (4.17) into equation (4.20) yields

$$\begin{aligned} Zy_2(s) &= \frac{T}{2} \frac{(1+z)}{(1-z)} \cdot \frac{2Tz}{(2+a_1T) - (4-2a_2T^2)z + (2-a_1T)z^2} \\ &= Z\left(\frac{1}{s(s^2 + a_1s + a_2)}\right) \end{aligned} \quad (4.21)$$

Proceeding in this manner, let a n th order differential equation be of the form

$$(D^{n-1} + a_1 D^{n-2} + \dots + a_{n-1}) y_{n-1}(t) = 1$$

$$y_{n-i}^{(0)} = y_{n-1}^{(0)} = \dots = y_{n-1}^{(n-2)}(0) = 0 \quad (4.22)$$

A recurrence relation can be constructed to be

$$Z\bar{y}_i(s) = \frac{T}{2} \left(\frac{1+z}{1-z} \right) Z\bar{y}_{i-1}(s) / \left[1 + a_1 T Z\bar{y}_{i-1}(s) \right] \quad (4.23)$$

$$2 \leq i < n$$

for $i=n$,

$$Z\bar{y}_n(s) \left[1 + a_n T (Z\bar{y}_{n-1}(s)) \right] = T \left[Zx(s) \right] \left[Zy_{n-1}(s) \right] \\ = \frac{T}{2} \left\{ x_0 Zy_{n-1}(s) \right\} \quad (4.24)$$

Initial conditions other than zero can be introduced at the last step.

C. Multiple-Integrator Substitution Program.

The multiple integrator substitution program casts the Laplace transform of the given differential equation into inverse powers of s by a suitable division and substitutes for them definite functions of z defined to be e^{-Ts} .

C. I. Halijak's Integrator Substitution Program.

This method uses trapezoidal convolution to generate z -transform of

$(1/s^n)$ and $(1/s^n)f$. The procedure is as follows:

$$\text{for } n = 1, \quad Z(1/s) = 1/(1 - z),$$

$$\begin{aligned} \text{for } n = 2, \quad Z(1/s^2) &= Z(1/s)Z(1/s)T - TZ(1/s), \\ &= Tz/(1 - z)^2, \end{aligned} \quad (4.25)$$

$$\begin{aligned} \text{for } n \geq 2, \quad Z(1/s^n) &= Z(1/s^{n-1})Z(1/s)T - Z(1/s^n)T/2, \\ &= Z\left(\frac{1}{s^{n-1}}\right) - \frac{T}{2} \frac{(1+z)}{(1-z)}, \\ &= Z\left(\frac{1}{s^2}\right) - \frac{T}{2} \left[\frac{(1+z)}{(1-z)} \right]^{n-2} \end{aligned}$$

Substituting Equation (4.25) into Equation (4.6) yields

$$Z\left(\frac{1}{s^n} f\right) \doteq \left[T^2 z / (1 - z)^2 \right] \left[\frac{T(1+z)}{2(1-z)} \right]^{n-2} (Z\bar{f} - 0.5f_0) \quad (4.26)$$

This method yields moderate accuracy approximations and is the basis of physically small computers.

SUMMARY

The result of an approximate computation of the solution of differential equation is affected by errors. These errors arise from different causes and affect final result in different ways.

Three types of errors that are most important are truncation error, round-off error, and accumulated error. Round-off error usually affects the last retained digit of the decimal representation; its effect can be minimized by retaining additional digits. Truncation error is due to discarded terms in an infinite series. Sometimes the remainder exceeds the sum of terms retained, thus making the calculated result meaningless. Therefore, an estimate of truncation error is essential. The importance of accumulated errors depends on rate of accumulation. If the accumulation error is unbounded, the solution becomes meaningless.

The search for Runge-Kutta type formulas is important. It seems that the method given in this report in deriving coefficients of Runge-Kutta formulas can be applied to investigate six and higher order formulas.

Trapezoidal convolution using the integral of the first two terms of Taylor series coincides with the average of right and left Riemann sum approximations. The truncation error is of order T^2 . Improved trapezoidal convolution using higher order Taylor Series terms seems worthy of further investigation.

Accuracy is not the only consideration for evaluating a computation process. The step size affects the solution accuracy as well as the cost; the more accurate solution is generally the more expensive. Frequently it is desired to have a solution within a certain accuracy with most economical

computation. The selection of optimal size depends on the method used and the problem solved.

A number of numerical transform techniques have been developed in the Z-transform language. Many of these exhibit difficulties for the solution of differential equation with non-zero initial conditions. The programs generated by Halijak have shown that proper reintroduction of initial conditions is required for better approximation.

APPENDIX 1

Error Formulas for an Interpolating Polynomial

Let the function $z(x)$ be defined on an interval containing the $q+1$ distinct points x_0, x_1, \dots, x_q . It is well known that among all polynomials in x of degree not exceeding q there exists exactly one polynomial $P(x)$ which satisfies the relations

$$P(x_i) = z(x_i) \quad i = 0, 1, 2, \dots, q \quad (\text{I} - 1)$$

The uniqueness of this interpolating polynomial follows from the fact that the difference of any two such polynomials is a polynomial of degree q which has $q+1$ zeros and therefore vanishes identically. Existence can be proved by exhibiting the polynomial explicitly, in the form

$$P(x) = L(x) \sum_{i=0}^q \frac{z(x_i)}{L^i(x_i) (x - x_i)} \quad (\text{I} - 2)$$

where $L(x) = (x - x_0)(x - x_1) \cdots (x - x_q)$

The error committed in this approximation can be estimated from the following Lemmas.

Lemma 1. Let $z(x)$ have a continuous derivative of order $q+1$ in J . Then for every point x in J there exists a point ξ in the smallest interval I containing both x and the points x_i ($i=1, 2, \dots, q$) such that

$$z(x) - P(x) = \frac{1}{(q+1)!} L(x) z^{(q+1)}(\xi) \quad (\text{I} - 3)$$

Lemma II. Let $z(x)$ satisfies the same hypothesis as in Lemma I. Then for every $x_k (k=0,1,2,\dots,q)$ there exists a number ξ such that

$$z'(x_k) - p'(x_k) = \frac{1}{(q+1)!} z^{(q+1)}(\xi) L'(x_k) \quad (\text{I} - 4)$$

REFERENCES

1. Tustin A.
A method of analysing the behavior of linear systems in terms of time series, Journal IEE, proceeding at the Convention on Automatic Regulators and Servomechanisms, vol, 94, part II-A, pp 130ff, May, 1947.
2. Madwed A.
Number Series Method of solving linear and non-linear Differential Equations, Rep. No. 6445-T-26, Instrumentation Lab., MIT; April, 1950.
3. Boxer R. and Thaler S.
A Simplified Method of Solving Linear and Nonlinear Systems, Proceedings IRE, vol. 44, no. 1, pp. 84-101, Jan. 1956.
4. Armarel S., Lambert L. and Millman J.
A Transformation Calculus for Obtaining Approximation Solutions to Linear Integro-Differential Equations with Constant Coefficients, Technical Report T-7/c, Columbia University Electronics Research Laboratory, August, 1955.
5. Halijak C. A.
Digital Approximation of the Solutions of Differential Equations Using trapexoidal Convolution, ITM-64, Bendix Systems Division, The Bendix Corporation, Ann Arbor, Michigan, August 26, 1960.
6. Halijak C. A.
Algebraic Methods in Numerical Analysis of Ordinary Differential Equations, Special Report No. 37, Kansas State University Bulletin.
7. Henrici P.
Discrete Variable Methods in Ordinary Differential Equations, John Wiley & Sons, Inc., New York. 1962
8. Kunz K. S.
Numerical analysis, McGraw-Hill Book Company, Inc., New York. 1942
9. Boxer R.
A Note on Numerical Transform Calculus, Proceeding IRE, vol. 45, pp. 1401-1406; October, 1957.
10. Kwan R. K.
On the Analysis of Closed-Loop systems by Means of Digital Computing Techniques. Master Thesis, McGill University, Montreal, Canada. 1963.

11. Hamming R. W.
Numerical methods for scientists and engineers, McGraw-Hill
book Company, 1962.

ACKNOWLEDGMENT

The author wishes to express his appreciation to Dr. Charles A. Halijak and Dr. Floyd W. Harris of the Department of Electrical Engineering for their tutelage and guidance which have implicitly contributed to this report.

SURVEY OF DIGITAL SOLUTION
OF DIFFERENTIAL EQUATIONS

by

ING-WEN HWANG

B.S.E.E., National Taiwan University, 1963

AN ABSTRACT OF A MASTER'S REPORT

submitted in partial fulfillment of the

requirements for the degree

MASTER OF SCIENCE

Department of Electrical Engineering

KANSAS STATE UNIVERSITY
Manhattan, Kansas

1966

There are a number of numerical methods for approximating solution of a differential equation. These methods have two things in common; the calculations are performed with discrete values and on a step-by-step basis. The purpose of this report is to review those methods that are frequently used on digital computers. All these methods are divided into two classes; the classical numerical techniques and the numerical transform techniques. Classical numerical techniques yield two types; one-step methods and multiple step methods. In one-step methods the value of y_n is solved from its previous step. In multiple step methods the value of y_n is solved from its several previous steps y_{n-1} , y_{n-2} , \dots , y_{n-q} . The virtue of one-step methods is that they are self starting whereas multiple step methods require more than one starting point that they are not self starting. However, multiple step methods are more easier in continuing a solution than one-step methods. The numerical transform techniques were developed recently by Tustin, Madwed, Boxer-Thaler, and Halijak. These methods use z-transform as a language. Three classes are discernable; Tustin program, single integration program, and multiple integration substitution program. Differences among them are the transition from integration to multiple integration.

Truncation error and accumulated relative error are important factors that affect availability of a method. It is necessary to estimate these errors so as not to make computation meaningless.