UNIFORMLY ACCURATE NUMERICAL SOLUTIONS TO DIFFERENTIAL EQUATIONS USING EXTRAPOLATION AND INTERPOLATION

Thesis for the Degree of Ph. D. MICHIGAN STATE UNIVERSITY RICHARD ALLAN ROGERS 1974



This is to certify that the

thesis entitled

"Uniformly Accurate Numerical Solutions to Differential Equations Using Extrapolation and Interpolation" presented by

Richard A. Rogers

has been accepted towards fulfillment of the requirements for

Ph.D. degree in Mathematics

Gerald D. Taylor/ Mary J. Winter

Major professor

Date July 25, 1974

O-7639

4.,.-: 4

Herald D. Zaylon Wenter



In this w

for solving ordir those methods that all powers of h a fixed integer. with such methods grid points below We develop the " combines extrapo ^{CCefficient} func produce a highly

finest mesh. Wh Melds uniform a æsh.

In Chapt ls hodified so a Problems. In ad ^{γ, p}ereyra's Sint boundary v

ABSTRACT

UNIFORMLY ACCURATE NUMERICAL SOLUTIONS TO DIFFERENTIAL EQUATIONS USING EXTRAPOLATION AND INTERPOLATION

By

Richard Allan Rogers

In this work we are concerned with numerical methods for solving ordinary differential equations. We consider those methods that have asymptotic error expansions involving all powers of h^q , where h is the steplength and q is a fixed integer. The process of extrapolation can be employed with such methods to obtain highly accurate solutions at grid points belonging to the coarsest mesh. In Chapter I we develop the "pullback interpolation method". This method combines extrapolation with Hermite interpolation of the coefficient functions for the asymptotic error expansion to produce a highly accurate solution at all grid points of the finest mesh. When q is 1 or 2 pullback interpolation yields uniform accuracy at all grid points of the finest mesh.

In Chapter II the pullback interpolation method is modified so as to be applicable to boundary value problems. In addition, an elementary proof of the stability Of V. Pereyra's finite difference scheme for solving two point boundary value problems is given.

In Chap:

equations with c

are shown to be

Because of the pacy obtained thre

for these proble:

A finite
second order dela
in Chapter III.
have an asymptot
steplength.

In Chapter III we consider difference differential equations with constant retardation. The methods of Chapter I are shown to be applicable to first order delay equations.

Because of the presence of the delay term, the uniform accuracy obtained through pullback interpolation is indispensible for these problems.

A finite difference scheme for directly solving second order delay equations is constructed and analyzed in Chapter III. The global discretization error is shown to have an asymptotic error expansion in even powers of the steplength.

UNIFORM TO DIFFEREN

in partial

 \mathfrak{D}^{6}

UNIFORMLY ACCURATE NUMERICAL SOLUTIONS TO DIFFERENTIAL EQUATIONS USING EXTRAPOLATION AND INTERPOLATION

Ву

Richard Allan Rogers

A THESIS

Submitted to
Michigan State University
in partial fulfullment of the requirements
for the degree of

DOCTOR OF PHILOSOPHY

Department of Mathematics

TO MY WIFE KAREL

I would in the following of the following and encouragement to Dr. Lee M. Son advice and guidan

ACKNOWLEDGEMENTS

I would like to thank my major professors,

Dr. Gerald D. Taylor and Dr. Mary J. Winter for their help
and encouragement. Also, I want to extend my appreciation
to Dr. Lee M. Sonneborn and Dr. Gerald D. Taylor for their
advice and guidance over the course of my graduate career.

Section 3

Section 4

Section 5

CHAPTER II: TWO F

Section 1

Section 2

Section 3

CHAPTER III: THI DII RE

Section ;

Section ; Section :

Section :

Section :

Section :

Section

BIBLIO

TABLE OF CONTENTS

	INTRO	מסטכיו	rion	1											
CHAPTER			BACK INTERPOLATION METHOD FOR VALUE PROBLEMS	10											
	Section	1:	Statement of the Problem	10											
	Section	2:	The Pullback Interpolation Method	16											
	Section	3:	The Pullback Method for $q=1$ and 2	27											
	Section	4:	Determination of $e_m^{\prime}(a)$	31											
	Section	5:	Numerical Results for Initial Value Problems	56											
CHAPTER	II: TWO	POIN	NT BOUNDARY VALUE PROBLEMS	77											
	Section	1:	The Problem and its Discretization	77											
	Section	2:	Stability	84											
	Section	3:	The Numerical Method	90											
CHAPTER	CHAPTER III: THE NUMERICAL SOLUTION OF DIFFERENCE DIFFERENTIAL EQUATIONS WITH CONSTANT RETARDATION														
			First Order Equations												
			Second Order Equations												
	Section	2.1:	Consistency	112											
			: Stability												
	Section	2.3	The Global Discretization Error	132											
	Section	2.4:	Implementing the Second Order Method	141											
	Section	2.5	Numerical Results for Second Order Equations	147											
	BIBLI	OGRA I	РНУ	152											

LIST OF TABLES

Table	1	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	14
Table	2	•	•	•	•	•	•		•		•		•	•	•	•	•				•	•	•	14
Table	3	•	•	•	•	•	•	•	•		•		•	•	•	•	•	•	•	•	•	•	•	4 5
Table	4	•	•	•		•	•	•	•	•		•	•	•	•	•	•	•	•	•	•	•	•	52
Table	5	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	63
Table	6	•	•	•	•	•	•	•		•	•	•	•	•	•	•	•	•	•	•	•	•	•	65
Table	7	•	•	•	•	•	•	•		•				•	•	•		•		•		•	•	67
Table	8	•	•	•	•	•	•	•	•		•	•	•	•	•	•	•	•	•	•	•	•	•	68
Table	9		•		•	•	•					•	•	•	•			•		•	•	•	•	70
Table	10	•	•		•	•	•		•	•	•	•	•	•	•	•	•		•	•			•	71
Table	11	•	•			•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•		•	72
Table	12	•	•		•	•	•		•			•	•	•	•	•	•	•	•	•	•	•	•	74
Table	13	•	•	•		•	•	•			•		•	•			•	•	•			•	•	105
Table	14	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•			•		•	•	107
Table	15	•		•		•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	148
Table	16					•	•																•	150

In this migues for obtainmitial value prosolutions by meanwill also considuonputed solutio:

Consider

$$x'(t) = f$$

(1)

$$x(a) = c$$

We will assume the

Conditions on f(

known and are giv

The major

are based on disc

inate solution to

ia,b]. We will o

wints are equall wint set can be

$$G = \{t_n = a\}$$

were the parameto

INTRODUCTION

In this section we will discuss some standard techniques for obtaining numerical solutions to first order initial value problems. The process of refining computed solutions by means of extrapolation will be explained. We will also consider the question of uniform accuracy of the computed solution.

Consider the first order initial value problem

$$x'(t) = f(t,x(t)),$$
(1)
$$x(a) = \alpha, \quad a \le t \le b.$$

We will assume that (1) has a unique solution $\varphi(t)$ which depends continuously on the initial condition $x(a) = \alpha$. Conditions on f(t,x(t)) which will guarantee this are well known and are given in Chapter I.

The majority of procedures for solving (1) numerically are based on discretization. As such, they yield an approximate solution to (1) on a discrete point set contained in [a,b]. We will only consider the case where the discrete points are equally spaced in [a,b]. In this case the discrete point set can be conveniently represented as a grid

$$G = \{t_n = a+nh: n=0, 1, ..., N\},\$$

where the parameter $h = \frac{b-a}{N}$ and is called the steplength.

The theo t(t) and the ap by $X(t_n,h)$ for

A computwhich takes the $X(t_{n+j},h)$ and to be a linear is referred to a.

Three of rule (sometimes of rule and Gragg's

Euler's r

 $X(t_{n+1},h) =$ $\sum_{i \neq a \text{ tion } (2) \text{ exhi}}$

is obtained given
imate solution at
at the preceeding
imate solution at
a one-step method

tition in (1). T

everything on the ste trying to comp

The theoretical solution to (1) will be denoted by $\phi(t) \quad \text{and the approximate numerical solution will be denoted}$ by $X(t_n,h)$ for each grid point t_n .

A computational method for determining $X(t_{n+k},h)$ which takes the form of a linear relationship between $X(t_{n+j},h)$ and $f(t_{n+j},X(t_{n+j},h))$, $j=0,1,\ldots,k$ is said to be a linear k-step method. The class of all such methods is referred to as the class of linear multistep methods.

Three of the simpler multistep methods are Euler's rule (sometimes called the Euler-Cauchy method), the trapezoid rule and Gragg's modified midpoint rule.

Euler's rule is summarized by the equation

(2)
$$X(t_{n+1},h) = X(t_n,h) + hf(t_n,X(t_n,h)), n=0,1,...,N-1.$$

Equation (2) exhibits how the approximate solution at t_{n+1} is obtained given that one has already obtained an approximate solution at t_n . Since information is required only at the preceding grid point in order to obtain an approximate solution at the next grid point, Euler's rule is a one-step method. To initialize or start the method requires one piece of information which is given by the initial condition in (1). That is, take $X(t_0,h)=x(a)=\alpha$. Since everything on the right hand side of (2) is known when we are trying to compute $X(t_{n+1},h)$, the relationship (2) is

This

te det

explicit and Euler's rule is an explicit multistep method.

The trapezoid rule is given by

(3)
$$X(t_{n+1},h) = X(t_n,h) + \frac{h}{2}[f(t_n,X(t_n,h)) + f(t_{n+1},X(t_{n+1},h))],$$

$$n=0,1,\ldots,N-1.$$

This is again a one-step method and is initialized by taking $X(t_0,h)=x(a)=\alpha$. However, the right hand side of (3) is not completely known. Indeed, when trying to compute the solution at t_{n+1} the term $f(t_{n+1},X(t_{n+1},h))$ is unknown because it involves the solution we are trying to compute. Whenever f(t,x(t)) is linear in x(t), and in some other special cases, equation (3) can be solved explicitly for $X(t_{n+1},h)$. However, in general this is not possible and the trapezoid rule is said to be an implicit multistep method because of this behavior. When (3) cannot be solved explicitly for $X(t_{n+1},h)$ a root finding procedure or functional iteration can be used to find the solution.

The midpoint rule is an explicit two-step method defined by

(4)
$$X(t_{n+2},h) = X(t_n,h) + 2hf(t_{n+1},X(t_{n+1},h)), n=0,1,...,N-2.$$

To compute $X(t_{n+2},h)$ we need to know both $X(t_n,h)$ and $X(t_{n+1},h)$. Thus, the midpoint rule requires two starting values $X(t_0,h) = \alpha$ and $X(t_1,h)$. The second of these can be determined in a variety of ways.

is t

use :

Secon to ma

proce

and

٠.

giver

(5)

(3)

tor a

ezoid

solut

(6)

The f

fixed

Gragg's modification of the midpoint rule [11,12] is twofold. First, to obtain the additional starting value use Euler's rule. That is,

$$X(t_1,h) = \alpha + hf(a,\alpha)$$
.

Second, at the grid point $t_N = b$ use a smoothing procedure to make the computed solution more stable. The smoothing procedure consists of averaging three computed solutions and is similar to a device originally employed by Milne and Reynolds [20,21]. Gragg's modified midpoint rule is formally given by

$$X(t_1,h) = \alpha + hf(a,\alpha);$$

(5)
$$X(t_{n+2},h) = X(t_n,h) + 2hf(t_{n+1},X(t_{n+1},h)), \quad n=0,1,...,N-1;$$

 $X(b,h) = \frac{1}{4}X(t_{N-1},h) + \frac{1}{2}X(t_{N},h) + \frac{1}{4}X(t_{N+1},h).$

The three numerical methods given by equations (2)

(3) and (5) have an important similarity. Each has an asymptotic error expansion that involves all powers of h^q for a fixed integer q. That is, for Euler's rule, the trapezoid rule, and Gragg's modified midpoint rule, the computed solution satisfies

(6)
$$X(t_n,h) = \varphi(t_n) + \sum_{k=1}^{\infty} e_k(t_n) h^{qk}$$
.

The functions $e_k(t)$ are independent of h and q is a fixed integer. For Euler's rule q=1 and for the trapezoid



rule

(6)

grid

7ati

(61)

The

beha

We s

C i

depe

(6)

mod:

Ste

0£ 8

traj

שודוו

(6)

rule and Gragg's modified midpoint rule q=2. The expansion (6) for Gragg's modified midpoint rule is only valid at $t_N = b$ while for the other two methods (6) is valid at all grid points t_n . The expansion (6) is valid only when $f \in C^{\infty}[[a,b] \times (-\infty,\infty)]$.

If f has only a finite number of continuous derivatives, a truncated version of (6) is valid, namely

(6')
$$X(t_n,h) = \varphi(t_n) + \sum_{k=1}^{M} e_k(t_n) h^{qk} + O(h^{q(M+1)}).$$

The notation $\mathfrak{G}(h^j)$ means that the function being suppressed behaves like a constant multiplied by h^j as $h \to 0$. Formally,

we say that $g(t) = O(h^j)$ if $\frac{|g(t)|}{h^j} \le C$ as $h \to 0$, where C is a constant. The length, M, of the expansion (6') depends on q and the number of continuous derivatives of f that exist.

The existence of asymptotic expansions of the form

(6) or (6') for Euler's rule, the trapezoid rule and Gragg's modified midpoint rule was originally proved by Gragg [11,12]. Stetter [31] and Pereyra [24] have also studied the existence of such expansions. Gragg's results for Euler's rule and the trapezoid rule are presented in more detail in Chapter I.

The existence of such expansions is important because numerical methods which have error expansions of the form

(6) or (6') are amenable to extrapolation. Basically, the



proc

q.1.00

as t

back

Arch

firs

cent

extr

latt

exce

has

solu

às f

ggre

à + H k=0,

711

With

solu

K+1 expa

process of extrapolation is a means of combining several computed solutions, each of low accuracy, in such a manner as to obtain a computed solution with high accuracy.

Extrapolation is not a recent development, dating back to at least 1654 when Huygens [14] used it to improve Archimedes polygonal approximation to π . Extrapolation was first systematically studied by Richardson [29] early in this century and has often been referred to as either Richardson extrapolation or the deferred approach to the limit, the latter being the title of Richardson's 1927 paper. An excellent survey article on extrapolation and its applications has been written by Joyce [16].

Extrapolation, when applied to obtaining an accurate solution to a differential equation, is most easily explained as follows: Let H>0 be a fixed basic steplength and suppose an accurate solution to (1) is desired at the point a+H. Define a sequence of steplengths $h_k = H/2^k$ $k=0,1,\ldots,K$ and grids

$$G_k = \{t_i^k = a+ih_k : i=0,1,...,2^k\}.$$

All grids contain a and a+H and the grids are nested with $G_k \subset G_{k+1}$ $\forall k$. On each grid G_k compute a numerical solution to (1) using a method which has an asymptotic error expansion of the form (6') with M>K. At a+H we have K+l computed solutions $X(t_N,h_k)$ for k=0,1,...,K.



Extra

of th

the f

const

(7)

an it

perfo in (

suit

at a

able

accu

ځ.

Eowe

poin

extr $a + \frac{E}{2}$ less

extx

თხ_{ნმ}

grid

Extrapolation is the process of forming a linear combination of these K+1 solutions in such a manner so as to eliminate the first K error terms of the expansion (6'). That is, constants $\mathbf{c_k}$ are determined so that

(7)
$$\sum_{k=0}^{K} c_k X(a+H, h_k) = \varphi(b) + O(H^{q(K+1)}).$$

Aitken [1] and Neville [22] independently devised an iterative scheme by means of which extrapolation can be performed without explicit computation of the constants c_k in (7). The convergence of this iterative scheme under suitable hypotheses was established by Gragg [11,12].

Note that in order to obtain $\mathfrak{O}(\mathtt{H}^{q\,(K+1)})$ accuracy at a grid point you must have K+1 computed solutions available to work with. Thus, extrapolation will yield $\mathfrak{O}(\mathtt{H}^{q\,(K+1)})$ accuracy only at the point a+H which is common to all grids \mathtt{G}_k . Extrapolation can be performed at other grid points. However, it will not yield $\mathfrak{O}(\mathtt{H}^{q\,(K+1)})$ accuracy at these points. For instance, using the fact that $\mathtt{G}_k \subseteq \mathtt{G}_{k+1}$ $\forall k$, extrapolation will yield $\mathfrak{O}(\mathtt{H}^{q\,(K+1)})$ accuracy at the midpoint $\mathtt{a} + \frac{\mathtt{H}}{2}$. At other grid points, extrapolation will yield even less accuracy.

Lindberg [19] has developed a method based on extrapolation and Lagrange interpolation which can be used to obtain $\mathcal{O}(H^{q(K+1)-1})$ accuracy all at grid points of the finest grid G_K .



pull

and

at a

four

and

Lind

inte

nune

pres

ordi

two

18;

Pere

scj.e

jas ;

Full ролис

Secti

In Chapter I we present what we have termed "the pullback interpolation method". It utilizes extrapolation and Hermite interpolation to obtain $\mathcal{O}(H^{q(K+1)})$ accuracy at all points of the finest grid when q=1 or 2. The first four sections of Chapter I are devoted to developing the method and establishing a theoretical basis for it. In Section 5 Lindberg's method is presented and compared with pullback interpolation both theoretically and numerically. Extensive numerical tests were performed and these results are also presented in Section 5 of Chapter I.

The focus in Chapter I is entirely on first order ordinary differential equations. In Chapter II we consider two point boundary value problems of the form

(8)
$$x''(t) = f(t,x(t),x'(t)),$$
$$x(a) = A, x(b) = B.$$

Pereyra [24,25,26,27,28] has developed a finite difference scheme which yields an approximate solution to (8) that has an asymptotic error expansion of the form (6') with q=2.

Pereyra's results are summarized in Chapter II and pullback interpolation is modified so as to be applicable to boundary value problems of the form (8). In addition, in Section 2 of Chapter II, we present a new proof of the

stabi

to Pe

of di

ation

are a

showr

Nume

deve

(9)

The .

briji exba stability of Pereyra's difference scheme. In comparison to Pereyra's proof, ours is considerably more elementary.

In Chapter III we consider the numerical solution of difference differential equations with constant retardation. First order equations

$$\dot{x}(t) = f(t,x(t),x(t-r))$$

are analyzed in Section 1 and pullback interpolation is shown to be a viable solution technique for these problems. Numerical results are presented to support this contention.

The remainder of Chapter III is devoted to the development and analysis of a finite difference method for directly solving second order equations of the form

(9)
$$\ddot{x}(t) + f(t,x(t),x(t-r),\dot{x}(t),\dot{x}(t-r)) = 0.$$

The approximate solution to (9) is shown to have an asymptotic expansion of the form (6') with q=2. A modification of pullback interpolation is shown to be applicable.

<u>Sect</u>

valu

(1)

ñe s of

ies:

3011

the

73](ya][

of 1

à s;

ۇچىز ئۇچىن ويري

san

ђе (

CHAPTER I

THE PULLBACK INTERPOLATION METHOD FOR INITIAL VALUE PROBLEMS

Section 1. Statement of the Problem

In this chapter we consider the first order initial value problem.

$$y'(t) = f(t,y(t))$$
(1)
$$y(a) = \alpha \qquad a \le t \le b.$$

We shall assume that f(t,y(t)) is a continuous function of t and satisfies a uniform Lipschitz condition with respect to its second argument. Under these assumptions it is well known (see Keller [17]) that (1) has a unique solution, $\varphi(t)$, which depends Lipschitz continuously on the initial data $y(a) = \alpha$. y(t) may be either a scalarvalued or a vector-valued function. If y(t) is vector-valued then f(t,y(t)) will be a vector-valued function of the variable t and the vector y(t) and (1) will be a system of first order differential equations. This case also arises when we reduce a m order differential equation to a system of m first order differential equations. The standard technique for accomplishing this can be found in Lambert [18]. The numerical methods to be considered for solving (1) will work for either the

scalar extrap this c of the the sa withou refine on a p

(1),]

Corlis

k=0,1,

 $G_{k} = \frac{2^{k}+1}{2^{k}+1}$ $\frac{2^{k}}{2^{k}} = \frac{2^{k}}{2^{k}}$

for th give e

ution,

sion

(2)

coeff:

hr fo

extrapolation and the pullback method, to be explained in this chapter, can be done independently for each component of the solution vector to obtain a refined solution. Since the same work is done for each component, one at a time, without using any results involving other components, the refinement process appears well suited to implementation on a parallel processing computer such as ILLIAC IV (see Corliss [5]).

Turning our attention to the numerical solution of (1), let h>O be a fixed basic steplength and for each k=0,1,...,K define steplengths $h_k = h/2^k$ and grids $G_k = \{t_i^k: t_i^k = a+ih_k, i=0,1,...2^k\}$. Each grid G_k contains 2^k+1 points; all grids contain the two points $t_0^k = a$ and $t_0^k = a+h$; and the grids are nested, that is, $G_k \subset G_{k+1}$ Vk.

In what follows we shall be concerned with methods for the numerical solution of (1) on the grids \mathbf{G}_k which give an approximation, $\mathbf{Y}(\mathsf{t}_i^k,\mathsf{h}_k)$, to the theoretical solution, $\phi(\mathsf{t}_i^k)$, such that the error has an asymptotic expansion

(2)
$$Y(t_i^k, h_k) = \varphi(t_i^k) + \sum_{j=1}^{M} h_k^{jq} e_j(t_i^k) + O(h_k^{(M+1)q})$$

for each $i=0,1,\ldots,2^k$ and for each $k=0,1,\ldots,K$. The coefficient error functions, $e_m(t)$, are independent of h_k , for each k, and q is a fixed integer (usually 1 or 2)

peculi of the

employ

exist.

f(t,y

with

solut:

and S

expan expan

of th

rule

gpose

funct

equat

(3)

The a the t

e]'∙.

there

funct

];;6]

peculiar to the method being employed. The exact length of the expansion (2) will depend on both the method being employed and the number of derivatives of $f(\cdot,\cdot)$ which exist. In general, (M+1)q continuous derivatives of f(t,y(t)) with respect to t are required for an expansion with M error terms. This is equivalent to the theoretical solution g(t) having (M+1)q+1 continuous derivatives.

As mentioned before, Gragg [11,12], Pereyra [24], and Stetter [31] have investigated the existence of such expansions. Examples of methods which yield such error expansions are Euler's method (q=1), the usual generalization of the trapezoid rule (q=2) and Gragg's modified midpoint rule (q=2). An important result obtained in each of the above mentioned studies is that the coefficient error functions $e_m(t)$ satisfy an inhomogeneous linear variational equation on $a \le t \le b$, of the form

(3)
$$e_{m}(t) - J(t)e_{m}(t) = b_{m}(\cdot)$$
$$e_{m}(a) = 0, \qquad m=1,...,M.$$

The arguments of the inhomogeneous terms, $b_m(\cdot)$, involve the theoretical solution $\phi(t)$, previous error functions e_1, \ldots, e_{m-1} , the function f(t,y(t)), and various derivatives thereof. The differentiability of the various error functions depends on the differentiability of f(t,y(t)). The left hand side of (3) is the Frechet derivative of (1),



cons

Alte

at t

side

dep∈

to t

grid

of t

rot

Coi

at as

sev

7a]

ಶಿಸ್ತ

: (t

Wit

50; tit

2C]

Pair

considered as a differential operator, operating on $e_m(t)$. Alternately, J(t) is the Jacobian of f(t,y(t)) evaluated at the theoretical solution, $\phi(t)$, of (1). The left hand side of (3) may be obtained formally by assuming y(t) depends on a parameter λ , differentiating (1) with respect to this parameter and setting $e_m = \frac{dy}{d\lambda}$ and $e_m = \frac{dy}{d\lambda}$.

If we compute the numerical solution of (1) on the grids G_k k=0,1,...,M using a method which has an expansion of the form (2), we may then employ extrapolation to obtain a solution $Y(a+h) = \phi(a+h) + \mathcal{O}(h^{(M+1)}q)$. However we are not able to obtain comparable accuracy at the intermediate points. For instance, using extrapolation, the solution at the midpoint satisfied $Y(a+\frac{h}{2}) = \phi(a+\frac{h}{2}) + \mathcal{O}(h^{Mq})$. And, as the following example illustrates, we cannot interpolate several extrapolated values to obtain a solution with equivalent accuracy.

Example 1:
$$y' = y^2$$
; $y(0) = .2$; $0 \le t \le 3$.

The theoretical solution to this problem is $\varphi(t) = \frac{1}{5-t}$. Using the trapezoid rule and extrapolation with h=1, M=3 we compute the solution, Y(t), at the points t=1,2, and 3 by resolving the same problem three times with initial conditions determined by the computed solution at t=1 and 2. The results are given in Table 1 below.

for

75å

(t,

Eer

the

t

0 l

2 3

1 2 .2 .2 .2 .4

igyat ieyes

A variety of interpolation schemes are possible for determining the solution at intermediate grid points. Two such are summarized in Table 2 below. L(t) is the Lagrange cubic interpolation polynomial for the data (t,Y(t)) given in Table 1 and H(t) is the quartic Hermite interpolation polynomial for the same data with the added condition $Y'(0) = (Y(0))^2 = (.2)^2 = .04$.

TABLE 1

t	φ (t)	Y (t)	φ(t) - Y(t)
0	.200 000 000	.200 000 000	0
1	.250 000 000	.250 000 000	О
2	.333 333 333	.333 333 330	3 x 10 ⁻⁹
3	.5000 000 000	.499 999 762	2.38×10^{-7}

TABLE 2

$\varphi(t) - H(t)$	H(t)	φ (t)-L(t)	L(t)	φ (t)	t
331 -1.736x10- ⁴	.222395831	-1.736x10 ⁻³	.223958320	.22222222	$\frac{1}{2}$
06 4.018x10 ⁻⁴	.285312506	1.339x10 ⁻³	.284375013	.285714285	<u>3</u> 2
34 -1.562X10 ⁻³	.401562434	-3.125x10 ⁻³	.403124923	.40000000	<u>5</u>

As an examination of Tables 1 and 2 quickly reveals, neither of the interpolation schemes gives accuracy that is comparable to that of the computed solution.

In the next section we will present and analyze the pullback interpolation method. It is based on a Hermite interpolation scheme for approximating the error functions $e_m(t)$ successively, beginning with the last term of the error expansion. The details are worked out for a general expansion of the form (2). However in practice q is usually either one or two and the reader may find it helpful to bear this fact in mind.

<u>Sec</u>

pro

we <u>v</u>+]

(4)

Eq.

DO.

38

!=

(4

Wi

Ă.

. .

Section 2. The Pullback Interpolation Method

Still assuming that our numerical method and the problem at hand are such that the expansion (2) is valid, we now take K=M. At the point $t=a+h \in G_{\overline{O}}$ we have the M+l computed solutions and error expansions

(4)
$$Y(t,h_k) = \varphi(t) + \sum_{j=1}^{M} h_k^{jq} e_j(t) + O(h_k^{(M+1)q}), k=0,1,...,M.$$

Equation (4) can be regarded as a system of M+1 polynomial equations in h. Since the h_k 's are distinct these polynomials are linearly independent and (4) may be solved as a matrix system for the unknown vector $(\phi(t), e_1(t), e_2(t), \ldots, e_M(t))^T$.

In matrix form (4) is given by

$$(4') \qquad \qquad AU = Y + E$$

where

$$A = \begin{bmatrix} 1, h_0^q, h_0^{2q}, \dots, h_0^{(M-1)q}, h_0^{Mq} \\ 1, h_1^q, h_1^{2q}, \dots, h_1^{(M-1)q}, h_1^{Mq} \\ \vdots & \vdots & \vdots & \vdots \\ 1, h_M^q, h_M^{2q}, \dots & h_M^{(M-1)q}, h_M^{Mq} \end{bmatrix}$$

is $(M+1) \times (M+1)$; and U,Y and E are $(M+1) \times 1$ vectors given by

$$U = \begin{bmatrix} \varphi(t) \\ e_1(t) \\ \vdots \\ e_M(t) \end{bmatrix}; \quad Y = \begin{bmatrix} Y(t, h_0) \\ Y(t, h_1) \\ \vdots \\ Y(t, h_M) \end{bmatrix}; \quad \text{and} \quad E = \begin{bmatrix} \varphi(h_0^{(M+1)q}) \\ \varphi(h_1^{(M+1)q}) \\ \vdots \\ \varphi(h_M^{(M+1)q}) \end{bmatrix}$$

respectively.

Define (A,I) to be the (M+1) \times 2(M+1) matrix obtained by adjoining the identity matrix to A. A standard theorem from linear algebra (see Cullen [6]) tells us that if (A,I) is row equivalent to (I,B) — denoted by (A,I) \sim (I,B) — then B = A⁻¹.

Since A is nonsingular we can find a sequence of elementary row operations such that $(A,I) \sim (D,C)$ where D is the diagonal matrix whose (j,j) entry is $h^{(j-1)q}$, $j=1,\ldots,M+1$. Moreover, since $\forall k h_k = h/2^k$, the row operations which accomplish this reduction are independent of h, in that each involves multiplication or division by a constant only. (D,C) can be further reduced to (I,A^{-1}) by dividing the j^{th} row by $h^{(j-1)q}$.

Multiplying both sides of (4') by A^{-1} we have $U = A^{-1}AU = A^{-1}Y + A^{-1}E$. The vector U is the precise solution to the system (4) and the vector $A^{-1}Y$ is the numerical solution we can obtain. The vector $A^{-1}E$ is the vector of errors for the numerical solution.

Now $\forall k h_k \leq h$, so we can write

$$\mathbf{E} = \begin{bmatrix} O(h_{O}^{(M+1)q)} \\ O(h_{1}^{(M+1)q)} \\ \vdots \\ O(h_{M}^{(M+1)q} \end{bmatrix} = \begin{bmatrix} O(h^{(M+1)q}) \\ O(h^{(M+1)q}) \\ \vdots \\ O(h^{(M+1)q}) \end{bmatrix}.$$

Since each component of the jth row of A^{-1} has a factor of $h^{-(j-1)q}$ and this is the only dependence of these entries on h, the jth component of $A^{-1}E$ is $O(h^{(M+1)q-(j-1)q}) = O(h^{(M-j+2)q})$. Thus our computed solution $A^{-1}Y = (\bar{\varphi}(t), \bar{e}_1(t), \bar{e}_2(t), \ldots, \bar{e}_M(t))^T$ will be equal to the precise solution $U = (\varphi(t), e_1(t), e_2(t), \ldots, e_M(t))^T$ with an error of magnitude

$$A^{-1}E = (O(h^{(M+1)q}), O(h^{Mq}), O(h^{(M-1)q}), \dots, O(h^{q}))^{T}.$$

It should be emphasized that what we actually obtain when computing a solution to (4) is the vector $\mathbf{A}^{-1}\mathbf{Y}$ which is only an approximation to the actual solution U. This can be conveniently summarized by saying the solution to (4), $(\varphi(t), e_1(t), \ldots, e_M(t))^T$, is known with accuracy $(\mathfrak{G}(\mathbf{h}^{(M+1)q}), \mathfrak{G}(\mathbf{h}^{Mq}), \ldots, \mathfrak{G}(\mathbf{h}^q))^T$.

Since $e_M(a) = 0$ and $\varphi(a)$ is known exactly; we can obtain $e_M'(a)$ exactly from (3). That is, $e_M(a) = 0$ implies that $e_M'(a) = b_M(a)$. To evaluate

b_M(a) required to us. These differential

interpolating pieces of date interpolating pieces of date into the interpolating inter

this chapter

Thus

 $h_{\mathbf{k}}^{\mathbf{Mq}}$

where $\beta_{M} \equiv \frac{1}{2} M = \frac{$

proc print in G and G 1 cor such t, nan

inis teak
solutions ar

 $b_M(a)$ requires knowledge of higher order derivatives of $\phi(t)$, J(t), $f(t,\phi(t))$ and the errors $e_m(t)$, $m=1,\ldots,M-1$, evaluated at the point t=a, all of which are available to us. These derivatives may be obtained by successively differentiating (1) and (3) with respect to t. The details will be worked out for specific methods later in this chapter.

Thus we know $e_M(a)$ and $e_M(a)$ exactly and we know $e_M(a+h)$ to $\mathcal{O}(h^q)$. Construct $P_M(t)$ the Hermite interpolating polynomial of degree 2 to the above three pieces of data. It is well known that the error in this Hermite interpolation is $\mathcal{O}(h^3)$. Thus, since $e_M(a+h)$ is known to $\mathcal{O}(h^q)$ we have $P_M(s) = e_M(s) + \mathcal{O}(h^3) + \mathcal{O}(h^q)$. Now $\forall k=1,\ldots,M$, $h_k < h$ implies $h_k^\gamma < h^\gamma$ where γ is any positive integer and therefore $\mathcal{O}(h_k^\gamma) < \mathcal{O}(h^\gamma)$. Thus $\forall k=1,\ldots,M$ and $\forall s \in [a,a+h]$ we have

(5)
$$h_{k}^{Mq} P_{M}(s) = h_{k}^{Mq} e_{M}(s) + O(h^{\beta})$$

where $\beta_{\mathbf{M}} \equiv \min((\mathbf{M}+1)q, \mathbf{M}q+3)$. Note that for $q \le 3$ $\beta_{\mathbf{M}} = (\mathbf{M}+1)q$.

Proceeding to the grid G_1 , let t be any point in $G_1\backslash G_0$ (since G_0 contains $2^0+1=2$ points and G_1 contains $2^1+1=3$ points, there is only one such t, namely t=a+h/2). Since the grids are nested this $t\in G_k$ $\forall k=1,\ldots,M$ and we have the M computed solutions and error expansions

Y(t,h (6)

Substituting of M linear unknown vecto

(7) Y(t,h_k)

because FM <

Solvin (c(t), e₁(t),...

(3(h^βM),

We now For $t = a+h \in$ and for t = a3(h (\$M-(M-1)q)

e_{M-1}(a) the Hermite in interpolating t ^{of inter}polatic

is [[a,a+h].

(6)
$$Y(t,h_k) = \varphi(t) + \sum_{j=1}^{M-1} h_k^{jq} e_j(t) + h_k^{Mq} e_M(t) + O(h_k^{(M+1)q})$$
for $k=1,\ldots,M$.

Substituting (5) into (6), we obtain the following system of M linearly independent equations in h for the unknown vector $(\phi(t), e_1(t), \dots, e_{M-1}(t))^T$

(7)
$$Y(t,h_k) - h_k^{Mq} P_M(t) = \varphi(t) + \sum_{j=1}^{M-1} h_k^{jq} e_j(t) + O(h_k^{(M+1)q}) + O(h^{\beta_M})$$

$$= \varphi(t) + \sum_{j=1}^{M-1} h_k^{jq} e_j(t) + O(h^{\beta_M}),$$

 $k=1,\ldots,M$

because $\beta_{M} \leq (M+1)q$ and $h_{k} \leq h$.

We now have the following information about $e_{M-1}(s)$. For $t=a+h \in G_0$ we know $e_{M-1}(t)$ with accuracy $\mathcal{O}(h^{2q})$ and for $t=a+h/2 \in G_1 \setminus G_0$ we know $e_{M-1}(t)$ with accuracy $\mathcal{O}(h^{M-(M-1)q})$. As before $e_{M-1}(a)=0$ and we can determine $e_{M-1}(a)$ exactly. Thus we can construct $P_{M-1}(s)$ the Hermite interpolation polynomial of degree three interpolating these four data points. Since the error of interpolation will be $\mathcal{O}(h^4)$ and $h_k \in \{a,a+h\}$,

Procee

..., P_{M-J}(s) f [†] j=0,1,...,**J**, j

nomial of degre

(9) h_k(M-j)q_{P_M}

The numbers \hat{p}_j

(10) β_{M-j} ≡ ;

and f_{M} is de

Let t Gull contains $^{\infty ints}$ there a $G_{C}^{K} \subset G^{K+1}$ A K f_k for k=J+2the following

in h:

(8)
$$h_{k}^{(M-1)q} P_{M-1}(s) = h_{k}^{(M-1)q} e_{M-1}(s) + O(h^{(M-1)q+4})$$

 $+ O(h^{(M+1)q}) + O(h^{\beta_{M}})$
 $= h_{k}^{(M-1)q} e_{M-1}(s) + O(h^{\beta_{M-1}}),$

where $\beta_{M-1} \equiv \min (\beta_M, (M-1)p+4)$. Here we have used the fact that $\beta_M \leq (M+1)q$. We note that for $q \leq 2$, $\beta_{M-1} = (M+1)q$.

Proceeding by induction, suppose $P_M(s)$, $P_{M-1}(s)$, ..., $P_{M-J}(s)$ for J < M-1 have all been constructed, where $\forall j=0,1,\ldots,J$, $P_{M-j}(s)$ is a Hermite interpolation polynomial of degree $2^{j}+1$ with

(9)
$$h_k^{(M-j)q} P_{M-j}(s) = h_k^{(M-j)q} e_{M-j}(s) + O(h^{\beta_{M-j}})$$
 $\forall s \in [a,a+h].$

The numbers β_{M-j} are defined by

(10)
$$\beta_{M-j} \equiv \min (\beta_{M-j+1}, (M-j)q + 2^{j} + 2), j=1,...,J,$$

and $\beta_{\mathbf{M}}$ is defined by (5). Note that $\beta_{\mathbf{M-j}} \leq \beta_{\mathbf{M-j+1}}$ \forall j.

Let t be any point in $G_{J+1} \setminus G_J$. Since G_{J+1} contains $2^{J+1}+1$ points and G_J contains $2^{J}+1$ points there are $2^{J+1}-2^J=2^J$ such points. Since $G_k \subset G_{k+1} \setminus K$ each $t \in G_{J+1} \setminus G_J$ is also an element of G_k for $k=J+2,J+3,\ldots,M$. Thus at each such t we have the following system of M-J linearly independent equations in h:

since BM-J = F the system (11)

we obtain the

(3(h⁵M-J), 3(h

From so Obtain knowled

accuracy o(h

 $\frac{k_{\text{now}}}{s(h^{(J+2)}q)} e_{M-J-1}$ (a

eM-J-1(a+h/2)

inductive hypo

of e M-J-1(t)
3[h M-J-(M-J-1

^eM-J-1(t) at

with varying a $e^{\dot{N}-1-1}(a) = 0$

(11)
$$Y(t,h_{k}) = \sum_{j=0}^{J} h_{k}^{(M-j)q} P_{M-j}(t)$$

$$= \varphi(t) + \sum_{j=1}^{M-J-1} h_{k}^{jq} e_{j}(t) + O(h_{k}^{(M+1)q}) + \sum_{j=0}^{J} O(h^{\beta_{M-j}})$$

$$= \varphi(t) + \sum_{j=1}^{M-J-1} h_{k}^{jq} e_{j}(t) + O(h^{\beta_{M-J}}), \quad k=J+1, J+2, \dots, M;$$

since $\beta_{M-J} \leq \beta_{M-J+1} \leq \cdots \leq \beta_{M-1} \leq \beta_{M} \leq (M+1)q$. Solving the system (11) for the unknown $(\phi(t), e_1(t), \dots, e_{M-T-1}(t))^T$,

we obtain the solution with accuracy

$$(O(h^{\beta}M-J), O(h^{\beta}M-J^{-q}), O(h^{\beta}M-J^{-2q}), \dots, O(h^{\beta}M-J^{-(M-J-1)q}))^{T}.$$

From solving (11) at each point in $G_{J+1}\setminus G_J$ we obtain knowledge of $e_{M-J-1}(t)$ at 2^J points with

accuracy $\mathfrak{O}(h$). From our work on G_O we know $e_{M-J-1}(a+h)$ with accuracy $\mathfrak{O}(h^{(M+1)}q-(M-J-1)q)=\mathfrak{O}(h^{(J+2)}q)$. From our work on $G_1\setminus G_O$ we know

 $e_{M-J-1}(a+h/2)$ with accuracy O(h). By our inductive hypothesis $\forall j \in J$, we have obtained knowledge of $e_{M-J-1}(t)$ at 2^j points in $G_{J+1}\backslash G_j$ with accuracy $O(h^{M-J}(M-J-1))$ $O(h^{M-J}(M-J-1))$. Thus we know

$$e_{M-J-1}(t)$$
 at $1 + \sum_{j=0}^{J} 2^j = 1 + 2^{J+1} - 1 = 2^{J+1}$ points

with varying accuracies. In addition we also have $\mathbf{e_{M-J-1}(a)} = \mathbf{0} \quad \text{and we can determine} \quad \mathbf{e_{M-J-1}(a)} \quad \mathbf{e_{x-J-1}(a)} = \mathbf{0}$

Thus we have k and we know ϵ interpolation ಮಿove data. T $\mathfrak{I}(h^{2^{J+1}+2})$ and $h_{k}^{(M-J-1)}$ Where βM-J-1 the error in ir 3(h (M-J-1)q+2^{J+} min(hk(M-J-1)q

 $^{\text{since}}$ $^{\text{S}}$ $M-J \leq$ (

Hence b ^{j=0,1},..., M-1

 $^{(9)}$ and (10) ar

on the expression for

the commen

Thus we have knowledge of $e_{M-J-1}(t)$ at $2^{J+1}+1$ points and we know $e_{M-J-1}(a)$. Construct $P_{M-J-1}(s)$ the Hermite interpolation polynomial of degree $2^{J+1}+1$ based on the above data. The error of interpolation will be

$$O(h^{2^{J+1}+2})$$
 and $\forall s \in [a,a+h]$,

(12)
$$h_k^{(M-J-1)q} P_{M-J-1}(s) = h_k^{(M-J-1)q} e_{M-J-1}(s) + O(h^{\beta_{M-J-1}}),$$

where $\beta_{M\!-\!J\!-\!1}$ is determined in the following manner:

the error in interpolation is $h_k^{(M-J-1)q}O(h^{2^{J+1}+2}) =$

 $O(h^{(M-J-1)q+2}^{J+1}+2)$ and the error in the given data is

$$\min (\{h_{k}^{(M-J-1)q}, h_{j}^{\beta}, h_{j}^{(M-J-1)q}, h_{k}^{(M-J-1)q}, h_{k}^{(M-J-1$$

since $\beta_{M-J} \leq (M+1)q$. Thus $\beta_{M-J-1} = \min \left(\beta_{M-J}, (M-J-1)q + 2^{J+1} + 2\right)$.

Hence by induction we can determine $P_{M-j}(s)$ for $j=0,1,\ldots,M-1$ such that $P_{M-j}(s)$ has degree $2^{j}-1$ and (9) and (10) are valid for $j=0,1,\ldots,M-1$.

On the last grid G_M we have the following expression for $\phi(t)$ for each $t \in G_M \backslash G_{M-1}$,

^{*}see the comment following (11)

trG_{M-1}, sin

^{there} exists

 $t \in G^{J+1} \setminus G^{J}$. the first con

^{the system} (

Since solution Y(

(14)

Befor

point out that of the iterat

extrapolation

(13)
$$Y(t, h_{M}) = \sum_{j=0}^{M-1} h_{M}^{(M-j)q} P_{M-j}(t)$$

$$= \varphi(t) + O(h_{M}^{(M+1)q}) + \sum_{j=0}^{M-1} O(h_{M}^{\beta M-j})$$

$$= \varphi(t) + O(h_{M}^{\beta 1})$$

since $\beta_1 \leq \beta_2 \leq \dots \leq \beta_M \leq (M+1)q$.

To obtain the final solution, Y(t), at all grid points in G_M , the finest grid, we proceed as follows. For t=a+h, Y(t) is the first component, $\bar{\phi}(t)$, of the computed solution to the system (4). For $t\in G_M\setminus G_{M-1}$, Y(t) is the solution of (13). For any $t\in G_{M-1}$, since $\bigcup_{k=0}^M T_k = G_M$ and the G_k 's are nested, there exists an index J, depending on t, such that $t\in G_{J+1}\setminus G_J$. The solution at this t, Y(t), will be the first component, $\bar{\phi}(t)$, of the computed solution to the system (11).

Since β_1 is the smallest of the β_j 's, the solution Y(t) constructed above satisfies

(14)
$$Y(t) = \varphi(t) + O(h^{1}).$$

Before we examine the β_j 's in detail we should point out that the pullback method does not require use of the iterative scheme of Aitken and Neville to perform extrapolation. In the pullback method extrapolation is

actually acsystems.

Ιt

for the sysmatrix for matrix of the matrix

Spensystem (4)

progressive

(15)

The matrix

Tis matrix

actually accomplished when we solve the various matrix systems.

It should also be pointed out that the matrices for the systems (11) are "nested" in the sense that the matrix for the Jth system is an easily obtained submatrix of the (J-1)st system. Thus we need only define the matrix for (4) and the others can be obtained by progressively deleting rows and columns of this matrix.

Specifically, recall that the matrix for the system (4) - the case j=0 - is given by

(15)
$$\begin{bmatrix} 1, & h_0^q, & h_0^{2q}, & \dots & h_0^{(M-1)q}, & h_0^{Mq} \\ 1, & h_1^q, & h_1^{2q}, & \dots & h_1^{(M-1)q}, & h_1^{Mq} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1, & h_M^q, & h_M^{2q}, & \dots & h_M^{(M-1)q}, & h_M^{Mq} \end{bmatrix}$$

The matrix for the system (7) - the case j=l - is given by

$$\begin{bmatrix} 1, & h_1^q, & h_1^{2q}, & \dots, & h_1^{(M-1)q} \\ 1, & h_2^q, & h_2^{2q}, & \dots, & h_2^{(M-1)q} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots \\ 1, & h_M^q, & h_M^{2q}, & \dots, & h_M^{(M-1)q} \end{bmatrix}$$

This matrix can be obtained from (15) by deleting the first row and the last column of (15).

For $g \in S$ is given by

1,

which can eas:

rows and the

Consecutive it is only necessarily necessarily

For 1

have developed vandermonde s

system.

Turnine the consection.

For general J the matrix of the system (11) is given by

$$\begin{bmatrix} 1, & h_{J+1}^{q}, & h_{J+1}^{2q}, & \dots, & h_{J+1}^{(M-J-1)q} \\ 1, & h_{J+2}^{q}, & h_{J+2}^{2q}, & \dots, & h_{J+2}^{(M-J-1)q} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1, & h_{M}^{q}, & h_{M}^{2q}, & \dots, & h_{M}^{(M-J-1)q} \end{bmatrix}$$

which can easily be obtained by deleting the first J+1 rows and the last J+1 columns of the matrix (15).

Consequently, when using the pullback technique it is only necessary to define one matrix.

For large M the Vandermonde matrix (15) is known to be ill-conditioned. Björck and Pereyra [2] have developed an efficient algorithm for solving such Vandermonde systems, which can be utilized for solving our system.

Turning our attention to the $\beta_{\mbox{\it j}}$'s, we shall examine the cases $q{=}1$ and $q{=}2$ in detail in the next section.

Section

<u>Lemma l</u>

Proof.

Hence o

with

x is

the v

this

(16)

guq

(17)

The

Section 3. The Pullback Method for q=1 and 2

We will need the following rather obvious lemma.

Lemma 1. a) Let $f(x) = -2x + 2^{x+1}$ then $\forall x \in [1, \infty)$, $f(x) \ge 2$.

b) Let $g(x) = -x + 2^{x+1}$ then $\forall x \in [1, \infty), g(x) \ge 2$.

<u>Proof.</u> a) $f'(x) = -2 + 2^{x+1} \ln 2 > -2 + 2^x \ge 0$ for $x \ge 1$. Hence on $[1, \infty)$ f(x) is a monotone increasing function with f(1) = 2 and therefore $f(x) \ge 2$.

b) $g(x) = f(x) + x \ge f(x) \ge 2$ by part a) since x is positive. \square

Our previous inductive argument has established the validity of equation (10) for j=1,...,M-1. From this and (5) we can conclude for q=2 that

$$\beta_{M} = \min((M+1)2, 2M+3) = (M+1)2$$
(16)
$$\beta_{M-j} = \min(\beta_{M-j+1}, (M-j)2 + 2^{j} + 2), \quad j=1,...,M-1;$$

and for q=1 that

$$\beta_{M} = \min(M+1, M+3) = M+1$$
(17)
$$\beta_{M-1} = \min(\beta_{M-j+1}, M-j+2j+2), \qquad j=1,...,M-1.$$

Theorem 1. For q=1 and 2 we have

$$\beta_{M-j} = (M+1)q, j=0,1,...,M-1.$$

Proof.

and (17)

(16) and

from ou

for q=

q=1 or

which e

the pul at ever

extrap

the ac

<u>Proof.</u> The proof is by induction on j. From (16) and (17) we have $\beta_M = (M+1)q$ and for j=1 we have

$$\beta_{M-1} = \min (\beta_{M}, (M-1)q + 2 + 2)$$

$$= \min ((M+1)q, (M-1)q + 2 + 2)$$

$$= (M+1)q.$$

Assume $\beta_{M-j} = (M+1)q$ where 1 < j < M-1. Then by (16) and (17)

$$\beta_{M-(j+1)} = \beta_{M-j-1} = \min(\beta_{M-j}, (M-j-1)q + 2^{j+1} + 2)$$

$$= \min((M+1)q, Mq-qj + 2^{j+1} - q + 2)$$

from our inductive assumption. By Lemma 1 $-qj + 2^{j+1} \ge 2$ for q=1 or 2 and obviously $-q + 2 \ge 0$. Thus for q=1 or 2

$$(M+1)q \le Mq + 2 \le Mq - qj + 2^{j+1} - q + 2$$

which establishes that $\beta_{M-(j+1)} = (M+1)q$. \square

Theorem 1 shows that for the cases q=1 and 2 the pullback interpolation method yields the same accuracy at every point of the finest grid as that obtained by extrapolation at the endpoint, a+h.

If M is large then the polynomials P_0 , P_1 , etc.; will be of high degree and "Runge's phenomenon" may destroy the accuracy of our computed solution. This can be

nomials. Exalimitation on to which we kneed only intaccuracy compaffect the overwhere.

Befor

PM-j we need

e'm(a) is giv

few more comm

pullback meth

an L-dimensi

and ej(t) a

the applicabl

solved L ti

Also a comple

must be compu

point out the

computers in

In th

There all the a

circumvented to some degree by using lower degree polynomials. Examining the proof of Theorem 1 we see that the limitation on the accuracy we can obtain is the accuracy to which we know e_{M-j} . Thus for large M and j we need only interpolate on enough points to guarantee accuracy comparable to that of our data. This will not affect the overall accuracy of our scheme and may avoid "Runge's phenomenon".

Before we can actually compute the polynomials P_{M-j} we need to know $e_{M-j}^{\dagger}(a)$. The method for determining $e_{m}^{\dagger}(a)$ is given in the next section but at this point a few more comments on the actual implementation of the pullback method are in order .

In the case where $y(t) = (y^1(t)...,y^k(t))$ is an k-dimensional vector, the functions $Y(t,h_k)$, $\varphi(t)$ and $e_j(t)$ are themselves k-dimensional. Consequently the applicable system - (4), (7) or (11) - must be solved k times, once for each component of the vectors. Also a complete set of interpolation polynomials P_{M-j} must be computed for each component. We again wish to point out the potential for using parallel processing computers in this case.

There are cases where y is -dimensional and not all the above work is required. Specifically, suppose

order differe order differe one is intere equation and

of the solution on [0,b] who the problem (equation on

If the

as the initia:

In the solved 1 time at h accurate polynomials of components of

the *L*-dimensional problem resulted from reducing an *L*th order differential equation to a system of *L* first order differential equations. Furthermore, suppose that one is interested only in the solution to the original equation and not in its various derivatives.

If the solution is desired on [0,h] then one need utilize the pullback technique on the first component of the solution vector only. If the solution is desired on [0,b] where b>h the standard procedure is to solve the problem (1) on [0,h], then solve the same differential equation on $[h,h+\bar{h}]$ using the computed solution at h as the initial condition for the new problem. This procedure may be repeated as many times as necessary to reach b.

In this latter case, the system (4) must be solved ℓ times since we need to know all the derivatives at h accurately. However the other systems and the polynomials P_{M-j} need to be computed only for the first components of the vectors.

Section 4. D

In or P_{M-j}, of the determine e_M is given in t carried out f

results, we wi are for the ca plifies the ad valid for, and ensional case,

In ord

Denote at the theoret

J(t) =

and define sym ^a≤t≤b, by

f^(k) (t, φ(t))

for k=2,3,...

£(k)

vative of f(t In the case who

Section 4. Determination of em(a)

In order to construct the Hermite polynomials, P_{M-j} , of the previous section it is necessary to first determine $e'_{M-j}(a)$. The method for computing $e'_{m}(a)$ is given in this section. The actual computations are carried out for $m=1,\ldots,4$ and examples are given.

In order to make effecient use of Gragg's [11,12] results, we will follow his notation. The examples given are for the case when y is 1-dimensional. This simplifies the actual computations. However, the theory is valid for, and will be presented for, the general k-dimensional case, $y=(y^1,\ldots,y^k)$.

Denote by J the Jacobian matrix of f, evaluated at the theoretical solution $\,\phi_{}$

$$J(t) = \frac{\partial f(t, \phi(t))}{\partial y} \qquad a \le t \le b$$

and define symmetric k-linear operators $f^{(k)}(t,\phi(t))$, $a \le t \le b$, by

$$\mathbf{f^{(k)}}(\mathsf{t},\varphi(\mathsf{t}))_{1}\cdots_{k} = \sum_{l_{1}=1}^{l}\cdots\sum_{k_{k}=1}^{l}\frac{\partial^{k}\mathbf{f}(\mathsf{t},\varphi(\mathsf{t}))}{\partial y_{1}\cdots\partial y_{k}}y_{1}^{l_{1}}\cdots y_{k}^{l_{k}}$$

for $k=2,3,\ldots$.

 $f^{(k)}(t, \varphi(t))y_1...y_k$ is the k^{th} Fréchet derivative of $f(t, \varphi(t))$ operating on the vectors $y_1...y_k$. In the case where y is 1-dimensional, J(t) is the

à!

(

ġ:

(

()

1

first partial with respect to y of f and $f^{(k)}(t, \varphi(t))y_1 \cdots y_k = \frac{\partial^{(k)} f(t, \varphi(t))}{\partial y^k} y_1 \cdots y_k$; that is, $f^{(k)}(t, \varphi(t)) y_1 \cdots y_k$

is the k^{th} partial of f with respect to y multiplied by y_1, \ldots, y_k . The y_j 's are functions of t on which $f^{(k)}(t, \varphi(t))$ operates.

Let

(18a)
$$e_{O}(t) \equiv \varphi(t)$$

and, for m=1,2,..., let $e_m(t)$ satisfy

(18b)
$$e_{m}'(t) = J(t)e_{m}(t) + a_{m}(t) + b_{m}(t)$$
$$e_{m}(a) = 0 \qquad a \le t \le b,$$

where

(18c)
$$a_{m}(t) = -\sum_{k=1}^{m} \alpha_{k} e_{m-k}^{(qk+1)}(t)$$

and

(18d)
$$\sum_{m=1}^{\infty} b_m(t) z^m = \sum_{k=2}^{\infty} \frac{1}{k!} f^{(k)}(t, \varphi(t)) \left(\sum_{m=1}^{\infty} e_m(t) z^m \right)^k.$$

The integer q and the constants $\alpha_{\mathbf{k}}$ are determined from a generating function

(18e)
$$A(z) = \sum_{k=0}^{\infty} \alpha_k z^{qk}$$

Gragg [11,12] has shown that both Euler's rule and the trapezoid rule have asymptotic error expansions of the form (2) with the coefficient functions determined by

(18). For Euler's rule the generating function is

(19)
$$A(z) = \sum_{k=0}^{\infty} \frac{1}{(k+1)!} z^k = \frac{e^z - 1}{z}.$$

For the trapezoid rule

(20)
$$A(z) = \frac{2}{z} \tanh\left(\frac{z}{2}\right)$$

$$= 1 - \frac{1}{3}\left(\frac{z}{2}\right)^2 + \frac{2}{15}\left(\frac{z}{2}\right)^4 - \frac{17}{315}\left(\frac{z}{2}\right)^6 + \frac{62}{2835}\left(\frac{z}{2}\right)^8$$

$$- \dots + \frac{(-1)^{n+1}2^{2n}(2^{2n}-1)}{(2n)!}B_n\left(\frac{z}{2}\right)^{2n-1} + \dots$$

where B_n is the nth Bernoulli number.

From (18a) $e_{O}(t) \equiv \varphi(t)$ and from (1) $\varphi'(t) = f(t, \varphi(t)), \varphi(a) = \alpha$; therefore we can find higher order derivatives of $e_{O}(t)$ at t=a by successively computing total derivatives of f and evaluating them at t=a:

(21)
$$e_0^{(p)}(a) = \varphi^{(p)}(a) = \frac{d^{p-1}f(t,\varphi(t))}{dt^{p-1}}\Big|_{t=a}, p=1,2,...$$

From (18b) $e'_1(t) = J(t)e_1(t) + a_1(t) + b_1(t)$ and from (18d) $b_1(t) \equiv 0$. Using (18c), we find $a_1(t) = -\alpha_1 e_0^{(qk+1)}(t)$. Thus

(22)
$$e_1^{(p)}(a) = \frac{d^{p-1}[J(t)e_1(t)]}{dt^{p-1}}\Big|_{t=a} -\alpha_1 e_0^{(qk+p)}(a).$$

The second term on the right of (22) is known from (21) and since J(t) is known, the first term on the right

can be computed because it involves only lower order derivatives of $e_1(t)$.

Proceeding inductively we assume that all derivatives of e_0, e_1, \dots, e_{j-1} at the point t=a can be evaluated, then from (18b)

$$e'_{j}(t) = J(t)e_{j}(t) + a_{j}(t) + b_{j}(t)$$

so that

(23)
$$e_{j}^{(p)}(a) = \frac{d^{p-1}[J(t)e_{j}(t)]}{dt^{p-1}}\Big|_{t=a} + a_{j}^{(p-1)}(t)\Big|_{t=a} + b_{j}^{(p-1)}(t)\Big|_{t=a}$$

Let us now consider each term on the right of (23) in turn.

The derivatives of $e_j(t)$ appearing in $\frac{d^{p-1}[J(t)e_j(t)]}{dt^{p-1}}$ are all of order less than p and can be evaluated at t=a successively from the lowest to the highest. Also J(t) and all its derivatives can be computed. Thus, the first term on the right of (23) is known.

From (18c) we see that $a_j(t)$ depends on various derivatives of the error functions $e_0, e_1, \ldots, e_{j-1}$ all of which are known at the point t=a by our inductive hypothesis. Consequently $a_j^{(p-1)}(t) \Big|_{t=a}$ is known.

Finally, $b_j(t)$ can be determined by collecting the proper coefficients of z^j on the right hand side of (18d). $b_j(t)$ will depend on various Frechet derivatives of f operating on various error functions. The Frechet derivatives can be computed. The error functions appearing as arguments of the Frechet derivatives are from the collection $e_0, e_1, \dots e_{j-1}$. This is easily seen by observing that the outer sum on the right hand side of (18d) begins with k=2, which precludes the possibility of $e_j(t)$ appearing as an argument of a Frechet derivative in the expression for $b_j(t)$.

Once $b_j(t)$ is determined it is necessary to find its (p-1)st derivative. The Frechet derivatives can be differentiated with respect to t and evaluated at t=a and the derivatives of the error functions are known at t=a by our inductive hypothesis. Thus, $b_j^{(p-1)}(t)|_{t=a}$ can be computed and equation (23) is valid.

Actually, equation (23) can be expressed entirely in terms of $e_0(t)$, J(t), the Frechet derivatives and various derivatives of these functions evaluated at t=a. This will be done later for the case q=1.

Theoretically, at least, we can obtain $e_m^{\prime}(a)$ for any m=0,1,.... Computationally, the larger m is, the more complicated the expression to be evaluated. As we

shall see, the computations of $e_1'(a)$, $e_2'(a)$ and $e_3'(a)$ are fairly easy to perform for q=1 or q=2. In addition, $e_A'(a)$ is reasonable.

We consider the cases q=1 and q=2 separately. In what follows we will use a superscript in parentheses to denote differentiation with respect to t. The sole exception to this notation will be that $f^{(k)}$ will continue to denote the k^{th} Frechet derivative of f.

Assuming the validity of (18), with q=1, (18c) and (18e) become

(18c) '
$$a_{m}(t) = -\sum_{k=1}^{m} \alpha_{k} e_{m-k}^{(k+1)}(t)$$

and

(18e)'
$$A(z) = \sum_{k=0}^{\infty} \alpha_k z^k$$

respectively.

From (18b), with t=a, we have $e_m^{(1)}(a) = a_m(a) + b_m(a)$. However, since $e_m(a) = 0$, $\forall m$, we can conclude from (18d) that $b_m(a) = 0$, $\forall m$. Thus,

(24)
$$e_m^{(1)}(a) = a_m(a)$$

For
$$m=1$$
 $a_1(t) = -\alpha_1 e_0^{(2)}(t)$, so that

(25)
$$a_1^{(p)}(a) = -\alpha_1 e_0^{(p+2)}(a)$$
 for any p.

In particular, we have

(26)
$$e_1^{(1)}(a) = -\alpha_1 e_0^{(2)}(a)$$
.

In order to determine $e_m^{(1)}(a)$ for m=2,3,4 we will need the second through fourth derivatives of $e_1(t)$ evaluated at t=a. To obtain these we proceed as follows: differentiating (18b) with m=1, we obtain

$$e_1^{(2)}(t) = \frac{d[J(t)e_1(t)]}{dt} + a_1^{(1)}(t) + b_1^{(1)}(t).$$

Now $b_1(t) \equiv 0$ and $a_1(t) = -\alpha_1 e_0^{(2)}(t)$. Therefore

(27a)
$$e_{1}^{(2)}(t) = \frac{d[J(t)e_{1}(t)]}{dt} - \alpha_{1}e_{0}^{(3)}(t)$$
$$= J^{(1)}(t)e_{1}(t) + J(t)e_{1}^{(1)}(t) - A_{1}e_{0}^{(3)}(t).$$

Since $e_1(a) = 0$ and $e_1^{(1)}(a) = -\alpha_1 e_0^{(2)}(a)$ we have

(27b)
$$e_{1}^{(2)}(a) = -\alpha_{1}J(a)e_{0}^{(2)}(a) - \alpha_{1}e_{0}^{(3)}(a)$$
$$= -\alpha_{1}[e_{0}^{(3)}(a) + J(a)e_{0}^{(2)}(a)].$$

Differentiating (27a) we have

(28a)
$$e_{1}^{(3)}(t) = \frac{d^{2}[J(t)e_{1}(t)]}{dt^{2}} - \alpha_{1}e_{0}^{(4)}(t)$$

$$= J^{(2)}(t)e_{1}(t) + 2J^{(1)}(t)e_{1}^{(1)}(t) + J(t)e_{1}^{(2)}(t)$$

$$-\alpha_{1}e_{0}^{(4)}(t)$$

by Liebnitz's rule. Using (26), (27b) and $e_1(a) = 0$, we can evaluate (28a) for t=a and obtain

(28b)
$$e_1^{(3)}(a) = -2J^{(1)}(a)\alpha_1 e_0^{(2)}(a) - J(a)\alpha_1 [e_0^{(3)}(a) + J(a)e_0^{(2)}(a)]$$

$$-\alpha_1 e_0^{(4)}(a)$$

$$= -\alpha_1 \{ [2J^{(1)}(a) + (J(a))^2]e_0^{(2)}(a) + J(a)e_0^{(3)}(a) + e_0^{(4)}(a) \}.$$

Differentiating (28a) we have

(29a)
$$e_1(t) = \frac{d^3[J(t)e_1(t)]}{dt^3} - \alpha_1 e_0^{(5)}(t)$$

$$= J^{(3)}(t)e_1(t) + 3J^{(2)}(t)e_1^{(1)}(t) + 3J^{(1)}(t)e_1^{(2)}(t)$$

$$+J(t)e_1^{(3)}(t) - \alpha_1 e_0^{(5)}(t).$$

Using (26), (27b), (28b) and $e_1(a) = 0$ we have

(29b)
$$e_{1}^{(4)}(a) = -3\alpha_{1}J^{(2)}(a)e_{0}^{(2)}(a)-3\alpha_{1}J^{(1)}(a)[e_{0}^{(3)}(a)+J(a)e_{0}^{(2)}(a)]$$

$$-\alpha_{1}J(a)\{[2J^{(1)}(a)+(J(a))^{2}]e_{0}^{(2)}(a)+J(a)e_{0}^{(3)}(a)$$

$$+e_{0}^{(4)}(a)\}-\alpha_{1}e_{0}^{(5)}(a)$$

$$= -\alpha_{1}\{[3J^{(2)}(a)+5J(a)J^{(1)}(a)+(J(a))^{3}]e_{0}^{(2)}(a)+\frac{1}{2}e_{0}^{(3)}(a)+\frac{1}{2}e_{0}^{(3)}(a)+\frac{1}{2}e_{0}^{(3)}(a)+\frac{1}{2}e_{0}^{(4)}(a)+e_{0}^{(5)}(a)\}.$$

Note that all expressions for the derivatives of e_1 evaluated at t=a involve only the derivatives of e_0 and the derivatives of J with respect to t. The derivatives of J with respect to t are a straightforward calculation and the derivatives of e_0 can be obtained by successively differentiating and evaluating

,

the original equation (1).

For m=2;
$$a_2(t) = -\alpha_1 e_1^{(2)}(t) - \alpha_2 e_0^{(3)}(t)$$
 so that

(30)
$$a_2^{(p)}(a) = -\alpha_1 e_1^{(p+2)}(a) - \alpha_2 e_0^{(p+3)}(t)$$
 for any p.

In particular using (24) and (27b) we have

(31)
$$e_2^{(1)}(a) = a_2(a) = -\alpha_1 e_1^{(2)}(a) - \alpha_2 e_0^{(3)}(a)$$

$$= \alpha_1^2 [e_0^{(3)}(a) + J(a) e_0^{(2)}(a)] - \alpha_2 e_0^{(3)}(a)$$

$$= (\alpha_1^2 - \alpha_2) e_0^{(3)}(a) + \alpha_1^2 J(a) e_0^{(2)}(a).$$

Now $e_m^{(1)}$ (a) for m=3,4 will require knowledge of the second and third derivatives of e_2 (t) evaluated at t=a. Differentiating (18b) with m=2, we obtain

(32)
$$e_2^{(2)}(t) = \frac{d[J(t)e_2(t)]}{dt} + a_2^{(1)}(t) + b_2^{(1)}(t)$$
.

From (18d) we have

(33)
$$b_2(t) = \frac{1}{2}f^{(2)}(t,\varphi(t))e_1(t)e_1(t)$$

implying that

$$b_2^{(1)}(t) = \frac{1}{2} \frac{df^{(2)}(t, \phi(t))}{dt} e_1(t) e_1(t) + f^{(2)}(t, \phi(t)) e_1(t) e_1'(t).$$

Note that $f^{(2)}(t, \varphi(t)) \equiv f_{yy}(t, \varphi(t))$ for the one dimensional problem. Since $e_1(a) = 0$ we can conclude that $b_2^{(1)}(a) = 0$.

Using (30), (31) and $e_2(a) = 0$ we have

$$e_{2}^{(2)}(a) = J(a)e_{2}^{(1)}(a) + a_{2}^{(1)}(a)$$

$$= -\alpha_{1}J(a)e_{1}^{(2)} - \alpha_{2}J(a)e_{0}^{(3)}(a) - \alpha_{1}e_{1}^{(3)}(a) - \alpha_{2}e_{0}^{(4)}(a)$$

$$= -\alpha_{1}(e_{1}^{(3)}(a) + J(a)e_{1}^{(2)}(a)) - \alpha_{2}(e_{0}^{(4)}(a) + J(a)e_{0}^{(3)}(a))$$

Hence, $e_2^{(2)}$ (a) can be expressed in terms of derivatives of e_0 (a) by using (27b) and (28b),

$$e_{2}^{(2)}(a) = \alpha_{1}^{2}\{[2J^{(1)}(a)+(J(a))^{2}]e_{0}^{(2)}(a)+J(a)e_{0}^{(3)}(a)+e_{0}^{(4)}(a) + (J(a))^{2}e_{0}^{(2)}(a)+J(a)e_{0}^{(3)}(a)\}-\alpha_{2}(e_{0}^{(4)}(a)+J(a)e_{0}^{(3)}(a)).$$

Collecting terms involving the various derivatives of $\mathbf{e}_{\mathbf{O}}$ (a), we have

(34)
$$e_2^{(2)}(a) = (\alpha_1^2 - \alpha_2) e_0^{(4)}(a) + (2\alpha_1^2 - \alpha_2) J(a) e_0^{(3)}(a) + 2\alpha_1^2 [J^{(1)}(a) + (J(a))^2] e_0^{(2)}(a).$$

Differentiating (32) again, we have

$$e_2^{(3)}(t) = \frac{d^2[J(t)e_2(t)]}{dt^2} + a_2^{(2)}(t) + b_2^{(2)}(t).$$

By differentiating (33) twice and using $e_1(a) = 0$, we find that $b_2^{(2)}(a) = f^{(2)}(a, \varphi(a))e_1^{(1)}(a)e_1^{(1)}(a)$. From

(30),
$$a_2^{(2)}(a) = -\alpha_1 e_1^{(4)}(a) - \alpha_2 e_0^{(5)}(a)$$
 and since $e_2(a) = 0$

we have

$$e_{2}^{(3)}(a) = 2J^{(1)}(a)e_{2}^{(1)}(a)+J(a)e_{2}^{(2)}(a)-\alpha_{1}e_{1}^{(4)}(a)-\alpha_{2}e_{0}^{(5)}(a)$$

$$+f^{(2)}(a,\phi(a))e_{1}^{(1)}(a)e_{1}^{(1)}(a).$$

Substituting for $e_2^{(1)}(a)$, $e_2^{(2)}(a)$, $e_1^{(4)}(a)$ and $e_1^{(1)}(a)$ from equations (31), (34), (29b) and (26) respectively, we have

$$e_{2}^{(3)}(a) = 2(\alpha_{1}^{2} - \alpha_{2})J^{(1)}(a)e_{0}^{(3)}(a) + 2\alpha_{1}^{(2)}J(a)J^{(1)}(a)e_{0}^{(2)}(a)$$

$$+(\alpha_{1}^{2} - \alpha_{2})J(a)e_{0}^{(4)}(a) + (2\alpha_{1}^{(2)} - \alpha_{2})(J(a))^{2}e_{0}^{(3)}(a)$$

$$+(2\alpha_{1}^{2} [J(a)J^{(1)}(a) + (J(a))^{3}]e_{0}^{(2)}(a)$$

$$+\alpha_{1}^{2} [3J^{(2)}(a) + 5J(a)J^{(1)}(a) + (J(a))^{3}]e_{0}^{(2)}(a)$$

$$+\alpha_{1}^{2} [3J^{(1)}(a) + (J(a))^{2}]e_{0}^{(3)}(a) + \alpha_{1}^{2}J(a)e_{0}^{(4)}(a)$$

$$+\alpha_{1}^{2} e_{0}^{(5)}(a) - \alpha_{2}e_{0}^{(5)}(a)$$

$$+\alpha_{1}^{2} f^{(2)}(a,\varphi(a))e_{0}^{(2)}(a)e_{0}^{(2)}(a).$$

Collecting on the derivatives of $e_{\Omega}(a)$, we have

(35)
$$e_{2}^{(3)}(a) = (\alpha_{1}^{2} - \alpha_{2}) e_{0}^{(5)}(a) + (2\alpha_{1}^{2} - \alpha_{2}) J(a) e_{0}^{(4)}(a)$$

$$+ [(5\alpha_{1}^{2} - 2\alpha_{2}) J^{(1)}(a) + (3\alpha_{1}^{2} - \alpha_{2}) (J(a))^{2}] e_{0}^{(3)}(a)$$

$$+ \alpha_{1}^{2} [3J^{(2)}(a) + 9J(a) J^{(1)}(a) + 3 (J(a))^{3}] e_{0}^{(2)}(a)$$

$$+ \alpha_{1}^{2} f^{(2)}(a, \varphi(a)) e_{0}^{(2)}(a) e_{0}^{(2)}(a).$$

With m=3 from (18c) we have

$$a_3(t) = -\alpha_1 e_2^{(2)}(t) - \alpha_2 e_1^{(3)}(t) - \alpha_3 e_0^{(4)}(t)$$
 so that

(36)
$$a_3^{(p)}(a) = -\alpha_1 e_2^{(p+2)}(a) - \alpha_2 e_1^{(p+3)}(a) - \alpha_3 e_0^{(p+4)}(a)$$

for any p. From (24) we have

$$e_3^{(1)}(a) = a_3(a) = -\alpha_1 e_2^{(2)}(a) - \alpha_2 e_1^{(3)}(a) - \alpha_3 e_0^{(4)}(a)$$

Substituting (34) and (28b) for $e_2^{(2)}$ (a) and $e_1^{(3)}$ (a) respectively, we have

$$e_{3}^{(1)}(a) = -\alpha_{1}(\alpha_{1}^{2} - \alpha_{2}) e_{0}^{(4)}(a) - \alpha_{1}(2\alpha_{1}^{2} - \alpha_{2}) J(a) e_{0}^{(3)}(a)$$

$$-2\alpha_{1}^{3}[J^{(1)}(a) + (J(a))^{2}] e_{0}^{(2)}(a)$$

$$+\alpha_{1}\alpha_{2}[2J^{(1)}(a) + (J(a))^{2}] e_{0}^{(2)}(a)$$

$$+\alpha_{1}\alpha_{2}J(a) e_{0}^{(3)}(a) + \alpha_{1}\alpha_{2}e_{0}^{(4)}(a) - \alpha_{3}e_{0}^{(4)}(a).$$

Combining like derivates of $e_{O}(a)$, this becomes

(37)
$$e_3^{(1)}(a) = (2\alpha_1\alpha_2 - \alpha_1^3 - \alpha_3)e_0^{(4)}(a) + 2(\alpha_1\alpha_2 - \alpha_1^3)J(a)e_0^{(3)}(a)$$

 $+[(\alpha_1\alpha_2 - 2\alpha_1^3)(J(a))^2 + 2(\alpha_1\alpha_2 - \alpha_1^3)J^{(1)}(a)]e_0^{(2)}(a).$

Differentiating (18b) with m=3, we have

$$e_{3}^{(2)}(t) = \frac{d[J(t)e_{3}(t)]}{dt} + a_{3}^{(1)}(t) + b_{3}^{(1)}(t).$$
From (18d), $b_{3}(t) = \frac{1}{6} f^{(3)}(t, \phi(t))e_{1}(t)e_{1}(t)e_{1}(t) + f^{(2)}(t, \phi(t))e_{1}(t)e_{2}(t)$ so that $b_{3}^{(1)}(a) = 0$ as

 $e_1(a) = e_2(a) = 0$. Substituting for $a_3^{(1)}(t)$ and using $e_3(a) = 0$, we have

$$e_3^{(2)}(a) = J(a)e_3^{(1)}(a)-\alpha_1e_2^{(3)}(a)-\alpha_2e_1^{(4)}(a)-\alpha_3e_0^{(5)}(a)$$
.

Substituting (37), (35) and (29b) into this equation yields

$$\begin{split} e_{3}^{(2)}(a) &= (2\alpha_{1}\alpha_{2} - \alpha_{1}^{3} - \alpha_{3})J(a)e_{0}^{(4)}(a) + 2(\alpha_{1}\alpha_{2} - \alpha_{1}^{3})(J(a))^{2}e_{0}^{(3)}(a) \\ &+ [(\alpha_{1}\alpha_{2} - 2\alpha_{1}^{3})(J(a))^{3} + 2(\alpha_{1}\alpha_{2} - \alpha_{1}^{3})J(a)J^{(1)}(a)]e_{0}^{(2)}(a) \\ &+ (\alpha_{1}\alpha_{2} - \alpha_{1}^{3})e_{0}^{(5)}(a) + (\alpha_{1}\alpha_{2} - 2\alpha_{1}^{(3)})J(a)e_{0}^{(4)}(a) \\ &+ [(2\alpha_{1}\alpha_{2} - 5\alpha_{1}^{3})J^{(1)}(a) + (\alpha_{1}\alpha_{2}[3\alpha_{1}^{3})(J(a))^{2}]e_{0}^{(3)}(a) \\ &- \alpha_{1}^{3}[3J^{(2)}(a) + 9J(a)J^{(1)}(a) + 3(J(a)]^{3}e_{0}^{(2)}(a) \\ &- \alpha_{1}^{3}f^{(2)}(a,\phi(a))e_{0}^{(2)}(a)e_{0}^{(2)}(a) + \alpha_{1}\alpha_{2}e_{0}^{(5)}(a) \\ &+ \alpha_{1}\alpha_{2}J(a)e_{0}^{(4)}(a) + \alpha_{1}\alpha_{2}[3J^{(1)}(a) + (J(a))^{2}]e_{0}^{(3)}(a) \\ &+ \alpha_{1}\alpha_{2}[3J^{(2)}(a) + 5J(a)J^{(1)}(a) + (J(a))^{3}]e_{0}^{(2)}(a) - \alpha_{3}e_{0}^{(5)}(a) \,. \end{split}$$

This can be rewritten as

(38)
$$e_3^{(2)}(a) = 2(\alpha_1\alpha_2 - \alpha_1^3 - \alpha_3)e_0^{(5)}(a) + (4\alpha_1\alpha_2 - 3\alpha_1^3 - \alpha_3)J(a)e_0^{(4)}(a)$$

$$+[5(\alpha_1\alpha_2 - \alpha_1^3)J^{(1)}(a) + (4\alpha_1\alpha_2 - 5\alpha_1^3)(J(a))^2]e_0^{(3)}(a)$$

$$+[3(\alpha_1\alpha_2 - \alpha_1^3)J^{(3)}(a) + (7\alpha_1\alpha_2 - 11\alpha_1^3)J(a)J^{(1)}(a) +$$

$$(2\alpha_1\alpha_2 - 5\alpha_1^3)(J(a))^3]e_0^{(2)}(a)$$

$$-\alpha_1^3f^{(2)}(a, \varphi(a))e_0^{(2)}(a)e_0^{(2)}(a).$$

Finally for m=4 we have

$$\begin{split} e_4^{(1)} &(a) = a_4(a) = -\alpha_1 e_3^{(2)} (a) - \alpha_2 e_2^{(3)} (a) - \alpha_3 e_1^{(4)} (a) - \alpha_4 e_0^{(5)} (a) \,. \\ &\text{Using (38), (35) and (29b), this can be written as} \\ e_4^{(1)} &(a) = 2 \left(\alpha_1^4 + \alpha_1 \alpha_3 - \alpha_1^2 \alpha_2 \right) e_0^{(5)} (a) + \left(3 \alpha_1^4 + \alpha_1 \alpha_3 - 4 \alpha_1^2 \alpha_2 \right) J(a) e_0^{(4)} (a) \\ &\quad + \left[5 \left(\alpha_1^4 - \alpha_1^2 \alpha_2 \right) J^{(1)} (a) + \left(5 \alpha_1^4 - 4 \alpha_1^2 \alpha_2 \right) (J(a))^2 \right] e_0^{(3)} (a) \\ &\quad + \left[3 \left(\alpha_1^4 - \alpha_1^2 \alpha_2 \right) J^{(2)} (a) + \left(11 \alpha_1^4 - 7 \alpha_1^2 \alpha_2 \right) J(a) J^{(1)} (a) + \\ &\quad + \left(5 \alpha_1^4 - 2 \alpha_1^2 \alpha_2 \right) J^{(2)} (a) + \left(11 \alpha_1^4 - 7 \alpha_1^2 \alpha_2 \right) J(a) J^{(1)} (a) + \\ &\quad + \left(5 \alpha_1^4 - 2 \alpha_1^2 \alpha_2 \right) J^{(2)} (a) e_0^{(2)} (a) \\ &\quad + \left(\alpha_2^2 - \alpha_1^2 \alpha_2 \right) e_0^{(5)} (a) + \left(\alpha_2^2 - 2 \alpha_1^2 \alpha_2 \right) J(a) e_0^{(4)} (a) \\ &\quad + \left(\left(2 \alpha_2^2 - 5 \alpha_1^2 \alpha_2 \right) J^{(1)} (a) + \left(\alpha_2^2 - 3 \alpha_1^2 \alpha_2 \right) (J(a))^2 \right] e_0^{(3)} (a) \\ &\quad - \alpha_1^2 \alpha_2 \left[3 J^{(2)} (a) + 9 J(a) J^{(1)} (a) + 3 (J(a))^3 \right] e_0^{(2)} (a) \\ &\quad + \alpha_1 \alpha_3 e_0^{(5)} (a) + \alpha_1 \alpha_3 J(a) e_0^{(4)} (a) \\ &\quad + \alpha_1 \alpha_3 \left[3 J^{(1)} (a) + \left(J(a) \right)^2 \right] + \left(J(a) \right) e_0^{(3)} (a) \\ &\quad + \alpha_1 \alpha_3 \left[3 J^{(2)} (a) + 5 J(a) J^{(1)} (a) + \left(J(a) \right)^3 \right] e_0^{(2)} (a) \\ &\quad - \alpha_2 e_0^{(5)} (a) \end{aligned}$$

Rearranging terms, we get

^{§tep 2} Comp

Step 3

e(3)
e(4)
e(4)

(39)
$$e_{4}^{(1)}(a) = (2\alpha_{1}^{4} - 3\alpha_{1}^{2}\alpha_{2} + 3\alpha_{1}\alpha_{3} + \alpha_{2}^{2} - \alpha_{4})e_{0}^{(5)}(a)$$

$$+ (3\alpha_{1}^{4} - 6\alpha_{1}^{2}\alpha^{2} + 2\alpha_{1}\alpha_{3} + \alpha_{2}^{2})e_{0}^{(4)}(a)$$

$$+ [(5\alpha_{1}^{4} - 10\alpha_{1}^{2}\alpha_{2} + 3\alpha_{1}\alpha_{3} + 2\alpha_{2}^{2})J^{(1)}(a) +$$

$$(5\alpha_{1}^{4} - 7\alpha_{1}^{2}\alpha_{2} + \alpha_{1}\alpha_{3} + \alpha_{2}^{2})(J(a))^{2}]e_{0}^{(3)}(a)$$

$$+ [3(\alpha_{1}^{4} - 2\alpha_{1}^{2}\alpha_{2} + \alpha_{1}\alpha_{3})J^{(2)}(a) + (11\alpha_{1}^{4} - 16\alpha_{1}^{2}\alpha_{2} +$$

$$5\alpha_{1}\alpha_{3})J(a)J^{(1)}(a) + (5\alpha_{1}^{4} - 5\alpha_{1}^{2}\alpha_{3} + \alpha_{1}\alpha_{3})(J(a))^{3}]e_{0}^{(2)}(a)$$

$$+ (\alpha_{1}^{4} - \alpha_{1}^{2}\alpha_{2})f^{(2)}(a, \phi(a))e_{0}^{(2)}(a)e_{0}^{(2)}(a).$$

To summarize, equations (26), (31), (37) and (39) can be used to determine $e_1^{(1)}(a)$, $e_2^{(1)}(a)$, $e_3^{(1)}(a)$ and $e_4^{(1)}(a)$ respectively. A computationally more effective recursive procedure for q=1 is given in Table 3 below.

TABLE 3

Step 1 Compute
$$e_0^{(1)}(a)$$
, $e_0^{(2)}(a)$, $e_0^{(3)}(a)$, $e_0^{(4)}(a)$, $e_0^{(5)}(a)$

Step 2 Compute $J(a)$, $J^{(1)}(a)$, $J^{(2)}(a)$,

Step 3 $e_1^{(1)}(a) = -\alpha_1 e_0^{(2)}(a)$
 $e_1^{(2)}(a) = J(a)e_1^{(1)}(a) - \alpha_1 e_0^{(3)}(a)$
 $e_1^{(3)}(a) = 2J^{(1)}(a)e_1^{(1)}(a) + J(a)e_1^{(2)}(a) - \alpha_1 e_0^{(4)}(a)$
 $e_1^{(4)}(a) = 3J^{(2)}(a)e_1^{(1)}(a) + 3J^{(1)}(a)e_1^{(2)}(a) + J(a)e_1^{(3)}(a)$
 \vdots
 $-\alpha_1 e_0^{(5)}(a)$

e4

simply the

as they bec

TABLE 3 continued

Step 4
$$e_2^{(1)}(a) = -\alpha_1 e_1^{(2)}(a) - \alpha_2 e_0^{(3)}(a)$$

 $e_2^{(2)}(a) = J(a) e_2^{(1)}(a) - \alpha_1 e_1^{(3)}(a) - \alpha_2 e_0^{(4)}(a)$
 $e_2^{(3)}(a) = 2J^{(1)}(a) e_2^{(1)}(a) + J(a) e_2^{(2)}(a) - \alpha_1 e_1^{(4)}(a)$
 $\vdots - \alpha_2 e_0^{(5)}(a) + f^{(2)}(a, \varphi(a)) e_1^{(1)}(a) e_1^{(1)}(a)$

Step 5
$$e_3^{(1)}(a) = -\alpha_1 e_2^{(2)}(a) - \alpha_2 e_1^{(3)}(a) - \alpha_3 e_0^{(4)}(a)$$

 $e_3^{(2)}(a) = J(a)e_3^{(1)}(a) - \alpha_1 e_2^{(3)}(a) - \alpha_2 e_1^{(4)}(a) - \alpha_3 e_0^{(5)}(a)$
.

Step 6
$$e_4^{(1)}(a) = -\alpha_1 e_3^{(2)}(a) - \alpha_2 e_2^{(3)}(a) - \alpha_3 e_1^{(4)}(a) - \alpha_4 e_0^{(5)}(a)$$

.

Table 3 lists the equations necessary for the computation of $e_1^{(1)}(a), \ldots, e_4^{(1)}(a)$ for a method which satisfies (18) with q=1. Of course, Table 3 can be continued. Each new error function, whose first derivative at t=a is to be computed, requires one additional computation at each step of the table and the addition of one more step to the table. The only computations in Table 3 which are not simply the result of plugging in formulae are those in steps 1 and 2. and the computation of various Fréchet derivatives as they become necessary.

To i

 $e_2^{(1)}(a), e_3^{(1)}$

in example 1

Example 2:

the solution

У

From equation

rule. In th

Step 1: e_O(1

(

()

0

e (

 e_{l}^{C}

Step 2:

J

_

To illustrate the use of Table 3 we compute $e_1^{(1)}(a)$, $e_2^{(1)}(a)$, $e_3^{(1)}(a)$ and $e_4^{(1)}(a)$ for the differential equation in example 1 using a numerical method for which q=1.

Example 2: Suppose we are using Euler's rule to compute the solution to

$$y'=y^2; y(0)=.2; 0 \le t \le 1.$$

From equation (19) we have that $\alpha_k = \frac{1}{(k+1)!}$ for Euler's rule. In this case Table 3 becomes

Step 1:
$$e_0^{(1)}(0) = (y(0))^2 = .040000$$

 $e_0^{(2)}(2) = 2y(0)y^{(1)}(0) = .016000$
 $e_0^{(3)}(0) = 2[y(0)y^{(2)}(0)+(y^{(1)}(0))^2] = .009600$
 $e_0^{(4)}(0) = 2[y(0)y^{(3)}(0)+3y^{(1)}(0)y^{(2)}(0)] = .007680$
 $e_0^{(5)}(0) = 2[y(0)y^{(4)}(0)+4y^{(1)}(0)y^{(3)}(0)+3(y^{(2)}(0))^2]$
 $= .007680$

Step 2:
$$J(0) = 2y(0) = .400000$$

$$J^{(1)}(0) = 2y^{(1)}(0) = .080000$$

$$J^{(2)}(1) = 2y^{(2)}(0) = .032000$$

Step 3: e₁ e (

Step 4: e₂

Step 5: e₃

Step 6: e4

Ιt

Steps 1 and Also note t

functions a

of the erro

Tur (18) is val

Step 3:
$$e_1^{(1)}(0) = -.008000$$

 $e_1^{(2)}(0) = -.008000$
 $e_1^{(3)}(0) = -.008520$
 $e_1^{(4)}(0) = -.009936$

Step 4:
$$e_2^{(1)}(0) = .002400$$

$$e_2^{(2)}(0) = .001810$$

$$e_2^{(3)}(0) = .006076 since $f^{(2)}(a, \phi(a)) = 2.$$$

Step 5:
$$e_3^{(1)}(0) = -.007950$$

 $e_3^{(2)}(0) = -.004882$

Step 6:
$$e_4^{(1)}(0) = .0017780$$

It should be emphasized that the computations in

Steps 1 and 2 are performed by the user not by the machine.

Also note that the signs on the derivatives of the error functions alternate in precisely the same manner as the signs of the errors for Euler's rule.

Turning our attention to the case q=2, we assume (18) is valid with q=2 so that (18c) and (18e) become

(18c)"

and

(18e)"

In

analogous

about the

b_l(t

b₂ (t

(40) b₃(t

b₄ (t

el (a)...

necessary

ħave

since e₄ (a

e₄(1) (a

(18c)"
$$a_m(t) = -\sum_{k=1}^{m} \alpha_k e_{m-k}^{(2k+1)}(t)$$

and

(18e)"
$$A(z) = \sum_{k=0}^{\infty} \alpha_k z^{2k}$$
.

In order to construct a table for the case q=2 analogous to Table 3 we will need to collect more information about the functions $b_m(t)$ for $m=1,\ldots,4$. From (18d) we have

$$b_{1}(t) = 0$$

$$b_{2}(t) = \frac{1}{2}f^{(2)}(t, \varphi(t))e_{1}(t)e_{1}(t)$$

$$(40) b_{3}(t) = \frac{1}{6}f^{(3)}(t, \varphi(t))e_{1}(t)e_{1}(t)e_{1}(t)+f^{(2)}(t, \varphi(t))e_{1}(t)e_{2}(t)$$

$$b_{4}(t) = \frac{1}{24}f^{(4)}(t, \varphi(t))e_{1}(t)e_{1}(t)e_{1}(t)e_{1}(t)+\frac{1}{2}f^{(3)}(t, \varphi(t))$$

$$e_{1}(t)e_{1}(t)e_{2}(t)+f^{(2)}(t, \varphi(t))e_{1}(t)e_{3}(t)$$

$$+\frac{1}{2}f^{(2)}(t, \varphi(t))e_{2}(t)e_{2}(t).$$

Since our goal is to produce a table for evaluating $e_1^{(1)}(a), \ldots, e_4^{(1)}(a)$ we examine what information will be necessary to determine $e_4^{(1)}(a)$. From (18b) with m=4 we have

$$e_4^{(1)}(a) = a_4(a)$$

since $e_4(a) = b_4(a) = 0$. From (18c)" we can conclude $e_4^{(1)}(a) = -\alpha_1 e_3^{(3)}(a) - \alpha_2 e_2^{(5)}(a) - \alpha_3 e_1^{(7)}(a) - \alpha_4 e_0^{(9)}(a)$.

In gener

e_m(a)

to deter $e_1^{(7)}$ (a)

Thus in

given be

function

zero. H

b₍₁₎

b₂(2)

p⁵(3)

b(4)

using t

In general, by (18b) we see that

$$e_m^p(a) = \frac{d^{p-1}[J(t)e_m(t)]}{dt^{p-1}} + a_m^{(p-1)}(t) + b_m^{(p-1)}(t); m=1,...,4.$$

Thus in order to determine $e_3^{(3)}(a)$ we need to know $b_3^{(2)}(a)$; to determine $e_2^{(5)}(a)$ we need $b_2^{(4)}(a)$; and to determine $e_1^{(7)}(a)$ we need $b_1^{(6)}(a)$. The necessary computations are given below. To simplify the notation, the arguments of the functions are suppressed.

Since $b_1(t) \equiv 0$ all derivatives of $b_1(t)$ are zero. From (40), $b_2(t) = \frac{1}{2}f^{(2)}(t,\phi(t))e_1(t)e_1(t)$ and therefore;

$$b_{2}^{(1)}(t) = \frac{1}{2} \frac{d(f^{(2)})}{dt} e_{1}e_{1} + f^{(2)}e_{1}e_{1}^{(1)}$$

$$b_{2}^{(2)}(t) = \frac{1}{2} \frac{d^{2}f^{(2)}}{dt^{2}} e_{1}e_{1} + 2(f^{(2)})^{(1)}e_{1}e_{1}^{(1)} + f^{(2)}e_{1}e_{1}^{(2)}$$

$$+ f^{(2)}e_{1}^{(1)}e_{1}^{(1)}$$

$$b_{2}^{(3)}(t) = \frac{1}{2} \frac{d^{3}f^{(2)}}{dt^{3}} e_{1}e_{1} + 3(f^{(2)})^{(2)}e_{1}e_{1}^{(1)} + 3\frac{d(f^{(2)})}{dt} e_{1}e_{1}^{(2)}$$

$$+ 3\frac{d(f^{(2)})}{dt} e_{1}^{(1)}e_{1}^{(1)} + 3f^{(2)}e_{1}^{(1)}e_{1}^{(2)} + f^{(2)}e_{1}e_{1}^{(3)}$$

$$b_{2}^{(4)}(t) = \frac{1}{2} \frac{d^{4}f^{(2)}}{dt^{4}} e_{1}e_{1} + 4\frac{d^{3}f^{(2)}}{dt^{3}} e_{1}e_{1}^{(1)} + 6\frac{d^{2}f^{(2)}}{dt^{2}} e_{1}^{(1)}e_{1}^{(1)}$$

$$+ 6\frac{d^{2}f^{(2)}}{dt^{2}} e_{1}e_{1}^{(2)} + 4\frac{d(f^{(2)})}{dt} e_{1}e_{1}^{(3)} + 12\frac{d(f^{(2)})}{dt} e_{1}^{(1)}e_{1}^{(2)}$$

$$+ 3f^{(2)}e_{1}^{(2)}e_{1}^{(2)} + 4f^{(2)}e_{1}^{(1)}e_{1}^{(3)} + f^{(2)}e_{1}e_{1}^{(4)}.$$

Using the fact that $e_1(a) = 0$ we can obtain

b₂(1)

b₂(2)

(41) b₂(3)

b₂(4)

υs

have

 $b_3^{(1)}(t) =$

 $b_{(2)}^{(2)}(t) =$

Since el

$$b_{2}^{(1)}(a) = 0$$

$$b_{2}^{(2)}(a) = f^{(2)}(a, \varphi(a)) e_{1}^{(1)}(a) e_{1}^{(1)}(a)$$

$$(41) \quad b_{2}^{(3)}(a) = 3[f^{(2)}(a, \varphi(a)) e_{1}^{(1)}(a) e_{1}^{(2)}(a)$$

$$+ \frac{df^{(2)}(t, \varphi(t))}{dt}|_{t=a} e_{1}^{(1)}(a) e_{1}^{(1)}(a)$$

$$b_{2}^{(4)}(a) = 6 \frac{d^{2}f^{(2)}(t, \varphi(t))}{dt^{2}}|_{t=a} e_{1}^{(1)}(a) e_{1}^{(1)}(a)$$

$$+12 \frac{df^{(2)}(t, \varphi(t))}{dt}|_{t=a} e_{1}^{(1)}(a) e_{1}^{(2)}(a)$$

$$+3f^{(2)}(a, \varphi(a)) e_{1}^{(2)}(a) e_{1}^{(2)}(a) +4f^{(2)}(a)$$

$$(a, \varphi(a)) e_{1}^{(1)}(a) e_{1}^{(a)}(a).$$

Using the definition of $b_3(t)$ given in (40), we have

$$b_{3}^{(1)}(t) = \frac{1}{6} \frac{df^{(3)}}{dt} e_{1}e_{1}e_{1}^{1}e_{2}^{1}f^{(3)}e_{1}e_{1}^{(1)}+(f^{(2)})^{(1)}e_{1}e_{2}$$

$$+ f^{(2)}e_{1}^{(1)}e_{2}^{1}+f^{(2)}e_{1}e_{2}^{(1)}$$

$$b_{3}^{(2)}(t) = \frac{1}{6} \frac{d^{2}f^{(3)}}{dt^{2}} e_{1}e_{1}^{1}+\frac{df^{(3)}}{dt} e_{1}e_{1}^{(1)}+f^{(3)}e_{1}e_{1}^{(1)}e_{1}^{(1)}$$

$$+\frac{1}{2}f^{(3)}e_{1}e_{1}^{(2)}+\frac{d^{2}f^{(2)}}{dt^{2}} e_{1}e_{2}^{1}+\frac{df^{(2)}}{dt} e_{1}^{(2)}+e_{1}^{(2)}e_{2}^{(1)}+f^{(2)}e_{1}^{(2)}e_{2}^{(1)}+f^{(2)}e_{1}^{(2)}e_{2}^{(1)}+f^{(2)}e_{1}^{(2)}e_{2}^{(2)}.$$

Since $e_1(a) = e_2(a) = 0$, at t=a we have

(42)

evaluatir

Step 1:

Step 2:

Step 3:

$$b_3^{(1)}(a) = 0$$
(42)
$$b_3^{(2)}(a) = 2f^{(2)}(a, \varphi(a))e_1^{(1)}(a)e_2^{(1)}(a).$$

We can now summarize the computations involved for evaluating $e_1^{(1)}(a), \ldots, e_4^{(1)}(a)$ when q=2.

TABLE 4

Step 1: Compute
$$e_0^{(1)}(a)$$
, $e_0^{(2)}(a)$,..., $e_0^{(8)}(a)$, $e_0^{(9)}(a)$,...
Step 2: Compute $J(a)$, $J^{(1)}(a)$,..., $J^{(5)}(a)$,...
Step 3: $e_1^{(1)}(a) = -\alpha_1 e_0^{(3)}(a)$
 $e_1^{(2)}(a) = J(a) e_1^{(1)}(a) - \alpha_1 e_0^{(4)}(a)$
 $e_1^{(3)}(a) = 2J^{(1)}(a) e_1^{(1)}(a) + J(a) e_1^{(2)}(a) - \alpha_1 e_0^{(5)}(a)$
 $e_1^{(4)}(a) = 3J^{(2)}(a) e_1^{(1)}(a) + 3J^{(1)}(a) e_1^{(2)}(a) + J(a) e_1^{(3)}(a)$
 $e_1^{(4)}(a) = \frac{3}{k=0} {4 \choose k} J^{(k)}(a) e_1^{(4-k)}(a) - \alpha_1 e_0^{(7)}(a)$
 $e_1^{(6)}(a) = \frac{4}{k=0} {5 \choose k} J^{(k)}(a) e_1^{(5-k)}(a) - \alpha_1 e_0^{(8)}(a)$
 $e_1^{(7)}(a) = \frac{5}{k=0} {6 \choose k} J^{(k)}(a) e_1^{(6-k)}(a) - \alpha_1 e_0^{(9)}(a)$
 \vdots

Step 4: b₂(2)
b₂(2)

b₂

Step 5: e

е

6

TABLE 4 continued

Step 4:
$$b_2^{(1)}(a) = 0$$

$$b_2^{(2)}(a) = f^{(2)}(a, \varphi(a)) e_1^{(1)}(a) e_1^{(1)}(a)$$

$$b_2^{(3)}(a) = 3[f^{(2)}(a, \varphi(a)) e_1^{(1)}(a) e_1^{(2)}(a)$$

$$+ \frac{df^{(2)}(t, \varphi(t))}{dt} \Big|_{t=a} e_1^{(1)}(a) e_1^{(1)}(a)]$$

$$b_2^{(4)}(a) = 6 \frac{d^2 f^{(2)}(t, \varphi(t))}{dt} \Big|_{t=a} e_1^{(1)}(a) e_1^{(1)}(a)$$

$$+ 12 \frac{df^{(2)}(t, \varphi(t))}{dt} \Big|_{t=a} e_1^{(1)}(a) e_1^{(2)}(a)$$

$$+ 3f^{(2)}(a, \varphi(a)) e_1^{(2)}(a) e_1^{(2)}(a)$$

$$\vdots$$
Step 5: $e_2^{(1)}(a) = -\alpha_1 e_1^{(3)}(a) - \alpha_2 e_0^{(5)}(a)$

$$e_2^{(2)}(a) = J(a) e_2^{(1)}(a) - \alpha_1 e_1^{(4)}(a) - \alpha_2 e_0^{(6)}(a)$$

$$e_2^{(3)}(a) = 2J^{(1)}(a) e_2^{(1)}(a) + J(a) e_2^{(2)}(a) - \alpha_1 e_1^{(5)}(a)$$

$$-\alpha_2 e_0^{(7)}(a) + b_2^{(2)}(a)$$

$$e_2^{(4)}(a) = \sum_{k=0}^{2} {3 \choose k} J^{(k)}(a) e_2^{(3-k)}(a) - \alpha_1 e_1^{(6)}(a)$$

$$-\alpha_2 e_0^{(8)}(a) + b_2^{(3)}(a)$$

$$e_2^{(5)}(a) = \sum_{k=0}^{3} {4 \choose k} J^{(k)}(a) e_2^{(4-k)}(a) - \alpha_1 e_1^{(7)}(a)$$

$$\vdots$$

$$-\alpha_2 e_0^{(9)}(a) + b_2^{(4)}(a)$$

TABLE 4 continued

Step 6:
$$b_3^{(1)}(a) = 0$$

 $b_3^{(2)}(a) = 2f^{(2)}(a, \varphi(a))e_1^{(1)}(a)e_2^{(1)}(a)$
...

Step 7:
$$e_3^{(1)}(a) = -\alpha_1 e_2^{(3)}(a) - \alpha_2 e_1^{(5)}(a) - \alpha_3 e_0^{(7)}(a)$$

 $e_3^{(2)}(a) = J(a) e_3^{(1)}(a) - \alpha_1 e_2^{(4)}(a) - \alpha_2 e_1^{(6)}(a) - \alpha_3 e_0^{(8)}(a)$
 $e_3^{(3)}(a) = 2J^{(1)}(a) e_3^{(1)}(a) + J(a) e_3^{(2)}(a)$
 $\vdots - \alpha_1 e_2^{(5)}(a) - \alpha_2 e_1^{(7)}(a) - \alpha_3 e_0^{(9)}(a) + b_3^{(2)}(a)$

Step 8:
$$e_4^{(1)}(a) = -\alpha_1 e_3^{(3)}(a) - \alpha_2 e_\beta^{(5)}(a) - \alpha_3 e_1^{(7)}(a) - \alpha_4 e_0^{(9)}(a)$$

.

Table 4 can be continued. Each new $e_m^{(1)}(a)$ to be computed will require two additional computations at each existing step of the table and the addition of two new steps to the table. One new step will contain the computations of $b_{m-1}^p(a)$ for p=1 and 2 and will preced the step where $e_{m-1}^{(1)}(a)$ is computed. The other new step will be the computation of $e_m^{(1)}(a)$.

Example 3: Compute $e_1^{(1)}(a), \ldots, e_4^{(1)}(a)$ for the equation $y' = y^2$; y(0) = .2; $0 \le t \le 1$ when the trapezoid rule is

used as the numerical method for solving the problem. The generating function for the coefficients in (18) for the trapezoid rule is given by (20). Specifically, $\alpha_1 = -\frac{1}{12}, \quad \alpha_2 = \frac{1}{120}, \quad \alpha_3 = \frac{-17}{315} \cdot \frac{1}{26} \quad \text{and} \quad \alpha_4 = \frac{62}{2835} \cdot \frac{1}{28}.$ If the computations in Table 4 are performed for this example the results are: $e_1^{(1)}(a) = .0008$; $e_2^{(1)}(a) = .000032$; $e_3^{(1)}(a) = 1.3867 \times 10^{-6}$; and $e_4^{(1)}(a) = 6.229 \times 10^{-8}$.

As is apparent, the computation of $e_m^{(1)}(a)$ for $m \ge 4$ becomes quite involved, particularly when q=2. However, there are applications (see Section 1 of Chapter 3) of the pullback method with $M \le 4$ which are more accurate than other methods currently available.

Section 5. Numerical Results for Initial Value Problems

In this section the pullback method and the method due to Lindberg [19] are compared theoretically and numerically. Since Lindberg's approach is conceptually and notationally quite different from ours, a discussion of his method is included. We will confine the discussion initially to the case M=4 and q=2 as this will suffice to point out the differences between the methods. Unfortunately, Lindberg's notation and that used here are in some instances in complete opposition to one another. The differences will be pointed out in footnotes.

Let h>0 be the basic steplength and for $k=0,1,\ldots,4$ define steplengths $h_k=h/2^k$ and grids $G_k=\{t_i^k=a+ih_k:i=0,\ldots,2^k\}$. In addition, let $h=h/2^4=h/16*$.

Assume that the numerical method being employed to solve the problem (1) is such that the expansion

(43)
$$Y(t_i^k, h_k) = \varphi(t_i^k) + e_1(t_i^k) h_k^2 + e_2(t_i^k) h_k^4 + e_3(t_i^k) h_k^6 + e_4(t_i^k) h_k^8 + O(h_k^{10})$$

is valid for each $i=0,1,\ldots,2^k$ and for each $k=0,1,\ldots,4$.

Lindberg's method is as follows. Number all grid points according to their order of occurrence in the finest

*In Lindberg's paper h is taken to be the basic steplength and $h_k=2^kh$ k=0,1,...,4. Thus the grids G_k and steplengths h_k are numbered in the opposite order.

grid, obtaining t_0 =a, t_1 ,..., t_{16} =a+h. At each t_i perform as many extrapolations as possible and denote the computed solution after n extrapolations as $Y_{n+1}(t_i, h)$. Thus at t_{16} we can perform 4 extrapolations obtaining $Y_1(t_{16}, h)$, $Y_2(t_{16}, h)$, $Y_3(t_{16}, h)$,..., $Y_5(t_{16}, h)$ where $Y_1(t_{16}, h)$ is the solution computed with the numerical method; i.e., $Y_1(t_i, h) = Y(t_i^4, h_4)$. At t_8 we can perform 3 extrapolations; at t_4 and t_{12} , 2 extrapolations and at t_2 , t_6 , t_{10} and t_{14} , 1 extrapolation.

According to Lindberg each $Y_{j}(t,h)$ satisfies the relationship

(44)
$$Y_{j}(t,h) = \varphi(t) + \sum_{v=j}^{\infty} \chi_{jv} e_{v}(t) h^{2v}$$

where for each j

(45)
$$\chi_{j\nu} = \prod_{p=1}^{j-1} \frac{2^{2p}-2^{2\nu}}{2^{2p}-1}.$$

The goal here is to define Y_5 at all points of the finest grid. Initially Y_5 is known only at the two points t_0 =a and t_{16} =a+h=a+16 \hat{h} . At each of these points form $Y_5(t,\hat{h})-Y_4(t,\hat{h})=-\chi_{44}e_4(t)\hat{h}^8+O(\hat{h}^{10})$ by (44). Via linear interpolation obtain an $O(\hat{h}^{10})*$ approximation to $-\chi_{44}e_4(t_8)\hat{h}^8$. Add this approximation to $Y_4(t_8,\hat{h})$ to obtain $Y_5(t_8,\hat{h})+O(\hat{h}^{10})$. At the next stage form the

^{*}This will be discussed in some detail later.

differences $Y_5(t,\hat{h})-Y_3(t,\hat{h})$ at the points $t=t_0,t_8$ and t_{16} . By (44), $Y_5(t,\hat{h})-Y_3(t,\hat{h})=-\chi_{33}e_3(t)\hat{h}^6-\chi_{34}e_4(t)\hat{h}^8$ +O(h^{10}). Using quadratic Lagrange interpolation obtain an $O(\hat{h}^3)$ approximation to $-\chi_{33}e_3(t)\hat{h}^6-\chi_{34}e_4(t)\hat{h}^8$ at $t=t_4$ and $t=t_{12}$ which is added to $Y_3(t,h)$ at these points to obtain $Y_5(t,\hat{h})$ for $t=t_4$ and $t=t_{12}$. However $Y_5(t,\hat{h})$, $t=t_4,t_{12}$, is only an $O(h^9)$ approximation to O(t). This process is continued until Y_5 has been defined at all t_i , $i=0,\ldots,16$. The result, according to Lindberg is $Y_5(t_i,\hat{h})=O(t_i)+O(\hat{h}^9)$ for $i=1,\ldots,7$ and $0,\ldots,15$ while $Y_5(t_i,\hat{h})=O(t_i)+O(\hat{h}^{10})$ for i=8 and 16.

Now the pullback interpolation method will yield $\mathfrak{O}(h^{1O})$ accuracy at all t_i 's and Lindberg's method yields $\mathfrak{O}(h^9)$ at most t_i 's. Since h=h/16 it looks as if Lindberg's method is of higher accuracy. This is not so for the following reasons: First, Lindberg presents the error analysis for extrapolation in terms of the smallest stepsize available, h. The 'O' notation is designed to be information—supressing and the remainder terms in extrapolation which are said to be $\mathfrak{O}(h^\beta)$ are actually of the form $\mathrm{Ch}^\beta g(t)$ where C and β are constants and g(t) is a continuous function of t. On a closed t interval g(t) will be bounded and we can write $\mathrm{Ch}^\beta g(t) \leq \mathrm{C_0} h^\beta = \mathfrak{O}(h^\beta)$. We could also write $\mathrm{Ch}^\beta g(t) = \hat{\mathrm{Ch}}^\beta g(t) \leq \hat{\mathrm{C_0}} h^\beta = \mathfrak{O}(h^\beta)$ by defining $\hat{\mathsf{C}} = 16^\beta \mathsf{C}$. The problem here is the fact that $\hat{\mathsf{C}}_\mathsf{O}$ will be

quite large compared to C_0 . When the error in extrapolation is expressed in terms of the largest or basic steplength, h, the constant C will be smaller than 1. If this same error is expressed in terms of the smallest stepsize, $\stackrel{\wedge}{h}$, the constant $\stackrel{\wedge}{C}$ will be of the form (45) and is larger than 1. In fact, when expressing the error in terms of the smallest stepsize we are unable to prove the convergence of the Aitken-Neville extrapolation scheme; indeed, an analysis of (45) shows that the constants $\mathcal{K}_{\mathbf{k}\mathbf{k}}^{\to\infty}$ as the stepsizes $h_{\mathbf{k}}^{\to0}$.

Secondly, the magnitude of the constant & further increased in Lindberg's analysis when interpolation is performed. Recall that Lindberg expresses the error in interpolation in terms of the smallest stepsize also. This means that at the first step when Lindberg is performing linear interpolation to $-\chi_{44}^{}e_{4}^{}(t)^{h}^{8}$ at the points t_0 and t_{16} he obtains the error $-\frac{(t-t_0)(t-t_{16})}{2!}\chi_{\Delta\Delta}e_{\Delta}^{(2)}(\xi)h^8 \quad \text{where} \quad t_0 < \xi < t_{16}. \quad \text{Lindberg}$ calls this error $O(h^{10})$; to do this one must interpret the maximum of $(t-t_0)(t-t_{16})$ as being $O(h^2)$. The maximum occurs at the midpoint t_{Q} and is easily seen to be $(\frac{h}{2})^2 = \frac{h^2}{4} = \frac{16^2 h^2}{4}$. To call this $O(h^2)$ further increases the size of the constant being suppressed. Note that in the error analysis for the pullback method, $(\frac{h}{2})^2$ is interpreted as $\frac{1}{4}h^2 = O(h^2)$ which causes the constant to decrease. When Lindberg claims an error of the form $\mathfrak{O}(h^{\beta})$ and the pullback method claims $\mathfrak{O}(h^{\beta})$, the inequality, $\mathfrak{O}(h^{\beta}) \leq \mathfrak{O}(h^{\beta})$, holds when the constants C_0 and c_0 are taken into account. Thus $\mathfrak{O}(h^{10})$ in actuality is smaller than $\mathfrak{O}(h^{9})$ and the pullback interpolation method yields a more accurate solution.

This increase in accuracy is due solely to the fact that we are using Hermite interpolation with the additional data $e_m^{(1)}$ (a). In fact, if the pullback interpolation method is changed so that Lagrange interpolation is used the results are identical to Lindberg's. In this case $O(h^{\beta}) \equiv O(h^{\beta})$.

Let's examine the differences between using Hermite and Lagrange interpolation when q=2 for arbitrary M.

For j=0, we know $e_M(t)$ at two points with accuracy $\mathcal{O}(h^2)$. Let $P_M(t)$ be the Hermite polynomial constructed as outlined and let $L_M(t)$ be the Lagrange interpolation polynomial for these two data points. We have $P_M(t) = e_M(t) + \mathcal{O}(h^3) + \mathcal{O}(h^2)$ and $L_M(t) = e_M(t) + \mathcal{O}(h^2) + \mathcal{O}(h^2)$ where the first ' \mathcal{O} ' term is the error in interpolation and the second is the error in the given data. Thus both $P_M(t)$ and $L_M(t)$ are $\mathcal{O}(h^2)$ approximations to $e_M(t)$. When $L_M(t)$ is used numerically, as it is in Lindberg's method, the solution at the midpoint of the interval is,

in most instances, less accurate than the solution which can be obtained by just performing extrapolation (see the examples concluding this section). This phenomenon, which Lindberg calls attention to in his paper, does not occur when $P_{\mathsf{M}}(\mathsf{t})$ is used.

For j=1, we know $e_{M-1}(t)$ at three points with accuracy $O(h^4)$. In this case, $P_{M-1}(t) = e_{M-1}(t) + O(h^4) + O(h^4) = e_{M-1}(t) + O(h^4)$ and $L_{M-1}(t) = e_{M-1}(t) + O(h^3) + O(h^4) = e_{M-1}(t) + O(h^3)$. Here it is the error in interpolation which determines the final accuracy so that there is a real gain when Hermite interpolation is used.

The situation is the same when j=2, that is, $P_{M-2}(t) = e_{M-t}(t) + O(h^6) + O(h^6) \quad \text{while} \quad L_{M-2}(t) = e_{M-2}(t) + O(h^5) + O(h^6). \quad \text{Again the error in interpolation dominates}$ and Hermite interpolation yields higher accuracy.

For any $j \ge 3$, it is the error to which we know $e_{M-j}(t)$ that dominates, since for $j \ge 3$ we have $2^j+1>2j+2$. The last inequality can be established by a proof similar to that of Lemma 1. Thus the advantages of using Hermite interpolation occur when j=0,1 and 2.

We now examine some examples of results obtained using the pullback interpolation method when numerically solving first order initial value problems. All numerical computations were done on the CDC 6500 installation at

Michigan State University. Calculations were done in single precision floating point arithmetic yielding 14 accurate decimal places.

The first example in this section is a continuation of Example 2, Section 4.

Example 4: Solve $y'=y^2$, y(0)=.2, $0\le t\le 1$ using Euler's method and pullback interpolation with M=4.

The basic stepsize is taken to be h=1 and, since M=4, the finest grid, G_4 , consists of 17 equally spaced points in [0,1]. That is,

$$G_A = \{t_i : t_i = 0 + i \ h/16, i = 0, 1, ..., 16\}.$$

The theoretical solution is y(t) = 1/(5-t) and the numerical results are reported in Table 5 below. Table 5 is constructed to exhibit the error in the solution computed by Euler's rule and the error in the solution computed by pullback interpolation at each grid point of G_4 . In each case the error is the computed solution minus the theoretical solution.

TABLE 5

i	Error using Euler	Error using pullback
0	O	0
1	-3.15×10^{-5}	$.01 \times 10^{-7}$
2	-6.53×10^{-5}	2.97×10^{-7}
3	-10.11×10^{-5}	8.47×10^{-7}
4	-13.93×10^{-5}	13.16×10^{-7}
5	-17.99×10^{-5}	14.87×10^{-7}
6	-22.31×10^{-5}	13.26×10^{-7}
7	-26.92×10^{-5}	9.17×10^{-7}
8	-31.82×10^{-5}	3.62×10^{-7}
9	-37.05×10^{-5}	-2.01×10^{-7}
10	-42.63×10^{-5}	-6.58×10^{-7}
11	-48.58×10^{-5}	-8.92×10^{-7}
12	-54.93×10^{-5}	-8.26×10^{-7}
13	-71.71×10^{-5}	-4.97×10^{-7}
14	-68.95×10^{-5}	-1.26×10^{-7}
15	-76.68×10^{-5}	54×10^{-7}
16	-84.96×10^{-5}	-3.30×10^{-7}

The results given in Table 5 can be quickly summarized by noting that Euler's method can guarantee only two accurate digits at all grid points of G_4 while the pullback method can guarantee 5 accurate digits.

Pullback interpolation based on Euler's rule with M=4 was also used to compute the solution to $y'=y^2$, y(0)=.2 with basic stepsizes of h=1/2 and h=1/4. In the first case the solution is computed on [0,1/2] and in the second case the solution is computed on [0,1/4]. The results are summarized in Table 6 below. The first two columns are the results for h=1/2 and the last two are the results for h=1/4. Table 6 is arranged to exhibit the correspondence between grid points for the two different mesh sizes.

TABLE 6

	h = 1/2	h = 1/4
i	Error using pullback	i Error using pullback
0	0	0 0
		1 .15 \times 10 ⁻¹⁰
1	1.96 x 10 ⁻¹⁰	2 11.93 x 10 ⁻¹⁰
		$3 33.93 \times 10^{-10}$
2	189.25 x 10 ⁻¹⁰	4 50.92 \times 10 ⁻¹⁰
		5 55.90 \times 10 ⁻¹⁰
3	532.79×10^{-10}	6 47.10 \times 10 ⁻¹⁰
		7 27.94×10^{-10}
4	817.32×10^{-10}	8 3.45 \times 10 ⁻¹⁰
		9 21.15 \times 10 ⁻¹⁰
5	905.78 x 10 ⁻¹⁰	10 40.66×10^{-10}
		11 50.02 x 10 ⁻¹⁰
6	777.95 x 10 ⁻¹⁰	$12 45.84 \times 10^{-10}$
		13 29.33 \times 10 ⁻¹⁰
7	486.78×10^{-10}	14 9.12 \times 10 ⁻¹⁰
		15 1.17×10^{-10}
8	110.79 x 10 ⁻¹⁰	$16 .40 \times 10^{-10}$
9	269.05 x 10 ⁻¹⁰	
10	571.77 x 10 ⁻¹⁰	
11	718.91 x 10 ⁻¹⁰	
12	657.64 x 10 ⁻¹⁰	
13	408.36 x 10 ⁻¹⁰	
14	108.33 x 10 ⁻¹⁰	
15	26.89 x 10 ⁻¹⁰	

 32.13×10^{-10}

16

As can be seen by examing Table 5 and Table 6 the smaller the basic steplength is the more accurate are the computed solutions. This is in agreement with results for exprapolation (see Lambert [18] and Gragg [11,12]).

Example 5: Solve y'=y, y(0)=1, $0 \le t \le 1$ using the trapezoid rule, and pullback interpolation with M=4. Using a basic steplength of h=1 we again obtain the solution at 17 equally spaced points in [0,1]. The theoretical solution of this equation is $y(t)=e^t$ and the error reported is the computed solution minus the theoretical solution. The results of Example 5 and the next example will be presented simultaneously in Table 7.

Example 6: Solve $y'=-\sin t$, y(0)=1, $0\le t\le 1$ using the trapezoid rule and pullback interpolation with M=4. The basic steplength is taken to be h=1 and the grids are the same as in Example 5. The theoretical solution is given by $y(t)=\cos(t)$.

TABLE 7

	Error in pullback	Error in pullback
i	for y'=y with h=l	for y'=-sin t with h=1
0	0	0
1	$.12 \times 10^{-10}$	3×10^{-13}
2	$.48 \times 10^{-10}$	7×10^{-13}
3	1.05×10^{-10}	9×10^{-13}
4	1.70×10^{-10}	7×10^{-13}
5	2.30×10^{-10}	3×10^{-13}
6	2.74×10^{-10}	1×10^{-13}
7	2.94×10^{-10}	0.0×10^{-13}
8	2.90×10^{-10}	2×10^{-13}
9	2.65×10^{-10}	3×10^{-13}
10	2.34×10^{-10}	3×10^{-13}
11	2.11×10^{-10}	$.2 \times 10^{-13}$
12	2.16×10^{-10}	1.2×10^{-13}
13	2.63×10^{-10}	2.5×10^{-13}
14	3.57×10^{-10}	3.6×10^{-13}
15	4.86×10^{-10}	2.9×10^{-13}
16	6.13×10^{-10}	-2.1×10^{-13}

The equations in examples 5 and 6 were solved in the same manner with a basic stepsize of h=1/2 to obtain solutions at 17 equally spaced points in [0,1/2]. The results of these computations are presented in Table 8.

TABLE 8

	Error in pullback	Error in pullback
i	for $y'=y$ with $h=1/2$	for $y'=-\sin t$ with $h=1/2$
0	0	0
1	0.0×10^{-13}	-0×10^{-14}
2	$.2 \times 10^{-13}$	-1 x 10 ⁻¹⁴
3	$.4 \times 10^{-13}$	-1×10^{-14}
4	$.6 \times 10^{-13}$	-1×10^{-14}
5	$.9 \times 10^{-13}$	-0×10^{-14}
6	1.0 x 10 ⁻¹³	-1×10^{-14}
7	$.7 \times 10^{-13}$	-0×10^{-14}
8	$.3 \times 10^{-13}$	-1×10^{-14}
9	4×10^{-13}	-1×10^{-14}
10	9×10^{-13}	-1 x 10 ⁻¹⁴
11	-1.4×10^{-13}	-1×10^{-14}
12	-1.7×10^{-13}	-2×10^{-14}
13	-1.8 x 10 ⁻¹³	-2×10^{-14}
14	-1.4×10^{-13}	-2×10^{-14}
15	9×10^{-13}	-3×10^{-14}
16	6×10^{-13}	-3×10^{-14}

As the results presented in Table 7 and Table 8

so vividly indicate, pullback interpolation will yield

uni form accuaracy at all grid points of the finest grid.

Lindberg in [19] presents numerical results ob tained when solving y'=y, y(0)=1, on [0,1] using the trapezoid rule and his interpolation method. The largest error Lindberg obtains is at t_5 and has magnitude 16×10^{-10} . Examining Table 7 we see that the largest error produced for this equation when using the pullback method is the error in extrapolation at the endpoint, $t_{16}=1$. The magnitude of this error is 6.13×10^{-10} .

In the next series of examples extrapolation,
Lindberg's method and pullback interpolation are compared
numerically.

Example 7: Solve $y'=y^2$, y(0)=1, $0 \le t \le 1$ using the trapezoid rule. Choosing a basic steplength of h=1 and M=3, the solution is computed first by using extrapolation as often as is possible at each grid point of the finest grid, $G_3 = \{t_i: t_i=0+i \text{ h/8}, i=0,1,\ldots,8\}$. Secondly, the solution is computed on G_3 by using Lindberg's method as described earlier in this section. Lastly, the solution is computed on G_3 using pullback interpolation. All errors are given as the computed solution minus the theoretical solution. The results are presented in Table 9.

TABLE 9

i	Error in the best extrapolated value	Error using Lindberg's method	Error using pullback
0	· O	0	0
1	1.7×10^{-6}	$.8 \times 10^{-10}$	4×10^{-10}
2	-1.0×10^{-8}	-26.7×10^{-10}	-3.0×10^{-10}
3	5.9×10^{-6}	$-30.3 \times 10^{-10*}$	-4.7×10^{-10}
4	3.7×10^{-10}	-5.9×10^{-10}	-1.9×10^{-10}
5	1.2 x 10 ⁻⁵	22.0×10^{-10}	3.7×10^{-10}
6	-5.3×10^{-8}	23.9×10^{-10}	$6.2 \times 10^{-10*}$
7	$1.9 \times 10^{-5*}$	-4.0×10^{-10}	$.8 \times 10^{-10}$
8	-1.5×10^{-10}	-1.5×10^{-10}	-1.5 x 10 ⁻¹⁰

^{*}greatest absolute error for the method

Example 8: Solve y'=-sin t, y(0)=1, 0\(\subseteq t \leq 1 \) using the trapezoid rule. We again take h=1 and M=3 and compute solutions using extrapolation, Lindberg's method, and Pullback interpolation. The results are given in Table 10.

TABLE 10

i	Error in the extrapolar		Error u	_	Error pullb	using ack
0	0		0			0
1	1.02 x	10 ⁻⁵	178.36	x 10 ⁻¹⁰	9.46	x 10 ⁻¹⁰
2	- 4.22 x	10 ⁻⁸	37.54	x 10 ⁻¹⁰	6.33	x 10 ⁻¹⁰
3	9.05 x	10 ⁻⁵	- 47.24	x 10 ⁻¹⁰	- 2.47	× 10 ⁻¹⁰
4	9.96 x	10 ⁻¹⁰	- 9.23	x 10 ⁻¹⁰	.37	x 10 ⁻¹⁰
5	24.62 x	10 ⁻⁵	31.70	x 10 ⁻¹⁰	4.92	x 10 ⁻¹⁰
6	-36.48 x	10 ⁻⁸	- 63.15	x 10 ⁻¹⁰	-17.54	x 10 ⁻¹⁰
7	46.76 x	10 ^{-5*}	-239.94	x 10 ^{-10*}	-62.65	× 10 ^{-10*}
8	.96 x	10-10	96	x 10 ⁻¹⁰	96	x 10 ⁻¹⁰

^{*}greatest absolute error for the method

Comparing the results in Example 7 and Example 8, we see that the solutions obtained with both Lindberg's method and the pullback interpolation method, at the intermediate grid points t_i , $i=1,\ldots,7$, are more accurate than the best extrapolated values. The pullback interpolation method is the most accurate of all, being almost a full decimal place better than Lindberg's method at all intermediate grid points.

As mentioned before, one drawback to Lindberg's method is that his results at the midpoint of the interval are no better, and in fact often worse, than the results

obtained by just performing extrapolation. This phenomenon does not occur when the pullback interpolation method is used, as a comparison of the errors at t_4 in Table 9 and Table 10 clearly demonstrates.

Examining Table 10 we note that the pullback method is dramatically less accurate at t_7 than it is at all other grid points. In an attempt to explain this behavior we computed the following example.

Example 9: Solve y'=-sin(t), y(0)=1 on $0\le t\le 1/2$. The only difference between this example and Example 8 is that we now take h=1/2 as the basic stepsize. The results are presented as

TABLE 11

i	Error in the best extrapolated value	Error in Lindberg's method	Error in pullback
0	0	0	0
1	6.36×10^{-7}	71.71×10^{-12}	3.89×10^{-12}
2	-6.62×10^{-10}	15.02×10^{-12}	2.55×10^{-12}
3	5.71×10^{-6}	-19.21×10^{-12}	-1.15 x 10 ⁻¹²
4	-3.92×10^{-12}	-3.87×10^{-12}	$.03 \times 10^{-12}$
5	1.58×10^{-5}	12.74×10^{-12}	1.99×10^{-12}
6	-5.89×10^{-9}	-25.85×10^{-12}	-7.54×10^{-12}
7	$3.07 \times 10^{-5*}$	-98.14 x 10 ^{-12*}	$-26.90 \times 10^{-12*}$
8	-1.1×10^{-13}	- .11 x 10 ⁻¹²	11 x 10 ⁻¹²

^{*}greatest absolute error for the method

Once again t_7 is the least accurate solution for all three methods. The only reasonable explanation which suggests itself is that the solution cost is changing very rapidly (relative to its behavior on the rest of the interval) between t_7 and t_8 . This behavior is compensated for by extrapolation at t_8 , but since no extrapolation is done at t_7 no correction is possible there.

From these examples it appears that pullback interpolation is quite sensitive to the accuracy to which we know the solution at the intermediate grid points and the type of equation we are solving. Thus it appears that it is the error in the original data and not the error in interpolation that is determining the overall accuracy of the method.

We next examine what happens when we solve the same equation on intervals [a,b], [b,c] and [c,d] using the last computed solution as initial data for the next interval.

Example 10: Solve $y'=y^2$, y(0)=1, $0 \le t \le 3/2$ by computing the solution on intervals of length 1/2 and using the computed solution at the endpoint of the previous interval as the initial data for the equation on the next interval. The method of solution on each interval will be the trapezoid rule and pullback interpolation with M=3. The results are presented as Table 12.

TABLE 12

[0,1/2]		[1/2,1]		[1,3/2]	
i	Error	i	Error	i	Error
0	0	0	0	0	0
1	01×10^{-12}	1	.56 x 10^{-12}	1	1.13×10^{-12}
2	55×10^{-12}	2	-1.20×10^{-12}	2	-4.39×10^{-12}
3	74×10^{-12}	3	-1.95×10^{-12}	3	-6.93×10^{-12}
4	.13 \times 10 ⁻¹²	4	$.27 \times 10^{-12}$	4	64×10^{-12}
5	1.59×10^{-12}	5	4.19×10^{-12}	5	11.24×10^{-12}
6	$2.21 \times 10^{-12*}$	6	$5.85 \times 10^{-12*}$	6	$16.07 \times 10^{-12*}$
7	$.94 \times 10^{-12}$	7	2.39×10^{-12}	7	5.04×10^{-12}
8	$.56 \times 10^{-12}$	8	1.13×10^{-12}	8	$.84 \times 10^{-12}$

Some deterioration in the accuracy of the computed solution for larger t can be observed in Table 12. This can be partly attributed to an accumulation of roundoff errors and partly to the sensitivity of the pullback method to the accuracy of the data it receives.

If we compare the first two columns of errors in Table 12 to the results given in Table 9, we see that in terms of greatest absolute error it is better to take a smaller basic steplength and solve the problem several times in succession. However, we should point out that this can be quite expensive computationally.

^{*}greatest absolute error

Also, note that the largest absolute error always occurs at the same relative position of the three grids, namely at t₆. This is further support for our contention that the nature of the equation and the accuracy of the data points used when interpolating are the factors which control the overall accuracy of the pullback method.

In each example computed using the trapezoid rule the largest error in the pullback method has occurred to the right of the midpoint. This is no coincidence. The trapezoid rule yields more accurate solutions at grid points nearer the initial point and consequently extrapolation will be more accurate in the first half of the interval. Also, when we perform Hermite interpolation we have an extra data point at the initial grid point. All of these factors combined make it reasonable to expect that the largest errors will be produced to the right of the midpoint. Thus, a solution computed using pullback interpolation should be most reliable in the first half of the interval on which the solution is computed.

Note that the largest error in Lindberg's method can occur in the first half of the interval, as an examination of Table 9 reveals. This is due to the fact that it is the error in interpolation which determines the overall accuracy of Lindberg's method.

In summary, pullback interpolation coupled with the trapezoid rule will yield highly accurate solutions at all

grid points. In fact, if the solutions obtained with the trapezoid rule and extrapolation are sufficiently accurate, pullback interpolation will yield uniform accuracy at all grid points.

Pullback interpolation coupled with Euler's rule is not recommended as a viable solution technique. Euler's rule is simply not accurate enough to enable pullback interpolation to operate effectively.

CHAPTER II

TWO POINT BOUNDARY VALUE PROBLEMS

Section 1. The Problem and Its Discretization

In this chapter we will consider two point boundary value problems of the form

$$x''(t)-f(t,x(t),x'(t)) = 0$$
(1)
$$x(a) = A, x(b) = B.$$

We assume that $f(t,x,x') \in C^1[[a,b] \times (-\infty,\infty) \times (-\infty,\infty)]$, f(t,y,z) is uniformly Lipschitz continuous in y and z, $0 < \varepsilon < \frac{\partial f}{\partial y}$ and $\left|\frac{\partial f}{\partial z}\right| \le K$ where K is a constant. Under these assumptions (see Keller [17]) problem (1) has a unique solution which we will denote by $\phi(t)$.

The continuous problem (1) will be denoted by

$$(1')$$
 $F(x) = 0.$

The operator $F(x) \equiv x''(t) - f(t,x(t),x'(t))$ maps the Banach space of twice continuously differentiable functions defined on [a,b] into C[a,b].

To obtain a discrete version of (1), let $n \ge 2$ be any natural number, define h=(b-a)/n and form the uniform mesh $\{t_k=a+ih\}_{i=0}^n$ in [a,b]. The discrete problem is then given by

$$\frac{X_{i+1}^{-2X_{i}^{+}X_{i-1}}}{h^{2}} - f(t_{i}, X_{i}, \frac{X_{i+1}^{-}X_{i-1}}{2h}) = 0, \quad i=1, \dots, n-1,$$

$$X_{0} = A, \quad X_{n} = B,$$

where we have introduced the notation $X_i = X(t_i)$. Problem (2) may be thought of as a nonlinear system of equations in E^{n-1} with the unknown being the vector $(X_1, \ldots, X_{n-1})^T$. The solution to (2) will be denoted by X(h) and we introduce the operator notation

$$(2') F_h(X) = 0$$

for problem (2).

In order to set up the correspondence between the continuous and discrete problems we will need to define a space discretization operator w_h . Thus, let z(t) be an arbitrary function defined on [a,b] and define w_h acting on z by

$$\omega_{h} z = \omega_{h}(z(t)) = (z(t_{1}), \dots, z(t_{n-1}))^{T}$$

$$\equiv (z_{1}, \dots, z_{n-1})^{T}.$$

Note that $w_n z$ is the vector in E^{n-1} whose components are obtained by evaluating z at the grid points.

Problem (1) and the discretization (2) have been studied by Pereyra [24,27]. Stetter [31] and Pereyra [26,28] have also studied the special case of (1) when x' is not present in the equation.

To formalize what is meant by convergence of the discrete solution X(h) to the continuous solution $\phi(t)$ we have the following definition (see Lambert [18] and Pereyra [28]).

<u>Definition 1</u>: We shall say that the discrete solution X(h) converges discretely to the theoretical solution $\phi(t)$ if and only if

(3)
$$\lim_{h\to 0} ||x(h) - \omega_h \varphi||_{(h)} = 0,$$

where $\|\cdot\|_{(h)}$ is the maximum norm on E^{n-1} .

The subscript (h) will be omitted from the norms throughout the remainder of this chapter. We should point out that Pereyra in [27] utilizes a different norm than the one we have used here, while in his earlier work [24] he uses the maximum norm.

Typically, discrete convergence depends on the two properties of consistency and stability of the discrete operator F_h in (2'). In formulating the definitions of these concepts we have followed the approach of Pereyra [25,28].

<u>Definition 2</u>: The operator F_h is said to be consistent of order p>0 if, for each solution $\phi(t)$ of (1) and for all $h\leq h_0$ we have

(4)
$$||F_h(\omega_h \varphi)|| = O(h^p).$$

<u>Definition 3</u>: The operator F_h is stable if for any pair of discrete functions U and W and for all $h \le h_0$ there exists a constant C > 0, independent of h, such that

(5)
$$\|U-W\| \le C \|F_h(U) - F_h(W)\|.$$

The standard result for discrete operators, that consistency and stability imply discrete convergence, is valid here. This theorem, in the context of linear multistep methods for ordinary differential equations, is originally due to Dahlquist [8,9]. For completeness we present Pereyra's statement of this result and briefly summarize his proof.

Lemma l (Pereyra []): If F_h is stable then it is locally invertible around $w_h \phi$ and the inverse mapping is locally Lipschitz continuous for all $h \leq h_O$.

<u>Proof:</u> Let $B_h = B(w_h \varphi, \rho)$, be the open ball of radius ρ centered at $w_h \varphi$, where $\rho > 0$ is independent of h. If $U, W \in B_h$ then stability implies that F_h is one-to-one on B_h , for otherwise we can violate (5). Therefore $F_h: B_h \to F_h(B_h)$ is a bijection, implying that F_h^{-1} exists on $F_h(B_h)$. If $X, Y \in F_h(B_h)$ then we can write (5) as

$$\|F_{h}^{-1}(X) - F_{h}^{-1}(Y)\| \le c \|X - Y\|.$$

Theorem 1. (Pereyra [25 28]): Assume (1) has a unique solution φ and F_h is stable on $B_h \equiv B(w_h \varphi, \rho)$ for all $h \le h_0$.

Let $F_{\mbox{$h$}}$ be consistent of order $\,p\,$ with $\,F.\,$ Then there exists an $\,\overline{h}_{\mbox{$O$}}>0\,$ such that:

- (a) $\forall h \leq \bar{h}_0$ there exists a unique solution X(h) for the discrete problem $F_h(X) = 0$, and
- (b) the solution satisfies

 (6) $||x(h) \omega_h \varphi|| = O(h^p)$.

Equation (6) can be summarized by saying the solutions are (discretely) convergent of order p.

<u>Proof:</u> By Lemma 1, F_h is a homeomorphism between its domain B_h and its range $R_h \equiv F_h(B_h)$ $\forall h \leq h_0$. Brouwer's Invariance of Domain Theorem implies F_h maps the interior of B_h onto the interior of R_h and the boundary onto the boundary.

If $V \in \partial B_h$, stability implies that $\frac{\rho}{C} \leq \|F_h(V) - F_h(\omega_h \phi)\|$ and therefore by letting V vary over ∂B_h we see that $B(F_h(\omega_h \phi), \frac{\rho}{C}) \subset R_h.$

Consistency implies $\|F_h(w_h\phi)\| = O(h^p)$ and therefore $\|F_h(w_h\phi)\| \to 0$ as $h \to 0$. Thus, there exists $\bar{h}_0 \le h_0$ such that $\|F_h(w_h\phi)\| < \frac{\rho}{C}$ $\forall h \le \bar{h}_0$ and $O \in B(F_h(w_h\phi), \frac{\rho}{C})$. But F_h is one-to-one and onto and therefore there exists a unique $X(h) \in B_h$ such that $F_h(X(h)) = 0$.

Stability implies that

$$\|\mathbf{x}(\mathbf{h}) - \mathbf{\omega}_{\mathbf{h}} \mathbf{\phi}\| \le C \|\mathbf{F}_{\mathbf{h}}(\mathbf{x}(\mathbf{h})) - \mathbf{F}_{\mathbf{h}}(\mathbf{\omega}_{\mathbf{h}} \mathbf{\phi})\| = C \|\mathbf{F}_{\mathbf{h}}(\mathbf{\omega}_{\mathbf{h}} \mathbf{\phi})\| = O(\mathbf{h}^{\mathbf{p}}). \quad \Box$$

Thus, in order to establish the discrete convergence of the numerical solution to the theoretical solution it suffices

to show the numerical method is both consistent and stable.

The usual technique for proving that F_h is consistent is to examine the local truncation error at each grid point t_i . The local truncation error $\tau_h(t_i)$ is obtained by applying F_h to the discretization of the theoretical solution and evaluating the result at t_i . Formally

for $i=1,\ldots,n-1$.

If we assume that f(t,y,z) has M total derivatives with respect to t on [a,b], then by expanding $F_h(w_h\phi)(t_i)$ in a Taylor series about $\phi(t_i)$ it can be shown that the local truncation error $\tau_h(t_i)$ satisfies

(8)
$$\tau_h(t_i) = F(\phi) \Big|_{t_i = 1} + \sum_{k=1}^{M} h^{2k} \{ \frac{2}{(2k+2)!} \phi^{(2k+2)} (t_i) - g_{2k}(\cdot) \} + O(h^{2M+2}),$$

where t_i is any grid point in [a,b]. The functions g_{2k} are obtained from

(9)
$$\sum_{k=1}^{M} h^{2k} g_{2k}(\cdot) = \sum_{k=1}^{M} \left\{ \frac{1}{k!} \frac{\partial^{k} f}{\partial z^{k}} \right\}_{t_{i}, \phi(t_{i}), \phi'(t_{i})}$$

$$\left\{ \sum_{j=1}^{M} \frac{h^{2j}}{(2j+1)!} \phi^{(2j+1)}(t_{i}) \right\}^{k},$$

by rearranging the right hand side in powers of h.

For instance,

$$g_2(\cdot) = \frac{1}{3!} \frac{\partial f}{\partial z} | (t_i, \varphi(t_i), \varphi'(t_i)) \varphi'''(t_i)$$

and

$$g_{4}(\cdot) = \frac{1}{5!} \frac{\partial f}{\partial z} |_{(t_{i}, \phi(t_{i}), \phi'(t_{i}))} \phi^{(5)}(t_{i})$$

$$+ \frac{1}{2!} \frac{\partial^{2} f}{\partial z^{2}} |_{(t_{i}, \phi(t_{i}), \phi'(t_{i}))} (\frac{1}{3!} \phi''(t_{i}))^{2}.$$

Because equations (8) and (9) are well-known we have not given the details for constructing them. The reader interested in more detail is referred to [27] or the work in Section 2.1 of Chapter 3.

By observing that $f(\phi) \equiv 0$, equation (8) implies that $f_h(X)$ is consistent of order 2.

Section 2. Stability

In this section we will investigate the stability of the discrete scheme (2). The fact that (2) is stable in the uniform norm has been proven by Pereyra [24]. The proof given by Pereyra uses the theory of monotone operators and is quite technical in nature. We present here a simple and direct proof of the fact that F_h is stable.

In what follows we will be considering vectors V which properly belong to E^{n+1} . However we will restrict these vectors to the (n-1)-dimensional subspace consisting of all vectors in E^{n+1} whose first and last components are given by the boundary conditions in (2). That is, all vectors will have the form $V = (A, V_1, \dots, V_{n-1}, B)^T$, and we will regard F_h as an operator on $V = (V_1, \dots, V_{n-1}) \in E^{n-1}$.

Let U and W be any two vector in E^{n-1} . For $2 \le i \le n-2$ the i^{th} component of $F_h(U)-F_h(W)$ is given by (10) $[F_h(U)-F_h(W)]_i$

$$= h^{-2} (U_{i+1} - 2U_{i} + U_{i-1}) - f(t_{i}, U_{i}, \frac{U_{i+1} - U_{i-1}}{2h})$$

$$- h^{-2} (W_{i+1} - 2W_{i} + W_{i-1}) + f(t_{i}, W_{i}, \frac{W_{i+1} - W_{i-1}}{2h})$$

$$= h^{-2} [(U_{i+1} - W_{i+1}) - 2(U_{i} - W_{i}) + (U_{i-1} - W_{i-1})$$

$$- [f(t_{i}, U_{i}, \frac{U_{i+1} - U_{i-1}}{2h}) - f(t_{i}, W_{i}, \frac{W_{i+1} - W_{i-1}}{2h})]$$

Using the Mean Value Theorem for continuous functions

of two variables (see Widder [33]) and regarding

$$f(t_i, v_i, \frac{v_{i+1}-v_{i-1}}{2h})$$
 as a function of v_i and $\frac{v_{i+1}-v_{i-1}}{2h}$,

we can obtain

(11)
$$f(t_{i}, U_{i}, \frac{U_{i+1}^{-U_{i-1}}}{2h}) - f(t_{i}, W_{i}, \frac{W_{i+1}^{-W_{i-1}}}{2h})$$

$$= f_{y}(t_{i}, \alpha_{i}, \beta_{i}) (U_{i}^{-W_{i}}) + f_{z}(t_{i}, \alpha_{i}, \beta_{i}) \cdot (\frac{U_{i+1}^{-U_{i-1}}}{2h} - \frac{W_{i+1}^{-W_{i-1}}}{2h})$$

where

(12a)
$$\alpha_{i} = W_{i} + \theta_{i} (U_{i} - W_{i})$$

and

(12b)
$$\beta_i = \frac{W_{i+1} - W_{i-1}}{2h} + \theta_i \left(\frac{U_{i+1} - U_{i-1}}{2h} - \frac{W_{i+1} - W_{i-1}}{2h} \right)$$

with $0 < \theta_i < 1$.

Substituting (11) into (10) we have

$$[F_{h}(U) - F_{h}(W)]_{i}$$

$$= h^{-2} (U_{i-1} - W_{i-1}) - 2h^{-2} (U_{i} - W_{i}) + h^{-2} (U_{i+1} - W_{i+1})$$

$$- f_{y}(t_{i}, \alpha_{i}, \beta_{i}) (U_{i} - W_{i}) + \frac{1}{2h} f_{z}(t_{i}, \alpha_{i}, \beta_{i}) (U_{i-1} - W_{i-1})$$

$$- \frac{1}{2h} f_{z}(t_{i}, \alpha_{i}, \beta_{i}) (U_{i+1} - W_{i+1})$$

$$= h^{-2} [(1 + \frac{h}{2} f_{z}(t_{i}, \alpha_{i}, \beta_{i})) (U_{i-1} - W_{i-1})$$

+
$$(-2-h^2f_y(t_i,\alpha_i,\beta_i))(U_i-W_i)$$

+ $(1-\frac{h}{2}f_z(t_i,\alpha_i,\beta_i))(U_{i+1}-W_{i+1})$].

For i=1 and i=n-1 we get similar expressions except the terms multiplying $(U_O - W_O)$ and $(U_n - W_n)$ are identically zero in the respective equations because $U_O \equiv W_O$ and $U_n \equiv W_n$ by our earlier convention.

If we write (13) for $i=1,\ldots,n-1$ as a matrix system we obtain

(14)
$$F_h(U) - F_h(W) = M_h(U, W) (U-W),$$

where $M_h(U,W)$ is an $(n-1) \times (n-1)$ matrix whose rows are given by

$$[M_{h}(U,W)]_{1} = h^{-2}(-2-h^{2}f_{y}(t_{1},\alpha_{1},\beta_{1}),1-\frac{h}{2}f_{z}(t_{1},\alpha_{1},\beta_{1}),0,...,0),$$

$$[M_{h}(U,W)]_{i} = h^{-2}(0,...,0,1+\frac{h}{2}f_{z}(t_{i},\alpha_{i},\beta_{i}),-2-h^{2}f_{y}(t_{i},\alpha_{i},\beta_{i}),$$

$$1-\frac{h}{2}f_{z}(t_{i},\alpha_{i},\beta_{i}),0,...,0)$$

for $i=2,\ldots,n-2$, and

$$[M_{h}(U,W)]_{n-1} = h^{-2}(0,...,0,1 + \frac{h}{2} f_{z}(t_{n-1},\alpha_{n-1},\beta_{n-1}),$$

$$-2-h^{2} f_{y}(t_{n-1},\alpha_{n-1},\beta_{n-1})),$$

With the nonzero entries for the ith row, $2 \le i \le n-2$, occurring in the (i-1)st, ith and (i+1)st positions. The arguments α_i and β_i are given by (12).

Taking norms in (14) we have

$$\|F_{h}(U) - F_{h}(W)\| = \|M_{h}(U, W)(U-W)\|$$
.

Multiplying both sides of this equation by $\frac{1}{\|U-W\|}$, using the linearity of $M_h(U,W)$, and noting that

 $\frac{U-W}{\|U-W\|}$ has norm one, we have for all $U \neq W$

$$\frac{\|F_{h}(U) - F_{h}(w)\|}{\|U - w\|} = \frac{\|M_{h}(U, w)(U - w)\|}{\|U - w\|}$$

$$= \|M_{h}(U, w) \frac{U - w}{\|U - w\|}\|$$

$$\geq \inf_{\|Z\|=1} \|M_{h}(U, w) Z\|.$$

Thus, if $\|M_h(U,W) Z\|$ is bounded below by a constant C, independent of h, for all vectors Z belonging to the unit ball in E^{n-1} we can write

$$\|F_{h}(U) - F_{h}(W)\| \ge C \|U - W\|,$$

which proves that F_h is stable.

To establish that $\|\mathbf{M}_{h}(\mathbf{U},\mathbf{W})\mathbf{Z}\|$ is bounded below we Proceed as follows:

let $h_0 = 2/K$ and define

$$a_i = \frac{h}{2} f_z(t_i, \alpha_i, \beta_i)$$

and

$$b_i = f_y(t_i, \alpha_i, \beta_i)$$

for each $i=1,\ldots,n-1$.

For any $h < h_0$, because of our hypotheses on f_v and f_z , we see that $|a_i| < 1$ and $\beta_i > \varepsilon > 0$ for all $i, 1 \le i \le n-1$.

In this notation the matrix $M_h(U,W)$ is given by

$$[M_h(U,W)]_1$$
 = $h^{-2}(-2-h^2b_1, 1-a_1, 0, ..., 0)$;
 $[M_h(U,W)]_i$ = $h^{-2}(0, ..., 0, 1+a_i, -2-h^2b_i, 1-a_i, 0, ..., 0)$
for $2 \le i \le n-2$; and

$$[M_h(U,W)]_{n-1}$$
 = $h^{-2}(0,...,0,1+a_{n-1},-2-h^2b_{n-1})$.

Now let $Z = (Z_1, ..., Z_{n-1})^T$ be any unit vector in E^{n-1} with respect to the maximum norm on E^{n-1} . This means that at least one component of Z has absolute value one and no component of Z has absolute value larger than one.

Suppose $Z_1 = 1$, then the first component of Mh (U,W)Z has absolute value

$$h^{-2} | -2 - h^{2}b_{1} + (1 - a_{1}) z_{2} | \ge h^{-2} (|2 + h^{2}b_{1}| - |1 - a_{1}| ||z_{2}|)$$

$$> h^{-2} (2 + h^{2}b_{1} - 2 ||z_{2}|)$$

$$\ge h^{-2} (2 + h^{2}b_{1} - 2)$$

$$= b_{1}$$

$$> \epsilon.$$

We have used the facts that $b_1 > 0$, $|a_1| < 1$ and $|z_2| \le 1$ in making the above estimates.

Similarly, if $Z_{n-1}=1$ then the last component of $M_h(U,W)Z$ has absolute value larger than ε .

Suppose $Z_i=1$ where $2 \le i \le n-2$, then the i^{th} component of $M_h(U,W)Z$ has absolute value

$$h^{-2} | (1+a_{i}) z_{i-1}^{-2-h^{2}b_{i}} + (1-a_{i}) z_{i+1} |$$

$$\geq h^{-2} (2+h^{2}b_{i}^{-} | (1+a_{i}) z_{i-1}^{-1} + (1-a_{i}) z_{i+1}^{-1} |)$$

$$\geq h^{-2} (2+h^{2}b_{i}^{-} (|z_{i-1}^{-1} + z_{i+1}^{-1}| + |a_{i}^{-1} | |z_{i-1}^{-1} - z_{i+1}^{-1} |))$$

$$\geq h^{-2} (2+h^{2}b_{i}^{-} (|z_{i-1}^{-1} + z_{i+1}^{-1}| + |z_{i-1}^{-1} - z_{i+1}^{-1} |))$$

$$\geq h^{-2} (2+h^{2}b_{i}^{-2})$$

$$\geq h^{-2} (2+h^{2}b_{i}^{-2})$$

$$\geq \epsilon \cdot$$

The next to last inequality is valid for the following reason.

Consider the expression

$$|\alpha+\beta|+|\alpha-\beta|$$
 with $|\alpha|\leq 1$, $|\beta|\leq 1$.

Squaring we obtain

(15)
$$(\alpha+\beta)^2+2|\alpha^2-\beta^2|+(\alpha-\beta)^2=2\alpha^2+2\beta^2+2|\alpha^2-\beta^2|$$
.

If $\alpha^2 = \beta^2$ then this is trivially ≤ 4 so without loss of generality we can assume $\alpha^2 > \beta^2$. In this case equation (15) becomes

$$(|\alpha+\beta|+|\alpha-\beta|)^2 = 2\alpha^2+2\beta^2+2\alpha^2-2\beta^2 = 4\alpha^2 \le 4$$
.

Thus for any two real numbers α and β with $|\alpha| \le 1$ and $|\beta| \le 1$ we have

$$|\alpha+\beta|+|\alpha-\beta|\leq 2.$$

Therefore, for any unit vector Z, some component of $M_h(U,W)Z$ has absolute value larger than ε . This implies that $\|M_h(U,W)Z\|>\varepsilon$ and consequently

$$\inf_{\parallel z\parallel=1} \|\mathbf{M}_{h}(\mathbf{U},\mathbf{W})\,\mathbf{z}\,\| \geq \, \, \boldsymbol{\varepsilon}.$$

With this our proof of stability is complete. In order to obtain the desired result we quite explicitly assumed that $0 < \varepsilon < \frac{\partial f}{\partial y}$. The formal hypothesis given by Pereyra in [24] is that $0 \le \frac{\partial f}{\partial y}$. However, he implicitly assumes that $\frac{\partial f}{\partial y}$ is bounded away from zero in his proof of stability. Thus our proof is as general as Pereyra's, and considerably more elementary.

Since the operator F_h is both consistent and stable we can apply Theorem 1 to conclude that the discrete scheme (2) has a unique solution, X(h), which converges discretely to the theoretical solution, $\varphi(t)$, of problem (1).

Section 3. The Numerical Method

In this section we will summarize the results of Pereyra [24,27] and Stetter [31] concerning the global discretization error for the method (2). Using the results obtained by these authors we will develop a pullback interpolation scheme for solving boundary value problems of the form (1).

The global discretization error for the numerical method (2) is defined to be

$$X(h) - \omega_h \varphi$$
.

Both Pereyra [24,27] and Stetter [31] have presented a general theorem which, when applied to our problem, becomes

Theorem 2. (Pereyra, Stetter):

Let F and F_h have M+l continuous Fréchet derivatives. Then for sufficiently small h the global discretization error satisfies

(16)
$$x(h) - w_h \varphi = \sum_{k=1}^{M} h^{2k} w_h (e_k(t)) + O(h^{2M+2}).$$

The functions $e_k(t)$ are independent of h and satisfy the linear two point boundary value problem

$$e_{k}^{"}(t)-f_{z}(t,\varphi(t),\varphi'(t))e_{k}^{'}(t)-f_{y}(t,\varphi(t),\varphi'(t))e_{k}(t)=b_{k}(\cdot),$$
(17)

$$e_k(a) = 0, e_k(b) = 0.$$

The functions b_k depend on previous error functions

 e_1, \dots, e_{k-1} , various derivatives of φ , and various Fréchet derivatives of F.

We will prove an analogous theorem in Chapter 3 and the reader interested in details can refer to either Section 2.3 of the next chapter or to the aforementioned papers.

To set up the numerical method, define stepsizes $h_k = \frac{b-a}{2^{k+1}} \quad \text{for} \quad k=0,1,\ldots,M \quad \text{and construct the uniform grids}$ $G_k = \{t_i^k = a + ih_k : i = 0,1,\ldots,2^{k+1}\} \subset [a,b].$

Using the numerical method (2), compute a solution vector $X(h_k)$ on each grid G_k .

The computation of $X(h_k)$ for each k requires us to solve a system of $2^{k+1}-1$ equations. These equations will be nonlinear whenever f is nonlinear in either x or x'. In the nonlinear case we must use a root finding procedure to solve the system (e.g. Newton's method). In the linear case, the matrix is tridiagonal and can be easily solved using an LU decomposition and back substitution (see Isaacson and Keller [15]).

Once the solution vectors $X(h_k)$ for $k=0,1,\ldots,M$ have been obtained, extrapolation can be performed to obtain an $\mathfrak{S}(h_0^{2M+2})$ approximation to $\mathfrak{P}(t)$ at $t=\frac{b-a}{2}$. This has been studied by Pereyra. As was the case for initial value problems, extrapolation at other grid points does not yield comparable accuracy.

The task of implementing a pullback interpolation scheme for obtaining $\mathfrak{S}(h_0^{2M+2})$ accuracy at all grid points of G_M is simplified in this case. This is due to the fact that from (17) we now have two pieces of information available to us concerning each of the error functions e_k , namely, $e_k(a) = e_k(b) = 0$. Recall that for initial value problems only one piece of information was readily available concerning e_k . Most of the effort involved in implementing pullback interpolation for initial value problems was expended in obtaining an approximation to e_k' at the initial point. However, by using the boundary conditions $e_k(a) = 0$ and $e_k(b) = 0$, we will have enough data points to enable Lagrange interpolation to yield accuracy comparable to that of our data (see the discussion in Section 5 of Chapter 1).

Thus a pullback interpolation scheme based on Lagrange interpolation can be devised for solving boundary value problems of the form (2).

At $t = \frac{b-a}{2} \equiv a+h_O$ we solve an $(M+1) \times (M+1)$ system to obtain the solution vector $(\phi(t), e_1(t), \dots, e_M(t))$ with accuracy $(\mathcal{O}(h_O^{2M+2}), \mathcal{O}(h_O^{2M}), \dots, \mathcal{O}(h_O^2))$. Thus we know $e_M(a) = e_M(b) = 0$ exactly and $e_M(a+h_O)$ with accuracy $\mathcal{O}(h_O^2)$. Construct $L_M(s)$, the second degree Lagrange interpolating polynomial to the above data. Then $L_M(s) = e_M(s) + \mathcal{O}(h_O^2) + \mathcal{O}(h_O^3)$. The first ' \mathcal{O} ' term is the error in the initial data and the second is the error in interpolation.

Notice that we are using a sharper order estimate for the error of interpolation than was utilized in Chapter 1. It is well known that when performing interpolation on equally spaced points the error is of order a power of the stepsize (see Isaacson and Keller [15]). However, in Chapter 1, when interpolating on the three points a, $a + \frac{h}{2}$, and a+h we referred to the error as being $\mathfrak{G}(h^3)$. Actually, this error is $\mathfrak{G}((\frac{h}{2})^3) < \mathfrak{G}(h^3)$, so we were quite generous in our error estimates.

In order to show that pullback interpolation for solving boundary value problems yields $\mathcal{O}(h_0^{2M+2})$ accuracy at all grid points of the finest grid, we must use the sharper estimate for the error in interpolation when interpolating to e_M . When interpolating to the other error functions we can again be 'generous' and interpret the error as being in terms of h_0 which will be larger than all the other stepsizes.

The details for constructing the pullback interpolation scheme for boundary value problems are nearly identical to those given in Chapter 1. The exceptions are knowledge about e_k' is replaced by $e_k(b) = 0$ and we have one more data point each time we interpolate. Therefore, pullback interpolation will yield $\mathfrak{G}(h^{2M+2})$ accuracy at all grid points of the finest grid when used in conjunction with the numerical method (2) for solving two point boundary value problems of the form (1).

CHAPTER III

THE NUMERICAL SOLUTION OF DIFFERENCE DIFFERENTIAL EQUATIONS WITH CONSTANT RETARDATION

Section 1. First Order Equations

In this section we will be concerned with difference differential equations of the form

(1)
$$\dot{x}(t) = f(t,x(t),x(t-r)), r>0, t>0.$$

Since r > 0, (1) is an equation of retarded type. For brevity, we will sometimes refer to (1) as a delay differential equation.

This equation and variants of it occur frequently in mathematics. For instance, Lord Cherwell is credited by Wright [34] for using the equation

$$\dot{x}(t) = -\alpha x(t-1)(1+x(t))$$

in the study of the application of probability methods to the theory of asymptotic prime number densities. Cunningham [7] uses a variant of this equation as a growth model to describe a fluctuating population of organisms under certain conditions. Cooke and Yorke [4] use a variant of (1) as an economic model for the growth of capital stock.

Because of the presence of the delay term x(t-r) in (1), it is necessary to know the solution at time t-r to obtain a solution at time t. Therefore to solve (1) on the interval (0,r], it is necessary to specify the solution

on [-r,0] by means of an initial function, $\varphi(t)$. A solution of (1) with initial value φ is a continuous function which agrees with $\varphi(t)$ for $t \in [-r,0]$ and satisfies (1) for t>0. The solution for t>0 will also be denoted by $\varphi(t)$.

The usual method of obtaining the theoretical solution to (1) is called "the method of steps". A continuous initial function $\phi(t)$ is specified on [-r,0] and the initial value problem

$$\dot{x}(t) = f(t,x(t),\varphi(t-r))$$
(2)
$$x(0) = \varphi(0), \qquad 0 \le t \le r$$

is solved. Denoting the solution to (2) by $\varphi(t)$, the process is repeated on intervals of length r to obtain a solution, $\varphi(t)$, defined on $[-r, \bullet)$.

In general, the solution $\varphi(t)$ will have jump discontinuities in various derivatives at the multiples of r, zero included. In terms of smoothness of the solution, the worst possible situation is that a jump discontinuity will occur in the (k+1)st derivative of $\varphi(t)$ at the point kr for $k=0,1,\ldots$. By adding conditions on $\varphi(t)$ for $t\in[-r,0]$ these jump discontinuities can be avoided to some extent. For instance, if we require that the left hand derivative of $\varphi(t)$ at t=0, $\varphi'(0-)$, exists and satisfies

$$\varphi'(O-) = f(O,x(O),\varphi(-r));$$

then the jump discontinuity at kr is in the (k+2)nd derivative for $k=0,1,\ldots$. By imposing similar conditions on the initial function we can force the jump discontinuities to occur at whatever derivative of the solution we desire.

For f to have n continuous total derivatives with respect to t on (0,r] it is necessary that $\varphi(t)$ be n times continuously differentiable on [-r,0]. In this case the solution will be n+1 times continuously differentiable on (0,r]. However, without the added conditions discussed above on the left hand derivatives of $\varphi(t)$ at t=0 there will still be a discontinuity in the first and all higher order derivatives of the solution at t=0.

Note that as we proceed to the right using the method of steps to solve (1) we gain differentiability of the solution provided we do not lose any differentiability of f. That is, if f is n times continuously differentiable with respect to each of its arguments on $(-\infty,\infty)$ and the initial function is n times continuously differentiable on [-r,0], then the solution to (1) will be n+k times continuously differentiable for $t \in ((k-1)r, kr]$, $k=1,2,\ldots$.

In order to insure the existence of a unique solution, $\varphi(t)$, to (1) which is continuous on $[-r,\infty)$ we will make the following assumptions (see Hale [13]):

- (i) the initial function $\varphi(t)$ is continuous on [-r, 0];
- (iii) f(t,u,v) satisfies uniform Lipschitz conditions
 with respect to u and v.

Under these hypotheses we can also show that the solution to (1) depends continuously on the initial function $\phi(t)$ (see Hale [13]).

In order to obtain a numerical solution of (1) we shall require more stringent conditions on f and the initial function. These conditions will be stated as they become necessary.

To solve (1) numerically, given an initial function $\varphi(t)$ on [-r,0], define $F(t,x(t)) \equiv f(t,x(t),\varphi(t-r))$ and numerically solve the initial value problem

(3)
$$\dot{x}(t) = F(t,x(t))$$
$$x(0) = \varphi(0), \qquad 0 \le t \le r.$$

To solve (3) numerically we select a basic stepsize h in such a manner that r is an integer multiple of h; i.e., r = Nh for some integer N > 0. The reason for this restriction on h will become apparent later.

As in Chapter 1 we define stepsizes $h_k = h/2^k$ and grids $G_k = \{t_i^k = ih_k : i=0,1,\ldots,2^k\}$ for k=0,1,...,M. We assume that the numerical method employed to solve (3) is such that an asymptotic expansion of the form

(4)
$$X(t_i^k, h_k) = \varphi(t_i^k) + \sum_{j=1}^{M} e_j(t_i^k) h_k^{jq} + O(h_k^{(M+1)q})$$

is valid at each t_i^k for every h_k . Here, $X(t_i^k, h_k)$ is the numerical solution to (3) at the grid point t_i^k obtained with stepsize h_k and $\phi(t_i^k)$ is the solution to (1) at t_i^k .

The existence of an expansion (4) will, in general, require the existence of (M+1)q continuous derivatives of the solution $\varphi(t)$. This wil be assured if we assume that f(t,u,v) has (M+1)q-1 continuous derivatives with respect to each of its arguments and that the initial function $\varphi(t)$ is (M+1)q-1 times continuously differentiable on [-r,0].

If the numerical method used to solve (3) is such that q=1 or 2, then we can proceed as in Chapter 1 to obtain a numerical solution X(t) which satisfies

$$X(t) = \varphi(t) + \varphi(h^{(M+1)q}).$$

This relationship will be valid for all t belonging to the finest grid, G_{M} .

The presence of jump discontinuities in derivatives of the solution at multiples of the delay r does not affect

the applicability of the pullback interpolation method. The reason for this is that the pullback interpolation method employs derivatives that are computed using the differential equation at the initial point. Since the differential equation is valid only to the right of the initial point all derivatives considered in Chapter 1 are right hand derivatives at the initial point. If f(t,x(t),x(t-r)) has M total derivatives with respect to t and $\phi(t)$ has M right hand derivatives at t=-r, the solution of our delay differential equation will have M right hand derivatives at zero and therefore at all other multiples of r also. These derivatives can be computed by differentiating the differential equation in the same manner as was done in Chapter 1.

Repeating the solution procedure for the problem,

$$\dot{x}(t) = F(t,x(t))$$

$$x(h) = X(h), h \le t \le 2h,$$

we can obtain an $O(h^{(M+1)q})$ solution to (1) at the points $\{h+ih_M: i=0,1,\ldots,2^M\}$.

If we repeat the above procedure N times we will have a computed solution, X(t), satisfying

(5)
$$X(t) = \varphi(t) + \varphi(h^{(M+1)q})$$

for all t such that

$$t \in \{jh+ih_{M}: i=0,1,\ldots,2^{M}, j=0,1,\ldots,N-1\} = \{ih_{M}: i=0,\ldots,N2^{M}\}$$

$$\subset [0,r].$$

Let's examine what happens when we now try to obtain a solution to (1) on [r,2r]. Suppose t is any grid point in [r,2r]; because $G_k \subset G_M$ $\forall k$, t can be represented as $t=r+ih_M$ where $0 \le i \le N2^M$. Since $t \in [r,2r]$, $t-r \in [0,r]$ and in order to solve (1) we must be able to evaluate $\phi(t-r)$. Using the above representation of t, we can write $t-r=r+ih_M-r=ih_M$ with $0 \le i \le N2^M$. This is a grid point in [0,r] and by (5) we have

(6)
$$\varphi(t-r) = X(ih_{M}) + O(h^{(M+1)q}).$$

It is easy to see that there is a 1-1 correspondence between grid points in [r,2r] and grid points in [0,r]. This fact and equation (6) emphasize the importance of the pullback method for obtaining numerical solutions to difference differential equations. No numerical method is any more accurate than the data it receives and the solution of (1) on [r,2r] requires the solution on [0,r] as data. The uniform accuracy guaranteed by (6) means that the entire process used to obtain the solution on [0,r] can be repeated to obtain a solution on [r,2r].

Because normal extrapolation without pullback interpolation produces a solution whose accuracy varies from grid point to grid point, it is not an effective procedure for this problem.

Extrapolation has been used to obtain solutions to problems of the form (1). Feldstein [10] has developed two algorithms based on Euler's rule that have (partial) asymptotic expansions of the form

(7)
$$X(t) = \varphi(t) + e(t)h + \varphi(h^2)$$
.

An expansion of the form (7) justifies the use of one extrapolation.

Cooke and List [3] have developed an algorithm which they call the "modified Euler-Feldstein" algorithm which can be used to solve delay equations with non-constant retardation. They indicate that a sketch of a proof that their method has an expansion of the form (7) has been supplied by John Hutchison. However, neither the proof nor its sketch is contained in [3].

When using the "Euler-Feldstein" or the "modified Euler-Feldstein" algorithm the basic stepsize h is taken to be quite small (usually r/16 or smaller) and the problem is solved on several grids with stepsizes $h/2^k$. Extrapolation is performed at all multiples of the basic stepsize. The retarded term x(t-r) is obtained by using the solution in the previous interval that was computed using the same stepsize with which one is now working.

This procedure appears to work well in practice but is computationally quite expensive since the basic

of the form (7) is not sufficient justification for using more than one extrapolation.

Of course, systems of delay equations can be solved numerically by generalizing the above procedure in the obvious manner. A more interesting problem is the equation

$$\dot{x}(t) = f(t,x(t),x(t-r_1),x(t-r_2),...,x(t-r_n))$$

where $r_i>0$ is a rational number for $i=1,\dots,n$. The same numerical procedure will work for this problem provided the basic stepsize is chosen in a manner which insures that the information required by the delayed terms is available. Without loss of generality, assume the delays r_i are ordered: $0 \le r_1 \le r_2 \le \dots \le r_n$. Let d be the least common multiple of the denominators of the r_i for $i=1,\dots,n$, then with h=1/d this equation may be solved in the same manner as (1). This choice of the basic stepsize insures that each r_i will be an integer multiple of the smallest stepsize employed and therefore $t-r_i$ will be a grid point at which the solution is known $\forall i$ whenever t itself is a grid point.

We close this section with two numerical examples.

Example 1: Solve $\dot{x}(t) = -x(t-1)$, $0 \le t \le 3$, numerically for the initial function $\phi(t) = t^2$ on [-1,0].

Here, the delay is r=1 and the theoretical solution is given by

$$\varphi(t) = \begin{cases} \frac{-(t-1)^3 - 1}{3} & 0 \le t \le 1 \\ \frac{(t-2)^4 + 4t - 9}{12} & 1 \le t \le 2 \\ \frac{-(t-3)^5 - 10t^2 + 65t - 96}{60} & 2 \le t \le 3 \end{cases}.$$

Note that $\varphi(t)$ is given by a polynomial on each subinterval.

To compute a numerical solution on [0,3], we first solve the equation numerically on [0,1] using the trapezoid rule, extrapolation and pullback interpolation with M=3. The finest grid G_3 contains nine equally spaced points in [0,1]. Denote the computed solution at these nine points by $X_1(t)$. Using $X_1(t)$ as the (discrete) initial function we next solve the same equation on [1,2] in the same manner to obtain $X_2(t)$. This procedure is repeated a third time to obtain $X_3(t)$ on [2,3].

The computations of $e_1'(t)$, $e_2'(t)$ and $e_3'(t)$ at the initial points 0,1 and 2 are greatly simplified for this problem. One factor contributing to this simplification is that the Jacobian J(t), is identically zero. The second factor is that all derivatives of the solution $\phi(t)$ at the points 0,1 and 2 can be expressed in terms of the derivatives of the initial function $\phi(t)=t^2$ at

t=-1 and the computed solution at smaller multiples of r=1. In fact, by noting that for any k=0,1,2,... and for any j>k we have $x^{(j)}(kr) = x^{(j-1)}((k-1)r)$ it is easily seen that $x^{(j)}(kr) = x^{(j-k-1)}(-r)$. While, if $j \le k$, then $x^{(j)}(kr) = x((k-j)r)$.

The numerical results are presented as Table 13. The error given is the computed solution minus the theoretical solution. Each column of Table 13 presents the errors at the nine equally spaced grid points t_0, t_1, \ldots, t_8 contained in the interval indicated by the heading of the column.

TABLE 13

i	Error on [0,1]	Error on [1,2]	Error on [2,3]
0	0.0	0.0×10^{-13}	-1.5×10^{-13}
1	0.0×10^{-13}	3 x 10 ⁻¹³	0.0×10^{-13}
2	0.0×10^{-13}	4×10^{-13}	.2 x 10 ⁻¹³
3	0.0×10^{-13}	5 x 10 ⁻¹³	$.3 \times 10^{-13}$
4	0.0 x 10 ⁻¹³	6×10^{-13}	$.4 \times 10^{-13}$
5	0.0×10^{-13}	8 x 10 ⁻¹³	$.5 \times 10^{-13}$
6	0.0×10^{-13}	-1.0 x 10 ⁻¹³	$.6 \times 10^{-13}$
7	o.o x 10 ⁻¹³	-1.3 x 10 ⁻¹³	$.7 \times 10^{-13}$
8	0.0×10^{-13}	-1.5 x 10 ⁻¹³	.8 x 10 ⁻¹³

The errors reported in Table 13 are almost entirely

due to the limitations of machine accuracy and the particular

routine used to compute the Hermite interpolating polynomials. Only the first twelve digits can be absolutely guaranteed to be accurate and to twelve decimal places all errors are zero! This is not particularly surprising, since the trapezoid rule and extrapolation yield extremely good results when the theoretical solution is a polynomial, as it is in this case.

Example 2: Solve $\dot{x}(t) = -x(t-1)$, $0 \le t \le 3$ numerically for the initial function $\phi(t) = e^t$ on [-1,0]. The differential equation is the same as that in Example 1. The theoretical solution for this initial function is given by

$$\varphi(t) = \begin{cases} -e^{(t-1)} + 1 + e^{-1} & 0 \le t \le 1 \\ e^{(t-2)} - (1 + e^{-1}) (t-1) & 1 \le t \le 2 \\ -e^{(t-3)} + (1 + e^{-1}) \frac{t^2}{2} - 2 (1 + e^{-1}) (t-1) & 2 \le t \le 3 \end{cases}.$$

In this case the solutions are not polynomials.

The same numerical method was used to compute the solution as was used in Example 1. The results are presented as Table 14.

TABLE 14

i	Error on [0,1]	Error on [1,2]	Error on [2,3]
0	0.0	-1.2 x 10 ⁻¹⁰	-7.2×10^{-10}
1	-5.3×10^{-10}	4.5×10^{-10}	-12.8×10^{-10}
2	-3.2×10^{-10}	3.0 x 10 ⁻¹⁰	-11.2 x 10 ⁻¹⁰
3	2.3×10^{-10}	-2.6×10^{-10}	-5.6×10^{-10}
4	$.7 \times 10^{-10}$	-1.2 x 10 ⁻¹⁰	- 7.0 x 10 ⁻¹⁰
5	-2.4×10^{-10}	2.1 x 10 ⁻¹⁰	-10.2 x 10 ⁻¹⁰
6	12.0 x 10 ⁻¹⁰	-12.7 x 10 ⁻¹⁰	4.7×10^{-10}
7	42.4×10^{-10}	-46.7 x 10 ⁻¹⁰	38.9×10^{-10}
8	-1.2 x 10 ⁻¹⁰	- 7.2 x 10 ⁻¹⁰	.3 x 10 ⁻¹⁰

Since the theoretical solution is not a polynomial, the computed solution is not as accurate as the computed solution in Example 1. However, we do obtain eight reliable decimal places at all grid points, which is comparable to the accuracy that extrapolation yields at the endpoints 1, 2 and 3.

Also, it should be noted how stable the numerical results are. There is very little deterioration in the accuracy of the computed solution as we move to the right. Thus, extrapolation and pullback interpolation appears to be a viable solution technique for first order delay equations.

Section 2. Second Order Equations

For the remainder of this chapter we will be working with the second order difference differential equation of retarded type

(8)
$$\ddot{x}(t) + f(t,x(t),x(t-r), \dot{x}(t),\dot{x}(t-r)) = 0, r > 0, t > 0.$$

By writing (8) as a system of two first order equations and imposing conditions (i), (ii) and (iii) of Section 1 or the more general conditions of Sansone [30] on this system we can insure the existence of a unique solution to (8) which depends continuously on the initial data (see Hale [13] and Norkin [23]).

In either case, f is assumed to be a continuous function of its arguments. For our purposes it is enough to assume the existence of a unique solution to (8) for t>0 whenever we are given a continuously differentiable function $\phi(t)$ on [-r,0].

To solve (8) theoretically, the method of steps can again be used. The analysis of the behavior of the solution to (8) is completely analogous to that given for first order equations in Section 1 of this chapter. Accordingly, we will not go into this in detail. When specific aspects of the theoretical solution become important they will be mentioned.

If (8) is written as a system we can of course use the method discussed in Section 1 to solve this system.

The rest of this chapter will be devoted to an investigation of a direct method for solving (8) which does not involve reducing it to a system of first order equations.

We introduce the notation

$$(8') F(x) = 0$$

for the continuous operator defined by (8) and we will refer to (8) as the continuous problem.

For ease in referring to partial derivatives we will refer to equation (8) as

$$\ddot{x}(t) + f(t, u, v, y, z) = 0.$$

In setting up a discrete analog of (8) we actually get two slightly different problems; one when we are trying to obtain the solution on [O,r] and the second when we seek a solution on an interval of the form [(L-1)r,Lr] with L>1. The version of the problem for the interval [O,r] is a special case of the version for the general interval [(L-1)r,Lr] and will be explained in the subsection on implementing the algorithm (see section 2.4). Thus until further notice we will concern ourselves with the discrete analog of (8) on the interval [(L-1)r,Lr] with L>1.

Let R>1 be any given natural number and define h=r/R. Form the uniform mesh $\{t_i=(L-1)r+ih\}_{i=-(R+1)}^R$.

We assume the solution has already been obtained at the grid points t_i for i=-(R+1),-R,...,0 and we denote this solution by $\phi(t_i)$.

Setting $X_i = X(t_i)$, the discrete problem corresponding to (8) is given by

$$\frac{X_{i+1}^{-2X_{i}^{+}X_{i-1}}}{h^{2}} + f(t_{i}, X_{i}, X_{i-R}, \frac{X_{i+1}^{-}X_{i-1}}{2h}, \frac{X_{i+1-R}^{-}X_{i-1-R}}{2h}) = 0,$$
(9)
$$i = 0, 1, \dots, R-1;$$

$$X_{-j} = \varphi(t_{-j}), j=0,1,...,R+1.$$

The discrete problem (9) will be denoted by

(9')
$$F_h(x) = 0.$$

Since X_{-j} is known for j=0,...,R+1, this may be thought of as a (in general) nonlinear system of equations in E^R with the unknown being the vector $(X_1,...,X_R)^T$. The solution to (9) will be denoted by X(h).

Since each equation in (9) taken in the order i=0,1,...,R-l introduces only one unknown, we can solve the system uniquely by forward substitution. In general, the non-linear nature of f will necessitate the use of a root finding procedure (e.g. Newton's method) to solve each equation.

We will adopt the same space discretization operator w_h that was utilized in Chapter 2. That is, for any continuous function z(t) on $[-r,\infty)$

$$w_h z(t) = (z(t_1), \dots, z(t_R))^T$$

$$\equiv (z_1, \dots, z_R)^T.$$

We will continue to use the definitions of discrete convergence, consistency and stability (definitions 1, 2 and 3 respectively) that were adopted in Section 1 of Chapter 2.

The proofs of Lemma 1 and Theorem 1 in Section 1 of Chapter 2 are valid in this setting. In fact, Pereyra mentions in [28] that these results are not truly tied up to the two-point boundary value problem.

In order to establish the existence and discrete convergence of a unique solution X(h) to (9) we will again show that (9) is consistent and stable.

In order to prove that F_h is stable we have to make some restrictive assumptions about the equation (8). These assumptions are necessary for our proof of stability to be valid. However, we believe that F_h is stable for a much larger class of functions f(t,u,v,y,z) than what we are able to establish.

Section 2.1. Consistency

Let $\varphi(t)$ denote the theoretical solution to the continuous problem, (8). If we apply the operator F_h defined by (9) to the discretization of $\varphi(t)$ we obtain what is known as the local truncation error, τ_h . Formally

$$\tau_{h}(t_{i}) = F_{h}(\omega_{h}\phi)(t_{i})$$

$$\equiv \frac{\phi(t_{i+1}) - 2\phi(t_{i}) + \phi(t_{i-1})}{h^{2}} + f(t_{i}, \phi(t_{i}), \phi(t_{i-R}), \frac{\phi(t_{i+1}) - \phi(t_{i-1})}{2h}, \frac{\phi(t_{i+1}) - \phi(t_{i-1})}{2h})$$

where t_i is any grid point i=1,...,R.

If $\tau_h(t_i) = \mathfrak{O}(h^p)$ Vi then we obviously have that our method is consistent of order p. Instead of investigating (10) directly we let $\mathbf{x}(t)$ be any sufficiently differentiable function and investigate $F_h(\mathbf{w}_h\mathbf{x})(t_i)$. If we make the "localizing" assumption that no previous truncation errors have been made, we may use Taylor's Theorem repeatedly to obtain an asymptotic expansion for $F_h(\mathbf{w}_h\mathbf{x})(t_i)$.

Let t_i be any grid point, i=1,...,R. If x(t) has M_1+2 derivatives at t_i we can write

(11)
$$x(t_{i+1}) = x(t_i) + hx'(t_i) + \frac{h^2}{2!}x''(t_i)$$

$$+ \sum_{j=3}^{M_1+1} \frac{h^j}{j!} x^{(j)} (t_i) + O(h^{M_1+2})$$

and

(12)
$$x(t_{i-1}) = x(t_{i})-hx'(t_{i}) + \frac{h^{2}}{2!}x''(t_{i})$$

 $+\sum_{j=3}^{M_{1}+1} \frac{(-1)^{j}h^{j}}{j!} x^{(j)}(t_{i}) + O(h^{M_{1}+2}).$

Using (11) and (12) we have

(13)
$$\frac{x(t_{i+1})-2x(t_{i})+x(t_{i-1})}{h^{2}}$$

$$= h^{-2} \left(h^{2}x''(t_{i}) + \sum_{k=2}^{K} \frac{h^{2k}}{2k!} x^{(2k)}(t_{i}) + O(h^{2[\frac{M_{1}}{2}]+2})\right)$$

$$= x''(t_{i}) + 2\sum_{k=2}^{K} \frac{h^{2k-2}}{2k!} x^{(2k)}(t_{i}) + O(h^{2K})$$

where $[\cdot]$ is the greatest integer function and $K = \left[\frac{M_1}{2}\right]$.

Expanding f(',',',y,') in a Taylor series about

$$x'(t_i)$$
 where $y = \frac{x(t_{i+1})-x(t_{i-1})}{2h}$ we obtain

(14)
$$f(\cdot, \cdot, \cdot, \cdot, y, \cdot) = f(\cdot, \cdot, \cdot, x'(t_{i}), \cdot) + f_{y}(\cdot, \cdot, \cdot, x'(t_{i}), \cdot) (y-x'(t_{i}))$$

$$+ f_{yy}(\cdot, \cdot, \cdot, x'(t_{i}), \cdot) \frac{(y-x'(t_{i}))^{2}}{2!}$$

$$+ \sum_{j=3}^{M_{2}} \frac{1}{j!} \frac{\partial^{j} f}{\partial y^{j}} | (\cdot, \cdot, \cdot, x'(t_{i}), \cdot) (y-x'(t_{i}))^{j}$$

$$+ O[(y-x'(t_{i}))^{M_{2}+1}],$$

where M_2+1 is the number of partial derivatives of f(t,u,v,y,z) with respect to y which exist.

Using (11), (12) and the definition of y, we obtain

(15)
$$y-x'(t_i) = \frac{1}{2h} \{x(t_{i+1})-x(t_{i-1})\}-x'(t_i)$$

$$= \frac{1}{2h} \{2hx'(t_i)+2\sum_{j=1}^{K} \frac{h}{(2j+1)!} (t_i)+O(h^{2K+1})\}-x'(t_i)$$

$$= \sum_{j=1}^{K} \frac{h^{2j}}{(2j+1)!} x^{(2j+1)} (t_i)+O(h^{2K}).$$

Substituting (15) in (14), we have

$$f(\cdot,\cdot,\cdot,y,\cdot) = f(\cdot,\cdot,\cdot,x'(t_{i}),\cdot) + f_{y}(\cdot,\cdot,\cdot,x'(t_{i}),\cdot) \cdot \frac{K}{\sum_{j=1}^{K} \frac{h^{2j}}{(2j+1)!} x^{(2j+1)} (t_{i}) + o(h^{2K})}{(2j+1)!} \times \frac{h^{2j}}{(2j+1)!} x^{(2j+1)} (t_{i}) + o(h^{2K})}^{2K} + \frac{h^{2j}}{k+3} \frac{h^{2j}}{h^{2k}} (\cdot,\cdot,\cdot,x'(t_{i}),\cdot) \cdot \frac{K}{j+1} \frac{h^{2j}}{(2j+1)!} x^{(2j+1)} (t_{i}) + o(h^{2k})^{k} + o((\sum_{j=1}^{K} \frac{h^{2j}}{(2j+1)!} x^{(2j+1)} (t_{i}) + o(h^{2k}))^{k} + o(h^{2K})^{k} + o$$

But $O(\{\sum_{j=1}^{K} \frac{h^{2j}}{(2j+1)!} \times (2j+1) (t_i) + O(h^{2K})\}^{M_2+1}) = O(\{O(h^2) + O(h^{2K})\}^{M_2+1}) = O(h^{2K}) + O(h^{2K})^{M_2+1}) = O(h^{2K}) + O(h^{2K})^{M_2+1} = O(h^{2K}), \text{ and for any } \ell=1, \ldots, M_2,$ we have, from the binomial theorem,

$$\left\{ \sum_{j=1}^{K} \frac{h^{2j}}{(2j+1)!} x^{(2j+1)} (t_{i}) + O(h^{2K}) \right\}^{\ell} = \left(\sum_{j=1}^{K} \frac{h^{2j}}{(2j+1)!} x^{(2j+1)} (t_{i}) \right)^{\ell} + O(h^{2K}).$$

Consequently we can write

(16)
$$f(\cdot, \cdot, \cdot, \cdot, y, \cdot) = f(\cdot, \cdot, \cdot, x'(t_{i}), \cdot) + f_{y}(\cdot, \cdot, \cdot, x'(t_{i}), \cdot) \cdot \sum_{j=1}^{K} \frac{h^{2j}}{(2j+1)!} x^{(2j+1)} (t_{i}) + \sum_{k=2}^{M_{2}} \frac{1}{k!} \frac{\partial^{k} f}{\partial y^{k}} (\cdot, \cdot, \cdot, x'(t_{i}), \cdot) \cdot \left\{ \sum_{j=1}^{K} \frac{h^{2j}}{(2j+1)!} x^{(2j+1)} (t_{i}) \right\}^{k} + O(h^{2K}).$$

Next, expand $f(.,.,x'(t_i),z)$ and its partial derivatives with respect to y as functions of z in Taylor series about $x'(t_i-r)$ and then set $z=\frac{x(t_{i+1-R})-x(t_{i-1-R})}{2h}$ to obtain $(17) \quad f(.,.,x'(t_i),z)=f(.,.,x'(t_i),x'(t_i-r)) \\ + f_z(.,.,x'(t_i),x'(t_i-r))(z-x'(t_i-r)) \\ + \sum_{k=2}^{M_3} \left\{\frac{1}{k!} \frac{\partial^k f}{\partial z^k} \middle| (.,.,x'(t_i),x'(t_i-r)) \right\} \\ (z-x'(t_i-r))^k \\ + \mathcal{O}((z-x'(t_i-r))^{M_3+1}).$

In the same manner, for $m=1,...,M_2$, we can obtain

(18)
$$\frac{\partial^{m} f}{\partial y^{m}} \Big|_{(\cdot, \cdot, \cdot, \cdot, x'(t_{i}), z)} = \frac{\partial^{m} f}{\partial y^{m}} \Big|_{(\cdot, \cdot, \cdot, x'(t_{i}), x'(t_{i}-r))} + \frac{\partial^{m+1} f}{\partial z \partial y^{m}} \Big|_{(\cdot, \cdot, \cdot, x'(t_{i}), x'(t_{i}-r))} (z-x'(t_{i}-r)) + \sum_{\ell=2}^{M_{3}} \frac{1}{\ell!} \frac{\partial^{m+\ell} f}{\partial z^{\ell} \partial y^{m}} \Big|_{(\cdot, \cdot, \cdot, x'(t_{i}), x'(t_{i}-r))} (z-x'(t_{i}-r))^{\ell} + O(\{z-x'(t_{i}-r)\}^{M_{3}+1});$$

where M_3+1 is the number of partial derivatives of f(t,u,v,y,z) with respect to z which exist.

Since $t_i-r \equiv t_{i-R}$, an expansion similar to (15) is valid for $z-x'(t_i-r)$. Introducing the notation $\frac{\partial^{\circ}f}{\partial y^{\circ}}$ for the function f we can, in direct analog to (16), rewrite (17) and (18) as

(19)
$$\frac{\partial^{m} f}{\partial y^{m}} | (\cdot, \cdot, \cdot, x'(t_{i}), z) = \frac{\partial^{m} f}{\partial y^{m}} | (\cdot, \cdot, \cdot, x'(t_{i}), x'(t_{i}-r))$$

$$+ \sum_{k=1}^{M_{3}} \left\{ \frac{1}{k!} \frac{\partial^{m+k} f}{\partial z^{k} \partial y^{m}} | (\cdot, \cdot, \cdot, x'(t_{i}), x'(t_{i}-r)) \right\}$$

$$\left\{ \sum_{j=1}^{K} \frac{h^{2j}}{(2j+1)!} x^{(2j+1)} (t_{i}-r) \right\}^{k}$$

$$+ O(h^{2K}),$$

for $m=0,1,\ldots,M_2$.

Combining equations (16) and (19), we have

$$f(\cdot,\cdot,\cdot,y,z) = f(\cdot,\cdot,\cdot,x'(t_{i}),x'(t_{i}-r))$$

$$+ \sum_{k=1}^{M_{3}} \frac{1}{k!} \frac{\partial^{k} f}{\partial z^{k}} \Big|_{(\cdot,\cdot,\cdot,x'(t_{i}),x'(t_{i}-r)) = 1}^{K} \frac{\sum_{j=1}^{K} \frac{h^{2j}}{(2j+1)!} x^{(2j+1)} (t_{i}-r)}^{k} \Big|_{j=1}^{K} \frac{\partial^{k} f}{\partial z^{k}} \Big|_{(\cdot,\cdot,\cdot,x'(t_{i}),x'(t_{i}-r)) = 1}^{K} \frac{\sum_{j=1}^{K} \frac{h^{2j}}{(2j+1)!} x^{(2j+1)} (t_{i})}{\sum_{j=1}^{K} \frac{h^{2j}}{(2j+1)!} x^{(2j+1)} (t_{i})} \Big|_{j=1}^{K} \frac{\partial^{k} f}{\partial z^{k} \partial y} \Big|_{(\cdot,\cdot,\cdot,x'(t_{i}),x'(t_{i}-r))}^{K} \frac{\sum_{j=1}^{K} \frac{h^{2j}}{(2j+1)!} x^{(2j+1)} (t_{i})}{\sum_{j=1}^{2} \frac{h^{2j}}{(2j+1)!} x^{(2j+1)} (t_{i})} \Big|_{j=1}^{K} \frac{\partial^{k} f}{\partial z^{k} \partial y} \Big|_{(\cdot,\cdot,\cdot,x'(t_{i}),x'(t_{i}-r))}^{K} \Big|_{j=1}^{K} \frac{h^{2j}}{(2j+1)!} x^{(2j+1)} (t_{i})^{k} \Big|_{j=1}^{K} \frac{h^{2j}}{(2j+1)!} x^{(2j+1)} (t_{i})^{k} \Big|_{j=1}^{K} \frac{\partial^{k} f}{(2j+1)!} x^{(2j+1)} (t_{i}-r)^{k} \Big|_{j=1}^{K} \frac{\partial^{k} f}{(2k+1)!} x^{(2j+1)} (t_{i})^{k} \Big|_{j=1}^{K} \frac{\partial^{k} f}{(2k+1)!} \Big|_{j=1}^{K} \frac{\partial^{k} f}{(2k+1)!}$$

Combining equations (13) and (20), and evaluating all derivatives at $(t_i, x(t_i), x(t_i-r), x'(t_i), x'(t_i-r))$ we can write $F_h(\omega_h x)(t_i)$ as

$$\begin{split} F_{h}\left(\omega_{h}^{\times}\right)\left(t_{i}\right) &= \frac{x\left(t_{i+1}\right)-2\left(t_{i}\right)+x\left(t_{i-1}\right)}{h^{2}} + f\left(t_{i},x\left(t_{i}\right),x\left(t_{i-R}\right),\right. \\ &= \frac{x\left(t_{i+1}\right)-x\left(t_{i-1}\right)}{2h}, \frac{x\left(t_{i+1-R}\right)-x\left(t_{i-1-R}\right)}{2h}) \\ &= x''\left(t_{i}\right)+2\sum_{k=2}^{K}\frac{h^{2k-2}}{(2k)!}x^{\left(2k\right)}\left(t_{i}\right)+O(h^{2K}) \\ &+ f\left(t_{i},x\left(t_{i}\right),x\left(t_{i}-r\right),x'\left(t_{i}\right),x'\left(t_{i}-r\right)\right) \\ &+ \sum_{k=1}^{M_{3}}\frac{1}{k!}\frac{\partial^{k}f}{\partial z^{k}}\left(\sum_{j=1}^{K}\frac{h^{2j}}{(2j+1)!}x^{\left(2j+1\right)}\left(t_{i}-r\right)\right)^{k} \\ &+ \sum_{k=1}^{M_{2}}\frac{1}{k!}\left[\sum_{k=0}^{M_{3}}\frac{\partial^{k+k}f}{\partial z^{k}\partial y^{k}}\left(\sum_{j=1}^{K}\frac{h^{2j}}{(2j+1)!}x^{\left(2j+1\right)}\left(t_{i}-r\right)\right)^{k}\right] \\ &- \left[\sum_{j=1}^{K}\frac{h^{2j}}{(2j+1)!}x^{\left(2j+1\right)}\left(t_{i}\right)\right]^{k}\right\} \\ &+ O(h^{2M_{2}+2}) + O(h^{2M_{3}+2}) + O(h^{2K}) \,. \end{split}$$

Define $M = \min(K-1, M_2, M_3)$; then, reindexing, we have

$$(21) \quad F_{h}(\omega_{h}x) (t_{i}) = F(x) (t_{i}) + 2 \sum_{k=2}^{M} \frac{h^{2k-2}}{(2k)!} x^{(2k)} (t_{i})$$

$$+ \sum_{k=1}^{M} \frac{1}{k!} \frac{\partial^{k} f}{\partial z^{k}} (\sum_{j=1}^{M} \frac{h^{2j}}{(2j+1)!} x^{(2j+1)} (t_{i}-r))^{k}$$

$$+ \sum_{k=1}^{M} \frac{1}{k!} \{ [\sum_{k=0}^{M} \frac{\partial^{k+k} f}{\partial z^{k} \partial y^{k}} (\sum_{j=1}^{M} \frac{h^{2j}}{(2j+1)!} x^{(2j+1)} (t_{i}-r))^{k} \}$$

$$[\sum_{j=1}^{M} \frac{h^{2j}}{(2j+1)!} x^{(2j+1)} (t_{i})^{k} \}$$

$$+ O(h^{2M+2}).$$

If we regroup (21) in powers of h we obtain

(21')
$$F_h(\omega_h^x)(t_i) = F(x)(t_i) + \sum_{j=1}^{M} h^{2j} g_{2j}(\cdot) |_{x,t_i} + O(h^{2M+2}),$$

where the arguments of the functions $g_{2j}(\cdot)$ involve various derivatives of x evaluated at t_i and t_i -r, and various partial derivatives of f(t,u,v,y,z) with respect to y and/or z evaluated at $(t_i,x(t_i-r),x'(t_i),x'(t_i-r))$. For instance

$$g_2(\cdot) = \frac{1}{12}x^{(4)}(t_i) + \frac{1}{6}(\frac{\partial f}{\partial z}x^{(3)}(t_i-r) + \frac{\partial f}{\partial y}x^{(3)}(t_i))$$

and

$$g_{4}(\cdot) = \frac{2}{6!} x^{(6)} (t_{i}) + \frac{1}{2!} \frac{1}{(3!)^{2}} \frac{\partial^{2} f}{\partial z^{2}} (x^{(3)} (t_{i}-r))^{2}$$

$$+ \frac{1}{5!} (\frac{\partial f}{\partial z} x^{(5)} (t_{i}-r) + \frac{\partial f}{\partial y} x^{(5)} (t_{k}))$$

$$+ \frac{1}{(3!)^{2}} \frac{\partial^{2} f}{\partial z \partial y} x^{(3)} (t_{i}-r) x^{(3)} (t_{i}) + \frac{\partial^{2} f}{\partial y^{2}} (x^{(3)} (t_{i}))^{2}).$$

To obtain the local truncation error for the discrete problem (9), by (10) we need only take x(t) to be the theoretical solution $\varphi(t)$ of problem (8) and use either equation (21) or (21'). Thus,

(22)
$$\tau_h(t_i) = F(\phi)(t_i) + \sum_{j=1}^{M} h^{2j} g_{2j}(\cdot) |_{\phi, t_i} + O(h^{2M+2})$$

from equation (21').

To establish the consistency of our discrete operator we need only note that $F(\phi) \equiv 0$. Consequently

 $\tau_h(t_i) = O(h^2)$ for every grid point t_i , i=1,...,R. Thus our discrete operator is consistent of order 2.

Section 2.2. Stability

In this section we investigate the stability of the discrete scheme (9). The first part of the analysis is done for the general method and is analogous to the work done in Section 2 of Chapter 2.

In investigating the stability of the discrete problem (9), it is convenient to view F_h as a function from E^R to E^R . The domain of F_h will consist of vectors $V = (V_1, \dots, V_R)^T \in E^R$. Of course, F_h continues to have the initial conditions

$$v_{-j} = \phi(t_{-j}), j=0,1,...,R+1.$$

When considering the j^{th} component of $F_h(V)$ we will suppress the (constant) first, third, and fifth arguments of f(t,u,v,y,z) and write

$$f(t_{j-1}, v_{j-1}, v_{j-1-R}, \frac{v_{j}-v_{j-2}}{2h}, \frac{v_{j-R}-v_{j-2-R}}{2h})$$

$$= f(\cdot, v_{j-1}, \cdot, \frac{v_{j}-v_{j-2}}{2h}, \cdot) .$$

Note that the jth component of $F_h(V)$ is obtained from the (j-1)st equation in (9) for j=1,...,R.

Let U,W be any two vectors in \mathbf{E}^R with $\mathbf{U} \neq \mathbf{W}$. The cases j=1 and j=2 are considered individually later. For $3 \leq \mathbf{j} \leq R$ the jth component of $\mathbf{F}_h(\mathbf{U}) - \mathbf{F}_h(\mathbf{W})$ is

given by

$$[F_{h}(Y) - F_{h}(W)]_{j}$$

$$= h^{-2} (U_{j-2} - 2U_{j-1} + U_{j}) + f(\cdot, U_{j-1}, \cdot, \frac{U_{j} - U_{j-2}}{2h}, \cdot)$$

$$-h^{-2} (W_{j-2} - 2W_{j-1} + W_{j}) + f(\cdot, W_{j-1}, \cdot, \frac{W_{j} - W_{j-2}}{2h}, \cdot)$$

$$= h^{-2} [(U_{j-2} - W_{j-2}) - 2(U_{j-1} - W_{j-1}) + (U_{j} - W_{j})]$$

$$+ f(\cdot, U_{j-1}, \cdot, \frac{U_{j} - U_{j-2}}{2h}, \cdot) - f(\cdot, W_{j-1}, \cdot, \frac{W_{j} - W_{j-2}}{2h}, \cdot) .$$

Proceeding as in Chapter 2, we use the Mean Value

Theorem for continuous functions of two variables (Widder [33])

to obtain

(24)
$$f(\cdot, U_{j-1}, \cdot \frac{U_{j}^{-U_{j-2}}}{2h}, \cdot) - f(\cdot, W_{j-1}, \cdot, \frac{W_{j}^{-W_{j-2}}}{2h}, \cdot)$$

$$= f_{u}(\cdot, \alpha_{j}, \cdot, \beta_{j}, \cdot) (U_{j-1}^{-W_{j-1}})$$

$$+ f_{v}(\cdot, \alpha_{j}, \cdot, \beta_{j}, \cdot) (\frac{U_{j}^{-U_{j-2}}}{2h} - \frac{W_{j}^{-W_{j-2}}}{2h});$$

where

(25a)
$$\alpha_{j} = W_{j-1} + \theta_{j} (U_{j-1} - W_{j-1})$$

and

(25b)
$$\beta_{j} = \frac{w_{j}^{-W}_{j-2}}{2h} + \beta_{j} (\frac{u_{j}^{-U}_{j-2}}{2h} - \frac{w_{j}^{-W}_{j-2}}{2h}),$$

with $0 < \theta_{i} < 1$.

Substituting (24) into (23) we have for $3 \leq j \leq R$

$$[F_{h}(U) - F_{h}(W)]_{j}$$

$$= h^{-2} (U_{j-2} - W_{j-2}) - 2h^{-2} (U_{j-1} - W_{j-1}) + h^{-2} (U_{j} - W_{j})$$

$$+ f_{u}(\cdot, \alpha_{j}, \cdot, \beta_{j}, \cdot) (U_{j-1} - W_{j-1})$$

$$+ \frac{1}{2h} f_{y}(\cdot, \alpha_{j}, \cdot, \beta_{j}, \cdot) (U_{j} - W_{j}) - \frac{1}{2h} f_{y}(\cdot, \alpha_{j}, \cdot, \beta_{j}, \cdot)$$

$$(\frac{U_{j-2} - W_{j-2}}{2h})$$

$$= h^{-2} [(1 - \frac{h}{2} f_{y}(\cdot, \alpha_{j}, \cdot, \beta_{j}, \cdot)) (U_{j-2} - W_{j-2})$$

$$+ (-2 + h^{2} f_{u}(\cdot, \alpha_{j}, \cdot, \beta_{j}, \cdot)) (U_{j-1} - W_{j-1})$$

$$+ (1 + \frac{h}{2} f_{y}(\cdot, \alpha_{j}, \cdot, \beta_{j}, \cdot)) (U_{j} - W_{j})].$$

For j=1 and 2 we can reason in the same manner to obtain

(27)
$$[F_h(U) - F_h(W)]_1 = h^{-2} (1 + \frac{h}{2} f_y(\cdot, \alpha_1, \cdot, \beta_1, \cdot)) (U_1 - W_1)$$

and

(28)
$$[F_h(U) - F_h(W)]_2 = h^{-2}[(-2 + h^2 f_u(\cdot, \alpha_2, \cdot, \beta_2, \cdot) (U_1 - W_1) + (1 + \frac{h}{2} f_y(\cdot, \alpha_2, \cdot, \beta_2, \cdot) (U_2 - W_2)].$$

The functions $\alpha_1, \beta_1, \alpha_2$, and β_2 are given by

(29a)
$$\alpha_1 = \varphi(t_0) \equiv W_0$$

(29b)
$$\beta_1 = \frac{W_1 - W_{-1}}{2h} + \theta_1 (\frac{U_1 - W_1}{2h}),$$

(30a)
$$\alpha_2 = w_1 + \theta_2 (U_1 - W_1)$$
,

and

(30b)
$$\beta_2 = \frac{W_2 - W_0}{2h} + \theta_2 (\frac{U_2 - W_2}{2h})$$

with $0 < \theta_1 < 1$ and $0 < \theta_2 < 1$.

Writing (26), (27) and (28) as a matrix system we obtain

(31)
$$F_h(U) - F_h(w) = M_h(U, w) (U-w)$$

where $M_h(U,W)$ is an $R \times R$ matrix whose rows are given by:

$$[M_h(U,W)]_1 = h^{-2}(1 + \frac{h}{2}f_Y(\cdot,\alpha_1,\cdot\beta_1,\cdot),0,0,\ldots,0);$$

(32)
$$[M_h(U,W)]_2 = h^{-2}(-2+h^2f_u(\cdot,\alpha_2,\cdot\beta_2,\cdot),1+\frac{h}{2}f_y(\cdot,\alpha_2,\cdot\beta_2,\cdot),$$

$$[M_{h}(U,W)]_{j} = h^{-2}(0,...,0,1 - \frac{h}{2}f_{y}(\cdot,\alpha_{j},\cdot\beta_{j},\cdot),$$

$$-2+h^{2}f_{u}(\cdot,\alpha_{j},\cdot,\beta_{j},\cdot),1 + \frac{h}{2}f_{y}(\cdot,\alpha_{j},\cdot,\beta_{j},\cdot),$$

$$0,...,0)$$

for 3 < j < R;

with the non-zero entries for $j \ge 3$ occurring in the (j-2)nd, (j-1)st and jth positions. The arguments α_j and β_j for $j=1,\ldots,R$ are given by the appropriate equation (25), (29) or (30). Note that since R = r/h the dimensionality of $M_h(U,W)$ depends on h.

Since $M_h(U,W)$ is not diagonally dominant, as it was for boundary value problems, we are not able to prove directly that $M_h(U,W)$ is bounded below. However, we do have the following equivalent formulation of the concept of stability.

Suppose $M_h(U,W)$ has an inverse $M_h^{-1}(U,W)$ that is bounded above in norm, independent of h. Then we can write (31) as

$$U-W = M_h^{-1}(U, W) (F_h(U) - F_h(W)).$$

Taking norms and using the fact that the norm of a product is less than or equal to the product of the norms, we have

where C is the bound on $\|M_h^{-1}(U,W)\|$. This inequality is precisely the definition of stability given in Section 1 of Chapter 2.

The norm inequality we have used is valid as long as we use the matrix norm induced by our vector norm. Since we are working with the maximum vector norm on $\mathbf{E}^{\mathbf{R}}$, the induced matrix norm is the maximum absolute row sum (see Isaacson and Keller [15]).

We are also unable to prove directly that $\|\mathbf{M}_h^{-1}(\mathbf{U},\mathbf{W})\|$ is uniformly bounded above. However, we are able to establish this fact with some additional hypotheses on f.

For notational convenience we will assume that r=1. This is a rescaling of the delay interval and has no effect on any of our results. We will also write the dimensionality of the matrix system (34) as n = R = 1/h.

The special case for which we will prove stability is given by the following assumptions:

(33)
$$\frac{\partial f}{\partial y} \equiv O,$$

$$-2 < -K < \frac{\partial f}{\partial u} < -\sigma < O,$$

where $\sigma \in K$ are positive constants.

Write $b_i = f_u(\cdot, \alpha_i, \cdot, \beta_i, \cdot)$. Then the assumptions (33) imply that b_i is negative for all i; in fact,

(34)
$$-\kappa < b_i < -\sigma$$
.

Rewriting the matrix $M_h \equiv M_h(U,W)$ given in (32) in terms of b_i and n and imposing the assumptions (33), we obtain

$$M_{h}=n^{2} \cdot \begin{bmatrix} 1 & , & 0 & , & 0 & , & 0 & , & \dots & 0 \\ -2+\frac{b_{2}}{n^{2}}, & 1 & , & 0 & , & 0 & , & & \dots \\ 1 & , & -2+\frac{b_{3}}{n^{2}}, & 1 & , & 0 & , & & \dots \\ 0 & , & 1 & , & -2+\frac{b_{4}}{n^{2}}, & 1 & , & & \dots \\ \vdots & \vdots & \ddots & \ddots & \ddots & \ddots & \ddots \\ \vdots & \ddots & \ddots & \ddots & \ddots & \ddots & \ddots \\ \vdots & \ddots & \ddots & \ddots & \ddots & \ddots & \ddots \\ \vdots & \ddots & \ddots & \ddots & \ddots & \ddots & \ddots \\ \vdots & \ddots & \ddots & \ddots & \ddots & \ddots & \ddots \\ \vdots & \ddots & \ddots & \ddots & \ddots & \ddots & \ddots \\ \vdots & \ddots & \ddots & \ddots & \ddots & \ddots & \ddots \\ \vdots & \ddots & \ddots & \ddots & \ddots & \ddots & \ddots \\ \vdots & \ddots & \ddots & \ddots & \ddots & \ddots & \ddots \\ \vdots & \ddots & \ddots & \ddots & \ddots & \ddots & \ddots \\ \vdots & \ddots & \ddots & \ddots & \ddots & \ddots & \ddots \\ \vdots & \ddots & \ddots & \ddots & \ddots & \ddots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots \\ \vdots & \ddots & \ddots & \ddots & \ddots &$$

Define

then

$$S_{h}^{M}{}_{h} = \begin{bmatrix} \frac{1}{2} - \frac{b_{2}}{4n^{2}}, & -\frac{1}{4}, & 0 & \dots & 0 \\ -\frac{1}{4}, & \frac{1}{2} - \frac{b_{3}}{4n^{2}}, & -\frac{1}{4}, & & & \\ 0 & & & \ddots & & \\ \vdots & & & \ddots & & \\ 0 & & & -\frac{1}{4}, & \frac{1}{2} - \frac{b_{n}}{4n^{2}}, & -\frac{1}{4} \end{bmatrix}$$

Write

$$(35) S_h^M = I_h - T_h$$

with

where I_h is the $n \times n$ identity matrix.

Because of (34), for n sufficiently large (h sufficiently small) we have

$$0 < \frac{1}{2} + \frac{b_i}{4n^2} < \frac{1}{2} - \frac{\sigma}{4n^2}$$

for i=2,...,n. Therefore the first through (n-1)st rows of T_h all have absolute row sum less than $1-\frac{\sigma}{4n^2}$.

All entries of the last row of T_h are positive, so the absolute row sum of the $n^{\mbox{th}}$ row is

$$1 - \frac{1}{4n^{2}} - \frac{1}{4n^{4}} \sum_{i=1}^{n-1} ib_{n+1-i}$$

$$= 1 - \frac{1}{4n^{2}} + \frac{1}{4n^{4}} \sum_{i=1}^{n-1} i|b_{n+1-i}|$$

$$\leq 1 - \frac{1}{4n^{2}} + \frac{1}{4n^{4}} K \sum_{i=1}^{n-1} i$$

$$= 1 - \frac{1}{4n^{2}} + \frac{Kn(n-1)}{8n^{4}}.$$

Let $\delta = 2 - K > 0$, then

$$\frac{1}{4n^{2}} - \frac{Kn(n-1)}{8n^{4}} = \frac{1}{8n^{2}} \left(2 - K(\frac{n-1}{n})\right)$$

$$> \frac{1}{8n^{2}} (2 - K)$$

$$= \frac{\delta}{8n^{2}}.$$

Therefore,

$$1 - \frac{1}{4n^2} - \frac{1}{4n^4} \sum_{i=1}^{n-1} ib_{n+1-i} \le 1 - \left(\frac{1}{4n^2} - \frac{Kn(n-1)}{8n^2}\right)$$

$$< 1 - \frac{\delta}{8n^2}.$$

Thus,

$$\|T_h\| \le \max (1 - \frac{\sigma}{4n^2}, 1 - \frac{\delta}{8n^2})$$

$$= 1 - \frac{1}{4n^2} \min (\sigma, \frac{\delta}{2}).$$

Consequently, for all n $||T_h|| < 1$. This implies by a well known theorem in linear algebra (see Varga [32] or Isaacson and Keller [15]) that $(I_h - T_h)^{-1}$ exists and moreover

(36)
$$\|(\mathbf{I}_{h}^{-\mathbf{T}_{h}})^{-1}\| \leq \frac{1}{1-\|\mathbf{T}_{h}\|}.$$

 ${\bf S}_{\bf h}$ is easily seen to be invertible by observing that it is row equivalent to ${\bf I}_{\bf h}.$ Therefore ${\bf M}_{\bf h}$ must be nonsingular and from (35) we have

(37)
$$M_h^{-1} = (I_h - T_h)^{-1} S_h$$

Taking norms in (37) and using the norm inequality (36), we can write

(38)
$$\|M_{h}^{-1}\| \leq \|(I_{h}^{-1}T_{h}^{-1})^{-1}\| \cdot \|S_{h}\|$$

$$\leq \frac{\|S_{h}\|}{1-\|T_{h}\|} \cdot$$
Since $\|T_{h}\| \leq 1 - \frac{1}{4n^{2}} \min (\sigma, \frac{\delta}{2})$, we have
$$1 - \|T_{h}\| \geq 1 - (1 - \frac{1}{4n^{2}}, \min (\sigma, \frac{\delta}{2}))$$

$$= \frac{1}{4n^{2}} \min (\sigma, \frac{\delta}{2}).$$

But

$$\|\mathbf{S}_{h}\| = \max \left(\frac{1}{4n^{2}}, \frac{1}{4n^{4}} \sum_{i=1}^{n} i\right)$$

$$= \max \left(\frac{1}{4n^{2}}, \frac{1}{4n^{4}} \frac{n(n+1)}{2}\right)$$

$$= \frac{1}{4n^{2}} \max \left(1, \frac{n+1}{2n}\right)$$

$$= \frac{1}{4n^{2}}.$$

Finally, using the above estimates in (38), we have $\|M_h^{-1}\| \leq \frac{\frac{1}{4n^2}}{\frac{1}{4n^2}\min\ (\sigma,\frac{\delta}{2})} = \frac{1}{\min\ (\sigma,\frac{\delta}{2})}.$

Since this is independent of n, and therefore of h, we can conclude that $\|M_h^{-1}\| \le C$; where C is a finite constant independent of h.

Thus F_h is stable and Theorem 1 of Chapter 2 allows us to conclude that (9) has a unique solution X(h) which converges discretely to $\phi(t)$.

Clearly, the assumptions (33) are required by our method of proof. We believe that there is nothing inherent to the operator F_h which would prevent a more general proof, but at this time we are unable to construct one.

In the next section we will return to the study of the general operator F_h , without the assumptions (33), and examine its global discretization error.

Section 2.3. The Global Discretization Error

In this section we investigate $\mathbf{E}(\mathbf{h}) \equiv \mathbf{X}(\mathbf{h}) - \omega_{\mathbf{h}} \varphi$, the global discretization error for our numerical method. Assuming the discrete convergence of our method, we know that $\mathbf{E}(\mathbf{h}) \to 0$ as $\mathbf{h} \to 0$. We seek more information as to the nature of $\mathbf{E}(\mathbf{h})$ as a function of \mathbf{h} . Ideally we would like $\mathbf{E}(\mathbf{h})$ to have an asymptotic expansion in even powers of \mathbf{h} . That is, we would like to show that $\forall \mathbf{t} \in [(\mathbf{L}-\mathbf{l})\mathbf{r},\mathbf{L}\mathbf{r}]$.

(39)
$$E(h) = \omega_h \left(\sum_{k=1}^{M} e_k(t) h^{2k} \right) + O(h^{2M+2})$$

with the functions $e_k(t)$ independent of h. Equivalently, we want to show how to determine $e_k(t)$ for $k=1,2,\ldots,M$ and $t \in [(L-1)r, Lr]$ so that

(40)
$$S(h) = E(h) - \omega_h \left(\sum_{k=1}^{M} e_k(t) h^{2k} \right) = O(h^{2M+2}).$$

If we write

$$u(h) = \sum_{k=1}^{M} e_k(t) h^{2k}$$

then we can rewrite (40) as

$$S(h) = E(h) - \omega_h(\mu(h)).$$

Let $I = F_h(\omega_h(\phi + \mu(h)))$, then since $F_h(X(h)) = 0$ we have $\|I\| = \|F_h(\omega_h(\phi + \mu(h))) - F_h(X(h))\|.$

Using the linearity of $\ \omega_h^{}$ and the stability of $\ F_h^{},$ we can write

$$||\mathbf{I}|| \ge \frac{1}{C} ||\mathbf{X}(\mathbf{h}) - \mathbf{w}_{\mathbf{h}} \mathbf{\varphi} - \mathbf{w}_{\mathbf{h}} \mathbf{u}(\mathbf{h})||$$

 $= \frac{1}{C} ||\mathbf{E}(\mathbf{h}) - \mathbf{w}_{\mathbf{h}} \mathbf{u}(\mathbf{h})||$
 $= \frac{1}{C} ||\mathbf{S}(\mathbf{h})||.$

Consequently, if we can show how to choose $e_k(t)$ for $k=1,\ldots M$ so that $\|I\|=\mathfrak{G}(h^{2K+2})$ then (40) will be valid and we will have the desired asymptotic expansion.

Theorem 1. Let $F_h(V)$ have M+1 continuous Frechet derivatives with respect to V. Then,

$$x(h) - w_h \varphi = \sum_{k=1}^{M} h^{2k} w_h e_k(t) + O(h^{2M+2})$$

where the vectors $w_h e_k(t)$ are solutions of

(41)
$$F_h'(\omega_h \varphi) \omega_h e_k = \omega_h B_k.$$

The vectors $w_h^B{}_k$ will be defined in the proof.

<u>Proof.</u> Expand I = $F_h(w_h(\phi+\mu(h)))$ in a Taylor series about the vector $w_h\phi$ to get

(42)
$$I = F_h(\omega_h \varphi) + F_h'(\omega_h \varphi) \omega_h u(h) + \sum_{k=2}^{M} \frac{1}{k!} F_h^{(k)}(\omega_h \varphi) (\omega_h u(h))^k + R_{M+1}$$

Here $F_h^{(k)}(\omega_h\phi)$ is the k^{th} Frechet derivative of F_h evaluated at $\omega_h\phi$. The remainder term R_{M+1} involves $[\mu(h)]^{M+1}$ which is $\mathcal{O}(h^{2M+2})$ and therefore R_{M+1} is $\mathcal{O}(h^{2M+2})$.

 $F_h\left(w_h\phi\right) \ \ \text{is the vector of local truncation errors}$ for the discrete operator $F_h. \ \ \text{By equation (22) each component}$ of $F_h\left(w_h\phi\right) \ \ \text{is given by}$

$$\tau_{h}(t_{i}) = F(\varphi)(t_{i}) + \sum_{j=1}^{M} h^{2j} g_{2j}(\cdot) |_{\varphi, t_{i}} + O(h^{2M+2}).$$

Using the space discretization w_h , we can write

(43)
$$F_h(\omega_h \varphi) = \omega_h \tau_h$$

$$= \omega_h F(\varphi) + \sum_{j=1}^{M} h^{2j} \omega_h (g_{2j}(\cdot) |_{\varphi}) + O(h^{2M+2}).$$

Here we have used the linearity of w_h . Also, the $O(h^{2M+2})$ term in (43) is understood to be a vector each of whose components is $O(h^{2M+2})$.

Noting that $F(\phi) \equiv 0$ and substituting (43) into (42), we obtain

(44)
$$I = \sum_{j=1}^{M} h^{2j} w_h g_{2j}(\cdot) \Big|_{\varphi} + F_h(w_h \varphi) w_h u(h) + \sum_{k=2}^{M} \frac{1}{k!} F_h^{(k)}(w_h \varphi) [w_h u(h)]^k + O(h^{2M+2}).$$

Using the definition of $\mu(h)$, it is immediate that the coefficient of h^2 in (44) is

$$w_h g_2 + F_h'(w_h \varphi) w_h e_1$$

In order to eliminate this term from the expansion (44), define $w_h B_1 = -w_h g_2$ and take $w_h e_1$ to be the solution of

$$F_h'(\omega_h \varphi) \omega_h e_1 = \omega_h B_1$$
.

With $w_h e_1$ so determined we can write (44) as

$$\begin{split} \mathbf{I} &= \sum_{j=2}^{M} h^{2j} \mathbf{w}_{h} \mathbf{g}_{2j} + \mathbf{F}' (\mathbf{w}_{h} \mathbf{\phi}) \, \mathbf{w}_{h} (\sum_{j=2}^{M} h^{2j} \mathbf{e}_{j}) \\ &+ \sum_{k=2}^{M} \frac{1}{k!} \, \mathbf{F}_{h}^{(k)} (\mathbf{w}_{h} \mathbf{\phi}) \left[\, \mathbf{w}_{h} \mathbf{u}(\mathbf{h}) \, \right]^{k} \, + \, \mathcal{O}(h^{2M+2}) \, . \end{split}$$

Collecting the terms involving h4 we have

$$h^4 w_h g_4 + F'(w_h \varphi) h^4 w_h e_2 + \frac{1}{2} F_h^{(2)}(w_h \varphi) h^2 w_h e_1 h^2 w_h e_1$$

Define $w_h B_2 = -w_h g_4 - \frac{1}{2} F_h^{(2)} (w_h \varphi) w_h e_1 w_h e_1$ and obtain $w_h e_2$ by solving

$$F_h'(\omega_h \varphi) \omega_h e_2 = \omega_h B_2$$

Substituting $\omega_h^e{}_2$ in the expression for I we can eliminate the terms involving $h^4.$

Suppose $w_h e_1, w_h e_2, \ldots, w_h e_J$ have all been determined as solutions to equations of the form (41), so that the expression for I is

$$I = \sum_{j=J+1}^{M} h^{2j} w_{h} g_{2j}(\cdot) |_{\varphi} + F'(w_{h} \varphi) w_{h} \mu(h)$$

$$+ \sum_{k=J+1}^{M} \frac{1}{k!} F_{h}^{(k)}(w_{h} \varphi) [w_{h} \mu(h)]^{k}$$

$$+ \sum_{k=J+1}^{JM} G_{2k}(\cdot) h^{2k} + O(h^{2M+2}).$$

The arguments of the functions $G_{2k}(\cdot)$ involve various error vectors $w_h e_j$ for $j=1,\ldots,J$ and various Frechet derivatives

$$F_h^{(j)}$$
 ($w_h \varphi$) for $j=2,...,J$. If
$$\sum_{j=2}^{J} \frac{1}{j!} F_h^{(j)} (w_h \varphi) [w_h u(h)]^j$$

is expressed as a power series in h^2 then for $k=J+1,\ldots,JM$ the function $G_{2k}(\cdot)$ can be obtained as the coefficient of h^{2k} in this expansion.

Collecting terms that include h^{2J+2} we have

$$w_h^g_{2J+2}h^{2J+2}+G_{2J+2}(\cdot)h^{2J+2}+F'(w_h^{\varphi})w_h^e_{J+1}.$$

Because the function G_{2J+2} includes only the error vectors $w_h e_1, \ldots, w_h e_J$ as arguments of the various Fréchet derivatives it contains, we are able to define $w_h^B_{J+1}$ as

$$w_h^B_{J+1} = -w_h^g_{2J+2} - G_{2J+2}$$

Therefore, we can determine $w_h e_{J+1}$ as the solution to

$$F_h'(^\omega_h^\varphi) \omega_h^e_{J+1} = \omega_h^B_{J+1}.$$

Hence by induction Theorem 1 is valid. □

Thus the global discretization error will have an expansion of the form (39) provided $F_h(V)$ has M+l continuous Fréchet derivatives and we can solve equations of the form (41).

In the finite dimensional case, the Fréchet derivative of a discrete differential operator $F_h(V)$ is just the Jacobian matrix of the operator considered as a vector-valued

function of the R-dimensional vector $V = (V_1, ..., V_R)^T$. For our particular difference operator defined by (9) the jth row of the Jacobian matrix is obtained by computing the partials with respect to the components of V of equation (9) with i=j-1.

We will denote the j^{th} row of $F_h'(V)$ by $[F_h'(V)]_j$. Introducing the notation

$$f^{k}(v) = f(t_{k}, v_{k}, v_{k-R}, \frac{v_{k+1}^{-v_{k-1}}}{2h}, \frac{v_{k+1-R}^{-v_{k-1-R}}}{2h})$$

and using it for the partial derivatives of ~f(t,u,v,y,z) , we can write $~F_h^{\,\prime}(\omega_h^{\,}\phi)~$ as

$$[F'_{h}(\omega_{h}\phi)]_{1} = h^{-2}(1 + \frac{h}{2}f_{y}^{\circ}(\omega_{h}\phi), 0, ..., 0);$$

$$[F'_{h}(\omega_{h}\phi)]_{2} = h^{-2}(-2 + h^{2}f_{u}^{1}(\omega_{h}\phi), 1 + \frac{h}{2}f_{v}^{1}(\omega_{h}\phi), 0, ..., 0);$$

and for $3 \le j \le R$,

$$[F'_h(\omega_h^{\varphi})]_{j} = h^{-2}(0, ..., 0, 1 - \frac{h}{2}f_y^{j-1}(\omega_h^{\varphi}), -2 + h^2f_u^{j-1}(\omega_h^{\varphi}),$$

$$1 + \frac{h}{2}f_y^{j-1}(\omega_h^{\varphi}), 0, ..., 0)$$

with the non zero entries in the j^{th} row for $3 \le j \le R$ occurring in the (j-2)nd, (j-1)st and j^{th} positions.

Since $F_h'(\omega_h^{\phi})$ is lower triangular, it will be nonsingular provided none of the diagonal entries vanish. For any j the $(j,j)^{th}$ entry of $F_h'(\omega_h^{\phi})$ is

$$h^{-2} + \frac{1}{2h} f_{v}(\cdot)$$
.

This will be non zero provided

$$f_{V}(\cdot) \neq \frac{-2}{h}.$$

If we assume that $f_y(t,u,v,y,z)$ is bounded, then as $h \to 0$ the left hand side of (45) is bounded while the right hand side diverges to $-\infty$. Thus for sufficiently small h, (45) will be valid and $F_h'(\omega_h \phi)$ will be nonsingular. Consequently, the equation (41) will be uniquely solvable.

Now let us consider the system of equations

(46)
$$F_h'(\omega_h \varphi) \omega_h e_k = 0$$

in more detail. The jth equation in (46) is given by

$$h^{-2} (1 - \frac{h}{2} f_{y}^{j-1} (\omega_{h} \varphi)) e_{k} (t_{j-2}) + h^{-2} (-2 + h^{2} f_{u}^{j-1} (\omega_{h} \varphi)) e_{k} (t_{j-1})$$

$$+ h^{-2} (1 + \frac{h}{2} f_{y}^{j-1} (\omega_{h} \varphi)) e_{k} (t_{j}) = 0.$$

Equivalently, we can write this as

(47)
$$\frac{e_{\mathbf{k}}(t_{j-2}) - 2e_{\mathbf{k}}(t_{j-1}) + e_{\mathbf{k}}(t_{j})}{h^{2}} + f_{\mathbf{y}}^{j-1}(\omega_{\mathbf{h}}\varphi) \frac{e_{\mathbf{k}}(t_{j}) - e_{\mathbf{k}}(t_{j-2})}{2h} + f_{\mathbf{u}}^{j-1}(\omega_{\mathbf{h}}\varphi) e_{\mathbf{k}}(t_{j-1}) = 0.$$

This equation is the discrete analog of

(48)
$$e_{k}''(t) + f_{v}(\phi) e_{k}'(t) + f_{u}(\phi) e_{k}(t) = 0$$

at
$$t=t_{j-1}$$
, where $f_y(\phi) = f_y(t,\phi(t),\phi(t-r),\phi'(t),\phi'(t-r))$.

Equation (48) is the linear variational equation associated

with the continuous operator F(x) = 0 and is often denoted by $F'(\phi)e_k = 0$.

By examining the first and second equations in (46) we see that $e_k(t_{-1}) = e_k(t_0) = 0$. In fact, since $E(h) = X(h) - \omega_h \phi$ and we have taken $X_{-j} = \phi(t_{-j})$ for all $j=0,1,\ldots,R+1$ we must have $e_k(t_{-j}) = 0$ $\forall k$ and $\forall j=0,\ldots,R+1$.

One method for obtaining (48) is to let $h \to 0$ in (47). If we do this, the points t_{-j} become dense in the interval [(L-2)r,(L-1)r] and we therefore must have $e_k(t) \equiv 0$ on [(L-2)r,(L-1)r]. This implies that $e_k'(t) \equiv 0$ on [(L-2)r,(L-1)r]. Since the solution to (48) will have a continuous first derivative, we see that the discrete problem (47) with the initial conditions $e_k(t_{-1}) = e_k(t_0) = 0$ corresponds to the continuous problem (48) with the initial conditions $e_k(t_0) = e_k'(t_0) = 0$.

Consequently, the vectors $\ w_h e_k$ in (41) are actually discretizations of differentiable functions e_k (t) which satisfy

$$e_{k}''(t) + f_{y}(\phi) e_{k}'(t) + f_{u}(\phi) e_{k}(t) = B_{k}(\cdot)$$

$$(48')$$

$$e_{k}(t_{0}) = e_{k}'(t_{0}) = 0.$$

To summarize, we have established that the global discretization error has an asymptotic expansion in even powers of h:

$$X(h)(t_i) = \varphi(t_i) + \sum_{k=1}^{M} e_k(t_i) h^{2k} + O(h^{2M+2}).$$

Moreover, the functions $e_k(t)$ are independent of h and satisfy $e_k(t_0) = e_k'(t_0) = 0$.

In the next section we will discuss the implementation of this method and the modifications necessary to obtain a solution on [0,r].

Section 2.4. Implementing the Second Order Method

In order to obtain a numerical solution to

(8)
$$\ddot{x}(t)+f(t,x(t),x(t-r),\dot{x}(t),\dot{x}(t-r)) = 0, r>0$$

on [0,r], some modifications of the discrete method are needed. The first modification is based on the fact that on [0,r] we have more information available to us concerning the behavior of the theoretical solution than we do later.

To obtain the theoretical solution to (8) on [0,r] we start with a given initial function $\varphi(t)$ whose derivative $\varphi'(t)$ is known on [-r,0]. The extra information we have in this case is knowledge of $\varphi'(t)$. This is incorporated into the discrete version of (8) as added initial conditions.

Let R>0 be any given natural number and define h=r/R. Construct the uniform grid $\{t_i=0+ih\}_{i=-R}^R$. Since we are given $\Phi(t)$ and $\Phi'(t)$ for $t\in[-r,0]$ we may define a discrete version of (8) for $t\in[0,r]$ by

$$\frac{X_{i+1}^{-2X_{i}^{+}X_{i-1}}}{h^{2}} + f(t_{i}, X_{i}, \phi_{i-R}, \frac{X_{i+1}^{-}X_{i-1}}{2h}, \phi'_{i-R}) = 0, \quad i=0,1,\ldots,R-1;$$

$$\begin{aligned} \phi_{i-R} &= \phi(t_{i-R}) &\equiv \phi(t_{i}-r), \\ \phi_{i-R}' &= \phi'(t_{i-R}) &\equiv \phi'(t_{i}-r), & i=0,1,\dots,R-1; \\ X_{-1} &= \phi(t_{-1}) &\equiv \phi(-h), \\ X_{O} &= \phi(t_{O}) &\equiv \phi(0). \end{aligned}$$

The other modification that we must make is based on the fact that the theoretical solution to (8) will, in general, not have a continuous second derivative at the origin.

Since the expression $\frac{X_{i+1}-2X_i+X_{i-1}}{h^2}$ for i=0 in

(49) is an approximation to the second derivative of the solution at t=0 we must take steps to insure that the second derivative exists and is continuous at t=0. However, we want to accomplish this in such a manner that the initial information X_0 , X_{-1} , ϕ_{i-R} and ϕ'_{i-R} for $i=0,1,\ldots,R-1$ is not altered.

One way of accomplishing this is to modify the given initial function $\varphi(t)$. Select $\varepsilon < 0$ such that $-h < \varepsilon < 0$ and define a new initial function $\overline{\varphi}(t)$ by using an interpolating polynomial in place of $\varphi(t)$ on $[\varepsilon,0]$. More precisely, let H(t) be the fourth degree Hermite interpolating polynomial which satisfies $P(\varepsilon) = \varphi(\varepsilon)$, $P'(\varepsilon) = \varphi'(\varepsilon)$, $P(0) = \varphi(0)$, $P'(0) = \varphi'(0)$ and $P''(0) = f(0,\varphi(0),\varphi(-r))$, $\varphi'(0),\varphi'(-r)$. Define a new initial function $\overline{\varphi}(t)$ on [-r,0] by

(50)
$$\bar{\phi}(t) = \begin{cases} \phi(t) & -r \leq t \leq \varepsilon \\ H(t) & \varepsilon \leq t \leq 0. \end{cases}$$

By construction $\bar{\phi}(t)$ will be continuously differentiable on [-r,0] and since $\bar{\phi}''(0) = f(0,\phi(0),\phi(-r),\phi'(0),\phi'(-r))$, the second derivative of the solution obtained with the initial

function $\bar{\phi}$ will be continuous at t=0. Now $\bar{\phi}(t)$ agrees with $\phi(t)$ for $-r \le t \le \varepsilon$ and also for t=0. This includes all points in [-r,0] at which initial information is needed in (49). Since, as mentioned earlier, the solution to (8) depends continuously on the initial function, the use of $\bar{\phi}(t)$ in place of $\phi(t)$ will not radically alter the character of the solution to (8). It will however make (49) a reasonable discretization for (8).

Note that, as we move to the right solving (8) with the method of steps, we gain differentiability of the solution (see the discussion for first order equations in Section 1 of this chapter). Thus the solution to (8) will have a continuous second derivative for all t>0 and the approximation $\frac{X_{1}+1^{-2X_{1}+X_{1}-1}}{h^{2}} \text{ will not cause any difficulties at other multiples of r.}$

The proof that the discrete problem (49) is consistent is a special case of Section 2.1. The only difference is that now we are using the actual value of $\phi'(t-r)$ instead of an approximation. This will simplify the expression for the local truncation error but does not alter the fact that the truncation error has an asymptotic expansion in even powers of h. Also, note that the proof of stability with the added hypotheses (33), and the analysis of the global discretization error given for (9) carry over directly to (49).

Using the appropriate discrete scheme, (9) or (49), we can compute a solution to (8) whose global discretization error has an asymptotic expansion in even powers of h. Since the coefficient error functions in this expansion are independent of h, extrapolation can be performed to obtain a more accurate solution.

Let R be any natural number and let $h_0 = r/R$ be the basic steplength. For each $k=1,\ldots,M$ define $h_k = h_0/2^k = \frac{r}{2^k R}.$ Define grids G_k by

$$\mathbf{G_k} = \{\mathbf{t_i^k} : \mathbf{t_i^k} = \mathbf{0} + \mathbf{ih_k}, \mathbf{i} = -2^k \mathbf{R}, -2^k \mathbf{R} + 1, \dots, 0, 1, \dots, 2^k \mathbf{R}\} \text{ for } k = 0, 1, \dots, M.$$

As before, the grids G_k are nested; $G_k \subset G_{k+1}$.

Given an initial function $\varphi(t)$ such that $\varphi(t)$ and $\varphi'(t)$ are continuous on [-r,0], select $\varepsilon < 0$ such that $-h_{M} < \varepsilon < 0$. With this ε define a new initial function, again denoted by φ , as in (50).

On each grid G_k compute the solution $X(h_k)$ to (49). Denoting the solution to (8) on [0,r] also by $\phi(t)$ we know from our previous work that

(51)
$$[X(h_k)]_i = \varphi(t_i^k) + \sum_{j=1}^M h^{2j} e_j(t_i^k) + O(h_k^{2M+2}), \text{ for } i=1,\ldots,2^k R.$$

If R=1, then we can use pullback interpolation as described in Chapter 1 to obtain a solution, X(h), which satisfies

(52)
$$[X(h)]_i = \varphi(t_i) + O(h_0^{2M+2})$$

for each $t_i \in G_M$. Actually the pullback interpolation method does not involve as much work in this case as it does for first order initial value problems. This is so because from (48') we have $e_j^!(0) = 0$ for $j=1,\ldots,M$. Therefore we do not have to go through all the work of constructing a table analogous to Table 4 of Chapter I to obtain $e_j^!(0)$.

If $R \ge 2$, then we can obtain uniform accuracy at all grid points of G_M without using any information about $e_j^!(0)$. The reason for this is that with $R \ge 2$ we have enough data points at each stage to do Lagrange interpolation with accuracy comparable to that of the data points. That is, when performing Lagrange interpolation to accomplish the pullback with $R \ge 2$, it is the accuracy to which we know the error functions that dominates (see the discussion in Section 5 of Chapter 1). Thus we do not need to do Hermite interpolation in this case and therefore the information about $e_j^!(0)$ is superfluous.

Once the solution X(h), satisfying (52), is obtained we can use the discrete method (9) to obtain a numerical solution on [r,2r]. This process, theoretically at least, can be repeated indefinitely on intervals of length r, to obtain a solution which is $O(h^{2M+2})$.

It should be noted however, that when using (9) to solve (8) on intervals of the form [(L-1)r, Lr], $L \ge 1$, the initial data is taken from the computed solution on the previous interval. Since the initial data will not be precise (the error in the initial data is $O(h^{2M+2})$) an accumulation or errors is unavoidable. Of course, this also occurs when any numerical method is used on successive intervals.

A variety of other schemes, based on (9), (48) and pullback interpolation, to solve (8) are possible. We could, for instance, use (49) and pullback interpolation on $[0,h_0]$ to obtain a solution. Then by modifying the initial concitions in (49) we could repeat this procedure on $[h_0, 2h_0]$. Repeating this R times we would have the solution on all of [0,r]. We could then follow the same format, using (9) in place of (49), to obtain a solution on [r,2r], etc.

We do not make a comparison of these schemes here.

In the next section we will merely exhibit some numerical examples computed using the first described solution technique.

Section 2.5. Numerical Results for Second Order Equations

This section consists of two examples. The first of these is

Example 3: Solve $\ddot{x}(t) = x(t-1)$, $0 \le t \le 1$ numerically for the initial function $\phi(t) = 12(t+1)^2 + 24$ on [-1,0]. The method of solution is the discretization (49), extrapolation and pullback interpolation with four grids. The steplengths employed are $h_k = 1/2^k$ for $k=0,\ldots,3$ and the grids are

$$G_k = \{t_i^k = 0 + ih_k : i=0,1,...,2^k\}$$

for k=0,...,3.

Note that the initial function $\varphi(t)$ satisfies $\varphi(-1) = \varphi''(0-)$, so the second derivative of the solution is continuous at the origin. Thus, the initial function does not need to be modified as in (50) for this example.

The theoretical solution is given by

$$\varphi(t) = t^4 + 12t^2 + 24t + 36$$

for $t \in [0,1]$. The numerical results presented in Table 15 are for the nine equally spaced points in [0,1] which belong to the finest grid G_3 . The error reported is once again the computed solution minus the theoretical solution.

TABLE 15

i	Error
0	0.0
1	0.0×10^{-12}
2	23×10^{-12}
3	23×10^{-12}
4	23×10^{-12}
5	23×10^{-12}
6	23×10^{-12}
7	0.0×10^{-12}
8	0.0×10^{-12}

The high accuracy of the numerical solution in this example is due to the fact that the theoretical solution is a low degree polynomial. Thus, extrapolation and consequently pullback interpolation are quite accurate in this case.

This accuracy is not to be expected for more general equations as the following example illustrates.

Example 4: Solve $\ddot{x}(t) + \frac{1}{2}x(t) - \frac{1}{2}x(t-\pi) = 0$ on $[0,\pi]$ for the initial function $\varphi(t) = 1$ -sin t, $-\pi \le t \le 0$. This example is from Norkin [23] and the theoretical solution is also given by $\varphi(t) = 1$ -sin t, $0 \le t \le \pi$. Thus, the second derivative of the solution is continuous at the origin and no modification of the initial function is required.

The method of solution is to use (49), extrapolation and pullback interpolation with four grids to obtain a numerical solution on $[0,\pi/4]$. Using this numerical solution as initial data, the discretization (9), extrapolation and pullback interpolation with four grids is employed to obtain a solution on $[\frac{\pi}{4},\frac{\pi}{2}]$. This process is repeated twice more to obtain numerical solutions on $[\frac{\pi}{2},\frac{3\pi}{4}]$ and $[\frac{3\pi}{4},\pi]$. On each interval of length $\pi/4$ the stepsizes used are $h_k = \frac{\pi}{4 \cdot 2^k}$ for k=0,...,3. Thus, on each of the four subintervals the finest grid contains nine equally spaced points which we number as t_0,\dots,t_8 .

The accuracies of the numerical results are reported as Table 16, with the error being computed in the usual manner.

TABLE 16

i	Error on $[0,\pi/4]$	Error on $[\pi/4,\pi/2]$
О	0.0	$.5 \times 10^{-8}$
1	15.9×10^{-8}	-1.0×10^{-6}
2	37.6×10^{-8}	-4.5×10^{-6}
3	30.2×10^{-8}	-6.4×10^{-6}
4	-1.6×10^{-8}	-6.6×10^{-6}
5	-33.2×10^{-8}	-6.7×10^{-6}
6	-40.0 x 10 ⁻⁸	-8.4×10^{-6}
7	-17.1×10^{-8}	-11.5×10^{-6}
8	$.5 \times 10^{-8}$	-11.9×10^{-6}
i	Error on $[\pi/2.3\pi/4]$	Error on $[3\pi/4,\pi]$
i O	Error on $[\pi/2, 3\pi/4]$ -11.9 x 10 ⁻⁶	Error on $[3\pi/4,\pi]$ -27.2 x 10 ⁻⁶
	_	
0	-11.9 x 10 ⁻⁶	-27.2×10^{-6}
o 1	-11.9×10^{-6} -13.7×10^{-6}	-27.2×10^{-6} -28.6×10^{-6}
0 1 2	-11.9×10^{-6} -13.7×10^{-6} -18.7×10^{-6}	-27.2×10^{-6} -28.6×10^{-6} -31.9×10^{-6}
O 1 2 3	-11.9×10^{-6} -13.7×10^{-6} -18.7×10^{-6} -21.2×10^{-6}	-27.2×10^{-6} -28.6×10^{-6} -31.9×10^{-6} -33.2×10^{-6}
O 1 2 3 4	-11.9×10^{-6} -13.7×10^{-6} -18.7×10^{-6} -21.2×10^{-6} -21.0×10^{-6}	-27.2×10^{-6} -28.6×10^{-6} -31.9×10^{-6} -33.2×10^{-6} -32.4×10^{-6}
O 1 2 3 4 5	-11.9×10^{-6} -13.7×10^{-6} -18.7×10^{-6} -21.2×10^{-6} -21.0×10^{-6} -20.6×10^{-6}	-27.2×10^{-6} -28.6×10^{-6} -31.9×10^{-6} -33.2×10^{-6} -32.4×10^{-6} -31.4×10^{-6}

The accuracy of the computed solution is considerably higher on the first subinterval than on the last three subintervals. Also the errors are steadily increasing at a very low rate over the last three subintervals. This suggests an accumulation of truncation errors.

This accumulation is to be expected for the following reason. For a linear equation, the discretizations (9) or (49) are similar to a well known linear multistep method for solving second order initial value problems (see Lambert [18]). The only difference is the manner in which the discretizations are initialized. The stability theory for linear multistep methods for solving second order equations predicts an accumulation of truncation errors. Thus, based on comparisons with similar methods, our numerical results are consistent with how our method should behave.

Note that the equation $\ddot{x}(t) + \frac{1}{2}x(t) - \frac{1}{2}x(t-\pi) = 0$ is not covered by our proof of stability. Yet the numerical results seem to indicate that the method is stable when applied to this equation. Consequently, we again reiterate our belief that stability is valid for a much larger class of functions than our proof indicates.

BIBLIOGRAPHY

BIBLIOGRAPHY

- 1. Aitken, A.C., On interpolation by iteration of proportional parts, Proc. Edinburgh Math. Soc., 2, 56-76 (1932).
- 2. Björck, A., and V. Pereyra, Solutions of Vandermonde systems of equations, Math. Comp., 24, 893-903 (1970).
- 3. Cooke, K.L., and S.E. List, The numerical solution of integro-differential equations with retardation, Department of Electrical Engineering, University of Southern California, Los Angeles, Technical Report No. 72-4 (1972).
- 4. Cooke, K.L. and J.A. Yorke, Equations modelling population growth, economic growth and gonorrhea epidemiology, preprint.
- 5. Corliss, G.F., <u>Parallel Rootfinding Algorithms</u>, Ph.D. thesis, Michigan State University, East Lansing, Michigan (1974).
- 6. Cullen, C.G., <u>Matrices and Linear Transformations</u>, Addison-Wesley Publishing Company, Reading, Mass. (1967).
- 7. Cunningham, W.J., A non-linear differential-difference equation of growth, Proc. Natl. Acad. Sci., U.S., 40, 709-713 (1954).
- 8. Dahlquist, G., Convergence and stability in the numerical integration of ordinary differential equations, Kungl. Tekniska Hogskolans Handlingar No. 130 (1959).
- 9. Dahlquist, G., A special stability problem for linear multistep methods, BIT, 3, 27-43 (1963).
- 10. Feldstein, A., <u>Discretization Methods for Retarded Ordinary</u>
 <u>Differential Equations</u>, Ph.D. thesis, University of
 California, Los Angeles (1964).
- 11. Gragg, W.B., Repeated Extrapolation to the Limit in the Numerical Solution of Ordinary Differential Equations, Ph.D. thesis, University of California, Los Angeles (1964).

- 12. Gragg, W.B., On extrapolation algorithms for ordinary initial value problems, SIAM J. Numer. Anal., 2, 384-403 (1965).
- 13. Hale, J., <u>Functional Differential Equations</u>, Springer-Verlag, New York, N.Y. (1971).
- 14. Huygens, C., <u>De Circuli Magnitudine Inventa</u>, Leiden (1654).
- 15. Isaacson, E., and H.B. Keller, <u>Analysis of Numerical Methods</u>, John Wiley and Sons, Inc., New York, N.Y (1966).
- 16. Joyce, D.C., Survey of extrapolation processes in numerical analysis, SIAM Review, 13, 435-490 (1971).
- 17. Keller, H.B., <u>Numerical Methods for Two-Point Boundary-Value Problems</u>, Blaisdell Publishing Company, Waltham, Massachusetts (1968).
- 18. Lambert, J.D., Computational Methods in Ordinary Differential Equations, John Wiley and Sons, London (1973).
- 19. Lindberg, G., A simple interpolation algorithm for improvement of the numerical solution of a differential equation, SIAM J. Numer. Anal., 9, 662-668 (1972).
- 20. Milne, W.E., and R.R. Reynolds, Stability of a numerical solution of differential equations, J. Assoc. Comput. Mach., 6, 193-203 (1959); Part II, J. Assoc. Comput. Mach., 7, 46-56 (1960).
- 21. Milne, W.E., and R.R. Reynolds, Fifth-order methods for the numerical solution of ordinary differential equations, J.Assoc. Comput. Mach., 9, 64-70 (1962).
- 22. Neville, E.H., Iterative interpolation, J. Ind. Math. Soc., 20, 87-120 (1934).
- 23. Norkin, S.B., <u>Differential Equations of the Second Order</u>
 with Retarded Argument, American Mathematical Society,
 Providence, Rhode Island (1972).
- 24. Pareyra, V., On improving an approximate solution of a functional equation by deferred corrections, Numer. Math., 8, 376-391 (1966).

- 25. Pereyra, V., Iterated deferred corrections for nonlinear operator equations, Numer. Math., 10, 316-323 (1967).
- 26. Pereyra, V., Accelerating the convergence of discretization algorithms, SIAM J. Numer. Anal., 4, 508-533 (1967).
- 27. Pereyra, V., Iterated deferred corrections for nonlinear boundary value problems, Numer. Math., 11, 111-125 (1968).
- 28. Pereyra, V., High order finite difference solution of differential equations, Seminar Notes, Stanford University, Palo Alto, California (1973).
- 29. Richardson, L.F., The deferred approach to the limit, I-single lattice, Trans. Roy. Soc. London, 226, 299-349 (1927).
- 30. Sansone, G., Teorema di esistenza di soluzioni per un sistema di equazioni funzionali differenziali, Ann. Mat. Pura Appl., 39, 65-67 (1955).
- 31. Stetter, H.J., Asymptotic expansions for the error of discretization algorithms for non-linear functional equations, Numer. Math., 7, 18-31 (1965).
- 32. Varga, R.S., <u>Matrix Iterative Analysis</u>, Prentice-Hall, Inc., Englewood Cliffs, New Jersey (1962).
- 33. Widder, D.V., Advanced Calculus, Prentice-Hall, Inc., Englewood Cliffs, New Jersey (1961).
- 34. Wright, E.M., A non-linear difference-differential equation, J. reine angew. Math., 194, 66-87 (1955).

REPORT DOCUMENTATION PAGE		READ INSTRUCTIONS BEFORE COMPLETING FORM
1. REPORT NUMBER	2. GOVT ACCESSION NO.	
4. TITLE (and Sublitle) UNIFORMLY ACCURATE NUMERICAL SOLUTIONS TO DIFFERENTIAL EQUATIONS USING EXTRAPOLATION AND INTERPOLATION.		5. TYPE OF REPORT & PERIOD COVERED
		6. PERFORMING ORG. REPORT NUMBER
7. Author(s)Richard Allan Rogers		8. CONTRACT OR GRANT NUMBER(s)
		AFOSR-72-2271
9. PERFORMING ORGANIZATION NAME AND ADDRESS Michigan State University Department of Mathematics East Lansing, Michigan 4882	24	10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS
11. CONTROLLING OFFICE NAME AND ADDRESS	F: - D 1	12. REPORT DATE
Air Force Office of Scienting 1400 Wilson Blvd.	August 1974	
Arlington, Virginia 22209		154
14. MONITORING AGENCY NAME & ADDRESS(if differen	t from Controlling Office)	15. SECURITY CLASS. (of this report)
		UNCLASSIFIED
	1	15a. DECLASSIFICATION/DOWNGRADING SCHEDULE
17. DISTRIBUTION STATEMENT (of the abstract entered	In Block 20, Il different from	m Report
18. SUPPLEMENTARY NOTES		
19. KEY WORDS (Continue on reverse side if necessary an	d identify by block number)	
20. ABSTRACT (Continue on reverse side if necessary and In this work we are co	ncerned with	numerical methods for We consider those

"pullback interpolation method". This method combines extrapolation with Hermite interpolation of the coefficient functions for the asymptotic error expansion to produce a highly accurate solution at all grid points of the finest mesh. When q is 1 or 2 pullback interpolation yields uniform accuracy at all grid points of the finest mesh.

In Chapter II the pullback interpolation method is modified so as to be applicable to boundary value problems. In addition, an elementary proof of the stability of V. Pereyra's finite difference scheme for solving two point boundary value problems is given.

In Chapter III we consider difference differential equations with constant retardation. The methods of Chapter I are shown to be applicable to first order delay equations. Because of the presence of the delay term, the uniform accuracy obtained through pullback interpolation is indispensible for these problems.

A finite difference scheme for directly solving second order delay equations is constructed and analyzed in Chapter III. The global discretization error is shown to have an asymptotic error expansion in even powers of the steplength.

