

National Library of Canada Bibliothèque nationale de Canada

Canadian Theses Service

Services des thèses canadiennes

Öttawa, Canada K1A DN4

## CANADIAN THESES

# THÈSES CANADIENNES

## NOTICE

The quality of this microfiche is heavily dependent upon the quality of the original thesis submitted for microfilming. Every effort has been made to ensure the highest quality of reproduction possible.

If pages are missing, contact the university which granted the degree.

Some pages may have indistinct print especially if the original pages were typed with a poor typewriter ribbon or if the university sent us an inferior photocopy.

Previously copyrighted materials (journal articles, published tests, etc.) are not filmed.

Reproduction in full or in part of this film is governed by the Canadian Copyright Act, R.S.C. 1970, c. C-30. Please read the authorization forms which accompany this thesis.

#### **AVIS**

La qualité de cette microfiche dépend grandement de la qualité de la thèse soumise au microfilmage. Nous avons tout fait pour assurer une qualité supérieure de reproduction.

S'il manque des pages, veuillez communiquer avec l'université qui a conféré le grade.

La qualité d'impression de certaines pages peut laisser à désirer, surtout si les pages originales ont été dactylographiées à l'aide d'un ruban usé ou si l'université nous a fait parvenir une photocopie de qualité inférieure.

Les documents qui font déjà l'objet d'un droit d'auteur (articles de revue, examens publiés, etc.) ne sont pas microfilmés.

La reproduction, même partielle, de ce microfilm est soumisé à la Loi canadienne sur le droit d'auteur, SRC 1970, c. C-30. Veuillez prendre connaissance des formules d'autorisation qui accompagnent cette thèse.

THIS DISSERTATION
HAS BEEN MICROFILMED
EXACTLY AS RECEIVED

LA THÈSE A ÉTÉ MICROFILMÉE TELLE QUE NOUS L'AVONS REÇUE



3/
•

National Library of Canada

Bibliothique nationale

CANADIAN INDIA ON INCROFICAL minus CABADAMAS Int incontent

NAME OF AUTHOR/NOW DE L'AUTEUR YULY - Shion Kuo
TITLE OF THESIS/TITHE DE LA THÈSE Equivalences Between Numerical Methods
for solving Differential Equations
F. C. J.
DEGREE FOR WHICH THESIS WAS PRESENTED!  GRADE POUR LEQUEL CETTE THÈSE PUT PRÉSENTÉE MOSTEY OF SCIENCE
YEAR THIS DEGREE CONFERNED/ANNÉE D'OUTENTION DE CE DEGRÉ 1983
NAME OF SUPERVISOR/MON DU DINECTEUR DE THÈSE R.D. RUSSELL
······································

Permission is hereby granted to the NATIONAL LIBRARY OF CANADA to microfilm this thesis and to lend or sell copies of the film.

The author reserves other publication rights, and neither the thesis ner extensive extracts from it may be printed or otherwise reproduced without the author's written permission.

L'enterisation est, par la précente, escerdée à la BALJOTHÉ

OUE NATIONALE DU CAMADA de miorefilmer estte thèse e
de prêter ou de vendre des exemplaires du film.

L'autour se récerve les autres draite de publication; el la thèce ni de lange entroite de colle-ci ne delvent être lagring en autrement maradules acon l'automation dealle de l'automatic

DATED/DATE Jud 1 83 SIGNED/SHORE

PERMANENT ADDRESS/MESIDENCE FIXE

# EQUIVALENCES BETWEEN NUMERICAL HETHODS FOR SOLVING DIFFERENTIAL EQUATIONS

by.

Yuh-Shiow Kuo

B.Sc. Fu-Jen Catholic University, 1978

#### THESIS SUBMITTED IN PARTIAL PULPILLMENT OF

THE REQUIREMENTS FOR THE DEGREE OF

HASTER OF SCIENCE in the Department

of

Mathematics

(C) Yuh-Shiow Kuo 1983

SIMON PRASER UNIVERSITY

July 1983

All rights reserved. This work may not be reproduced in whole or in part, by photocopy or other means, without permission of the author.

#### APPROVAL

Name: Yuh-Shiow Kuo

Degree: MASTER OF SCIENCE

Title of thesis: EQUIVALENCES BETWEEN NUMERICAL BETHODS FOR

SCLVING DIFFERENTIAL EQUATIONS

Examining Committee:

Chairman: B.S. Thomson

R.D. Russell Senior Supervisor

R.W. Lardner

E.M. Shoemaker

J. Paine
External Examiner
Visiting Assistant Professor
Department of Mathematics
Simon Fraser University

Date Approved: August 3, 1983

# PARTIAL COPYRIGHT LICENSE

I hereby grant to Simon Fraser University the right to lend my thesis or dissertation (the title of which is shown below) to users of the Simon Fraser University Library, and to make partial or single copies only for such users or in response to a request from the library of any other university, or other educational institution, on its own behalf or for one of its users. I further agree that permission for multiple copying of this thesis for scholarly purposes may be granted by me or the Dean of Graduate Studies. It is understood that copying or publication of this thesis for financial gain shall not be allowed without my written permission.

Title of Thesis/Dissertation:

Equival	inces Cet	ween Nameric	al Methedo
- C for	Solving	Pifferential	al Methodo
-	J	7	0
Z.			
2			

Author:

(signature)

Yuh-Shion Kuo.

(name)

(date)

#### ABSTRACT

A brief description of various types of boundary value problems and initial value problems is given. Existence and uniqueness theorems and stability properties of them are considered. One of our main purposes is to survey the well-known numerical methods such as initial value approaches, finite difference methods, and finite element methods for solving the problems. Them, various equivalences between the methods are shown. Lastly, the high order finite difference methods considered by Doedel and by Lynch-Rice are discussed from a different point of view. A relationship between the high order finite difference methods and collocation methods is presented.

Comparison of operation counts and numerical results for Doedel's methods, Lynch-Rice's methods, and collocation methods using B-splines and Gauss points is treated.

#### ACEBOULEBGEBERTS

I am grateful to Dr. R. D. Russell for accepting me as his student and for suggesting the topics of this thesis. Because of his kind and patient guidance, I could finish my thesis and take care of my little one at the same time.

Thanks to Dr. John Paine for his kindly help in part of the thesis.

I also want to say thanks to the Mathematics Department S.F.U. for offering enough financial support and to Mrs. Kathy Hannes and Mrs. Sylvia Holmes for their kindness and help in arranging Taships without which my studies would not have been possible.

# TABLE OF CONTENTS

Approval	•••••••••••••••••••••••••••••••••••••••	ii
Abstract	*******************	iii
acknowledgements	••••••	iv
List of Tables	****************	<b>v</b> ii
1. Introduction	***********	1
2. IVPs and BVPs		5
	s for ODEs	
_	••••••	
	***************************************	
	iqueness Theory	
	- · ·	
	<b>y</b>	
;	Solving BVPs	
	roaches	
3.1.1. Superpositi	ion	19
	oting	
3.1.3. Multiple Sh	booting	25
3.1.4. Stabilized	Harch	31
3.1.5. Invariant I	Imbedding	34
3.2. Pinite Difference	e Methods	37
. 3.2.1. A Simple So	cheme for a Second Order Problem	38
3.2.2. One Step Sc	chemes :Trapezoidal Rule and	, 4
Midpoint Rule	***************	39
3.2.3. Runge-Kutta	Schemes	44
3.3. Pinite Element Me	ethods	45
	Het hods	

3.3.2. Galerkin Methods48
3.3.3. Least Squares Hethod50
3.3.4. Hitz Method
4. Equivalences Between These Hethods54
4.1. Collocation and Finite Difference Bethods54
4.1.1. Collocation, Trapezoidal, and Hidroint Rules54
4.1.2. Collocation and Trapezoidal Rule56
4.2. Simple Ritz and Finite Difference
4.3. Collocation and Implicit Runge-Kutta59
4.4. Collocation and Multiple Shooting
4.5. Multiple Shooting and the Box Scheme
4.6. Invariant Imbedding and Hultiple Shooting68
5. Pinite Differences for Solving High Order Differential
5.1. Construction of the High Order Pinite Difference Approximation
5.1.1. The Approximations for Interior Subintervals77
5.1.2. The Approximation of Initial Conditions and Boundary Conditions
5.2 The Order of Consistency91
5.3. Stability of the Schemes94
5.4. Improved Order with Particular Choice of z101
5.5. An Equivalence Between Finite Difference and Collocation Methods
5.6. Work Estimates110
5.7. Experimental Results
6. Conclusion
5 8 9 7 1 9 4 C P C

# LIST OF TABLES

TABLE		•			•		,	PAGE
5.1	Operation	Counts	••••	• • • • • • •	• • • • •	• • • • • • • • •	• • • • • • • • • •	117
5. 2	Mumerical	results	for	Example	5.9.	• • • • • • • •	• • • • • • • •	120
5.3	Numerical	Results	for	Example	5.10	******	• • • • • • • •	121
5. 4	Numerical	Results	for	Brample	5.11		• • • • • • • • •	122

## 1. Introduction

For solving boundary value problems (BVPs) for ordinary differential equations (ODEs), the most common numerical methods are: initial value approaches; finite difference methods; and finite element methods. One of the areas of considerable interest to numerical analysts is the relationships between these methods. Once some equivalence has been found, properties of a method (e.g. the property of being well-conditioned for a multiple shooting method) can thus be shown to apply to the equivalent ones.

The main purposes of this thesis are to show equivalences between some well-known numerical methods and to give a clear view of the high order finite difference methods of Doedel [9] and Lynch-Rice [16].

In Chapter 2, general forms for ODEs are given to help maintain basic understanding. Before introducing numerical methods, existence and uniqueness theorems for the solutions of general initial value problems (IVPs) and linear BVPs are mentioned. Since it is extremely difficult to establish existence and uniqueness theorems for general BVPs, there are only some for restricted BVPs in special cases. One of them that relates to second order BVPs is mentioned in Section 2.2. Basic stability properties of Ortega [18] for linear IVPs and BVPs are provided. It is shown in Lentini-Osborne-Russell [15] that the

well-conditioning of a BVP is related to two bounding quantities: one involving the boundary conditions and the other involving the Green's function. Their result is described also. Having well-posed BVPs in hand, we consider numerical methods for solving them. In Chapter 3, to prepare the groundwork needed in the subsequent Chapters, the well-known methods are introduced, viz superposition, simple shooting, multiple shooting, and invariant imbedding for initial value approaches; trapezoidal rule, midpoint rule, and Runge-Kutta schemes for finite difference methods; collocation, Galerkin, least square, and Ritz methods for finite element methods. When presenting the methods, some equivalences between them are easily seen as they are introduced, e.g. equivalence between discrete Galerkin and collocation and equivalence between discrete least square and collocation.

Since finite difference methods involve unknowns which correspond directly to approximate solutions at mesh points, in Chapter 4 equivalences between finite difference methods and some other methods are emphasized. They include the trapezoidal rule and collocation; the midpoint rule and collocation; Runge-Kutta methods and collocation; the Box scheme and multiple shooting; multiple shooting and collocation; multiple shooting and invariant imbedding; and the simple Ritz and finite difference methods. While all of them are probably known, some are not in the literature. Although in theory, two methods are shown to be equivalent mathematically, i.e. they have the same

solution set, in practice their properties (e.g. order of accuracy) can be different. This difference is particularly important from a computational point of view.

An investigation was made for high order finite difference methods considered by Doedel [9] and by Lynch-Rice [16]. The main difference between their methods is that Doedel's methods include noncompact approximations while Lynch-Rice do not. In Chapter 5/ the construction of high order finite difference methods derived by Doedel is presented. Order of consistency and stability of the schemes are also discussed. For a general n-th order linear differential equation, choice of a set of auxiliary points which gives one higher order of accuracy is given. For the special n-th order differential equation H=D", formulae are provided to evaluate the uniquely determined right-hand-side coefficients of the approximations, and to find the location of the special auxiliary points which give the order of accuracy as high as possible. An obvious equivalence between these methods and collocation methods is given in Section 5.5. This equivalence appears to have not been previously observed and not shown in either Doedel [9] or in Lynch-Rice [16].

Work estimates for the finite difference methods for general n-th order differential equations are provided.

Comparison of computational work of Doedel's schemes and of Lynch-Rice's schemes is offered for general second order differential equations. Lynch-Rice [16] also compared their methods with five other methods, but the comparison here is

based on the Langrange polynomial interpolation basis functions rather than the set of basis functions they considered. From the operation counts, we conclude that Doedel's methods are more efficient than Lynch-Rice's schemes for same orders of accuracy if ore takes as few auxiliary points as possible. In Section 5.8, numerical examples show that Doedel's schemes with one auxiliary point are competive with Lynch-Rice's schemes. Hence, methods using one auxiliary point are the most efficient of the high order finite difference methods. Comparing some numerical results of the finite difference methods with that of collocation methods using B-splines and Gauss points, one finds that the finite difference methods can require a large number of subintervals to achieve any significant accuracy.

The last Chapter of the thesis is Chapter 6, the conclusion.

#### 2. IVPs and BVPs

We begin with a brief account of some of the basic prerequisites: general IVPs and BVPs, existence and uniqueness theorems, and problem stability.

# 2.1. Standard Problems for ODEs

In this section, some standard problems which correspond to the most common forms of BVPs are presented. Since ways to numerically solve BVPs can be closely connected with methods for solving IVPs, and since the theory of IVPs is closely related to that for BVPs, a treatment of IVPs is considered also.

#### 2.1.1. IVPs

The general IVP can be written as a first order system

(2.1a) 
$$y'(t) = f(t, y(t))$$
 t>a

$$(2.1b) y(a) = \alpha$$

where  $y(t) = (y(t), y(t), \dots, y(t))$  is the unknown

function, 
$$f(t,y) = (f(t,y), f(t,y), \dots, f(t,y))$$
 is the

nonlinear right hand side and  $\phi$  is a known n-vector of initial conditions which completely determes y(t).

A high order ODE

(m) 
$$(t) = f(t, y, y^{*}, ..., y^{*})$$

can be converted to the first order form (2.1a), by letting

The ODE has the equivalent form

Similarly, a system of high order ODEs can be reduced to a set of first order equations in this way.

In the simpler case where the IVP is linear, (2.1) is simplified to

(2.2a) 
$$y'(t) = \lambda(t)y(t) + f(t)$$
 t>a

$$(2.2b) y(a) = \emptyset$$

where A(t) is an n x n matrix and f(t) is an n-vector valued

function of t.

The linear system (2.2a) is called homogeneous if f(t)=0, and inhomogeneous otherwise.

#### 2.1.2. BVPs

Unlike IVPs, solutions to BVPs are not completely determined by initial information; the information is given at two or more points which normally correspond to the boundary of some physical region of interest. One basic form is the linear two point BVP

(2.3a) 
$$y'(t)=\lambda(t)y(t)+f(t)$$
 astsb

where y(t) and f(t) are n-vectors, h(t) is an n x n matrix, a and b are finite or infinite and g(t) is a constant n-vector,  $B_2$  and  $B_b$  are n x n matrices corresponding to n boundary conditions.

For (2.3) to have a unique solution, it is necessary but not sufficient that these boundary conditions be linearly independent, i.e. that the matrix  $(B_a \mid B_b)$  have n linearly independent columns or simply rank  $(B_a \mid B_b) = n$ .

BC of the general form (2.3b) are called nonseparated BC since each involves information about y(t) at both end points. If rank  $(B_d)$  <n or rank  $(B_b)$  <n, then the BC are called partially separated. The BC are called generated if they can be simplified

to

$$C_{\mathbf{a}} \underbrace{\mathbf{y}(\mathbf{a})}_{\mathbf{1}} = \mathbf{x}$$

$$C_{\mathbf{b}} \underbrace{\mathbf{y}(\mathbf{b})}_{\mathbf{0}} = \mathbf{x}$$

where  $C_{a}$  is a p x n matrix ,  $C_{b}$  is an (n-p)x n matrix , and  $p=rank(B_{a})$  .

A general linear multipoint BVP consists of the ODE (2.3a) and multipoint BC

(2.3c) 
$$\sum_{j=1}^{r} B y(j) = \alpha$$

where  $B_1, \ldots, B_r$  are n x n matrices,  $\alpha$  is an n-vector and  $a=f, < f_1 < \ldots < f_r = b$ .

A nonlinear two point BVP can normally be expressed in the form

$$(2.4a) \qquad y^{\perp}(t) = f(t, y(t)) \qquad a \le t \le b$$

(2.4b) 
$$g(y(a),y(b))=0$$

where g = (g, ..., g) and 0 is the zero n-vector.

A nonlinear ath order (scalar) BVP normally has the form

(a) 
$$(m-1)$$
  
(2.5a) y (t) = f(t,y(t),y'(t),...,y (t)) a \le t \le b

(2.5b) 
$$g(y(a),...,y(b),...,y(b)) = 0$$

where (2.5h) involves m-vectors g and 0 corresponding to the BCs.

In the linear case, (2.5a,b) simplifies to

(a) 
$$m-1$$
 (j)  
(2.6a)  $y$  (t) =  $a$  (t)  $y$  (t) +  $f$  (t)  $a \le t \le b$   
 $j=1$  j

As in the IVP case,  $\{2.5a,b\}$  and  $\{2.6a,b\}$  can be converted to the first order systems  $\{2.4a,b\}$  and  $\{2.3a,b\}$ , respectively, where the unknown solution is  $y(t) = (y,y^1,...,$ 

The most general BVP we consider involves a system of ODEs of different orders with multipointaBC which is called a mixed order system and can be written as

(a) 
$$(m-1)$$
  $(m-1)$   $(m-1)$ 

where 
$$y(t) = (y(t), ..., y(t))$$

# 2.2. Existence and Uniqueness Theory

Before introducing numerical methods for solving the above problems, existence and uniqueness theorems of the solutions of the problems are given in the following.

The 2.1: Let f(t,y) be continuous on  $D=\{(t,y): a \le t \le b, \|y-a\| \le R\}$ , where  $\|\cdot\|$  is some vector norm, and satisfy a Lipschitz condition with respect to y on D, that is, there is a constant K, such that for any (t,y) and (t,z) in D

$$\left\|\underbrace{f(t,y)-f(t,z)}\right\| \leq K\left\|\underbrace{y-z}\right\|.$$

If  $\|f(t,y)\| \le H$  on D and c=min{b-a,R/H}, then (2.1) has a unique solution for a \( \le t \le a + c \) (see Keller [12]).

Unfortunately, since it is extremely difficult to provide a result like the above for general BVPs, there are only existence and uniqueness theorems restricted for BVPs in special cases.

Many of them relate to an important class, the second order BVP

(2.8a) 
$$y''(t) = f(t,y,y')$$
  $a \le t \le b$ 

(2.8b) 
$$y(a) = \alpha$$
,  $y(b) = \alpha$ 

Por instance: (Bailey-Shampine-Waltman [5] )

Thu. 2.2: Suppose that  $f(t,y,y^*)$  is continuous on D=[a,b]  $x (-\infty,\infty)$  and satisfies there a Lipschitz condition, i.e. there exist constants L and K such that for every  $(t,y,y^*)$  and  $(t,z,z^*)$  in D

$$\left\|f(t,y,y')-f(t,z,z')\right\| \leq K \left\|y-z\right\| + L \left\|y'-z'\right\|.$$
If b-a<2d(L,K) where

then (2.8) has a unique solution. This result is the best possible.

While for general linear BVPs, when the problem is expressed in terms of an associated IVP, a general theorem is possible:

Thu. 2.3 Assume A(t) and f(t) are continuous in (2.3a). The BVP (2.3a,c) has a unique solution if and only if the matrix

(2.9) 
$$0:=B + \sum_{j=2}^{r} B Y (j)$$

is nonsingular, in which case the solution is

$$y(t)=Y(t)Q = (x-\sum_{j=2}^{r} BY(j)$$

$$= x-\{u\}f(u)du + \hat{y}(t)$$

where  $\hat{y}(t) = Y(t)$   $\begin{cases} t & -1 \\ Y(u) f(u) du, \text{ and } Y(t) \text{ is the fundamental} \end{cases}$  solution of (2.3a) which satisfies Y'(t) = A(t)Y(t) astished and Y(a) = I (see Keller [12]).

<u>Proof</u>: Let Y(t) be the fundamental solution that satisfies the above conditions. By direct substitution, it can be shown that

$$\frac{y(t)=Y(t)\left[\int_{a}^{b}+\int_{a}^{t}\frac{-1}{(u)f(u)du}\right]=Y(t)\int_{a}^{b}+y(t)}{a}$$

is the unique solution of (2.3). To satisfy the BC (2.3c), must be chosen such that

$$Q \bigwedge_{i} + \sum_{j=2}^{r} B Y {\binom{5}{j}} \int_{a}^{7_{j}} -1 Y (u) f(u) du = \infty$$

Hence (2.3a,c) has a unique solution if and only if Q is nonsingular.

# 2.3. Problem Stability

In this section, we discuss the stability properties of IVPs and BVPs. In general, when computing a quantity y from data t by some numerical method H, a problem is called unstable or ill-conditioned if "small" changes in the data t produce "large" changes in the solution y even if the method H is executed with no rounding error. The method H is called numerically unstable if small rounding errors introduced when using H produces large errors in the solution y even when the data are exact. As a rule, one should not try to compute numerically unstable quantities, and one should not use numerically unstable methods (Pranklin [10]).

Definition: Consider the initial value problem

(2. 10a) 
$$y'(t) = f(t,y)$$
 t>a  
(2. 10b)  $y(a) = a$ 

where 
$$y(t) = (y_1(t), \dots, y_n(t))$$

A solution y(t) is called <u>stable</u> (with respect to change in the initial conditions y(a)) if given any  $\xi > 0$ , there is a  $\delta > 0$  so that any other solution y(t) of (2.10a) for which

$$\|\mathbf{y}(\mathbf{a}) - \widehat{\mathbf{y}}(\mathbf{a})\| \le \delta$$

satisties

(2.11) ||y(t)-y(t)|| ≤ & for all toa.

The solution y(t) is asymptotically stable if, in addition to (2.11)

(2.12)  $\|\underline{y}(t) - \hat{y}(t)\| \rightarrow 0 \text{ as } t \rightarrow \infty$ 

and y(t) is relatively stable if, instead of (2.11)

(2.10c) y'(t)=Ay(t) t>a
stability can be fairly easily characterized (see Ortega [18]).

The 2.4 The unique solution of (2.10c,b) is stable if and only if all eigenvalues of A have nonpositive real part and any eigenvalues with zero real part belongs to a 1x1 Jordan block. Purthermore, y(t) is asymptotically stable if and only if all eigenvalues of A have negative real part (Ortega [18] and Pranklin [11]).

Thu. 2.5 If s is the set of eigenvalues of A with maximal real part, then the solution y(t) of (2.10c,b) is relatively stable if and only if  $\alpha$  has a component in the direction of a principal vector for some  $\lambda_i$  es, and this vector has the maximal degree

associated with od (Ortega [18]).

when A is a function of t, the eigenvalues of the Jacobian of A(t) can change sign, so it is only possible to give a characterization of stability corresponding to Thm. 2.4 which is local.

Now, consider the stability theory for BVPs. For simplicity, only the two point linear BVP (2.3a,b) is considered.

Assume (2.3 a,b) has a unique solution y(t) and also, the matrices  $B_a$  and  $B_b$  in (2.3b) are scaled such that

then (2.3 a,b) is stable if for any  $\xi > 0$ , there exists a  $\delta > 0$  such that the following is satisfied:

if max  $\{ | \delta B_{\mu} | , | \delta B_{\mu} | , | \delta A_{\mu} | \} \le \delta$  there is a  $b^* > a$  (independent of  $\xi$  and  $\delta$ ) such that for any  $b \ge b^*$  for which both (2.3 a,b) and the BVP (2.3 a),

(2.3b') 
$$\hat{B} \hat{y}(a) + \hat{B} \hat{y}(b) = \hat{\alpha}$$

where  $\hat{B} = B + \delta B$ ,  $\hat{B} = B + \delta B$  and  $\hat{a} = a + \delta a$  are well-posed, then a a a b b b

the respective solution  $\mathbf{y}(\mathbf{t})$  and  $\hat{\mathbf{y}}(\mathbf{t})$  satisfy

If we express the solution y(t) of (2.3 a,b) as

$$(2.16) y(t) = Y(t) s + v(t) a \le t \le b$$

where Y(t) is the n x n fundamental solution matrix which satisfies

$$(2.17) Y' = \lambda(t) Y = astSb$$

$$(2.17a)$$
  $Y(a) = I$ 

and v(t) is a particular solution of (2.3 a). Substituting (2.16) into the BC (2.3 b), we find that s is required to satisfy

(2.18) 
$$Q_{S} = \widehat{A}$$
where 
$$Q = \begin{bmatrix} B + B & Y & (b) \end{bmatrix}$$

$$\widehat{A} = \underline{A} - B & Y & (a) - B & Y & (b)$$

Let  $\tilde{y}(t)$  solve the differential equation (2.3 a) subject to the perturbed BC

(2. 19) 
$$\hat{B}_{a} y (a) + \hat{B}_{b} y (b) = \hat{\alpha}$$

where 
$$\delta B := \hat{B} - B$$
,  $\delta B := \hat{B} - B$ , and  $\delta c := \hat{a} - c$  are "small".

Prom Thm.2.3, since the nonsingularity of problem (2.3 a,b) is equivalent to the nonsingularity of the matrix Q defined in (2.18) which is the special case of (2.9) with r=2,  $f_1=a$ , and  $f_2=b$ , it is tempting to take the condition number of Q cond(Q):=||Q|| ||Q||

as an indication for the condition number of the BVP. This quantity however, turns out to be rather misleading at times. The reason is that Q contains the effects of the fundamental matrix Y(t) which is the solution of (2.17) subject to the initial values (2.17a), rather than the BC (2.3b). Thus, if the IVP behaves very differently than the BVP, Q may be

ill-conditioned even when the problem (2.3a,b) is not.

In general, write  $\hat{y}(t)$  as  $\hat{y}(t) = \hat{y}(t)$   $\hat{s} + \hat{y}(t)$   $\hat{s} + \hat{y}(t)$ 

and define \$5:=\$-s. We get

$$Q \cdot \delta s = \delta \alpha := \alpha - \alpha + \delta B \gamma (a) + \delta B \gamma (b)$$

Nov, the relevant quantity is

 $\hat{y}(t)-y(t)=Y(t)\delta s$ , not  $\delta s$  alone,

so an indication of the condition of the problem (2.3 a,b) is

the number

(2.20) 
$$k := \max_{1 \le a \le t \le b} |Y(t)Q|$$

rather than cond (Q).

The solution  $y^*(t)$  of (2.3a,b) can also be expressed as  $y^*(t) = Y(t)Q \propto + \begin{cases} b \\ G(t,s)f(s) ds \end{cases}$ 

where Q is in (2.18), G(t,s) is the Green's function for (2.3a,b).

The perturbations in the BCs give perturbed solutions which are related to  $|y^+(a)|$ ,  $|y^+(b)|$ , and hence to G(t,s). Therefore, conditioning of (2.3a,b) is also related to the boundedness of the Green's function G(t,s)

(Lentini-Osborne-Russell [15]).

It is obvious that  $k_1$  doesn't depend on the particular choice of fundamental matrix Y. If  $\Phi(t)$  is any n x n fundamental matrix satisfies (2.17), then there is a nonsingular n x n matrix P such that  $\Phi(t) = Y(t) P$ . Hence,

Therefore, to estimate the condition number of (2.3 a,b), it is tempting to use the bound

$$\max_{\mathbf{t}} \left| \underline{\Phi}(\mathbf{t}) \right| \left| \begin{bmatrix} \mathbf{B} \underline{\Phi}(\mathbf{a}) + \mathbf{B} \underline{\Phi}(\mathbf{b}) \end{bmatrix} \right| =$$

If we choose  $\xi(a) = I$ , i.e.  $\xi(t) = Y(t)$ , then this bound will frequently be misleading. Indeed, if Y(t) increases exponentially as t increases, then this bound is approximately cond (Q). To obtain a more realistic estimate,  $\xi(t)$  must be properly scaled. In particular, let

$$\underline{\Phi}(t) = \{\varphi_1(t), \dots, \varphi_n(t)\}$$
 be a fundamental solution matrix

max 
$$\| \underline{\Phi}(t) \| = 1$$
 and ast  $\leq b$   $\| \underline{\Phi}(t) \| = 1$  and the second se

Then the condition number of (2.3 a,b) can be approximated by

$$k = k$$
  $(b) := \| \begin{bmatrix} B & \Phi(a) + B & \Phi(b) \end{bmatrix}^{-1} \| .$ 

Hence the result follows:

Suppose the BVP (2.3 a,b) has a unique solution, then  $_3$ (2.3 a,b) is well-conditioned if  $k_1$  (b) =0(1) in (2.23) as b-> $\infty$ , and ill-conditioned if  $1/k_1$  (b) =0(1) as b-> $\infty$ .

Example: (Lentini-Osborne-Russell [15])

. Consider the BVP

$$y^{+}+y^{\pm}0$$
,  $y(0)=0=y(b)$ ,

then after some computation,  $k(b) = (1+\cos b)/\sin b$ ,

which means that for b away from multiple of  $\pi$ , the solution is not sensitive to small changes in the BCs. However, when b gets close to a multiple of  $\pi$ , the problem is unstable.

Stability properties of BTPs can also be indicated by the condition of the problems.

After investigating the existence and uniqueness of the solution of given problems and the stability of the problems, we then consider numerical methods for solving these problems.

## 3. Numerical Methods for Solving BVPs

Most of the interesting equations which occur in practice require that their solutions be obtained by numerical means. In this chapter, some well-known numerical methods for solving BVPs are discussed. First, initial value techniques, then finite difference methods, and finally finite element methods are discussed.

## 3.1 Initial Value Approaches

Initial value techniques play an important role in the numerical solution of BVPs. The basic process is to solve BVPs by solving IVPs with some arbitrary initial conditions, then find the solutions by satisfying the given BCs. In this section, several initial approaches such as superposition, shooting, stabilized march, and invariant imbedding are considered.

## 3.1.1. Superposition

Consider the linear BVP

(3.1a) 
$$y'(t) = \lambda(t) y(t) + f(t)$$
  $a \le t \le b$ 

(3.1b) 
$$\begin{array}{ccc}
B & y & (a) + B & y & (b) = & \\
a & & b & \\
\end{array}$$

Due to the linearity, the solution y(t) of (3.1) can be written

as

(3.2) 
$$y(t) = y(t;s) = Y(t)s + y(t)$$
where  $Y(t) = Y(t;a)$  is the fundamental solution matrix satisfying

$$Y'(t) = \lambda(t) Y(t)$$
  $a \le t \le b$ 

$$Y(a) = I$$

and v(t) is a particular solution of (3.1).

Here, s is to be determined so that the BC (3.1 b) are satisfied, so

$$\overset{=}{\underset{a}{\overset{=}{\bigcirc}}} \begin{bmatrix} \mathbf{I} (a) \overset{=}{\underset{b}{\overset{=}{\bigcirc}}} \mathbf{v} (a) \end{bmatrix} + \overset{=}{\underset{b}{\overset{=}{\bigcirc}}} \begin{bmatrix} \mathbf{I} (b) \overset{=}{\underset{b}{\overset{=}{\bigcirc}}} \mathbf{v} (b) \end{bmatrix}$$

OF

(3.3) 
$$Q_{S} = \hat{Q}$$
where  $Q = B + B Y (b)$  and a b
$$\hat{A} = A - B Y (a) - B Y (b)$$

If Q is nonsingular, then s can be obtained from (3.3) and so the solution y(t) is constructed. The above procedure is called the method of superposition. Generally, initial value methods which involve solving the ODE over [a,b] as an IVP, as the above does, are called shooting methods (see section 3.1.2).

when Q is formed, it is often ill-conditioned. Sometimes, this can be corrected by scaling, but realize that Q is nonetheless still ill-conditioned. For instance, consider the problem

$$y^{+}=10000y$$
  $y(0)=1$ ,  $y(1)=e$ 

A short computation gives

$$Q=B + B Y (b) = \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} + \begin{pmatrix} 0 & 0 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} \cosh 100 & (\sinh 100) / 100 \\ 100 & (\sinh 100) & \cosh 100 \end{pmatrix}$$

$$= \begin{pmatrix} 1 & 0 \\ \cosh 100 & \frac{(\sinh 100)}{100} \end{pmatrix} \approx \begin{pmatrix} 1 & 0 \\ 100 & 100 \\ \frac{e}{2} & \frac{e}{200} \end{pmatrix}$$

Q is ill-conditioned although by scaling (multiply row # by ), we get a well-conditioned problem.

When the BCs are partially separated, a reduced superposition method as decribed below can be used.

Suppose the partially separated BCs can be written as

(3.4a) 
$$C y(a) = d$$

(3.4b) D y (a) +D y (b) =
$$a$$

(3.4b) Dy(a)+Dy(b)=d a b 2 where C is a p x n matrix of rank p, D and D are q x n

matrices with p+q=n, and  $\alpha$  and  $\alpha$  are p- and q- vectors.

Write the solution of (3.1 a) and (3.4) as

$$y(t) = y(t;s) = \overline{Y}(t) \overline{s} + v(t)$$

where  $\overline{Y}(t)$  is an n x q matrix of fundamental solutions sastisf ying

(3.5a) 
$$\overline{Y}^{\dagger}(t) = \lambda(t) \overline{Y}(t)$$
 and

(3.5b) 
$$C \overline{Y}(a) = [0]$$

and the particular solution v(t) satisfies

(3.6) 
$$C v(a) = d$$
.

The q-vector's is determined so that the q boundary conditions

(3.4 b) are satisfied, i.e.

where  $\overline{Q} = [D \overline{Y}(a) + D \overline{Y}(b)]$  is a q x q matrix and b

$$\vec{a} = d - D \cdot v \cdot (a) - D \cdot v \cdot (b)$$
 a q-vector.

We need n linearly independant initial conditions to determine v(t) and each column of  $\overline{Y(t)}$ , while (3.5 b) or (3.6) only supply with p conditions. Hence we augment  $C_d$  by a q x n matrix G such that

Then, require the initial condition

$$\hat{B}_{\mathbf{a}} \mathbf{v}(\mathbf{a}) = \begin{pmatrix} \mathbf{a} \\ \mathbf{b} \\ \mathbf{0} \end{pmatrix} \qquad \hat{B}_{\mathbf{a}} \mathbf{v}(\mathbf{a}) = \begin{pmatrix} \mathbf{0} \\ \mathbf{I} \\ \mathbf{q} \end{pmatrix}$$

where I is the q x q identity matrix.

If we partition

where P has p columns and P has q columns,

then we have  $\underline{\underline{v}}(a) = \underline{\underline{F}} d$  and  $\overline{\underline{I}}(a) = \underline{\underline{F}} d$ .

After s is obtained, the solution y(t) is determined.

Although the superposition method is conceptually simple and works in many instances, it frequently gives very ill-conditioned problems even when the BVP is well-conditioned.

It has two major drawbacks: (Scott-Watts [23])

- 1. due to the finite word length used by computers, the solutions may lose their numerical independence. The resulting matrix problem in (3.3) may be so poorly conditioned that s cannot be determined accurately,
- 2. related to the finite word length of the computer, a loss of significance can occur even if the linear combination vector s has been computed accurately. This will normally occur if the fundamental solution Y(t) is large compared to the desired solution.

To overcome these difficulties, one can use multiple shooting (see section 3.1.3) or the stabilized march method (see section 3.1.4) to get a very well-conditioned matrix (Mattheij [17]) since they maintain stability by restricting the intergrations to smaller intervals, or by keeping the solutions nearly mutually orthogonal thus guaranteeing their independence over the entire interval. Also, normalizing the vectors at the initial point of each of the subintervals controls the growth of solutions (Scott-Watts [23]).

# 3.1.2. Simple Shooting

Simple shooting for solving

$$(3.7a) \quad y'=f(t,y) \quad a\leq t\leq b$$

(3.7b) 
$$g(y(a),y(b))=0$$

is the following: First guess the unknown initial values at a,

say  $\underline{s}$ . If we denote  $\underline{y}$  (t;  $\underline{s}$ ) as the solution of (3.7) subject to the initial conditions  $\underline{y}$  (a;  $\underline{s}$ ) =  $\underline{s}$ , then the problem reduces to finding a solution  $\underline{s}^{+}$  to a system of n nonlinear algebraic equations:

$$P(s) = 0$$
 where  $P(s) := g(s, y(b;s))$ .

One can solve the above system by Newton's method i.e. given  $s_o$ , then solve

(3.8) 
$$J(s) \Delta s = -P(s)$$
,  $s = s + \Delta s$ ,  $k = 0, 1, ...$ 
 $k = k$ 
 $k = k$ 

where J(s) = 2F/3s .

Define

$$B = \frac{3q(s; y(b, s))}{3s}$$

$$B = \frac{\partial g(\underline{s}; \underline{Y}(b,\underline{s}))}{\partial \underline{Y}(b,\underline{s})}.$$

By (3.7a),

$$\left(\frac{3\tilde{Z}}{3\tilde{A}}\right)_{i} = \left(\frac{3\tilde{Z}}{3\tilde{A}_{i}}\right)_{i} = \frac{9\tilde{Z}}{3\tilde{Z}} = \frac{9\tilde{Z}}{3\tilde{Z}} = \frac{9\tilde{Z}}{3\tilde{Z}} = \frac{9\tilde{Z}}{3\tilde{Z}}$$

and by the definition of y(a:s),

$$\frac{\partial y(a;g)}{\partial g} = \frac{\partial g}{\partial g} = I$$

Then

$$Y(t) = \frac{\partial Y(t;\underline{s})}{\partial \underline{s}}$$
 satisfies

(3.9a) 
$$Y'(t) = \underbrace{i(t, y(t; \underline{s}))}Y(t)$$
  $\underline{a \le t \le b}$ 

$$(3.9b)$$
  $Y(a) = I$ 

where

$$J(s) = \frac{\partial f(s)}{\partial s} = \frac{\partial g(s, f(b; s))}{\partial s}$$

$$= \underbrace{3q(\underline{s},\underline{y}(b;\underline{s}))}_{\partial \underline{s}} \underbrace{\frac{3\underline{s}}{\partial \underline{y}(b;\underline{s})}}_{\partial \underline{y}(b;\underline{s})} \underbrace{\frac{3\underline{y}(b;\underline{s})}{\partial \underline{y}(b;\underline{s})}}_{\partial \underline{y}(b;\underline{s})} \underbrace{\frac{3\underline{y}(b;\underline{s})}{\partial \underline{y}(b;\underline{s})}}_{\partial \underline{y}(b;\underline{s})}$$

$$= B + B \underline{y}(b) .$$

Thus, J(s) for the nonlinear case is like Q in the linear case.

For this method and for superposition, one can use the multiple shooting method to prevent the fundamental solution I from becoming numerically dependent or unbounded. This and alternative methods are described below.

## 3.1.3. Multiple Shooting

The multiple shooting method is a generalization of the shooting method which is designed to avoid the build-up of errors arising from the computation of fundamental solutions over large intervals (Keller [12]).

This is done by dividing [a,b] into J subintervals, solving IVPs involving fundamental solutions and a particular solution over each subinterval independently, and then taking the final solution as an appropriate combination of these solutions which satisfies the boundary conditions and is continuous across interior points connecting subintervals.

#### A. Linear case:

Multiple shooting for solving

(3. 10a) 
$$y^{+}(t) = \lambda(t) y(t) + f(t)$$
  $a \le t \le h$ 

takes the following general form:

divide [a,b] into J subintervals [t ,t ] (1≤j≤J) where

The fundamental solution  $Y_j$  (t) and particular solution  $v_j$  (t) on  $(t_j, t_{j+1})$  are obtained by solving

(3. 12a) Y'(t) = A(t) Y(t) 
$$t \le t \le t$$

j j j+1

for some given n x n matrix P , and

(3.13a) 
$$\forall'$$
 (t)=1(t)  $\forall$  (t)+f(t) t  $\leq$ t $\leq$ t  $\sim$ j  $\gamma$ +1

for some given vector v 1≤j≤J

Then constant vectors  $c_1, \ldots, c_J$  are determined so that on  $\{t_j, t_{j+1}\}$  the approximation solution y(t) defined by

(3.14) y(t) := Y(t)c + V(t)  $1 \le j \le J$ 

is continuous and satisfies the BCs.

Requiring the BC,

For the "standard" multiple shooting case, F = I and V = 0

j "j ~

(1≤j≤J) (and therefore c = y(t)).

~ j ~ j

The simplified form for this method is

# B. Nonlinear case:

Consider the nonlinear BVP

$$\underbrace{y'(t) = f(t, y(t))}_{a \le t \le b}$$

$$\underbrace{g(y(a), y(b)) = 0}_{a}$$

As in case A, [a,b] is divided into J subintervals [t,

t ]  $1 \le j \le J$ , and the IVPs j+1

(3.18a) 
$$y'(t)=f(t,y(t))$$
  $t \le t \le t$   $f'(t)=f(t,y(t))$ 

(3.18b) 
$$y$$
 (t) =c  $1 \le j \le J$   $\sim j$ 

are solved independently.

The constants  $c_j$ ,  $1 \le j \le J$ , are determined such that y(t) which is defined by

$$y(t) := y(t)$$
 on  $[t,t]$   $1 \le j \le J$ 
 $y(t) := y(t)$  on  $[t,t]$ 

satisfies continuity and the BCs, i.e.

and

$$g(y(a),y(b))=g(c,y(b))=0$$

Then the matrix form

(3.19) 
$$\phi(c) := \begin{cases} y(t;c) - c \\ 1 & 2 & 1 \\ 2 & 1 & 2 \end{cases} = 0$$

$$y(t;c) - c \\ 1 & 2 & 1 & 2 \\ 1 & 3 & 3 - 1 & 3 \\ 1 & 3 & 3 - 1 & 3 \\ 1 & 3 & 3 - 1 & 3 \\ 1 & 3 & 3 - 1 & 3 \\ 1 & 3 & 3 - 1 & 3 \\ 2 & 3 & 3 & 3 \\ 3 & 3 & 3 & 3 \\ 4 & 3 & 3 & 3 \\ 5 & 3 & 3 & 3 \\ 6 & 3 & 3 & 3 \\ 7 & 3 & 3 & 3 \\$$

where  $\tilde{z} = (c_1, \dots, c_j)^T$  and  $y_j(t; c_j)$  is a solution to (3.18).

We can solve (3.19) by Newton's method, i.e. given c, let

$$i$$
  $i$   $i$   $i+1$   $i$   $i$ 
 $J(c) = -\phi(c), c = c + 2c$   $i=0,1,...$ 

where the Jacobian matrix  $J(c) = \frac{\partial \phi}{\partial c}$ .

Define

$$g(\tilde{\lambda}(\cdot,\tilde{c})) = \frac{3\tilde{\lambda}(\tilde{a};\tilde{c})}{3\tilde{\lambda}(\tilde{a};\tilde{c});\tilde{\lambda}(\tilde{p};\tilde{c})}$$

$$d^{p}(\lambda(\cdot,c)) = \frac{9\lambda(p;c)}{9\lambda(p;c)}$$

$$\frac{9\lambda(p;c)}{9\lambda(p;c)}$$

From (3.18 a), we have

and (3.18 b) gives

Then Y (t;c) = 
$$\stackrel{\circ}{\text{j}}$$
 (t;c)

j  $\stackrel{\circ}{\text{j}}$   $\stackrel{\circ}{\text{j}}$  satisfies

$$Y'(t;c) = \frac{\partial f(t,y(t;c))}{\partial y} Y(t;c) \qquad t \le t \le t$$

$$j \sim j \qquad j \qquad j \rightarrow 1$$

and

$$J(c) = \begin{bmatrix} g_{1}(y_{1};c) & g_{2}(y_{1};c) & g_{3}(y_{1};c) & g_{4}(y_{1};c) & g_{4}(y_{1};$$

Thus, to solve  $\phi(c) = 0$  by Newton's method, each iteration involves solving a linear multiple shooting problem.

Indeed, to get a stable problem, several multiple shooting codes select the shooting points by attempting to satisfy  $\|Y_j^*(t_{j+1}^*)\| \le k$  for some constant k. In which case cond(M) =  $\|H\|\|H^{-1}\| \le (k+1), (k_1+k_2J)$  where k, and k<sub>2</sub> are given in (2.20) and (2.21).

The two sources of error using multiple shooting are

1. approximating the fundamental solution matrices and
particular solutions, and

2. solving the linear system of algebraic equations.

## 3.1.4. Stabilized March

The stabilized march procedure, like multiple shooting, is designed to maintain numerical linear independence of fundamental and particular solutions. While with multiple

shooting the fundamental solution Y; (t) and the particular solution Y; (t) are computed independently on each subinterval, this is not done with the stabilized march. The stabilized march also intends to economize on the multiple shooting method by reducing the number of fundamental solution components which must be computed in the same way that reduced superposition economizes on superposition (Scott-Watts [21]).

The procedure starts with  $Y_i$  (a) and  $y_i$  (a) satisfying (3.5b) and (3.6) given. Suppose we are given a fundamental solution  $Y_j$  (t) and a particular solution  $y_j$  (t) satisfying (3.5) and (3.6) for  $t \ge t_j$ . Let  $Y_j$  ( $t_j$ ) = $P_j$ , where  $P_j$  is a n x q matrix of rank q. When the solutions are becoming linearly dependent, say for  $t = t_{j+1}$ , then a factorization

$$(3.20) \qquad \mathbf{f} \quad (t) = \mathbf{P} \\ \mathbf{j} \quad \mathbf{j+1} \quad \mathbf{j+1} \begin{bmatrix} 0 \\ P \\ \mathbf{j+1} \end{bmatrix} = \begin{bmatrix} 1 & 1 & 2 \\ P & 1 & P \\ \mathbf{j+1} & \mathbf{j+1} \end{bmatrix} \begin{bmatrix} 0 \\ P \\ \mathbf{j+1} \end{bmatrix}$$

is performed, where P , P , P are n x p, n x q, and j+1 j+1 j+1

q x q matrices. Let

The process of computing a fundamental solution  $Y_i$  (t) and initial values  $F_i$  and  $y_i^o$  at  $t_i$  is continued for i=j+2, until eventually the point t=b is reached. As for reduced superposition, the computed solution is

(3.22) 
$$y(t) := Y(t) c + v(t) t \le t \le t$$
 $y(t) := Y(t) c + v(t) t \le t \le t$ 
 $y(t) := Y(t) c + v(t) t \le t \le t$ 

If  $b=t_{J+1}$ , then requiring y(t) to satisfy continuity and the BCs gives

Y (t ) c +v (t ) =Y (t ) c +v (t ) 
$$1 \le j \le J-1$$
  
j j+1  $^{\prime}$  j +1 j+1  $^{\prime}$  j+1  $^{\prime}$  j+1

Hultiplying by 
$$\begin{bmatrix} 0 \\ 2 \\ E \\ i+1 \end{bmatrix}$$
, then

$$\begin{bmatrix} 0 \\ p \\ j+1 \end{bmatrix} c = \begin{bmatrix} 0 \\ 1 \\ q \end{bmatrix} c - \begin{bmatrix} 0 \\ 2 \\ E \\ j+1 \end{bmatrix} (\forall (t) - \forall j+1) .$$

The last q BCs are D y(a)+D y(b)=D (P c + w)+D (Y c + a) b A = a 1 1 1 b A = a

(b) = d. In matrix form, we have

The stabilized march is not amenable to parallel processing like multiple shooting. However, it does compute a smaller solutions set and the multiple shooting points can be easily selected automatically. It also has the advantage that rapidly decaying fundamental solution information is only needed in the particular solution, viz one may be able to adapt an initial value solver such that it ignores these rapidly decreasing components in all integrations except that for this one solution vector (Lentini-Osborne-Russell [15]).

# 3.1.5. Invariant Imbedding

For simplicity, I only consider the BVP with linear separated BC

(3.24a) 
$$y'(t) = H(t) y(t) + h(t)$$
  $a \le t \le b$ 

$$(3.24b) B_{\mathbf{a}} \mathbf{y}(\mathbf{a}) = \mathbf{d}, B_{\mathbf{b}} \mathbf{y}(\mathbf{b}) = \mathbf{0}$$

Suppose (3.24) can be reformulated as

(3.25a) 
$$(\underbrace{\overset{\mathbf{u}}{\mathbf{v}}}(t)) = \begin{bmatrix} \lambda(t) & B(t) \\ C(t) & D(t) \end{bmatrix} (\underbrace{\overset{\mathbf{u}}{\mathbf{v}}}(t)) = (\underbrace{\overset{\mathbf{f}}{\mathbf{v}}}(t))$$

with BC

(3.25b) 
$$(K,K)$$
  $(\overset{\overset{\circ}{u}}{(a)})=\overset{\overset{\circ}{u}}{(a)}$   $(K,K)$   $(\overset{\overset{\circ}{u}}{(b)})=\overset{\overset{\circ}{u}}{(b)}$ 

where g(t), f(t), x are p-vectors, y(t), g(t), y are q-vectors, x are p x p matrices, y are y are

p x q matrices, C(t), K are q x p matrices, and D(t), K  $_{2}$ 

are q x q matrices.

The key idea in invariant imbedding is to replace a two point BVP by a set of initial value problems. One standard way is to first express the solution of (3.25 a) in the form (3.26) u(t)=R(t)v(t)+x(t)

where R(t) is a p x q matrix and x(t) is a p-vector. Substitute (3.26) into (3.25 a) to get

$$\underbrace{\mathbf{x}^{1} + \mathbf{R}C\mathbf{R} + \mathbf{R}D - \mathbf{A}\mathbf{R} - \mathbf{B}}_{\mathbf{x}^{2}} \underbrace{\mathbf{x}^{1} + \mathbf{R}C\mathbf{x} + \mathbf{R}g - \mathbf{A}\mathbf{x} - \mathbf{f}}_{\mathbf{x}^{2}} = 0$$

If we require that the coefficients of y(t) vanish in the first set above, and that (3.26) satisfies the BC (3.25 b) with the coefficient of y(a) set to zero,

then the solutions are found by solving three IVPs

(3.27a) 
$$\begin{cases} R' = [A(t) - RC(t)]R - RD(t) + B(t) \\ KR(a) + K = 0 \\ 0 & 1 \end{cases}$$

and

(3. 27c) 
$$\begin{cases} v' = [D(t) + C(t) R]v + C(t) x + g(t) \\ [KR(b) + K]v(b) = g - Kx(b) \\ 2 & 3 & 2 \end{cases}$$
 for a  $\leq t \leq b$ 

The first equation of (3.27 a) is called a matrix Riccati equation (See Reid [19]).

Suppose Y(t) is a fundamental solution matrix of (3.24 a), that is

$$Y^*=H(t)Y$$
  $a\le t\le b$   $\det(Y(a))\ne 0$ 

where 
$$H(t) = \begin{pmatrix} A(t) & B(t) \\ C(t) & D(t) \end{pmatrix}$$
.

If Y(t) can be partitioned into submatrices as H(t) can, i.e.

$$Y(t) = \begin{pmatrix} Y_{0}(t) & Y_{1}(t) \\ 0 & 1 \\ Y_{1}(t) & Y_{1}(t) \end{pmatrix} \} q$$

$$\underbrace{\frac{2}{p} \quad \frac{3}{q}}$$

and Y(t) satisfies

a) Y (t) if nonsingular on [a,b]

then the solution R(t) of (3.27 a) may be taken as R(t) =

 $Y_1(t)Y_3^{-1}(t)$  (Keller-Lentini [13]).

After R(t), x(t), and y(t) have been found, y(t) is easily found as a linear combination of them, and the solution y(t) of the BVP (3.24) is merely

$$\underline{y}(t) = (\frac{\underline{u}(t)}{\underline{v}(t)}) .$$

Even though invariant imbedding has the disadvantage of giving nonlinear initial value problems (for the matrix Riccati equation), it may well overcome the following two related difficulties of multiple shooting and stabilized march:

- when the related IVPs are unstable, short subintervals are necessary, and
- in the presence of rapidly growing/decreasing solutions,
   scaling can be a constant problem.

# 3.2. Pinite Difference Methods

Choose a mesh  $\Pi$ :  $a=t_1 < t_2 < \dots < t_N < t_{N+1} = b$ . The basic idea of the method involves finding approximate solution values at these mesh points  $t_j$  by the following:

- form a set of algebraic equations for the approximate solution values by replacing derivatives with difference quotients in the differential equations and the boundary conditions,
- 2. solve the resulting system of equations for the approximate solution. (Keller [12])

I start with a simple example.

3.2.1. A Simple Scheme for a Second Order Problem
Consider the scalar BVP

(3.29a) Lu(t) =-u''+a (t) u'+a (t) u=b(t) 
$$0 \le t \le 1$$

$$(3.29b)$$
  $u(0) = \alpha$  ,  $u(1) = \beta$ 

where a (t), a (t) and b(t) are continuous functions on 1

[0,1].

Take a uniform mesh for this problem ,i.e.  $t_i = (i-1) h$ , i=1,2,...,N+1, h=1/N.

Assume the exact solution of (3.29) exists. A mesh function

The differential equation (3.29 a) is approximated by

(3.30a) Lu(t) = L<sub>u</sub> = - 
$$\frac{u}{i+1}$$
  $\frac{-2u}{i}$  + a (t)  $\frac{u}{i+1}$   $\frac{-u}{i-1}$  + a (t)  $\frac{u}{i+1}$   $\frac{-u}{i+1}$   $\frac{-u}{i+1}$  + a (t)  $\frac{u}{i+1}$   $\frac{-u}{i+1}$  + a (t)  $\frac{u}{i+1}$  + a (t

and the BCs (3.29 b) give

(3.30b) 
$$u = 4$$
,  $u = 6$ 

If (3.30) is written in matrix form, then the following tridiagonal system is obtained:

$$\begin{bmatrix}
2/h^{2}+a & (t ), & -1/h^{2}+a & (t )/2h \\
0 & 2
\end{bmatrix}$$

$$-1/h^{2}-a & (t )/2h, & 2/h^{2}+a & (t ), & -1/h^{2}+a & (t )/2h \\
1 & 3
\end{bmatrix}$$

$$-1/h^{2}-a & (t )/2h, & 2/h^{2}+a & (t ), & -1/h^{2}+a & (t )/2h \\
1 & 3 & 1 & 3
\end{bmatrix}$$

$$-1/h^{2}-a & (t )/2h, & 2/h^{2}+a & (t ) \\
1 & 3 & 0 & 3
\end{bmatrix}$$

$$\begin{bmatrix}
u \\ 2 \\ u \\ 3 \\ . & .
\end{bmatrix}$$

$$\begin{bmatrix}
b(t) + ol(1/h^{2}+a & (t)/2h) \\
2 & b(t) \\
3 & .
\end{bmatrix}$$

$$b(t) \\
3 & .
\end{bmatrix}$$

$$b(t) \\
4 & .
\end{bmatrix}$$

$$b(t) + (3(1/h^{2}-a & (t)/2h) \\
b(t) + (3(1/h^{2}-a & (t)/2h) \\
1 & .
\end{bmatrix}$$

If  $a_0(t)$  is positive, the approximate mesh function  $\{u_j\}_{j=1}^{N+1}$  exists, the matrix is positive definite, and the process of Gauss elimination without pivoting for tridiagonal matrices is extremely simple and efficient and stable (Ascher-Russell [3]).

# 3.2.2. One Step Schemes : Trapezoidal Rule and Midpoint Rule

Por a given first order system of DEs, e.g. y'=f(t,y), if one integrates both sides of the equations, then the left hand

side of the integration directly gives the solutions of the DB.

Most finite difference methods involve converting to first order

systems and then selecting a discretization.

Consider the linear first order system

(3.31a) Ly (t) = y' (t) - 
$$\lambda$$
 (t) y (t) =  $b$  (t) a \le t \le b

(3.31b) B [
$$y(a),y(b)$$
]=B  $y(a)+B$   $y(b)=0$ 

where A(t), B, and B are n x n matrices.

a b

The two simplest one step schemes, which use only information about the approximate solution at  $t_i$  to obtain the approximate solution at  $t_{i+1}$ , are the trapezoidal method and midpoint method.

On a mesh  $\pi$ , a numerical solution  $\{y_i\}_{i=1}^{N+1}$  is sought where  $y_i$  is to approximate componentwise the exact solution  $y_i(t)$  at  $t=t_i$  and is required to satisfy the BCs

$$B[y,y]=By+By=0$$
 $-1-H+1$ 
 $a-1$ 
 $b-H+1$ 

The trapezoidal method is defined as

and the midpoint method or the Box scheme is defined by

$$L_{R} = \frac{Y}{-i+1} - \frac{1}{2} (t) (y + y) = b(t)$$

$$i + 1/2 - i + 1 - i - i + 1/2$$

where h = t - t and t := t + h / 2,  $1 \le i \le n$ . i + 1 + i = i + 1 / 2 = i = i For the nonlinear problem

(3.32a) 
$$y(t) = y(t) - f(t;y) = 0$$
  $a \le t \le b$ 

(3.32b) 
$$g(y(a),y(b))=0$$
,

the trapezoidal rule is given by

(3.33) 
$$g(y,y)=0$$
  
 $\sim 1 - 1 + 1$ 

1≤i≤#

and the midpoint rule is given by (3.33) and

again 
$$h = t - t$$
 and  $t := t + h / 2$  ,  $1 \le i \le N$  .

i i+1 i i+1/2 i i

For the nonlinear problem (3.32), it is not difficult to form the difference schemes (3.33), (3.34), and (3.35). The difficulty is in solving the resulting n(N+1) nonlinear algebraic equations, where n is the order of the first order system(3.32a) and N is the number of subintervals. For instance, n=5 and N=200 gives more than 1000 equations. We consider Newton's method to solve the nonlinear problems. From (3.33) and (3.34),

$$H_{\pi_{1}} = \frac{\frac{y-y}{2}}{h} - \frac{1}{2} [f(t,y) + f(t,y)] = 0$$

$$H_{\frac{y}{2}} = \frac{\frac{y}{3} - \frac{y}{2}}{\frac{h}{2}} - \frac{1}{2} \left[ \frac{f(t, y) + f(t, y)}{2 - 2} \right] = 0$$

(3.36) ...

$$H y = \frac{y}{h} - \frac{y}{h} - \frac{1}{2[f(t,y)]} + f(t,y) = 0$$

$$g(y,y) = 0$$

$$H = H - H + 1$$

Letting the system of equations (3.36) be  $\mathbf{r}(\mathbf{y}) = 0$  where

$$\overset{\mathbf{F}}{\sim} ( \overset{\mathbf{H}}{\mathsf{m}} \overset{\mathbf{y}}{\mathsf{n}} , \overset{\mathbf{N}}{\mathsf{m}} \overset{\mathbf{y}}{\mathsf{n}} , \dots ) \overset{\mathbf{T}}{\mathsf{n}} , \text{ then}$$

$$\frac{\partial \underline{r}(t,\underline{y})}{\partial \underline{y}} = 
\frac{\partial \underline{r}(t,\underline{y})}{\partial \underline{y}} + 
\frac{\partial \underline{r}(t,\underline{y})}{\partial \underline{r}} + 
\frac{\partial \underline{r}(t,\underline{r})}{\partial \underline{r}} + 
\frac{\partial \underline{r}(t,\underline{$$

Newton's method becomes solving for w from the following

$$(3.37) \qquad \frac{\partial F}{\partial Y} = J \cdot \begin{bmatrix} w \\ 2 \\ \vdots \\ w \\ 2 \end{bmatrix} = -F(y)$$

If written component-wise, (3.37) gives

(3.38) 
$$\frac{v - w}{h} - 1/2[\lambda(t) w + \lambda(t) w] = -h_{\pi} y$$
i

and

(3.39) B w +B w =-g(y,y)  

$$a^{-1}$$
 b^n+1  $\sim$  1  $\sim$  H+1

where {y} are known values from a former iteration,

~i i=1

(3.40)
$$B = \frac{\partial g}{\partial y} (y, y),$$

$$a = \frac{\partial g}{\partial y} (x, y),$$

$$b = \frac{\partial g}{\partial y} (y, y),$$

$$b = \frac{\partial g}{\partial y} (y, y),$$

and the next iterate is given by y :=y +w , i=1,...,N+1.

~i ~i ~i ~i

The system (3.38), (3.39), and (3.40) is a linear systems of equations for the correction vector  $\{\mathbf{w}_i\}_{i=1}^{N+1}$ , which looks like a trapezoidal discretization of some linear problems. In each iteration we have performed two operations in succession, discretization and linearization.

Let  $y^{(m)}(t)$  be an appropriately smooth function satisfying  $y^{(m)}(t_i) = y_i^m$ , i = 1, ..., N+1. If we first linearize the differential problem

$$H y = y' - f(t, y) = Dy - f(t, y),$$

Newton's method is :

given y, solve

$$(3.41) \qquad \frac{3}{3} \underbrace{y}_{\underline{y}} = - \underbrace{y}_{\underline{y}}^{\underline{m}} ,$$

and let y = y + y, n=0,1,2,...Since  $\frac{\partial y}{\partial y} = \frac{\partial y}{\partial y}$ ,

(3.41) is

$$\widetilde{D^{M-9\widetilde{\xi}}(f'\widetilde{\lambda})}\widetilde{a}=\widetilde{a}_{i}-\widetilde{9\widetilde{\xi}(f'\widetilde{\lambda})}\widetilde{n}$$

Now, if discretization (trapezoidal rule) is applied to solve these equations, from linearizing the differential equation, we get (3.38) and (3.39). Thus, the two operations of linearization and discretization are commutative for the trapezoidal scheme. The iterative scheme where we first linearize and then discretize is called quasilinearization.

# 3.2.3. Runge-Kutta Schemes

The Runge-Kutta methods involve using only information about the approximate solution at  $t_i$  to obtain the approximate solution at  $t_{i+1}$ . The general form of a <u>k-stage Runge-Kutta</u> scheme for (3.32) on a mesh is

$$(3.42) y = y + h \sum_{j=1}^{k} \beta f$$

$$\sim i+1 \sim i i j=1 j \sim i$$

(3.44) 
$$g(y, y) = 0$$

The method is called <u>explicit</u> if  $\alpha_{j|}=0$   $j\geq 1$  and <u>implicit</u> otherwise. Both the trapezoidal rule and midpoint rule are implicit Runge-Kutta schemes. The first is 2-stage and the second is 1-stage.

## 3.3. Finite Element Methods

Like the finite difference methods, finite element methods attempt to find approximate solutions of a BVP at a discrete set of points by satisfying the BC and ODE simultaneously throughout the interval. However, when solved by finite element methods, differential equations need not be converted to a first order system.

Consider the scalar problem

(3.45a) Ly(t):= y (t) - 
$$\sum_{j=0}^{n-1} a_j$$
 (j)  $a \le t \le b$ 

(3.45b) 
$$\sum_{l=0}^{m-1} \{b \ y \ (a) + c \ y \ (b) \} = r$$
 1\le j\le m.

For finite element methods, the approximate solution is a spline function  $s(t) \in P_{K,\pi,1}$ , where  $P_{K,\pi,1}$  is a collection of spline functions which are of order k (degree less than or equal to k-1) in each subinterval of  $\pi$ :  $a=t_1 < t_2 < ... < t_N < t_{N+1} = b$  and are --Ith order continuous at every mesh point. The advantages of selecting a spline space Pr 77 are that high order methods result and local basis representations produce banded matrix equations. For convenience, assume that  $s(t) \in P_{\kappa,\pi,1}^{\sigma}$  , the subspace of P K. T. 1 consisting of spline functions which satisfy the BC (3.45b). So the unknown solution parameters correspond to some representation for s(t). Letting the approximate solution be  $s(t) = \sum_{i=1}^{m} \alpha_i \Psi_i$  (t) where  $\{\Psi_i(t)\}$  is a basis of  $P_{\kappa,\pi,1}$ , H=dim(  $P_{K,\pi,1}^{\circ}$ ), we determine  $\alpha_{j}$ , j=1,...,5 by requiring s(t) to satisfy the differential equations in one of several natural ways. The basic types of finite element methods are collocation, Galerkin, least squares, and Ritz methods which are described below.

#### 3.3.1 Collocation Methods

For solving (3.45), the collocation solution s(t) is required to satisfy the differential equation (3.45a) exactly at H points  $\{z_i\}_{i=1}^M$  (called the collocation points) in [a,b], i.e.

(3.46) Ls<sub>1</sub>(z)=L(
$$\sum_{j=1}^{H} \alpha_j \exists_j \in \mathbb{Z}$$
)  $1 \le i \le \exists_j \in \mathbb{Z}$ 

Thus collocation requires the residual r(t):=Ls(t)-q(t) be set to zero at 4 points. In matrix form, (3.46) is

$$(3.47)$$
  $Cd = q$ 

where 
$$C:=(c)$$

$$ij i, j=1$$

$$g:=(q(z), \dots, q(z))$$

$$\alpha:=(\alpha, \dots, \alpha)$$

$$1$$

Suppose the collocation points  $\mathbf{z}_{\mathbf{K}}$  are in the jth subinterval and can be expressed as

$$z = 7 = 7 + \frac{1}{2}\rho$$

$$k \quad ji \quad j+1/2 = 2 i$$

Let P (t) be the Gauss Legendre polynomial of degree k. If k

P are chosen to be the zeros of P (t), 1≤i≤1, and to
i

satisfy  $-1 < \rho < \rho < \dots < \rho < 1$ , then z are called the Gauss

points. If  $\rho$  is the (i-1) st zero of [P (t)+P (t)]/(t-1),

 $2 \le i \le 1$ , and  $-1 = \binom{2}{2} < \ldots < \binom{2}{i} < 1$ , then z are called the Radau points. In the case that  $\binom{2}{i}$  is the (i-1)st zero of  $\binom{2}{i}$  (t),  $\binom{2}{i} < 1 \le i \le 1-1$ , and  $-1 = \binom{2}{2} < \ldots < \binom{2}{i} = 1$ , then z are called the Lobatto points.

#### 3.3.2. Galerkin Methods

Por the Galerkin method, the approximate solution  $\overline{s}(t) = \sum_{j=1}^{M} \overline{a_j} \psi_j(t)$  is determined so that the differential equation (3.45a) is satisfied in the sense that

(3.48) 
$$\int_{a}^{b} L\overline{s}(t) \psi(t) dt = \int_{a}^{b} q(t) \psi(t) dt \qquad 1 \le i \le 1.$$

That is, one requires the residual  $\vec{r}(t) = L\vec{s}(t) - q(t)$  to satisfy

$$\int_{a}^{b} \overline{r}(t) f(t) dt = 0 \qquad \text{for all } f(t) \in P^{\circ}_{k, \pi, 1}.$$

In matrix form, (3.48) is  $G_{\overline{a}} = \overline{q}$ ,

where 
$$3:=(j)$$
 =  $(\int_{a}^{b} L\psi(t)\psi(t)dt)$  if i, j=1

$$\overline{q} := (\overline{q}_{1}, \dots, \overline{q}_{H})^{T}$$

$$= (\int_{a}^{b} q(t) + (t) dt, \dots, \int_{a}^{b} q(t) + (t) dt)^{T}$$

Unless the BVP is extremely simple, the elements of G and  $\vec{q}$  must be approximated by a numerical quadrature. If the quadrature rule has the form

(3.50) 
$$\int_{a}^{b} f(t) dt = \sum_{i=1}^{Q} w_{i} f(z_{i}) ,$$

the resulting discrete Galerkin method solution s (t) =

 $\Sigma \stackrel{*}{a} \psi$  (t) is obtained from the system of equations j=1 j

(3.51) 
$$G \stackrel{*}{\bowtie} = \stackrel{*}{q}$$

where  $G = (g)$   $= (\sum_{i \neq i} x_i + (z_i) + (z_i)$ 
 $i \neq i$   $i \neq$ 

The discrete Galerkin equations (3.51) can be written as

If Q=M and B is nonsingular, then the collocation and discrete Galerkin methods are equivalent. Compared with the collocation method, the Galerkin method has the disadvantage that the integral coefficients must be evaluated and the advantage that, for the same order of convergence, smoother spline functions can be used.

# 3.3.3. Least Squares Method

The Least squares method is to find  $\hat{s}(t) = \sum_{j=1}^{M} \hat{a} \psi(t)$ 

such that 
$$E(s) = \min_{s \in P} E(s)$$
 where  $k, \pi, 1$ 

$$E(s) = E(\alpha, \dots, \alpha) := \int_{a}^{b} [Ls(t) - q(t)]^{2} dt$$

By setting 
$$\frac{\partial E}{\partial x_i} = 0$$
  $1 \le i \le B$ ,

we find this is equivalent to requiring

$$\int_{a}^{b} [L\hat{s}(t) - q(t)] L\psi(t) dt = 0 \qquad 1 \le i \le n$$

OL

(3.53) 
$$\int_{a}^{b} [L\widehat{s}(t) L \psi(t)] dt = \int_{a}^{b} q(t) L \psi(t) dt$$
 15i5H.

In matrix form, (3.53) is

(3.54) 
$$\widehat{\exists} \hat{x} = \widehat{q}$$
 where

$$\widehat{H} = (h) \qquad H \qquad = (\int_{a}^{b} L \psi(t) L \psi(t) dt) \qquad H \qquad i, j=1$$

$$\widehat{q} = (\widehat{q}_{1}^{2}, \dots, \widehat{q}_{N}^{2})^{T} \qquad = (\int_{a}^{b} L \psi(t) q(t) dt, \dots, \int_{a}^{b} L \psi(t) q(t) dt)^{T} \qquad \widehat{\alpha} = (\widehat{\alpha}_{1}^{2}, \dots, \widehat{\alpha}_{N}^{2})^{T} \qquad .$$

As with the Galerkin method, the discrete least squares method involves evaluating the integral coefficients in (3.55) by a numerical quadrature rule of the form (3.50).

The solution is 
$$\hat{s}^*(t) = \sum_{j=1}^{n} \hat{a}^*(t)$$
 if

(3.56) 
$$C DC\widehat{A} = C Dq$$
where D,C,and  $q$  are as in (3.52), (3.47) and  $\widehat{A} = (d_1, \dots, d_n)$ .

Clearly, if Q=M and the collocation matrix is nonsingular, the discrete least squares and collocation methods are equivalent.

#### 3.3.4. Ritz Method

Consider the problem

(3.57a) Ly(t) = 
$$\sum_{i=0}^{n/2} (-1)^{i} D (0^{i} (t) D y(t)) = q(t)$$

$$i-1$$
  $i-1$   $(3.57b)$  D y (a) = D y (b) = 0  $1 \le i \le n/2$ 

where Dy(t):=y'(t) and the (smooth) coefficient functions satisfy  $\sigma$  (t)  $\geq 0$   $1 \leq i \leq m/2-1$  and  $\sigma$  (t)  $\geq 0$  for  $a \leq t \leq b$ .

The operator L is called a <u>self-adjoint</u> operator, having the property that it satisfies

$$\int_{a}^{b} Lu(t) \forall (t) dt = \int_{a}^{b} u(t) L \forall (t) dt$$

for any  $u, v \in C_{\bullet}^{m}[a,b]$ , the space of functions in  $C_{\bullet}^{m}[a,b]$  which satisfy the BC (3.57b).

The Ritz method involves choosing an approximate subspace and letting the approximate solution be the function which 'minimizes the variational formulation I (u) for (3.57) over that subspace. Here

$$I(u) := \begin{cases} b & \frac{\pi}{2} \\ & \sum_{i=0}^{b} \int_{2i}^{\pi} (t) (D u(t)) -2q(t) u(t) \} dt. \end{cases}$$

The Ritz solution  $\widetilde{s}(t) = \sum_{j=1}^{H} \overrightarrow{\lambda} \gamma(t)$  satisfies

(3.58) 
$$I(\widehat{s}) = \min_{0} I(s)$$

$$s \in P$$

$$k_{n} = 1$$

Setting 
$$\frac{\partial I(s)}{\partial \alpha} = 0$$
 and  $\frac{\partial^2 I}{\partial \alpha^2} > 0$   $1 \le j \le h$ 

the matrix form is

$$(3.59)$$
  $\widetilde{G}_{\widetilde{d}} = \widetilde{q}$  where

$$\widetilde{G} = (\widetilde{g})$$

$$ij \ i, j=1$$

$$= (e(\psi(t), \psi(t)))$$

$$i, j=1$$

$$e(u, v) := \int_{a}^{b} (\sum_{i=0}^{m/2} (i) D u(t) D v(t)) dt$$

$$= \int_{a}^{b} Lu(t) v(t) dt$$

and  $\widetilde{q}$  is in (3.49).

It is clear that for solving (3.57), the Galerkin and Ritz methods are mathematically equivalent, at least if k,l≥m. However, the discrete Ritz method for which the integral coefficients (3.60) are approximated by using a quadrature, is generally different from the discrete Galerkin method (3.51) and has the advantage of preserving the matrix symmetry in (3.59).

while certain methods have their particular advantage in special cases, we usually only consider the collocation method since it appears to be generally the most efficient and since software for this method has been developed. For a comparison, see Russell-Varah [20].

#### 4. 'Equivalences Between These Bethods

In this Chapter, various equivalences between the methods mentioned in Chapter 3 are presented.

# 4.1. Collocation and Finite Difference Methods

#### 4.1.1. Collocation, Trapezoidal, and Midpoint Rules

Consider the BVP (3.7). The collocation methods using 2

Lobatto points and 1 Gauss point relate to the trapezoidal rule

and the midpoint rule (or Box scheme), respectively.

In particular, when solving the BVP (3.7) by collocation with approximate solution s(t) in  $P_{1,|\Pi|,1}$ , if the Gauss points are the collocation points then we have

(4.1) 
$$s'(t) = f(t), s(t)$$
  $1 \le i \le n$   
 $i+1/2 = i+1/2$ 

(4.2) 
$$g(s(a), s(b)) = 0$$
  
where t = (t +t )/2  
i+1/2 i i+1

Since s(t) = y,  $1 \le i \le N+1$ , the Newton form of the interpolating polynomial s(t) can be expressed as

$$s(t) = y + (y - y)(t-t)/h$$
  
 $\sim i - i + 1 \sim i i$ 

so

$$s'(t) = (y -y)/h$$
 -  $i+1 \sim i$  i

(i) Case I:

At the Gauss point t

$$i+1/2$$

$$s'(t) = (y - y)/h$$
  
 $\sim i+1/2 \sim i+1 \sim i i$ 

$$s(t) = y + (y - y)(t - t)/h$$
  
 $\sim i+1/2 \sim i \sim i+1 \sim i i+1/2 i i$ 

$$= (y + y)/2$$

$$\sim i+1 \sim i$$

Now (4.1) gives

$$N_{x}y = (y -y)/h - f(t , 1/2 (y + y )) = 0$$
 $N_{x}i \sim i+1 \sim i \sim i+1/2 \sim i \sim i+1 \sim$ 

and (4.2) gives g(s(y),s(y))=0,

which are the midpoint rule (3.33) and (3.35), one of the most widely used finite difference methods.

(ii) Case II:

When collocating at Lobatto points t , and t , i +1

$$\overset{s'}{\sim} \overset{(t)}{\sim} \overset{=f}{\sim} \overset{(t)}{\sim} \overset{s(t)}{\sim} \overset{i.e.}{\sim}$$

(4.3) 
$$(y - y)/h = f(t, y)$$
  
 $\sim i+1 \sim i$   $i \sim i$ 

and

$$s^{i}(t)=f(t,s(t)), i.e.$$
  
 $\sim i+1 \sim i+1 \sim i+1$ 

From (4.3) and (4.4), we get

$$H_{T_{-i}} = \{y - y\}/h - 1/2[f(t,y)+f(t,y)]=0$$

which is the trapezoidal rule in (3.34).

## 4.1.2. Collocation and Trapezoidal Rule

In the previous subsection, it has been shown that there is an equivalence between the trapezoidal rule and collocation method for  $P_{2,\,\pi,\,1}$ . Here, the same result is obtained for  $P_{3,\,\pi,\,2}$ .

Suppose s(t) = y

and

$$s'(t) = y' \qquad 1 \le i \le N+1$$

$$c \qquad i \qquad i$$

The Newton form of the interpolating polynomial s(t) gives

$$s(t) = y + y'(t-t) + (y - y - y'h)(t-t)^{2}/h^{2}$$
  
 $\sim i \sim i i \sim i+1 \sim i i i i$ 

Hence

$$s'(t) = y' + 2 (y - y - y'h) (t-t)/h^2,$$
  
 $\sim i \sim i + 1 \sim i \sim i i i$ 

so

and

Since s(t) satisfies the differential equation at collocation points, i.e.

$$s'(t)=f(t,s(t))$$
 and  $s'(t)=f(t,s(t))$ ,  
 $i+1 \sim i+1 \sim i+1 \sim i$ 

from above

$$2(y - y)/h - f(t,y) = f(t,y)$$
  
 $\sim i+1 \sim i \sim i \sim i+1 \sim i+1$ 

Hence, we have

$$M_{\pi} y = (y - y)/h - 1/2[f(t,y)+f(t,y)] = 0$$
  
 $M_{\pi} y = (y - y)/h - 1/2[f(t,y)+f(t,y)] = 0$ 

which is the trapezoidal rule in (3.34).

# 4.2. Simple Ritz and Pinite Difference

Consider the Bitz method for the simple second order BVP

$$(4.5a) -y''(t) + f(t)y(t) = q(t) \qquad a \le t \le b$$

$$(4.5b)$$
  $y(a) = 0$ ,  $y(b) = 0$ 

with  $\widetilde{s}(t) \in P_{a, \pi, 1}$ . For simplicity, consider a uniform mesh, i.e.  $h_i = (b-a)/N$  for  $1 \le i \le N$ . Let the piecewise linear B-spline basis functions which satisfy (4.5) be

$$B (t) = \begin{cases} (t-t)/h & t \in [t, t] \\ j-1 & j-1 \end{cases}$$

$$B (t) = \begin{cases} (t-t)/h & t \in [t, t] \\ 2 \le j \le N \end{cases}$$

$$0 & otherwise$$

(For the construction and evaluation of B-splines, see de Boor [6] [7], and Ascher-Russell [2]).

Then  $\widetilde{G} = (\widetilde{g})$   $\widetilde{q}$  in (3.59) can be shown to be if i, j=1, i

$$(4.6) \quad \widetilde{g}_{ij} = \begin{cases} \frac{2}{h} + \int_{t}^{t} \sigma(t) \left(\frac{t-t}{h}\right)^{2} dt + \int_{t}^{t} \sigma(t) \left(\frac{t-t}{h}\right)^{2} dt \\ \frac{-1}{h} + \int_{t}^{t} \sigma(t) \left(\frac{t-t}{h}\right) \left(\frac{t-t}{h}\right) dt \quad j=i\pm 1 \\ 0 \quad \text{otherwise} \end{cases}$$

for  $2 \le i, j \le N$  and

$$(4.7) \quad \widetilde{q} = \begin{cases} t \\ i \\ t \end{cases} q(t) \left(\frac{t-t}{h}\right) dt + \begin{cases} t \\ i+1 \end{cases} q(t) \left(\frac{t-t}{h}\right) dt \quad 2 \le i \le N .$$

In this case (3.59) can be expressed as

$$\begin{cases} -1/h + \int_{t-1}^{t} \sigma(t) \left( \frac{t-t}{h} \right) \left( \frac{t-t}{h} \right) dt \} y + \{2/h \} \\ + \int_{t-1}^{t} \sigma(t) \left( \frac{t-t}{h} \right)^{2} dt + \int_{t-1}^{t} \sigma(t) \left( \frac{t-t}{h} \right)^{2} dt \} y \\ + \{-1/h + \int_{t-1}^{t} \sigma(t) \left( \frac{t-t}{h} \right) \left( \frac{t-t}{h} \right) dt \} y \\ + \{-1/h + \int_{t-1}^{t} \sigma(t) \left( \frac{t-t}{h} \right) \left( \frac{t-t}{h} \right) dt \} y \\ + \{-1/h + \int_{t-1}^{t} \sigma(t) \left( \frac{t-t}{h} \right) \left( \frac{t-t}{h} \right) dt \} y \\ + \{-1/h + \int_{t-1}^{t} \sigma(t) \left( \frac{t-t}{h} \right) dt + \int_{t-1}^{t} \sigma(t) \left( \frac{t-t}{h} \right) dt \} y \\ + \{-1/h + \int_{t-1}^{t} \sigma(t) \left( \frac{t-t}{h} \right) dt + \int_{t-1}^{t} \sigma(t) \left( \frac{t-t}{h} \right) dt \} y \\ + \{-1/h + \int_{t-1}^{t} \sigma(t) \left( \frac{t-t}{h} \right) dt + \int_{t-1}^{t} \sigma(t) \left( \frac{t-t}{h} \right) dt \} y \\ + \{-1/h + \int_{t-1}^{t} \sigma(t) \left( \frac{t-t}{h} \right) dt + \int_{t-1}^{t} \sigma(t) \left( \frac{t-t}{h} \right) dt \} y \\ + \{-1/h + \int_{t-1}^{t} \sigma(t) \left( \frac{t-t}{h} \right) dt + \int_{t-1}^{t} \sigma(t) \left( \frac{t-t}{h} \right) dt \} y \\ + \{-1/h + \int_{t-1}^{t} \sigma(t) \left( \frac{t-t}{h} \right) dt + \int_{t-1}^{t} \sigma(t) \left( \frac{t-t}{h} \right) dt \} y \\ + \{-1/h + \int_{t-1}^{t} \sigma(t) \left( \frac{t-t}{h} \right) dt + \int_{t-1}^{t} \sigma(t) \left( \frac{t-t}{h} \right) dt \} y \\ + \{-1/h + \int_{t-1}^{t} \sigma(t) \left( \frac{t-t}{h} \right) dt + \int_{t-1}^{t} \sigma(t) \left( \frac{t-t}{h} \right) dt \} y \\ + \{-1/h + \int_{t-1}^{t} \sigma(t) \left( \frac{t-t}{h} \right) dt + \int_{t-1}^{t} \sigma(t) \left( \frac{t-t}{h} \right) dt \} y \\ + \{-1/h + \int_{t-1}^{t} \sigma(t) \left( \frac{t-t}{h} \right) dt + \int_{t-1}^{t} \sigma(t) dt + \int_{t-1}^{t} \sigma$$

#### 2≤i≤#

If the trapezoidal rule is used to approximate the coefficients in (4.6) and (4.7), then (4.8) becomes

$$-y$$
 /h +{2/h +h0(t)}y -y /h=hq(t)  
i i i+1 i

Dividing by h, we have

$$-\frac{y^{-2y+y}}{\frac{i+1}{h^2}} + f(t) y = q(t),$$

which is identical to a finite difference scheme for solving (4.5) (see Varah [22]).

#### 4.3. Collocation and Implicit Bunge-Kutta

when DEs are solved by most finite difference methods, they are converted to first order system. To relate collocation methods with Runge-Kutta methods, we consider the first order nonlinear DE

$$y' = f(t,y)$$

The collocation schemes for (4.9) are

$$y(t)=y$$
,  $y'(t)=f(t,y(t))$   
 $\sim i \sim i$   $\sim ij \sim ij \sim ij$ 

where t are collocation points which satisfy ij

and  $\binom{1}{j}$  are canonical points in [0,1]  $0 \le \binom{1}{1} < \binom{1}{2} < \dots < \binom{k}{k} \le 1$ .

Let f = f(t, y(t)), j=1,...,k, and express y in terms ij = ij = ij

of interpolation to the values y , f ,...,f , i.e. wi will wik

$$y(t) = y + h \sum_{i=1}^{k} f \phi_{i} \left(\frac{t-t}{h}\right)$$

where  $\phi$  (t) j=1,...,k are polynomials of degree at most k

on [0,1] determined by interpolation conditions

(4.10) 
$$\phi_{j}(0) = 0$$
  $\phi_{j}(\rho_{j}) = \delta_{j1}$   $1 = 1, ..., k$ 

here  $\delta$  denotes the Kronecker delta function.

Let

(4.11) 
$$\beta_{j} = \phi_{j}(1)$$
  $\alpha_{j1} = \phi_{j1}(\beta_{j})$ 

Then we get the equivalent implicit Runge-Kutta method

$$(4.12) y = y + h \sum_{j=1}^{k} (j f 1 \le i \le N)$$

$$f = f(t : y + h \sum_{j=1}^{k} (j f j \le i \le N) j = 1, ..., k$$

$$f = f(t : y + h \sum_{j=1}^{k} (j f j \ge i \le N) j = 1, ..., k$$

where  $\beta$ ,  $\alpha$  are given in (4.11) (Ascher-Weiss [4]).

Not every RK scheme is equivalent to a collocation scheme.

But, among the most accurate RK schemes (using Gauss, Radau,

and Lobatto points), the most important are in fact equivalent

to collocation schemes.

When k=1, 
$$\rho_1 = 1/2$$
, then by (4.10) and (4.11)
$$\rho_1 = \phi_1(1) = 1$$

$$\alpha_1 = \phi_1(2) = 1/2$$

$$\alpha_{11} = \phi_1(2) = 1/2$$

so (4.12) gives the midpoint rule, and the equivalence has been shown in Section 4.1.1.

when 
$$k=2$$
,  $herefore = 0$ ,

The equivalent trapezoidal rule was treated in Section 4.1.1 as

well.

## 4.4. Collocation and Multiple Shooting

We now consider the collocation using monomial basis functions for solving the differential equation

(4. 13a) Ly(t); =y (t) - 
$$\sum_{j=0}^{n-1} \sigma_j$$
 (t) y (t) =q(t) a \le t \le b

with the separated boundary conditions

(4.13b) 
$$\sum_{l=1}^{n} b y \quad (a) = 0 , \qquad \sum_{l=1}^{n} c y \quad (b) = 0$$

$$1 \le j \le n/2.$$

We consider collocation at Gaussian points with  $s(t) \in P_{K,\Pi,m}$  where  $k \ge 2\pi$ , and we assume that the order  $\pi$  of the DE is even.

In general, the spline solution is determined by two types of constraints, continuity conditions and discretization equations (collocation equations and BCs).

On [t ,t ], the local monomial basis considered has the i i+1

form

$$\left\{\begin{array}{c} \frac{j+1}{j+1} \\ \frac{j}{(j-1)!} \end{array}\right\} = 1$$

The collocation approximation s(t) can be written as

(4.14) 
$$s(t) = \sum_{j=1}^{n} z_{j-1} + \sum_{j=1}^{n} w_{j-1} + \sum_{j=1$$

where 
$$z = (z_1, ..., z_n) = (s(t_1), ..., s_n(t_n))$$
.

The scaling in the first sum and h in the second sum are only introduced for later notational convenience. Now, both continuity conditions and discretization conditions must be satisfied. For the continuity conditions

$$\frac{(r-1)}{s} + \frac{(r-1)}{(t)} = \frac{1}{s} + \frac{(r-1)}{i+1}$$
 1\left\{\frac{1}{2}\text{\$\frac{1}\text{\$\frac{1}{2}\text{\$\frac{1}{2}\text{\$\frac{1}{2}\text{\$\frac{1}{2}\text{\$\frac{1}{2}\text{\$\frac{1}\text{\$\frac{1}\text{\$\frac{1}\text{\$\frac{1}{2}\text{\$\frac{1}\text{\$\frac{1}\text{\$\

$$(4.15) z = 8 z + 0 w$$

$$(4.15) \gamma_{3+1} i \gamma_{3} i \gamma_{3}$$

B = (B ) is an m x m upper triangular matrix with i rj.
entries

and D = (D ) is an m x (k-m) matrix with entries i ri

(4.17) i 
$$m+1-r$$
  
D = h /( $m+j-r$ )! =0 (h) .

The collocation conditions in (t,t) give i + 1

(4.18) Ls(t) = h 
$$\sum_{i \neq j} \frac{ij}{(n+j-1)!} L[\frac{(t-t)^{n+j-1}}{h}]$$

where to are the collocation points in (t ,t ) for ir i+1

15r5k-m.

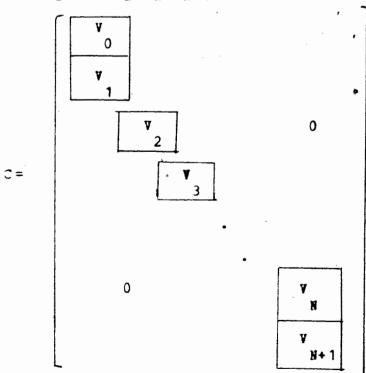
Writing (4.18) in matrix form, we have

where 
$$q = (q(t), \dots, q(t))$$
 $\sim i$ 
 $i,k-m$ 

and  $G = (G^{-})$  is a  $(k-m) \times (k-m)$  matrix with entries i rj

(4.21) 
$$G = \frac{(r)}{rj} - \sum_{i=1}^{n} G(t)h \frac{(r)}{(n+j-1)!}.$$

Thus the collocation matrix C corresponding to the unknowns



The blocks  $V_i$  (1 $\leq i \leq B$ ) of the matrix consisting of collocation and continuity equations are  $k \times (k+m)$  and have the structure

$$(4.22) \qquad \forall = \begin{bmatrix} \frac{1}{1} & \frac{1}{1} & \frac{1}{1} \\ \frac{1}{1} & \frac{1}{1} & \frac{1}{1} \end{bmatrix} k - a$$

where I is the m x m identity matrix.

and substitute this into the continuity equation (4.15) to obtain

$$(4.23) z = \begin{bmatrix} z + g \\ i & i \end{bmatrix}$$

where

The coefficient matrix for  $\{z_i\}$  corresponding to the BC i=1

(4.13) and (4.23) has the form

$$\mathbf{c}' = \begin{bmatrix} \mathbf{v} & \mathbf{I} & \mathbf{I} & \mathbf{I} \\ -\mathbf{v} & \mathbf{I} & \mathbf{I} \\ -\mathbf{v} & \mathbf{I} & \mathbf{I} \end{bmatrix}$$

Let {\vec{n} \{t;t\)} be the set of linearly independent j i j=1

solutions of Ly=0 subject to the initial conditions

(1-1)   
 
$$(t;t)=\delta$$
 ,  $1\le j,l\le n$ , and let y be the solution to  $j$  i i  $jl$ 

Ly=3 with z(y(t)) given, where  $z(y(t)):=(y(t),y^*(t),...)$ (s-1) T

,y (t )) . Then we have

(4.25) 
$$z(y(t)) = f(t), t)z(y(t))$$
  
 $i+1$   $i+1$   $i$ 

where % (t;t) is the fundamental matrix % (t;t) = (z(% (t;t) i i % 1 i

)),...z = (1 + (t;t)). If s(t) is the collocation approximation

of y(t), then

$$\left\| z(s(t_{i+1})) - z(y(t_{i+1})) \right\| = 0(h_{i}) - c(h_{i})$$

From the analysis of collocation,  $z(s(t)) = \int_{i}^{\infty} z(s(t))$ 

and 
$$z(y(t)) = \sum_{i=1}^{n} z(y(t)) + O(h)$$
. From (4.25), since

$$z(y(t))$$
 was arbitrary,  $\int_{1}^{2} = \pi(t, t) + O(h)$ , i.e.

is an approximation to the fundamental matrix

multiple shooting matrix (Ascher-Russell [3]).

As stated before, after doing condensation, the collocation methods give multiple shooting like matrices which have the advantage of being well-conditioned.

#### 4.5. Multiple Shooting and the Box Scheme

We now consider the relationship between multiple shooting and the Box scheme. Suppose the mesh is uniform, and h is small. The multiple shooting matrix (3.17) gives

OF

(4.26) -Y (t ) c +Y (t ) c = 
$$\forall$$
 (t ) - $\forall$  (t ) i i+1 ^i i+1 \( \times i \) i+1 ^i i+1 \( \times i \) i+1 \( \times i \) 1 \( \times i \) 1 \( \times i \)

If we set Y (t ) =I-(hA )/2 , 
$$1 \le i \le J-1$$
,  $i = i = 1/2$ 

then using a Taylor expansion and the fact that Y'=AY,

$$7 (t) = Y (t) + hY'(t) + O(h^{2})$$

$$= I - (hA) / 2 + hAY (t) + O(h^{2})$$

$$= I - (hA) / 2 + hAY (t) + O(h^{2})$$

Similarly, if v (t ) is a particular solution satsfying Fi i

Hence (4.26), gives

$$(4.27) = -(I + (h + 1)/2) c + (I - (h + 1)/2) c = hf + 0(h^2) .$$

$$i + 1/2 = i + 1/2 = i + 1/2 = i + 1/2$$

Since y is the multiple shooting approximation solution,

$$y := Y (t) c + v (t) = (I + (hA))/2)c$$
  
 $i i i i i i i$   $i - 1/2$   $v i$ 

and c = y + 0 (h), (4.27) gives  $\sim i \sim i$ 

v i+1/2

$$-(I+(h\lambda))/2)y+(I-(h\lambda))/2)y=hf+0(h)$$
  
 $i+1/2$   $\sim i$   $i+1/2$   $\sim i+1/2$ 

which can be written as

$$\frac{y - y}{-i + 1 - i} = \lambda \qquad \frac{y + y}{(-i + 1 - i)} + f \qquad +0 (1)$$

$$= \lambda \qquad y \qquad +f \qquad +0 (1)$$

$$= \lambda \qquad y + f \qquad +0 (1)$$

$$= \lambda \qquad 1 + 1/2 - i + 1/2 - i + 1/2$$

Thus, the Box scheme gives a matrix which is a discrete approximation to a multiple shooting matrix.

#### 4.6. Invariant Imbedding and Multiple Shooting

equivalence between invariant imbedding and the box scheme in the sense that a specific algorithm for solving the difference equations is valid if and only if an appropriate imbedding is valid. But the equivalence was not obvious. Recently, Lentini-Osborne-Russell [15] presented an easier way of getting a close relationship between multiple shooting and invariant imbedding. The Keller-Lentini [14] result is a special case of their presentation. When solving a BVP with separated BC, the relationship between factorizations of the multiple shooting matrix and invariant imbedding formulations of the BVP are shown in [15]. The equivalence is described below.

Consider the problem (3.24). In (3.17), converting each block [-Y]  $(t_{i+1})$   $P_{i+1}$  ] to  $[-P_{i+1}^{-1}]$   $Y_i(t_{i+1})$  I gives M similar to the "standard" multiple shooting matrix. Therefore, it is sufficient to consider only the "standard" multiple shooting matrix. Suppose the fundamental solutions  $Y_i$  (t) and particular solutions  $\psi_i$  (t) at  $t_{i+1}$  can be partitioned as

$$\frac{p}{q} = \begin{bmatrix}
0 & 1 \\
Y & Y \\
i & i
\end{bmatrix} p \\
\frac{2}{i} & \frac{3}{i} \\
\frac{7}{i} & \frac{7}{i}
\end{bmatrix} q , v(t) = \begin{bmatrix}
1 \\
v \\
i
\end{bmatrix} p \\
\frac{2}{v} \\
\frac{7}{i}
\end{bmatrix} q$$

For convenience, arrange the multiple shooting matrix so that it

(#.29) Ty=b.

Now, consider the factorization of (4.28). If  $K_0$  is nonsingular, multiply the first block by  $K_0^{-1}$  and perform the first step of elimination, then we have

Let 
$$V(t) = \begin{bmatrix} 1 \\ V(t) \\ 1 \\ 3 \\ V(t) \end{bmatrix} = V(t) \begin{bmatrix} -1 \\ k & k \\ 0 & 1 \\ -1 \end{bmatrix}$$
 be the column of

complementary function solutions and z = 0

The second step of elimination gives

$$\begin{bmatrix}
I & -P(t) & & & \\
P & 1 & & \\
3 & & \\
0 & V(t) & 0 & I \\
& & 1 & 2 & q \\
& & & & \\
\end{bmatrix}
\begin{bmatrix}
u(t) & x(t) \\
& 1 & \\
& v(t) & \\
& & & \\
& & & \\
\end{bmatrix}$$

$$\begin{bmatrix}
x(t) \\
x(t) \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\$$

where from section 3.1.5, R(t) = v(t)v(t)is the Riccati

matrix satisfying (3.27a) and x(t) satisfies (3.26). The next step gives

where

$$\begin{pmatrix}
4.30 & 1 & 1 & 1 \\
-Y & -Y & E & (t) & 2 & 2 \\
2 & 2 & 2 & 2 \\
-Y & -Y & E & (t) & 2 & 3 \\
-Y & -Y & E & (t) & 2 & 3 & -1
\end{pmatrix} = Y (t) \begin{pmatrix} -R & (t) & 2 & 2 \\ 2 & 3 & 2 & 3 \\ -Y & 2 & 2 & 2 & 2 \end{pmatrix}$$

$$= \underbrace{Y}_{2} \underbrace{(t)}_{3} \underbrace{Y}_{1} \underbrace{(t)}_{2} \underbrace{\begin{bmatrix} -1 \\ K & K \\ 0 & 1 \\ -1 \end{bmatrix}}_{1} \underbrace{\begin{bmatrix} 3 & 2 \\ (Y + Y & R(t)) \\ 1 & 1 \end{bmatrix}}_{1} -1$$

$$= \underbrace{\mathbf{Y}}_{1} \underbrace{(t)}_{1} \underbrace{\begin{pmatrix} -1 \\ K & K \\ 0 & 1 \\ -1 \end{pmatrix}}_{1} \underbrace{\begin{pmatrix} 3 & 2 & -1 \\ (\mathbf{Y} + \mathbf{Y} & R & (t)) \\ 1 & 1 & 1 \end{pmatrix}}_{1}$$

$$\begin{bmatrix} -I \\ 3 & 2 & -1 \\ 1 & 3 & 1 & 1 & 1 & 2 & 3 \\ 1 & 3 & 1 & 1 & 1 & 2 & 3 \\ \end{bmatrix} = v (t) = \begin{bmatrix} 1 \\ v (t) \\ 2 & 3 \\ \hline v (t) \\ 2 & 3 \end{bmatrix}.$$

Since  $\begin{bmatrix} K & K \\ 0 & 1 \end{bmatrix} \stackrel{\forall}{=} (t) = 0$ ,  $\stackrel{\forall}{=} (t)$  is a transformed set of

complementary function solutions; therefore, R(t) = v (t) v (3 2 3 2

t) - The process is continued until the stage that

 $\mathbf{v}$  (t ) is singular or partial pivoting is performed on the -i i+1

purpose, introduce a permutation M' of  $\widetilde{M}$  which preserves the zero structure in (4.28) and leaves the differential equation invariant (i.e. if y'=Hy is modified by the change of variables z=2y, then  $z^*=\widehat{H}z$  where  $\widehat{H}=PHP^{-1}$ ). Then P,Q in  $M^*=P\widetilde{M}Q$  must have the form

where  $T_p$ ,  $T_q$ ,  $T_N$  are identity matrices of size  $p \times p$ ,  $q \times q$ , and  $N \times N$ , respectively, and  $q_i = p_i$  are  $N \times N$  permutation matrices which satisfy

(see Keller-Lentini[13], Lentini-Osborne-Russell[15]). Using the same argument in (4.30), the theorem follows: 

Thm.4.1 Suppose that the BVP (3.25) has a unique solution. 
Assume  $K_0$  is nonsingular. Suppose further that the LU factorization for PHQ exists where

Here the N x N matrices  $p_{j}$ ,  $q_{j}$  occur (n-i) times and satisfy (4.31). Then the factorization reduces (4.28) to the block upper triangular form

(4.33) x 1 N+1 K-K R (t, 3 2

Here 
$$(\widetilde{\underline{u}}(t)) := q (\underline{u}(t)) \hat{v}(t) = q Y(t) \begin{bmatrix} R(t) \\ \widetilde{\underline{v}}(t) \end{bmatrix}$$
,  $\widetilde{\underline{v}}(t) = q Y(t) \begin{bmatrix} P(t) \\ I \end{bmatrix}$ ,  $\widetilde{\underline{v}}(t) = q Y(t) \begin{bmatrix} P(t) \\ I \end{bmatrix}$  (1 $\leq j \leq i$ ), where  $\widetilde{\underline{R}}(t)$  and  $\widetilde{\underline{x}}(t)$  are the solutions to the invariant imbedding equations corresponding to (3.27a,t) for the variables  $\widetilde{\underline{v}}(t)$ . Proof: see Lentini-Osborne-Pussell [15].

The factorization of  $\overline{\mathbf{M}}$  to (4.33), when P=Q=I, can be interpreted as a forward elimination corresponding to finding solutions  $\overline{\mathbf{M}}$ (t) and  $\underline{\mathbf{M}}$ (t) to the IVP (3.27a), (3.27b). The back surstitution on (4.33) to find  $\underline{\mathbf{M}}$  (j=N,...,0) then corresponds to finding the solution to (3.27c). For  $\mathbf{q}$   $\neq$ I, it corresponds to changing the invariant imbedding formulation at  $\mathbf{t}$ = $\mathbf{t}_{i+1}$ . When P,Q have any number of adjacent blocks of permutation matrices which are different, the matrix factorization corresponds to finding solutions for different imbedding formulations (Lentini-Osborne-Russell [15]). The equivalence between multiple shooting and invariant imbedding is therefore shown, and the meller-Lentini [14] result concerning the equivalence between invariant imbedding and the Box scheme follows using basically the same argument as in this section.

# 5. Finite Differences for Solving High Order Differential Equations

In this chapter, the construction of finite difference methods which give high-accuracy approximations to the solution of a high order linear differential equation My=f subject to linear BOs is investigated.

Define a mesh  $\pi$ :  $a=t_0 < t_1 < \dots < t_T = b$  and for any function w(t) on  $\pi$ , let  $w_j = w(t_j)$ . At mesh points,  $u_j$  is the estimate of y and u satisfies  $H_h u = \widehat{f}$ , together with appropriate BCs, where  $H_h u$  is a linear combination of values of u at stencil points (adjacent mesh points) and  $\widehat{f}$  is a linear combination of values of f at f auxiliary points close to the stencil points.

The construction is based on a local collocation procedure with polynomials, which is equivalent to the method of undetermined coefficients.

In section 1, the description of the discretization approximation is presented in the first part and some examples are given. In the second part the description of the finite difference approximations to boundary conditions follows. In section 2, the order of the truncation error is given when the location of the auxiliary points is independent of M. Stability of the scheme is discussed in section 3. And in section 4, methods with higher order of accuracy obtained by Doedel [9] and by Lynch-Rice [16] using special auxiliary points are presented

and followed by examples. An obvious equivalence between the high order finite difference methods and collocation methods is shown in section 5. Section 6 contains a comparison of the computational effort of Doedel's schemes and Lynch-Rice's schemes. This was not done in either [9] or [16]. Comparison with collocation methods is performed. Numerical results for Doedel's methods, Lynch-Rice's methods, and collocation methods are provided for comparison.

# 5.1. Construction of the High Order Pinite Difference Approximation

In this section, the construction of the nigh order finite inference methods is presented. We consider the interior sutintervals first.

## 5.1.1. The Approximations for Interior Subintervals

consider the ath order linear differential equation

(5.1) 
$$\forall y (t) = y \quad (t) + \sum_{k=0}^{n-1} a_k(t) y \quad (t) = f(t) \quad a \le t \le b$$

and the mesh a=t < t < ... < t = b.

Let  $u_j$  be the approximate solution of (5.1) at  $t=t_j$ , and in the subinterval  $\{t_{j-r_j}, t_{j+s_j}\}$  where  $r_j$  and  $s_j$  are positive constants. Let the difference operator  $H_h$  be

The right hand side of the approximation equation  $H_h u_j = \widetilde{f_j}$  is

$$\widetilde{f} = \sum_{j=1}^{n} e f(z)$$

where {d } and {e } are known coefficients and
j,i j,i

We will show how {i } and {e } are constructed below.

j,i j,i

For simplicity, we frequently omit the subscript j, e.g.  $z_{j,j}$ ,  $d_{j,j}$ ,  $e_{j,j}$ ,  $r_{j}$  and  $s_{j}$  become  $z_{j}$ ,  $d_{j}$ ,  $e_{j}$ ,  $r_{j}$  and  $s_{j}$ . Then the finite difference approximations to (5.1) at mesh points have the form

(5.2) 
$$\text{Mu} = \sum_{i=-r}^{s} du = \sum_{i=1}^{n} ef(z) = \widehat{f} \qquad r \le j \le J-s$$

where 
$$h = \max h$$
,  $h = t - t$ .

 $j \quad j \quad j = 1$ 

Since there are J-s-r+1 equations in (5.2) and n BCs, and the number of unknowns is J+1, one requires that r+s≥n and also incorporates more constraints if necessary.

The coefficients  $d_1$  and  $e_2$  are determined so that the approximation is exact on  $P_L$ , the space of all polynomials of degree at most L. i.e. if  $\mathbf{v}^1$  (t) (0  $\leq$  1  $\leq$  L) form a basis for  $P_L$  then  $d_1$ ,  $e_2$ , are made to satisfy the equations

(5.3) 
$$\sum_{i=-r}^{s} d u (t) - \sum_{i=1}^{n} e H u (z) = 0$$

$$1 = 0, \dots, L$$

The system (5.3) is homogeneous in  $d_{i}$ ,  $e_{i}$ . Therefore, in addition to (5.3), we take some convenient normalization equation such as one of

(5.4) 
$$e = 1$$
1

(5.4)  $\sum_{i} |e_{i}| = 1$ , or

(5.4)  $\sum_{i} |e_{i}| = 1$ 

To uniquely determine the r+s+m+1 unknowns d and e

from (5.3) and (5.4), I must be at least equal to r+s+m-1.

BCs for u are obtained in a similar way, they are treated in the next section.

Let p(t) be the polynomial in  $P_{r+S+m-1}$  which interpolates the solution y(t) and satisfies

write p(t) as a linear combination of the basis functions

(5.7) 
$$p(t) = \sum_{k=0}^{r+s+n-1} c w (t) .$$

Then (5.5), (5.6), and (5.7) give

$$r+s+m-1 \qquad k$$

$$\sum_{k=0}^{\infty} c w (t_{k-}) = u \qquad -r \le i \le s$$

with the operator # as in (5.1).

One way to find d; and e; for a general set of z; can be the following: Evaluate the determinant by expanding in terms of the last column and compare it with (5.2), introduce a normalizing factor E, then d; and e; are given by

where cof [.] is the cofactor of the given element in D

and a convenient normalizing factor E can be chosen as

(5.9) 
$$E = -\sum_{i=1}^{n} cof [f(z)]$$

with  $\frac{0}{M}y(t) = y$  (t) ...

If the w (t) are chosen so that

$$\frac{1}{t}(t) = \delta \qquad 0 \le 1 \le r + s$$
(5. 10)

then from (5.8) and (5.9), d and e can be calculated as

 $r+s+1 \qquad r+s+n-1$   $m+i+1 \qquad r+s+n-1$   $r+s+1 \qquad r+s+n-1$   $m+i+1 \qquad i-1 \qquad i-1$   $r+s+1 \qquad r+s+n-1$   $m+i+1 \qquad m+i+1$ 

15i Sm,

and

If n=1, then E=-1 and e=1.

In satisfy (5.10), w (t),  $0 \le l \le r + s + n - 1$ , can be chosen as

1≤1≤#-1

If r=0 and s=n, then (5.2) is the higher-order difference approximation with identity expansions (HODIE) considered by Lynch-Rice [16].

0 0 0 0

In the following, we use d ,e ,M ,f for the i i i

coefficients and the operators when M=D .

#### Example 5.1

\_ 2 1 0 Consider H=D +a (t) D+a (t) for 0≤t≤1.

Let r=s=m=1, suppose the mesh points are equally spaced

$$0 = \frac{(t-t)(t-t)}{2h^2}$$

$$\frac{1}{w(t)} = \frac{(t-t)(t-t)}{-h^2},$$

$$\frac{2}{v(t)} = \frac{(t-t)(t-t)}{\frac{1}{2h^2}},$$

so 
$$\frac{1}{1} = \frac{1}{2} \sqrt{(z_{1})} = \frac{1}{h^{2}-a} (t_{1})/2h$$
,

$$1 = 1 + (z) = -2/h^2 + a + (t),$$

$$\frac{2}{1} = 4 (z_{i}) = 1/h^{2} + a (t_{i})/2h ,$$

and (5.2) becomes

$$\frac{1}{j+1} + \frac{1}{j} + \frac{1}{j-1} + \frac{1}{j} + \frac{1}{j+1} + \frac{1}{j} + \frac{1}{j}$$

which is the usual divided difference approximation for

2 1 0 D +a (t) D+a (t).

#### Example 5.2

Now, consider M=D+a (t). Let r=0, s=1, m=2,

z = t, and z = t. Choose the basis functions 1 j 2 j+1

 $v^{2}$  (t) = (t-t) (t-t) /h<sup>2</sup>, j +1 j

so

$$\frac{1}{2} = \frac{1}{2} = \frac{1}$$

$$e = (MW^{2}(z))/E=1/2$$
,

then (5.2) becomes

$$(\dot{u} - u)/h + (a \hat{u} + a u)/2 = (f + f)/2$$
  
 $j+1$   $j$   $j$   $j+1$   $j+1$ 

which is the trapozoidal rule.

If m=1 and z=t is taken, then 1 j+1/2

$$d = Mw (z) = 1/h + a /2 .$$
1 1 j j+1/2

The difference approximation is therefore

$$(u -u)/h +a$$
  $(u +u)/2=f$   
 $j+1$   $j$   $j+1/2$   $j$   $j+1$   $j+1/2$ 

waich is the Box scheme.

5.1.2. The Approximation of Initial Conditions and Boundary Conditions

Approximation to initial or boundary conditions are required to complete the difference scheme. Here, only separated BCs are considered. Initial conditions can be considered as a special case of separated BCs.

Suppose two BCs for (5.1) are given by

(k) 
$$k-1$$
  
1 1 (i)  
(5.15a)  $3 \gamma(a) = \gamma(a) + \sum_{i=0}^{\infty} b \gamma(a) = \alpha k < n$ ,

(5.15b) 
$$B_{b} y (b) = y (b) + \sum_{i=0}^{2} b_{i} y (b) = \beta k < n.$$

The construction of a finite difference approximation to (5.15a) resembles that for (5.1) so we have

(5.16) 
$$B \quad u = \sum_{h=0}^{S} du = e \alpha + \sum_{i=1}^{m} e f(z) .$$

If the basis functions w (t)  $(0 \le 1 \le s + m)$  of P s + m

are chosen as in (5.10) for r=0, then the coefficients d and e can be calculated as

ງ≤i≤s

1≤i≤a

where 3 can be chosen as

unless a=0, in which case, set E=1.

The BC (5.15b) is treated similarly. The finite difference approximation to (5.15b) is

(5.21) B 
$$u = \sum_{h, p \in J} du = e \beta + \sum_{i=1}^{n} e f(z)$$

and the coeficients d ,e are i i

```
(5.22)
                        -r≤i≤0
(5.23)
(5.24)
                                                 (z )
i+1
```

and E can be chosen as

# Example 5.3

Suppose the BCs for Example 5.1 are given by

(5.26a) 
$$B_0 y(0) = y'(0) + b_0 y(0) = 0$$

(5.26b) 
$$B_1 y (1) = y^* (1) + b_1 y (1) = \beta$$

Let n=1 and s=1 for (5.26a) and

$$w = (t-t)/h$$
,

 $t = (t-t)/h$ ,

 $t = (t-t)/h$ ,

 $t = (t-t)/h$ ,

 $t = (t-t)/h^2$ 

so that

$$E = 8 v (z) = 2/h^2$$

(5.27b) 
$$e = -B \cdot w^2 \cdot (0) / E = h/2$$

The approximation to (5.26a) is

If r=1 and n=1 at the right end point, the approximation to (5.26b) is similarly given by

=1/h+a (2)+a (2)h/2

Letting

by (5.22), (5.23), (5.24), and (5.25)

$$E = H \quad \forall \quad (z) = 2/h^2$$

### 5.2 The Order of Consistency

The local truncation error of (5.2) is defined as  $(j=H_h, y_j - f_j)$  where  $y_j = y(t_j)$ , y(t) is the exact solution of (5.1) subject to appropriate BCs. If (j-1), as (j-1), then the finite difference approximation to (5.1) is said to be consistent. If there is a constant c and p is the largest integer such that

$$|7| \le ch^p$$
, as  $h_j \rightarrow 0$ 

then the difference approximation is said to be consistent of

order p.

Let y(t) be the exact solution of (5.1) subject to appropriate initial or boundary conditions and assume that y(t) is unique. Taking a Taylor expansion of y(t) at  $b_j$ , we get

for some b & [t ,t ], g & [t ,z ]
i j j+i i j i

Since (5.2) is exact on P and (t-t)  $\in$  P  $0 \le k \le L$ , the

quantity in square brackets is zero for each k. If we make the assumption that  $h/c\le h \le h$  ( $1\le j\le J$ ) for some c (such j

a family of meshes is called quasiumiform), then from
the calculation of d , e in (5.11), (5.12), and
i i

(5.13) with the basis functions w (t) in (5.14a), (5.14b),

it follows that there are constants c ,c ,and c which are independent of h such that

and 
$$\begin{vmatrix} d & \leq c & h & -r \leq i \leq s \\ e & \leq c & 1 \leq i \leq s \end{vmatrix}$$

$$\begin{vmatrix} e & \leq c & 1 \leq i \leq s \\ 1 & \leq c & 1 \leq i \leq s \end{vmatrix}$$

$$\begin{vmatrix} h & (z-t) & \leq c & h & z \leq z \leq z \\ 1 & 1 & 1 & 1 \end{vmatrix}$$

Recalling that L is at least r+s+m-1, we have the following theorem.

factor 2 be given by (5.11), (5.12), (5.13), respectively.

Assume that 870, and r+s2n. If the mesh is quasiuniform and h is small enough then at least n+1 of the d; are monzero and the order of consistency of the finite difference approximation (5.2) is greater than or equal to r+s+m-m (Doedel [8],[9] and Lynch-Rice [16]).

#### Brample 5.4

In Example 5.1, r=1, s=1, n=1, and n=2; therefore, the order of the scheme is at least r+s+n-n=1+1+1-2=1.

In fact, it is of order 2 as we will see in Section 5.4.

Example 5.5

In Example 5.2, the order of the scheme is 0+1+2-1=2, which is the order of the trapezoidal rule.

Now, consider the approximation for BCs. The truncation error of (5.16) is

$$7 = B \quad y(a) - e d - \sum_{i=1}^{n} e f(z_i)$$

where z aré collocation points in [t ,t ].

The truncation error of (5.21) is

$$7 = B$$
  $y(b) - e$   $3 - \sum_{i=1}^{R} e f(z)$ 

where z are the collocation points in [t ,t].

i J-r J

We have the following similar result to Thm.5.1.

Thu.5.2 Assume that  $E_0 \neq 0$  in (5.20) ((5.25)). If n=0 let  $n \geq k_1$  ( $r \geq k_2$ ). If n > 0, let  $n \geq k_2$  ( $n \geq k_3$ ). Then the order of consistency of the finite difference approximation to BCs (5.16) ((5.21)) is at least equal to  $n \geq k_3 + 1$  ( $n \geq k_3 + 1$ ) (Doedel [8][9]).

It follows that the order of consistency of (5.28), (5.29) in Example 5.3 is at least 1+1-1+1=2.

## 5.3. Stability of the Schenes

It is mentioned in Section 2.1.4 that one should not use numerically unstable methods. Once consistency of a new numerical method is established, for convergency it is necessary to establish stability. In order to guarantee convergence, we will now examine the stability of the schemes described, and in order to apply Kreiss! theory [14], we will assume that the mesh is uniform.

Consider the BVP (5.1) and the BCs

(5.30b) B (b) y (b) = 
$$\sum_{i=0}^{k} b_{i}$$
 (i)  $\sum_{k=0}^{k} b_{i}$  (b)  $\sum_{k=0}^{k} b_{i}$  (b)  $\sum_{k=0}^{k} b_{i}$  (b)  $\sum_{k=0}^{k} b_{i}$  (c)  $\sum_{k=0}^{k} b_{i}$  (d)  $\sum_{k=0}^{k} b_{i}$  (d)  $\sum_{k=0}^{k} b_{i}$  (d)  $\sum_{k=0}^{k} b_{i}$  (d)  $\sum_{k=0}^{k} b_{i}$  (f)  $\sum_{k=0$ 

(5.31) Lu = 
$$\sum_{i=1}^{s} d_i u = \sum_{i=1}^{s} e_i f(z_i) = f_i r \le j \le J - s$$
  
h j i=-r j, i j+i i=1 j, i j, i j

(5.32a) B (a) 
$$u = \sum_{k=0}^{k} d$$
 (a)  $u = e$  (a) b (a) h, k 0 i=0 k, i i k, 0 k

n (a)  
k  
+ 
$$\sum_{j=1}^{n} e_{j}(a) f(z_{j}(a)) = b_{j}(a)$$
  $1 \le k \le n$ 

(5.32b) B (b) 
$$u = \sum_{k=-r}^{0} d$$
 (b)  $u = e$  (b) b (b) h, k J  $i = -r$  (b) k, i , J+i k, 0 k

\*\* (b) k + 
$$\sum_{j=1}^{n} e_{j}(b) f(z_{k,j}(b)) = b'(b) n + 1 \le k \le n.$$

If the approximation is not compact, then r+s-n extra difference equations are required to match the number of equations and the number of unknowns. Suppose these extra equations involve the differential equation and are given by

(5.33a) Lu = 
$$\sum_{k=1}^{\infty} d_k u = \sum_{k=1}^{\infty} e_k f(z_k) = f(x_k)$$

$$k \le j \le r-1 \quad k \ge 0$$

$$0$$

(5.33b) Lu = 
$$\sum_{k=1}^{\infty} \frac{k}{u} = \sum_{k=1}^{\infty} \frac{k}{(x^{2})} = \widehat{f}$$
h j i=-r j, i J+i k=1 j, k j, k j

J-s+1≤j≤J-n+k

Letting u = (u, ..., u), then (5.31), (5.32), and (5.33)

can be expressed as

where f is the appropriate (J+1)-vector, and L is a h

(J+1)x(J+1) matrix.

Let 
$$7 = (7 (a), ..., 7 (a), 7 (..., 7 , 7 (b), ..., 7 (b), ...,$$

be the vector of the truncation errors with

$$7 (a) = B (a) y - b (a) 1 \le k \le n$$
  
 $k \quad h, k \quad 0 \quad k \quad 0$ 

$$7 = L y - \hat{f} \qquad k \le j \le J - n + k$$

$$j \qquad h \qquad j \qquad 0 \qquad 0$$

and 
$$7(b) = B(b) y - b(b) n + 1 \le k \le n$$
  
 $k \quad h, k \quad J \quad k \quad 0$ 

The finite difference scheme (5.34) is said to be stable if for all sufficiently small h, L exists and satisfies h

L Sc for some constant c independent of h.

Here for any (J+1)-vector  $g=(g,\ldots,g)$ , and for any  $(J+1)\times(J+1)$  matrix  $\lambda$ ,

5 ac

Let R denote the translation operator, i.e. R u =u . i i+1

Define D = (R-I)/h, and let Int u be the polynomial of

degree≤i which interpolates u at u ,u ,...u .

0 1 i

For a compact schese, we have

Thm.5.3 (Kreiss [14]) Assume the homogeneous problem corresponding to (5.1) and (5.30) only has the trival solution. If r+s=m, then there exists a constant k such

that for all solutions of (5.31) and (5.32) an a priori estimate

$$\left\|\begin{array}{c} \mathbf{n} \\ \mathbf{D} \mathbf{u} \right\| \leq \mathbf{K} \left( \left\|\mathbf{u}\right\| + \left\|\widetilde{\mathbf{f}}\right\| + \sum_{k=0}^{n-1} \left\|\mathbf{b}\right\| \right)$$

holds: (Here we define f;=0 for i=0,1,...r-1 and i=J-s+1,...
.,J). If the Equations (5.31) and (5.32) are consistent,
then these equations have, for every f and b and all
sufficiently small h, a unique solution u, and there is a
constant k such that

$$\|\mathbf{u}\| \leq k \left(\left\|\widetilde{\mathbf{f}}\right\| + \sum_{k=0}^{n-1} \left\|\widetilde{\mathbf{b}}_k\right\|\right)$$
.

Furthermore, the interpolated function Int a converges to

the solution y of the differential equation. i.e.

$$\begin{array}{c|c} 1in & | nt u-y | = 0 \\ h \to 0 & n \end{array}$$

Ĺ

We now consider the case when r+s>n. For later purpose,

write (5.31) in the form

(5.31°) 
$$L u = S(h) D u + \sum_{k=0}^{n-1} q D u = \hat{f}$$
  
h j + j-r k=0 k + j-r j

where 5(h) denotes a uniformly bounded difference operator of the form

$$S(h) = \sum_{k=0}^{r+s-h} S(h) R$$

and q are linear combinations of d k , j,i

For example, in example 5.1, S(h) = 1 + a h/2,

Let x = D , j = 0, 1, ..., J-n, then (5.31) can be written as j + j

(5.35) 
$$S(h) x = g j = r,...J - s$$
.

In most applications (5.35) has constant coefficients, i.e. we can write it as

where c' are constants. In this case, define the

characteristic polynomial c(t) associated with (5.31) where

If (5.34) is not compact, then we need r+s-n characteristic polynomials associated with the extra BCs (5.33). Let them have the form

(5.36) 
$$c$$
 (t) =  $\sum_{i=0}^{j} c$  t  $k \le j \le r-1, J-s+1 \le j \le J-n+k$   $i = 0$   $k \ge 0$ .

Also consider the homogeneous difference equations

(5.37a) 
$$\sum_{i=0}^{N} c v = 0 j=r,r+1,...J,J+1,...$$

with BCs

an d

(5.38) 
$$\sum_{i=0}^{N} c v = 0 j=J-s, J-s-1,..., 0,-1,...$$

with BCs

(5.38b) 
$$\sum_{i=0}^{j} c \quad \forall \quad =0 \quad J-s+1 \le j \le J-n+k$$

$$i=0 \quad j, \forall \quad =1 \quad 0$$

$$sup \quad |\forall \quad | \le constant \quad .$$

$$s \le j < po \quad |J-j|$$

Then Kreiss [14] has also shown:

Thm.5.4 Suppose the homogeneous problem corresponding to (5.1) and (5.30) only has the trivial solution. Assume the difference scheme (5.34) is consistent and all roots to of the characteristic equation c(t)=0 satisfy |t; | > 1. If the difference scheme is not compact, also suppose that the difference equations (5.37a,b) and (5.38a,b) have only the trivial

small h and the difference scheme is stable.

#### Example 5.6

Consider the problem

(5.39a) 
$$y'(t) = f(t)$$
  $0 \le t \le 1$   
(5.39b)  $y(0) = 0$   
Let  $r = n = s = 1$  and  $z = t$ 

then

hence S(h)=R and c(t)=t. Therefore, c(t)=0 has no root on the unit circle. Since c =0 and c =1, (5.37a) and (5.38a)

only have the trivial solution and by (5.33a) and w, w in example 5.2, d =-1/h, d =1/h, a =1.

0,0 0,1 0,0

Hence (5.37b) only has the trivial solution, and so does the homogeneous problem corresponding to (5.39). Since noncompact approximations to (5.39) with m=1 are always consistent (see Doedel 9), the scheme is stable.

Having some roots of c(t)=0 lie on the unit circle does not necessarily imply that the finite difference approximation is unstable. By numerical experience, it has been shown that such

approximations may lead to stable schemes (Doedel [9]).

# 5.4. Improved Order with Particular Choice of z

In (5.2), if the z are chosen properly, higher order

accuracy can be obtained. Doedel [9] only considered the choice of such z for which one higher order accuracy is

obtained. The details of this analysis are presented below.

0 1 r+s+n-1
Let w ,w ,...,w be a basis of P defined in r+s+n-1

(5.14a) and (5.14b). This linearly independent set can be extended to form a basis of P by adding a polynomial r+s+m

w (t) EP which vanishes at the mesh points. The

extra polynomial can be of the form

$$r+s+m$$
  $m-1$   $r+s$   
 $w$   $(t) = \prod_{k=1}^{m-1} (t-x) \prod_{k=0}^{m-1} (t-t)$   
 $k=1$   $k$   $k=0$   $j-r+k$ 

where x are in [t ,t ], and satisfy  $|x-t| \le ch$ , k j

 $1 \le k \le m-1$ . If we expand y(t) in terms of w (t), i.e.

$$y(t) = \sum_{k=0}^{r+s+n} c w (t) + O(h)$$

then the truncation error can be written in the form

$$-\sum_{i=1}^{n} \frac{k}{i} \frac{\sum_{j=0}^{n} s_{j} + n - n + 1}{i}$$

The quantity between square brackets vanishes for

 $0 \le k \le r+s+n-1$ , and since v (t )=0 (-r\leq i \leq s), 7 becomes j+i

Hence, it is clear that if the z are chosen so that

0 r+s+m r+s+m

H v (z)=0 1≤i≤m where v (t) is in P and
i r+s+m

r+s+m satisfies w (t)=0 0≤i≤s, then an extra order of

consistency can be obtained.

# Example 5.7

In Example 5.1, let

and z be the root of H w (t)=6(t-t) i.e. z =t,

1 j 1 j

which is what we chose for z in Example 5.1. Then, as stated in Example 5.4, the order of consistency of the scheme is 1+1=2.

However, if we treat the z; in (5.2) as unknowns, one could expect that for the special operator N<sub>h</sub>, higher orders up to r+s+2n-n can be achieved. It is shown in Lynch-Rice [16] that such z's exist and they also offered the special choice of z; for the case N=N°. This will be discussed next.

for N=D<sup>N</sup>=M<sup>0</sup>. For the general case of the variable coefficient operator N and sufficiently small h, it can be shown that there is a set of e's and a unique set of m auxiliary points z; with t<sub>j-r</sub> <z<sub>1</sub> <... <z<sub>m</sub> <t<sub>j+5</sub> such that the high order scheme is exact on P<sub>2m+n-1</sub> (see Lynch-Rice [16]). However, it is not clear how these z's can be found. Since their positions are problem dependent, it would not be practical to construct a high order scheme which is exact on P<sub>2m+n-1</sub> for a general N. In the following, a generalized result of Lynch-Rice [16] is shown. The basic process is the same as theirs, but r and s are no longer restricted to 0 and n, respectively. We now find the special location of z's which would give the order of consistency as high as possible.

On the jth subinterval, since the  $w^{1}(t)$  of (5.14a) are in  $P_{n}$ , their n-th derivatives are constants. When applying them to (5.3), we have

(5.40) 
$$\frac{0}{d} - \frac{n!}{s} \qquad \sum_{k=0}^{n} e^{-t} = 0.$$

$$i \qquad \uparrow \uparrow \qquad (t - t) \quad j=1 \quad j \quad k = -r, k \neq i \quad j+i \quad j+k$$

If (5.4c) is used, (5.40) gives

which means the operator  $H_h$  is n! times the usual divided difference approximation to  $H^0 = D^n$ . Thus

i.e.,  $H_h^o y(t_j)$  is the nth derivative of the unique polynomial in  $P_n$  which interpolates the values  $y(t_{j+i})$  at  $t_{j+i}$ ,  $i=r,\ldots,s$ .

By Taylor's Theorem, y(t) can be expressed as

(5.42) 
$$y(t) = \sum_{i=0}^{n-1} \frac{i}{j-r} + \int_{1-r}^{t} \frac{t}{(n-1)!} dx$$

(Goldberg[11]).

Substituting (5.42) into (5.41), since the nth divided difference of an element of  $P_{n-1}$  is zero, we get

where B (t;x) is the nth divided difference

g[t ,...,t ;x] with respect to t of n j-r j+s

$$g_{n}(t;x) = (t-x) = \begin{cases} (t-x) & n-1 \\ (t-x) & /(n-1)! & \text{if } t>x \\ 0 & \text{otherwise} \end{cases}$$

Hence B (t ;.) is the (n-1)st degree polynomial B-spline n j

with joints at the stencil points. Therefore, the

truncation error is

(5.44) 
$$7 = H y(t) - \hat{f}$$

where E D y is the quadrature error in using f as an h approximation to the intergral of n!B (t;x)D y. Let

If e are chosen such that

(5.46) 
$$e = n!$$
  $\begin{cases} t & i-1 \\ & B & (\bar{t}; x) \neq (x) dx & i=1,...,n \end{cases}$ ,  $j-r$ 

then since 
$$\sum_{i=1}^{n} v_{i}(t) = 1$$
 and  $\int_{t}^{t} t dx = 1/n!$ ,

we obtain  $\sum_{i=1}^{n} e^{-i}$ . But for the e's in (5.46), and any

[t ,t ], define the inner product j-r j+s

(5.47) 
$$(u,v) = \begin{cases} t & - \\ j+s & - \\ t & n \\ j-r \end{cases}$$

Let b, b, ... with b eP; be the normalized orthogonal

polynomials with respect to this inner product (call them the B-spline orthogonal polynomials). Each  $b_i$  has i distinct real zeros in  $(t_{j+r},t_{j+s})$  (call them B-spline Gauss points). If the maxiliary points  $z_i$  are the B-spline Gauss points for b, since Gauss quadrature is exact on  $P_{2m-1}$  for Gauss points, the high order finite difference approximation is exact on  $P_{2m+n-1}$ .

A sequence of orthonormal polynomials can be generated by a 3-term recurrence relation as follows:

$$\widetilde{b} (t) = 0,$$

$$-1$$

$$\widetilde{b} (t) = 1,$$

$$0$$

$$(5.48)$$

$$\widetilde{b} (t) = (t-B)\widetilde{b} (t) - c\widetilde{b} (t) \quad i=0,1,2,...$$

$$i+1 \quad i \quad i \quad i-1$$

where 
$$B = \begin{cases} t \\ j+s & 2 \\ xb & (x)B & (\overline{t};x)dx/s \end{cases}$$
,
$$t & i & n & j & i \\ j-r & & & i & n & j & i \end{cases}$$

$$c = \begin{cases} 0 & \text{if } i=0 \\ s / s & \text{if } i \neq 0 \\ i & i-1 \end{cases}$$

$$s = \begin{cases} t & \text{j+s } ^2 \\ & \text{b } (x) B (\overline{t}; x) dx \\ t & \text{i } n \text{ j} \end{cases}$$

and b (t) 
$$= \widetilde{b}$$
 (t) /s i=0,1,2,...

Since the z are roots of b (t)=0, it is obvious that

$$\sum_{i=1}^{n} eb(z)=0$$

and because (b ,b ) is positive, 7 in (5.44) is not zero

for polynomials in P, that is, the approximation is not 2m+n

exact on P . Hence the high order scheme has order at 2m+n

most 2m. If only j of the m Gauss points are used, then the high order scheme is of order m j.

#### Example 5.8

Consider the operator H=D . Let m=2,r=s=1.

Por convenience, consider a uniform mesh. Then

$$B (\bar{t}; x) = 
 \begin{cases}
 (x-t)/2h^2 & t < x < t \\
 j-1 & j-1 & j
 \end{cases}$$

$$(t -x)/2h^2 & t < x < t \\
 j+1 & j & j+1, \\
 0 & otherwise$$

and. (5.45) gives

From (5.46), we get

and by (5.48), we have

Hence the scheme

$$(u -2u +u )/h^2 = f(t +h/\sqrt{6})$$
  $j=1,2,...,J-1$   
 $j-1$   $j$   $j+1$   $j$ 

for y = f(t) (with appropriate BCs) is of order 4.

# 5.5. An Equivalence Between Pinite Difference and Collocation Methods

In this section, an equivalence between the finite difference methods and collocation methods is presented.

Consider the case n=1. Let  $v^1$  (t) in (5.10) be the set of basis functions for our collocation method. Then the collocation solution s(t) is

$$s(t) = \sum_{l=0}^{r+s} c \cdot v \cdot (t) \cdot t$$

$$1 = 0 \cdot 1$$
Since  $v \cdot (t) = \delta \quad 0 \le 1, k \le r+s, j=0,...,J$ ,
$$j-r+k \quad 1k$$

s(t )=c. If we write the collocation solution at j-r+1, ⇒1; j-r+l j-r+l  $s(t) = \sum_{i=0}^{r+s} v(t)u = i-r+1$ j-r+1 i=-r The collocation method for one collocation point and w (t) in (5.10) requires u to satisfies \_\_\_\_**j\*i** \_\_\_\_\_ (t) u = f(z)j+i On the other hand, from (5.11) -r≤i≤s The left sum in (5.2) is  $\sum$  du =  $\sum$  Hw (z) u i=-r i j+i i=-r 1 j+i and right sum in (5.2) gives -**乏** e f(z)=f(z) By (5.50) and (5.51) we have

 $\sum_{i=-r}^{\infty} Hv (z)u = f(z)$ 

which from (5,49) is collocation with one collocation point.

# 5.6. Work Estimates

In this section, the computational aspects of the higher order finite difference schemes are considered. Since the HODIE methods of Lynch-Rice [16] are similar to the above methods (call them Doedel's methods), comparison of the computational work for the same global accuracy is made. Comparison of efficiency with that of collocation methods using B-spline and Gaussian points is also considered.

find the coefficients by solving the linear algebraic system (5.3) instead of calculating the determinents in (5.11) and (5.12). This is how the HODIE methods find the coefficients. If the basis functions are chosen such that they satisfy (5.10), then (5.3) gives

(5.52) 
$$d = \sum_{i=1}^{n} e \, \text{Mw} \, (z) \quad -r \le i \le s$$

(5.53) 
$$\sum_{j=1}^{n} e^{j} + e^{j} = 0$$
  $r+s+1 \le 1 \le r+s+n-1$ 

for the interior subintervals.

If the normalization equation is chosen such that e=1, then (5.53) becomes

$$(5.54)$$
 Ae=-b

After the e are determined from (5.54), d can be found easily from (5.52), and the finite difference approximations at mesh points are the solutions u of (5.2).

Lynch-Rice [16] found that to obtain e's it is computationally more efficient to use a different set of basis functions ([16],p.363). As a result, (5.52) is no longer valid for the basis they considered, and so, one has to solve a system of r+s+1 algebraic equations to find the d;. Their set of basis functions is therefore less efficient for evaluating d;, and so, only basis functions satisfying (5.10) are considered here.

Consider the problem (5.1) subject to (5.30a,b) at a uniform partition  $t_k = kh$ , k = 0, ..., J, and a general set of m auxiliary points. Let MT be a multiplication/division time and F be a function evaluation time. Suppose that the values  $v^{1}(z_{i})$ ,0≤i≤n,0≤l≤r+s+m-1, 1≤j≤m, have been previously computed and stored (they do not depend on the subintervals), then the setup. time for A and b in (5.54) is ((n-1)(n-1)n+n(n-1)) ff+nmf. The function evaluation is for a general set of z: . If z; are at mesh points for some i, function values at these points can be stored beforehand and do not need to be recomputed. If (5.54) is solved by Gaussian elimination without pivoting, the solution time is  $[(m^2+m-3)(m-1)/3]$ HT. In (5.52), for a fixed i, it takes null time to evaluate  $Hv^{r+1}(z_i)$  ( since the  $a_k(z_i)$  have already been calculated in (5.54)). Since e = 1, it takes another (m-1) NT time to find  $d_{\frac{1}{2}}$  . The total time required to evaluate  $d_{\frac{1}{2}}$  ,  $-r \leq i \leq s$ and  $e_i$  ,  $1 \le i \le n$ , is therefore

#### $[nn(n-1)+(n^2+n-3)(n-1)/3+(r+s+1)(nn+n-1)]$

For simplicity, we only consider the case where there are the same number of auxiliary points in the consecutive subintervals involved in each row of (5.34). To get an order of r+s+m-n for a complete scheme, in (5.32a) (or (5.32b)), one has to take i from 0 to at least v=r+s-n+n<sub>K</sub>(a)-1 (or -v'=-(r+s-n+n<sub>K</sub>(b)-1 to 0) for each k. v (or v') is nonpositive only when r+s=n and n<sub>K</sub>(a)=0,1 (or n<sub>K</sub>(b)=0,1)-(recall that r+s≥n). For the case n<sub>K</sub>(a)=0 (or n<sub>K</sub>(b)=0), no approximation equation is needed since we have been given the exact solution at a (or b). If n<sub>K</sub>(a)=1 (or n<sub>K</sub>(b)=1), set v=1 (or v'=1) and pick m-1 auxiliary points in [t, t<sub>1</sub>] (or [t<sub>J-1</sub>, t<sub>J</sub>]). As in (5.52) and (5.53), e; in (5.19) can be calculated more efficiently by solving the system of mequations

(5.55) e B v (a) + 
$$\sum_{i=1}^{k} e^{ik}$$
 (z) =0, k= $v+1$ ,... $v+m$ , 0 k, a j=1 j j

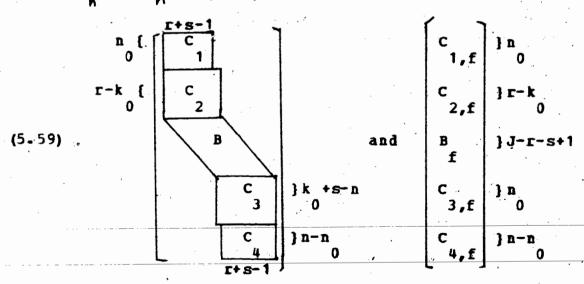
and d in (5.17) can be evaluated by

Then e in (5.24) can be calculated by solving

and d in (5.22) can be evaluated from

(5.58) 
$$d = e B \quad w \quad (b) + \sum_{i=1}^{r+i} e H w \quad (z) \quad -v' \le i \le 0.$$

Using the same argument as before, when  $n_{\kappa}(a) = 1$  (or  $n_{\kappa}(b)$ =1), there are  $2(n-1)+n(n-1)-(n-1)+(n-1)(n^2+n-3)/3+2(n+n)$ multiplications to approximate the k-th boundary condition. If v>0 ( $v^{+}>0$ ), there are  $(n_{V}(a)+1) = +na^{2}+n(a^{2}+3a-1)/3+(v+1)$   $(n_{V}(a)+1)$ +n+m) (or  $(n_{K}(b)+1) = +n = 2+n = (n^{2}+3n-1)/3+(v^{2}+1)^{2}(n_{K}(b)+n+m)$ ) multiplications to approximate k-th boundary condition. Denote the number of multiplications required for approximating a BC by MTK. The number of function evaluations for approximating BCs also depends on  $n_{K}(a)$  and  $n_{K}(b)$ , e.g. if all  $n_{K}(a)$  are the same and not equal to 0 or 1 for  $k=1,...n_0$ , then only nmP is needed for the first no BCs. Denote the number of function evaluations for approximating the first  $n_o$  BCs by  $P_d$  and the last  $n-n_o$  by  $P_h$ . In (5.33a) (or (5.33b)), i has to go from 0 to at least r+s (or from - (r+s) to 0) for each j. Since  $n_{\kappa}(a)$  (or  $n_{\kappa}(b)$ ) is less than n,  $r+s-n+n_{K}(a)-1$  (or  $r+s-n+n_{K}(b)-1$ ) is at most r+s-2. Hence,  $L_h$  and  $f_h$  in (5.34) have the form



where B is a trapezoid of width r+s+1 and height J-r-s+1. Depending upon  $n_{K}(a)$  and  $n_{K}(b)$ , some elements of  $C_{4}$  and  $C_{4}$  may

be zero. Note that one has to take different  $r-k_0$  sets of auxiliary points for  $C_2$  and different  $k_0+s-n$  sets of auxiliary points for  $C_3$ , otherwise, some rows of  $C_2$  or  $C_3$  will be identical and  $L_h$  will be singular. The setup time for  $L_h$  and  $f_h$  is therefore

C,C F+ 
$$\sum_{n=1}^{\infty}$$
 HT b k=n+1 k

where  $q=n\pi(m-1)+(m^2+m-3)(m-1)/3+(r+s+1)(nm+m-1)$ . If Gaussian elimination without pivoting is used to solve (5.34), then the solution time is

$$-1-n+k$$
)+(r+s+1-n)(2n+2r+2s-2n-1)+  $\sum_{j=1}^{n}$  j(j+2)]HT

for getting an upper triangular matrix and

r-k

u
[
$$\Sigma$$
 j+u(j-r-s)+  $\Sigma$  (u+j)+  $\Sigma$  (r+s-j) ]HT for back

j=1 j=1

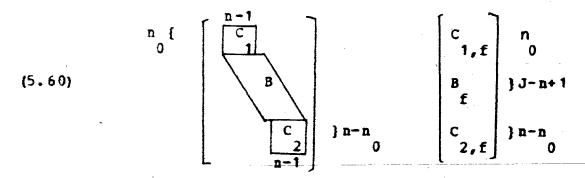
substitution, where u=s+1+k -n
0 0

Note that the work estimates calculated are based on an nax

(r+s-1) matrix  $C_1$  and an  $(n-n_0)$  x (r+s-1) matrix  $C_3$ . Since some elements of  $C_1$  and  $C_3$  may be zero, the actual work required should be no more than the work calculated from the formula.

Now consider the HODIE methods. Recall that in the HODIE methods r=0 and s=n. Hence, to compare with Doedel's method which gives the same global accuracy of order r+s+m-n, one must select p=r+s+m-n auxiliary points in each subinterval [t;,t;+n] for a HODIE method.

Again, let  $e_1 = 1$  and assume all values of  $w^2(z_j)^{(i)}$  are already stored. For approximating BCs, consider the case when p points are collocated. Then in (5.32a) (or (5.32b)), i goes from 0 to  $n_K(a) = 1$  (or  $n_K(b) = 1$ ). As before, if  $n_K(a) = 0$ , no approximation is needed. If  $n_K(a) = 1$ , then one finds  $d_1, d_1$  from  $e_0, e_1, \ldots, e_{p-1}$ . Let  $F_d^*$  and  $F_b^*$  be the function evaluation time for approximating BCs at a and b, respectively, and let  $HT_K^*$  be the multiplication time for each BC. Since the HODIE methods are compact, no extra equations are required.  $L_h$  and  $f_h$  of (5.34) then have the form



where B is a trapezoid of width n+1 and height J-n+1.
Hence, setup time for (5.34) is

If Gaussian elimination without pivoting is used to solve (5.34), since some elements of  $C_1$  and  $C_2$  may be zero, the total time for the second part of the implementation is at

most 
$$\begin{bmatrix} n & \sum_{j=0}^{n-1} (n-j) + n & (n+2-n) & (J-n+1-n) + (n+2-n) & (n+n-3) \\ 0 & j=0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

n-3
+ ∑ (n-j+2) (n-j) ]HT to get an upper triangular matrix and j=1

$$[(n-n)(n-n+1)/2 + (n-n+1)(J-n+1) + (2n-n-1)n/2]$$
HT to

obtain the solutions by back substitution.

The comparison below is done for second-order differential equations subject to Dirichlet BCs since they are the most important case and are simple. Consider the problem

$$Hy(t) = y''(t) + a(t)y''(t) + a'(t)y(t)$$
  $a \le t \le b$ 

$$y(a) = 0 = y(b)$$

For this case, n=2, n=1, n(a)=0=n(b), (L)=(L)

Table 5.1 for three different orders of accuracy (4,6, and 8).

The data for collocation using B-splines and Gauss points is derived from Russell-Varah [20].

Table 5.1

*	order					
r+s-n=	4	6	8			
0	(J-1) (82HT+12F)	(J-1) (186#T+18F)	(J-1) (354HT+24F)			
1	(J-1) (60HT+9F)	(J-1) (144NT+15F)	(J-1) (284HT+21F)			
2	(J-1) (42NT+6P)	(J-1) (110HT+12F)	(J-1) (226HT+18F)			
3	(J-1) (27HT+3P)	(J-1) (83HT+9P)	(J-1) (1798T+15P)			
4		(J-1) (60HT+6P)	(J-1) (140HT+12F)			
collo- cation	J (328T+4P)	J(828T+6P)	J (164HT+8P)			

The first row of Table 5.1 gives the computational work for the HODIE methods (Lynch-Rice [16]), and the first five rows give operation counts for five more general different Doedel's schemes. Lynch-Rice [16] picked t<sub>j+1</sub> ∈ [t<sub>j</sub>,t<sub>j+2</sub>] as one of the auxiliary points. Since the central mesh point of an odd-number-point difference operator is a zero of every odd-degree generalized B-spline orthogonal polynomial, one higher order of accuracy is obtained. In this case, one would expect that, for the same order of accuracy, the counts of the HODIE regular case in Lynch-Rice [16] is smaller than that of the HODIE methods we consider here. From Table 5.1, it is obvious that as r+s increases, the operation count for a given

order decreases. It seems on the basis of operation counts that one should pick as few auxiliary points as possible and let r+s be large. This view is supported in the numerical examples given in the next section.

If we compare the first row with the other rows of Table 5.1, Doedel's methods are more efficient than the HODIE methods for large r+s-n. By comparing the data in Table 5.1 with that of Table 9-2 in Lynch-Rice [16], we conclude that Doedel's methods are comparable even to the HODIE Gauss-type case when r+s-n is large enough. Collocation with B-splines is also competitive in work for second-order differential equations.

#### 5.7. Experimental Results

From Table 5.1, we have seen that if the number of auxiliary points is chosen as small as possible, Doedel's methods are much more efficient than the HODIE methods.

Numerical experiments have been run to support both the theorems in the previous sections and the above conclusion. All computations were performed on the SPU IBH 3033 using double precision arithmetic. In each experiment, the mesh considered is equal spaced. For the HODIE methods, auxiliary points are chosen such that  $z_{j,i} = t_{j+1} + (i-1)h/J$ . For Doedel's methods, only one auxiliary point is used, it being chosen as the midpoint of  $[t_{j+1}]$ ,  $j=0,\ldots,J-s$ . Note that the matrix  $L_h$  obtained when  $r=r_p$ ,  $s=r_0+s_0$ , is identical with the matrix  $L_h^s$  obtained when r=0,  $s=r_0+s_0$ ,

and  $z_i^*$  in  $[t_{j+r_0+s_0}]$  are at the same location corresponding to  $z_i$  in  $[t_{j-r_0}, t_{j+s_0}]$ .

Since we consider the case n=1, if the order of consistency is required to be greater than one for Doedel's methods then r+s must be greater than n, and hence r+s-n extra finite difference equations are needed. In the experiments, they are divided into two sets. If r+s-n is even, half of them are set to approximate the first s+1 solutions  $u_0, \ldots, u_s$ ; if r+s-n is odd, then (r+s-n+1)/2 equations are defined for  $u_0, \ldots, u_s$ . In both cases, the others are set for the last s+1 approximate solutions  $u_{J-S}$ , ...  $u_J$ . In  $\{t_0, t_5\}$ , the auxiliary point involved in the i-th equation of these extra equations is taken to be the (i+1)st mesh point. In  $\{t_{J-S}, t_s\}$ , the auxiliary point is defined to be the reflection of the corresponding one in  $\{t_0, t_s\}$ .

In the following tables, numerical results are shown for a number of cases. Orders of consistency considered are 2, 4, 6, and 8. The results using Doedel's methods with m=1 are the rows marked s=. The rows marked m= are the results using the HODIE methods. For collocation methods, numbers of collocation points used are 2,3, and 4 for order 4,6,and 8, respectively. The results which use collocation methods with B-splines, Gaussian points, and uniform meshes (COLLO) are given in the last columns. The notation .220-1 stands for .220 10 In the first column, (o) means the order of the schemes. J is the number of subintervals.

# Example 5.9

#### Consider the DE

$$y'' + y - 2y = 2(1 - 6t) e$$
  
 $y(0) = 0 = y(1)$ .

The solution of this problem is  $y(t)=2t(1-t)e^{-t}$ . The problem has been used by Doedel [8] [9]. The results are given in Table 5.2.

_		*	77533	718			
(0)		J=8	J=16	J=32	J=64	J=128	COLLO
	<b>n=</b> 2	. 220-1	.558-2	-140-2	.351-3	. 877-4	
2	s=3	. 137-1	. 348-2	. 876-3	-219-3	. 548-4	. •
4	R=4	.564-4	. 354-5	. 223-6	. 139-7	. 862-9	J=64
	s=5	-537-4	.623-5	. 440-6	-283-7	·179-8	.735-9
6	<b>m</b> =6	-608-7	. 958-9	. 158-10	xx	XX	J=16
	s=7	-826-6	. 123-7	- 187-9	.307-11	. 306-13	-800 <del>-</del> 11
8	<b>n</b> =8	. 398-10	XX	XX	XX	XX	J=8
	s=9		. 212-10	. 626-13	EE	II	.119-12

Table 5.2

xx Contaminated by roundoff.

#### Example 5.10

Let the DE be given by y\*-4y=4cosh(1)

y(0)=0=y(1).

The solution to this problem is y=cosh(2t-1)-cosh(1). The problem has been used by Lynch-Rice [16] (note that they had a mistake in [16], f(t)=4cosh(1), not f(t)=2cosh(1)). The results are given in the following table.

<u> Table' 5.3</u>

(0)		J=8	J=16	J=32	J=64	J= 128	COLLO
2	n=2	.795-2	. 198-2	- 496-3	-124-3	.310-4	Paking in the
	s=3	. 496-2	- 124-2	.310-3	-775-4	. 194-4	•
4	R=4	. 324-4	.203-5	. 127-6	.794-8	-491-9	J=64
	s=5	. 324-4	.359-5	-251-6	.162-7	. 102-8	.168-9
6	<b>n=</b> 6	.751-7	- 118-8	. 181-10	XX	II	J= 16
	s=7	-942-6	- 148-7	-229-9	-364-11	. 282-13	.360-11
8	n=8	.103-9	XX	XX	· XX	XX	J=8
	s=9		- 552-10	.144-12	XX	XX	.458-13

The last example shown is

# Example 5.11

The DE  $y^{n-400}y=400\cos^{2}\pi t+2\pi^{2}\cos^{2}\pi t$ y(0)=0=y(1)

 $-20 \ 20t \ -20 \ -20t \ -20$  has solution y=e e /(1+e )+e /(1+e )-cos<sup>2</sup> $\pi t$ .

The results are shown below.

Table 5.4

(0)		J=8	J=16	J=32	J=64	J=128	COLLO
2	<b>n=</b> 2	<b>-700+0</b>	. 149+0	.296-1	.665-2	. 156-2	
	<b>s=</b> 3	. 724-1	. 435-1	. 153-1	.370-2	.928-3	
	<b>n</b> =4		. 128-1	.723-3	-427-4	. 253-5	ÿ J=16
- 4	s=5	. 270-1	.740-2	-512-3	.522-4	. 431-5	.219-3
	<b>m</b> =6	. 265~1	.736-3	. 105-4	-502-6	-230-8	J=16
6	s=7	.208-1	. 364-2	. 137-3	. 191-5	-269-7	.182-5
8	<b>n=8</b>	.380-2	. 273-4	. 974-7	.356-,9	XX	J=16
	s≖9		. 173-2	. 274-4	.117-6	. 234-9	-891-8

From the results, the HODIE methods and Doedel's methods are quite competitive with each other. Comparing their results with that of COLLO, one requires larger values of J for the finite difference methods to achieve comparable accuracy. Taking the ratios of adjacent numbers in each row of the above tables, it is evident that both the HODIE methods and Doedel's methods give orders of consistency as predicted and there is no numerical instability. Since the implimentation of the HODIE methods involves solving an (m-1) x (m-1) matrix for each row of (5.34), the execution time of high order HODIE methods is much longer than of Doedel's methods for the same order of accuracy. From Table 5.1, we see that, when m=1, Doedel's methods are much cheaper than the HODIE methods.

#### 6. Conclusion

General forms of BVPs and TVPs have been given, existence and uniqueness theorems for solutions of IVPs and BVPs have been provided. Stability properties of IVPs and two point BVPs have been discussed. Some well-known numerical methods for solving BVPs have been presented.

numerical methods for solving differential equations. Some of them hold only in special circumstances. e.g. in Section 5.5, only one collocation point was considered, and there are many cases which have not been taken into account. Consequently, one remaining task would be to find more relationships between these methods (e.g., the HODIE methods and collocation when m>1)

We have discussed the high order finite difference methods thoroughly. As with finite element methods, when solved by the finite difference methods differential equations need not be converted to first order systems. Though one can get superconvergence using general B-spline and Gauss points, it is not practical for general n-th order differential equations. For the case H<sup>0</sup>=D<sup>n</sup>, one can find the B-spline Gauss points and the right-hand-side coefficients easily by using the formulae discussed before. While for general H, the location of the general B-spline Gauss points depends on H and the mesh points. Thus, it would be difficult to have a practical code for the

However, the methods have been shown to be computationally efficient. Operations counts and numerical results have shown that, to find the approximations by Doedel's methods more efficiently, one should use schemes which only involve one auxiliary point.

# References

- Ascher U. & Christiansen J. & Russell B.D., A Collocation Solver for Mixed Order Systems of Boundary Value Problems, Math. Comp., No. 33, 1979.
- 2. Ascher U. & Russell R.D., Evaluation of B-splines for Solving Systems of BVPs, Technical Report 77-14, Nov., 1977.
- 3. Ascher U. & Russell R.D., Numerical Solutions of Boundary Value Problems, manuscript , 1981.
- 4. Ascher U. & Weiss R., Collocation for Singular Perturbation Problems I: Pirst Order Systems with Constant Coefficients, To appear in SINUM.
- Bailey P.B. & Shampine L.F. & Waltman P.E., Monlinear Two Points Boundary Value Problems, Academic Press, New York, 1966.
- 6. DE Boor C.R., On Calculating with B-splines, J. of Approximation Theory v.6, pp50-62, 1972.
- 7. DE Boor C.R., Package for Calculating with B-splines, SIAM J. Numer. Anal. Vol 14, No. 3, June 1977.
- 8. Doedel E.J., The Construction of Finite Difference Approximations to Ordinary Differential Equations, SIAN J. Humer. Anal. Vol. 15, No. 3, June, 1978.
- 9. Doedel E.J., Difference Methods for Ordinary Differential Equations with Applocations to Parabolic Equations, Ph.D. Thesis UBC Jan., 1976.
- Pranklin J. M., Matrix Theory, Prentice-Hall Inc. Englewood Cliffs, N.J., 1968.
- 11. Goldberg R.R., Methods of Real Analysis, Xerox College Publishing, 1964.

- 12. Keller H.B., Numerical Methods for Two Points Boundary Value Problems, Blaisdell Pub. Co., Waltham, Mass., 1968.
- 13. Keller H.B. & Lentini H., Invariant Imbedding, the Box Scheme and an equivalence Between them, SINUM No. 19, 1982.
- 14. Kreiss H-O, Difference Approximations for Boundary and Eigenvalue Problems for Ordinary Differential Equations, Math. Comp. No. 26, v. 19, July, 1972.
- 15. Lentini H. & Osborne H. & Russell, R.D., The Close Relationships Between Methods for Solving ODE BVPs, manuscript, Jan. 1983.
- 16. Lynch R.B. & Rice J.R., A High-Order Difference Method for Differential Equations, Nath. Comp. V.34, No. 150, pp 333-372, 1980.
- 17. Hattheij R.H.H., The Stability of LU-Decompositions of Block Tridiagonal Matrices, Manuscript, Mathematisch Institut, Katholicke Universiteit, Mijmegan, The Metherlands, 1982.
- 18. Ortega J.H., Humerical Analysis, 2nd ed., Academic Press, N.Y. & London, 1972.
- 19. Reid W.T., Riccati Differential Equations, Academic Press N.Y. & London, 1972.
- 20. Russell R.D. & Varah J.H., A Comparison of Global Methods for Linear Two-Point Boundary Value Problems, Math. Comp. v.29, No.132, 1975.
- 21. Scott H.R. & Watts H.A., Computational Solution of Linear Two-Point BVPs via Orthonormalization, SINUM v. 14, No. 1, 1977.
- 22. Varah J.H., A Comparison of Some Humerical Methods for Two Point Boundary Value Problems, Math. Comp. v. 28, No. 127, July, 1974.
- 23. Scott H.R. & Watts H.A., SUPORT- A Computer Code for Two-Point BVPs via Orthonormalization, SAND75-0198 Sandia

Laboratories, Albuquerque, New Mexico, 1975.