

1989

# Boundary Approximation Methods For Some Free And Moving Boundary Problems

David Gerald Meredith

Follow this and additional works at: <https://ir.lib.uwo.ca/digitizedtheses>

---

## Recommended Citation

Meredith, David Gerald, "Boundary Approximation Methods For Some Free And Moving Boundary Problems" (1989). *Digitized Theses*. 1845.

<https://ir.lib.uwo.ca/digitizedtheses/1845>

This Dissertation is brought to you for free and open access by the Digitized Special Collections at Scholarship@Western. It has been accepted for inclusion in Digitized Theses by an authorized administrator of Scholarship@Western. For more information, please contact [tadam@uwo.ca](mailto:tadam@uwo.ca), [wlsadmin@uwo.ca](mailto:wlsadmin@uwo.ca).



National Library  
of Canada

Bibliothèque nationale  
du Canada

Canadian Theses Service

Service des thèses canadiennes

Ottawa, Canada  
K1A 0N4

## NOTICE

The quality of this microform 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.

Reproduction in full or in part of this microform is governed by the Canadian Copyright Act, R.S.C. 1970, c. C-30, and subsequent amendments.

## AVIS

La qualité de cette microforme 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.

La reproduction, même partielle, de cette microforme est soumise à la Loi canadienne sur le droit d'auteur, SRC 1970, c. C-30, et ses amendements subséquents.

**BOUNDARY APPROXIMATION METHODS FOR SOME FREE  
AND MOVING BOUNDARY PROBLEMS**

by

**David G. Meredith**

**Department of Applied Mathematics**

**Submitted in partial fulfilment  
of the requirements for the degree of  
Doctor of Philosophy**

**Faculty of Graduate Studies  
The University of Western Ontario  
London, Ontario  
August 1989**

**© David G. Meredith 1989**



National Library  
of Canada

Bibliothèque nationale  
du Canada

Canadian Theses Service    Service des thèses canadiennes

Ottawa, Canada  
K1A 0N4

The author has granted an irrevocable non-exclusive licence allowing the National Library of Canada to reproduce, loan, distribute or sell copies of his/her thesis by any means and in any form or format, making this thesis available to interested persons.

The author retains ownership of the copyright in his/her thesis. Neither the thesis nor substantial extracts from it may be printed or otherwise reproduced without his/her permission.

L'auteur a accordé une licence irrévocable et non exclusive permettant à la Bibliothèque nationale du Canada de reproduire, prêter, distribuer ou vendre des copies de sa thèse de quelque manière et sous quelque forme que ce soit pour mettre des exemplaires de cette thèse à la disposition des personnes intéressées.

L'auteur conserve la propriété du droit d'auteur qui protège sa thèse. Ni la thèse ni des extraits substantiels de celle-ci ne doivent être imprimés ou autrement reproduits sans son autorisation.

ISBN 0-315-51757-3

Canada

## **ABSTRACT**

Numerical methods for a class of free and moving boundary problems are considered. The class involves the solution of Laplace's equation on a domain which is changing shape with time. The position of the boundary is described by an evolution equation. With the time fixed, a boundary approximation method is employed to solve the potential problem. The boundary location at the next time is determined from the evolution equation using standard techniques and the process is repeated.

Two boundary methods are examined. Both are characterized by representing the approximate solution of the potential problem as a series of known basis functions, chosen from a complete set of particular solutions to the Laplace equation. In the first approach, the parameters, to be determined from the boundary data, appear linearly in the trial solution. The basis functions are closely related to the well studied harmonic polynomials and this permits an extensive analysis of the linear method. In particular, convergence of the method is demonstrated and some estimates on the degree of convergence are derived. In the second approach, the parameters appear nonlinearly. This approach is new, but derives from classical results on complex rational approximation and may be interpreted as an acceleration of the convergence of the linear technique.

The linear method is applied to a number of electrochemical machining examples and performs well for relatively smooth boundaries. A nonlinear approach is tested on the inverse machining problem with excellent results.

Both the linear and nonlinear methods are applied to several challenging examples of Hele-Shaw flow. In all instances, the nonlinear scheme outperforms the linear. The ease of programming, efficiency and concomitant accuracy of the nonlinear scheme make it an attractive choice for the numerical integration of a class of free and moving boundary problems.

## **ACKNOWLEDGEMENTS**

I would like to express my appreciation to Dr. H. Rasmussen for his patience and support over the course of this work.

I have had numerous fruitful discussions with Jake Greydanus and I thank him, in particular, for lending a helping hand with the use of the word processor.

I thank Michel Pettigrew for hours of enjoyable discussion on many topics and for sharing some of his boundless enthusiasm for the subject of mathematics.

A special thanks to Penny for help with many of the splendid figures in the thesis.

# TABLE OF CONTENTS

	Page
CERTIFICATE OF EXAMINATION	ii
ABSTRACT	iii
ACKNOWLEDGEMENTS	iv
TABLE OF CONTENTS	v
LIST OF TABLES	vi
LIST OF FIGURES	vii
LIST OF APPENDICES	x
CHAPTER 1 - INTRODUCTION	1
1.1 General Introduction	1
1.2 Approximation Methods	2
1.3 Objectives	3
1.4 Outline of Thesis	3
CHAPTER 2 - THE BOUNDARY APPROXIMATION METHOD	8
2.1 Introduction	8
2.2 Best Linear Approximation to Boundary Data	11
2.3 The Choice of Norm	15
2.4 Existence of a Best Approximation	17
2.5 Completeness of the Set of Basis Functions	18
2.6 Degree of Convergence	24
2.7 Characterization of, and Algorithms for, Computing the B.A.	30
2.8 Best Approximation by Nonlinear Functions	35
CHAPTER 3 - ELECTROCHEMICAL MACHINING PROBLEMS	43
3.1 Introduction	43
3.2 Review	47
3.3 Steady State ECM	49
3.4 Unsteady ECM	57
3.5 The Inverse ECM Problem	68
3.6 Discussion	77
CHAPTER 4 - HELE-SHAW FLOW	78
4.1 Introduction	78
4.2 Review	81
4.3 Cusping Solution: Linear Approximation	85
4.4 Saffman Finger: Linear Approximation	112
4.5 Hele-Shaw Flows: Nonlinear Approximation	125
4.6 Stability	153
CHAPTER 5 - CONCLUDING REMARKS	156
APPENDIX A.1 THE LOGARITHMIC TERM	158
APPENDIX A.2 THE LEAST SQUARES ALGORITHM	161
APPENDIX A.3 EXACT SOLUTIONS OF SOME HELE-SHAW PROBLEMS	164
REFERENCES	167
VITA	171

## LIST OF TABLES

Table	Description	Page
3.3.1	The Ratio $m:n$	53
3.3.2	Convergence for Increasing $n$	54
3.3.3	Effects of Proximity of Singularity	55
3.3.4	The Coefficients ( $n=30, p=15$ )	57
3.5.1	The Parameters ( $n=9$ )	73
3.5.2	Errors in Computed Cathode Coordinates	74
4.3.1	The Ratio $m:n$	93
4.3.2	Convergence	95
4.3.3	Convergence	96
4.3.4	Collocation	97
4.3.5	Comparison of Exact and Computed Values of $g(0,t)$	105
4.4.1	Comparison of Cusping and Saffman Solutions at $t=0.0$	115
4.4.2	Convergence	116
4.4.3	Convergence	117
4.4.4	Comparison of Exact and Computed Values of $g(0,t)$	122
4.5.1	Convergence - Cusping Profile	126
4.5.1(b)	The Nonlinear Parameters	127
4.5.2	Convergence - Saffman Profile	127
4.5.2(b)	The Nonlinear Parameters	128
4.5.3	Comparison of Exact and Computed Values of $g(0,t)$	137
4.5.4	Comparison of Exact and Computed Values of $g(0,t)$	145



## LIST OF FIGURES

Figure	Description	Page
2.1.1	Definition Sketch - Physical Plane	9
2.2.1	Transformed Plane ( $\zeta = \text{EXP}(-iz)$ )	12
2.5.1	Simply Connected Case	20
2.5.2	Doubly Connected Case	21
2.6.1	Semi-infinite Domain - Physical Plane	25
3.3.1	Definition Sketch	44
3.4.1	Time Evolution of ECM ( $\alpha = 1, a = 0.25$ )	61
3.4.2	Time Evolution of ECM ( $\alpha = 2, a = 0.25$ )	62
3.4.3	Time Evolution of ECM ( $\alpha = 3, a = 0.25$ )	63
3.4.4	Time Evolution of ECM ( $\alpha = 2.5, a = 0.5$ )	64
3.4.5	Time Evolution of ECM ( $\alpha = 2, h(x) = \text{EXP}(-4 \sin^2(x/2))$ )	67
3.5.1	Inverse ECM ( $\alpha = 3.0, a = 0.71$ )	75
3.5.2	Field Lines ( $\alpha = 3.0, a = 0.71$ )	76
4.3.1	Definition Sketch	86
4.3.2	Exact Cusping Behaviour	90
4.3.3(a)	Cusping Profiles Using a Linear Eulerian Approximation ( $n=15, m=60$ )	99
4.3.3(b)	Comparison of Exact and Computed Profiles ( $n=15, m=60$ )	100
4.3.4(a)	Cusping Profiles Using a Linear Eulerian Approximation ( $n=30, m=60$ )	102
4.3.4(b)	Comparison of Exact and Computed Profiles ( $n=30, m=60$ )	103
4.3.5	Error Growth ( $n=15, 30, 45$ )	104

	<b>Page</b>
4.3.6(a) Cusping Profiles Using a Linear Lagrangian Approximation (n=30,m=60)	107
4.3.6(b) Comparison of Exact and Computed Profiles (n=30,m=60)	108
4.3.7 Reduced $h_{max}$ $t=0.31$	109
4.3.8 Smoothed Cusping Profiles	111
4.4.1 Exact Saffman Finger	113
4.4.2(a) Saffman Profiles Using a Linear Eulerian Approximation (n=10,m=20)	118
4.4.2(b) Comparison of Exact and Computed Profiles (n=10,m=20)	119
4.4.3(a) Saffman Profiles Using a Linear Eulerian Approximation (n=20,m=40)	120
4.4.3(b) Comparison of Exact and Computed Profiles (n=20,m=40)	121
4.4.4 Error Growth (n=10,20,25)	123
4.5.1(a) Location of the Singularities in the Nonlinear Approximation of the Cusping Problem - $t=0.25$	129
4.5.1(b) Location of the Singularities in the Nonlinear Approximation of the Saffman Finger - $t=0.50$	130
4.5.2(a) Cusping Profiles Using a Nonlinear Lagrangian Approximation (n=3,m=29)	131
4.5.2(b) Comparison of Exact and Computed Profiles (n=3,m=29)	132
4.5.3(a) Cusping Profiles Using a Nonlinear Lagrangian Approximation (n=5,m=15)	133
4.5.3(b) Comparison of Exact and Computed Profiles (n=5,m=15)	134
4.5.4 Error Growth (n=3,4,5)	135
4.5.5(a) Saffman Profiles Using a Nonlinear Lagrangian Approximation	138

	<b>Page</b>
4.5.5(b) Comparison of Exact and Computed Profiles (n=3,m=38)	139
4.5.6(a) Saffman Profiles Using a Nonlinear Lagrangian Approximation (n=5,m=38)	140
4.5.6(b) Comparison of Exact and Computed Profiles (n=5,m=38)	141
4.5.7 Error Growth (n=3,4,5)	142
4.5.8 The Particle Trajectories Between t=0.0 and t=1.6	144
4.5.9 Local Error t=0.00(0.20)2.00	146
4.5.10(a) Saffman Profiles Using a Nonlinear Lagrangian Approximation (n=5,m=38)	147
4.5.10(b) Comparison of Exact and Computed Profiles (n=5,m=38)	148
4.5.11 Saffman Profiles Using a Nonlinear Eulerian Approximation (n=5,m=40)	149
4.5.12 Error Growth - Cusping Case	151
4.5.13 Error Growth - Saffman Finger	152

## LIST OF APPENDICES

Appendix		Page
Appendix A.1	The Logarithmic Term	158
Appendix A.2	The Least Squares Algorithm	161
Appendix A.3	Exact Solutions of Some Hele-Shaw Problems	164

The author of this thesis has granted The University of Western Ontario a non-exclusive license to reproduce and distribute copies of this thesis to users of Western Libraries. Copyright remains with the author.

Electronic theses and dissertations available in The University of Western Ontario's institutional repository (Scholarship@Western) are solely for the purpose of private study and research. They may not be copied or reproduced, except as permitted by copyright laws, without written authority of the copyright owner. Any commercial use or publication is strictly prohibited.

The original copyright license attesting to these terms and signed by the author of this thesis may be found in the original print version of the thesis, held by Western Libraries.

The thesis approval page signed by the examining committee may also be found in the original print version of the thesis held in Western Libraries.

Please contact Western Libraries for further information:

E-mail: [libadmin@uwo.ca](mailto:libadmin@uwo.ca)

Telephone: (519) 661-2111 Ext. 84796

Web site: <http://www.lib.uwo.ca/>

# **CHAPTER 1**

## **Introduction**

### **1.1 General Introduction**

Free and moving boundary problems comprise a wide class of mathematical models in the applied sciences. In broad terms, a boundary value problem is free if it involves the solution of a partial differential equation on some domain for which at least a portion of the boundary is unknown and is to be determined as part of the solution. If, in addition, the model is a time dependent one, with the boundary continuously changing, the problem is called a moving boundary problem, or MBP for brevity.

This work is concerned with the numerical integration of a subclass of free and moving boundary problems. Numerical schemes will be presented and applied to several interesting examples from the class. Each member of this class possesses essentially the same mathematical description. The field is governed by a potential equation and the evolution of the free boundary is described by a nonlinear partial differential equation.

The list of MBP whose field equation is the Laplace equation is long and includes: Hele-Shaw flows, flows in a porous medium, electrochemical machining and electroforming, injection moulding of plastics, Rayleigh-Taylor and Kelvin-Helmholtz instabilities, and nonlinear water waves.

In general, analytic solutions to MBP are rarely available and numerical schemes usually difficult to employ owing to the arbitrary and ever-changing nature of the boundary. As far as the numerical integration is concerned, it will be assumed that we may proceed in a step-wise manner, first calculating the solution to a potential problem for a given time and on a known region and then advancing the boundary position according to the evolution equation<sup>1</sup>. The process is applied repeatedly. In this prescription, the evolution equation usually reduces to the solution of a system of ordinary differential equations. The bulk of the effort then, resides with the manner in which the potential portion of the problem is treated. It is the method of solution adopted for this portion which characterizes the overall scheme and is central to this work.

## 1.2 Approximation Methods

In the past many different numerical techniques have been implemented according to the above step-by-step format. The potential problem has been solved by finite differences, finite elements and various boundary approximation techniques. The latter methods have not figured prominently in the class of MBP, with the exception of the boundary element method, which has seen extensive use.

Approximation methods involve the approximation of the true solution  $\phi$  by a function  $\phi_n$  which depends on a finite number of suitably chosen parameters. It is the choice of function  $\phi_n$  and the manner in which the parameters are determined which properly defines the scheme. Thus we may choose

$$\phi_n = \phi_n(b_1, b_2, \dots, b_n; \mathbf{x})$$

such that, for arbitrary choice of parameters  $b_j$ , one of the following options holds:

---

<sup>1</sup> A notable exception to this division of labour arises in variational inequality formulations. See Elliott and Ockendon (1982), for example.

(i)  $\phi_n$  satisfies the differential equation exactly

(ii)  $\phi_n$  satisfies the boundary conditions exactly

(iii)  $\phi_n$  satisfies neither the differential equation nor the boundary conditions exactly.

The  $b_j$  are then chosen so that  $\phi_n$  approximates, in some norm, the boundary conditions or the solution of the differential equation or both.

The first option is called a boundary approximation method and is the choice of scheme to be employed in this work. Options (ii) and (iii) are referred to as interior and mixed methods, respectively. (See Collatz, L. (1960) or the more recent work of Gottlieb, D. and Orszag, S. (1977) for example.) The latter authors refer to such schemes under the general heading of spectral methods, wherein they define a spectral method to be any scheme whereby the approximate solution to the boundary value problem is represented as a linear combination of known basis functions of the independent variables only. The definition only encompasses linear approximating methods and ignores the rather important area of nonlinear approximation. What is more, the term spectral method usually implies application to regular geometries. The basis functions are then chosen from a standard orthonormal set. Finally, these same techniques have often been referred to as methods of weighted residuals (see, for example Ames (1965)). We shall stick to the more general term, approximation methods.

In any event, each of the approximations (i), (ii) or (iii) provides a closed form expression representing the solution at all points of the given domain.

Each of the boundary, interior and mixed methods is properly characterized by formulating them as problems in functional approximation. For example, let  $D \subset R^2$  represent the domain with boundary  $\partial D$ . Let  $X(D)$  and  $X(\partial D)$  be two normed linear spaces



of functions defined on each of  $D$  and  $\partial D$ . Let  $L$  be an elliptic operator and  $B$  a given differential operator. With given functions  $l \in X(D)$  and  $b \in X(\partial D)$  we have the following boundary value problem:

$$L\phi(\mathbf{x}) = l(\mathbf{x}) \quad , \quad \mathbf{x} \in D \quad (1.2.1)$$

$$B\phi(\mathbf{x}) = b(\mathbf{x}) \quad , \quad \mathbf{x} \in \partial D \quad (1.2.2)$$

Define  $X(\bar{D})$ ,  $\bar{D} = D \cup \partial D$  to be the normed linear space of functions  $\phi$  defined on  $D \cup \partial D$  such that  $L\phi \in X(D)$  and  $B\phi \in X(\partial D)$ . The type of approximation problem is characterized by selecting a manageable subset  $M \subset X(\bar{D})$  and choosing functions  $\phi_n$  from this subset which minimize the error functional

$$E(\phi_n) = \alpha \|L\phi_n - l\|_D + \beta \|B\phi_n - b\|_{\partial D} \quad (1.2.3)$$

$\alpha, \beta \geq 0$ , where  $\|\cdot\|_D$  and  $\|\cdot\|_{\partial D}$  are suitable norms for the linear spaces  $X(D)$  and  $X(\partial D)$ .

Boundary methods correspond to  $\alpha = 0$ , interior methods to  $\beta = 0$  and mixed methods to both  $\alpha, \beta \neq 0$ .

Now, as already mentioned, the solution to the MBP will involve two steps:

- (a) the solution of a potential problem on a known but irregularly shaped domain,
- (b) the solution of the evolution equation for determination of the boundary position

at the next time step. The approximate solution of step (a) by way of a boundary method is most reasonable in light of the fact that

(1) at least a portion of the boundary is not likely to be coincident with a coordinate direction let alone coincide with a regularly spaced finite difference mesh, and

(2) the shape of the domain changes with time presenting a new potential problem on a new domain every time step.

### **1.3 Objectives**

This work began as an investigation of a numerical method used by Rienecker and Fenton (1981) and Fenton and Rienecker (1982) in their calculations of nonlinear water waves. Their scheme (essentially a linear boundary approximation method) is presented under the guise of a Fourier technique. Properties such as convergence are therefore accepted as a matter of course, relying on the extensive knowledge of the convergence properties of Fourier series. And their excellent computational results provide solid evidence of numerical convergence.

But there is much more that can be said if the method is properly formulated as a method in the theory of best approximations. This then, is the first goal of this work. Convergence of the sequence of best approximations is established and some error estimates are obtained.

Over the course of these investigations it was found that the practical applications of the linear boundary method were somewhat limited. This was largely due to the poor conditioning of the least squares matrices for the choice of basis functions used and is typical of this type of numerical scheme. To effectively deal with this, a related nonlinear approximation scheme was developed which accelerates the convergence the linear method. This then, is the second objective of this thesis, to numerically investigate the effectiveness of the nonlinear scheme and compare the performance of the two methods on practical problems.

### **1.4 Outline of Thesis**

In chapter 2, the approximation methods for the potential problem are formulated as extremal problems within the theory of best approximations. The convergence in norm of the linear boundary method is established and several error estimates are derived. Practical means for computing the best linear approximation are discussed. The nonlinear

approach does not lend itself readily to the same scrutiny and we suffice to establish the existence of a best approximation and to suggest a method for its determination. In both the linear and nonlinear cases, the possibility of approximation follows from classical results of Runge and Walsh.

Chapter 3 is intended to assess the performance of the linear method on a typical member of our class. Examples of both steady and unsteady electrochemical machining are considered. The treatment of the time-dependent portion of the moving boundary problem is discussed at this stage. The unsteady machining problem is a physically stable process. This means that any irregularities in the solution behaviour can be confidently attributed to the numerical scheme, thereby making this an attractive test choice.

On the other hand, the inverse free boundary problem of electrochemical machining is not well posed and can present subtle difficulties when its solution is attempted numerically. In chapter 3 we present a scheme which is accurate and widely applicable, provided a sufficient amount of data is available. In the event that this is not the case, a nonlinear approximation scheme has been proposed, the results of which appear quite promising.

In chapter 4, both a linear and nonlinear approximation method are applied to the unstable problems of Hele-Shaw flow. Such problems present a considerable challenge to any approximation method, as the physical instabilities of the flow pattern will feed on the accruing computational errors. This makes it difficult to isolate the origins of any irregularities that might arise in the solution behaviour. For example, is the observed phenomenon to be interpreted as normal physical development, or is it simply a manifestation of numerical instabilities inherent in the chosen approximation scheme?

Two Hele-Shaw examples are examined in this chapter. The two problems share almost identical initial conditions and yet the flow pattern which develops in each case is quite different. Thus, not only are the flows unstable, but the problems are ill-posed. This means that the computational errors committed enter the solution as small

perturbations and may well lead to a result quite different from the exact solution. Both the linear and nonlinear schemes are applied to each of these examples and both an Eulerian and a Lagrangian description of the moving surface are experimented with. Analytic solutions are available for comparison purposes.

We find that both numerical methods are able to simulate the flows for short periods of time, but only the nonlinear method is able to accurately follow the exact solution for large times.

All computations have been performed in single precision on the University of Western Ontario's CDC Cyber 170/835. The machine uses a sixty bit word together with a forty-eight bit mantissa, representing roughly fourteen digits of accuracy.

All plots of the free surface have been generated by straight line interpolation between data points.

## CHAPTER 2

### The Boundary Approximation Method

#### 2.1 Introduction

The class of MBP considered in this work is described by

$$\nabla^2 \phi(\mathbf{x}, t) = 0 \quad , \quad \mathbf{x} \in D(t) \subset R^3 \quad (2.1.1)$$

$$B\phi(\mathbf{x}, t) = b(\mathbf{x}) \quad , \quad \mathbf{x} \in \partial D(t) \quad (2.1.2)$$

$$f_t(\mathbf{x}, t) = -\nabla\phi(\mathbf{x}, t) \cdot \nabla f(\mathbf{x}, t) \quad , \quad \mathbf{x} \in \Gamma_1(t) \quad (2.1.3)$$

$$f(\mathbf{x}, 0) = f_0(\mathbf{x}) \quad , \quad \mathbf{x} \in \Gamma_1(0) \quad (2.1.4)$$

$\Gamma_1(t)$  is that portion of the boundary of  $D$  that is changing in time according to the evolution equation (2.1.3).  $f(\mathbf{x}, t) = 0$  describes the position of  $\Gamma_1(t)$ .

For the purposes of the present chapter, the domain in question is restricted to a two-dimensional subset of the complex  $z$ -plane ( $z = x + iy$ ). Furthermore, it is assumed that the solution  $\phi(x, y, t)$  is periodic in  $x$  of period  $2\pi$  and defined on the infinite periodic "strip"

$$h(x, t) \leq y \leq g(x, t) \quad , \quad -\infty < x < \infty$$

Both  $h$  and  $g$  are periodic functions of  $x$ .  $D(t)$  corresponds to one period of the strip, as shown in figure (2.1.1). That is, for given  $t$ ,

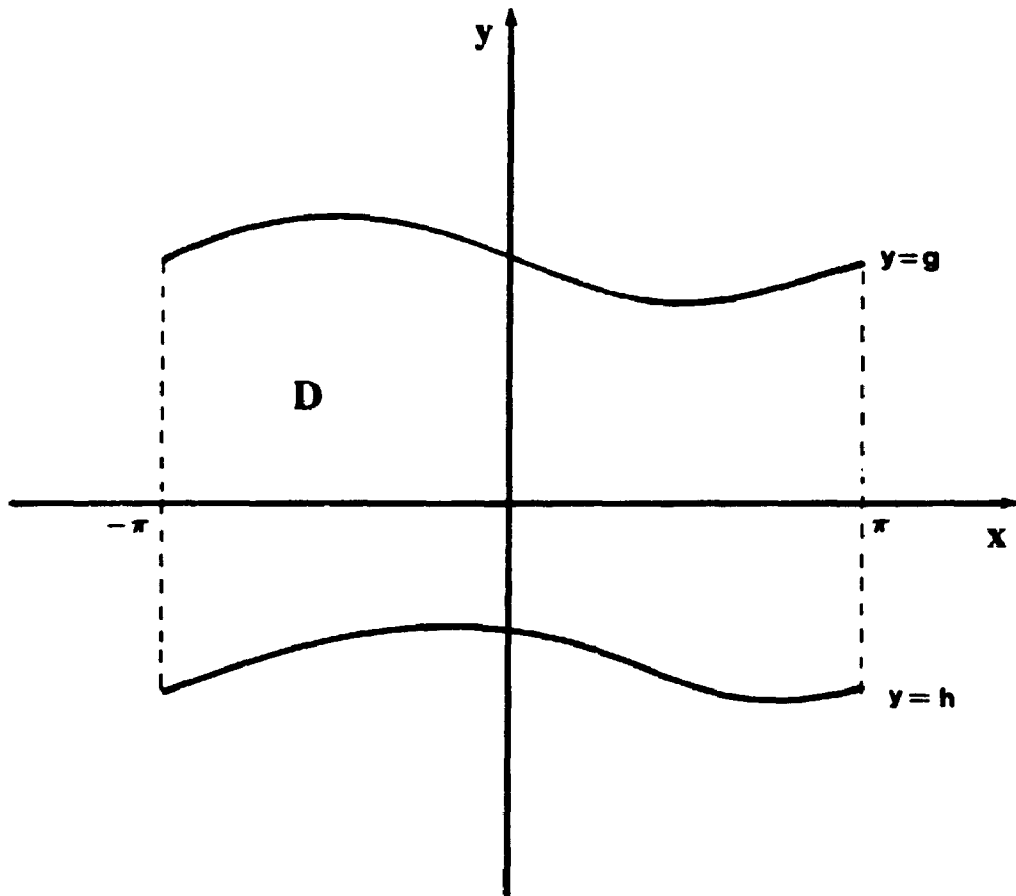


Fig. 2.1.1 Definition Sketch - Physical Plane

$$D(t) = \{(x, y) \mid h(x, t) < y < g(x, t) \quad , \quad -\pi < x < \pi\}$$

Fixing the time  $t_0$ , let us consider the problem posed by equations (2.1.1) and (2.1.2), where  $\partial D(t_0)$  is assumed known.

$$\nabla^2 \phi(\mathbf{x}) = 0 \quad , \quad \mathbf{x} \in D \quad (2.1.5)$$

$$B\phi(\mathbf{x}) = b(\mathbf{x}) \quad , \quad \mathbf{x} = (x, h(x)) \quad (2.1.6)$$

$$\mathbf{x} = (x, g(x))$$

$$\phi(-\pi, y) = \phi(\pi, y) \quad (2.1.7)$$

The approximation method is applied to this situation.

As outlined in section 1.2, the boundary approximation proceeds by first selecting trial functions  $\phi_n = \phi_n(b_1, \dots, b_n; \mathbf{x})$  which satisfy the field equation, in this case (2.1.5). The parameters  $b_i$  are then chosen so that  $\phi_n$  approximates, in some norm, the boundary conditions (2.1.6) and (2.1.7). Now, we distinguish between two types of approximation, namely

- (i) linear approximation, and
- (ii) nonlinear approximation.

In the first case, the trial solutions are represented by linear combinations of known functions  $u_j(\mathbf{x})$ ,

$$\phi_n = \sum_{j=1}^n b_j u_j(\mathbf{x}).$$

In the second case, the trial solutions are nonlinear functions of the unknown parameters and typically take the form

$$\phi_n = \sum_{j=1}^n b_j \gamma_j(a_1, \dots, a_m, \mathbf{x})$$

where the unknowns are the vectors  $\mathbf{b} \in R^n$  and  $\mathbf{a} \in R^m$ . In both cases, the boundary approximation method requires that the functions  $u_j$  and  $\gamma_j$  are solutions to equation (2.1.5). As we shall see, the case of linear approximation can be developed within the

comfortable framework of normed linear spaces. This point of view allows for an exhaustive analysis of the linear boundary approximation method, which is largely unavailable for the nonlinear case. For this reason, the bulk of chapter 2, namely sections 2.2 through 2.7, is devoted to the case of linear approximation. Section 2.8 will outline the nonlinear case.

Although much less is known about nonlinear approximation, it is often the case that such approximations provide superior results to those of the corresponding linear case (viz polynomial versus rational approximation in one-dimensional approximation problems). In some of the applications that follow in the later chapters, we will witness the excellent results that can be obtained with nonlinear approximation, but usually at the cost of greater computational effort.

## 2.2 Best Linear Approximation to Boundary Data

The linear approximation to the boundary value problem (2.1.5) - (2.1.7) is formulated below. The section concludes with an outline of the important aspects of the approximation and identifies later sections where each of these aspects is discussed in turn.

The methods of the present chapter represent, in their simplest terms, the approximation of functions  $\phi(x, y)$  from a certain class by other known functions. For much of the discussion which follows, the properties of the class of functions involved are best dealt with in a transformed plane. In particular, we perform a mapping of the domain in question to a region of the complex plane where the conditions of periodicity are an intrinsic property of the domain involved.

Consider then, a conformal mapping of the strip,

$$\zeta = e^{-z} \quad , \quad h(x) \leq y \leq g(x) \quad , \quad -\infty < x < \infty .$$



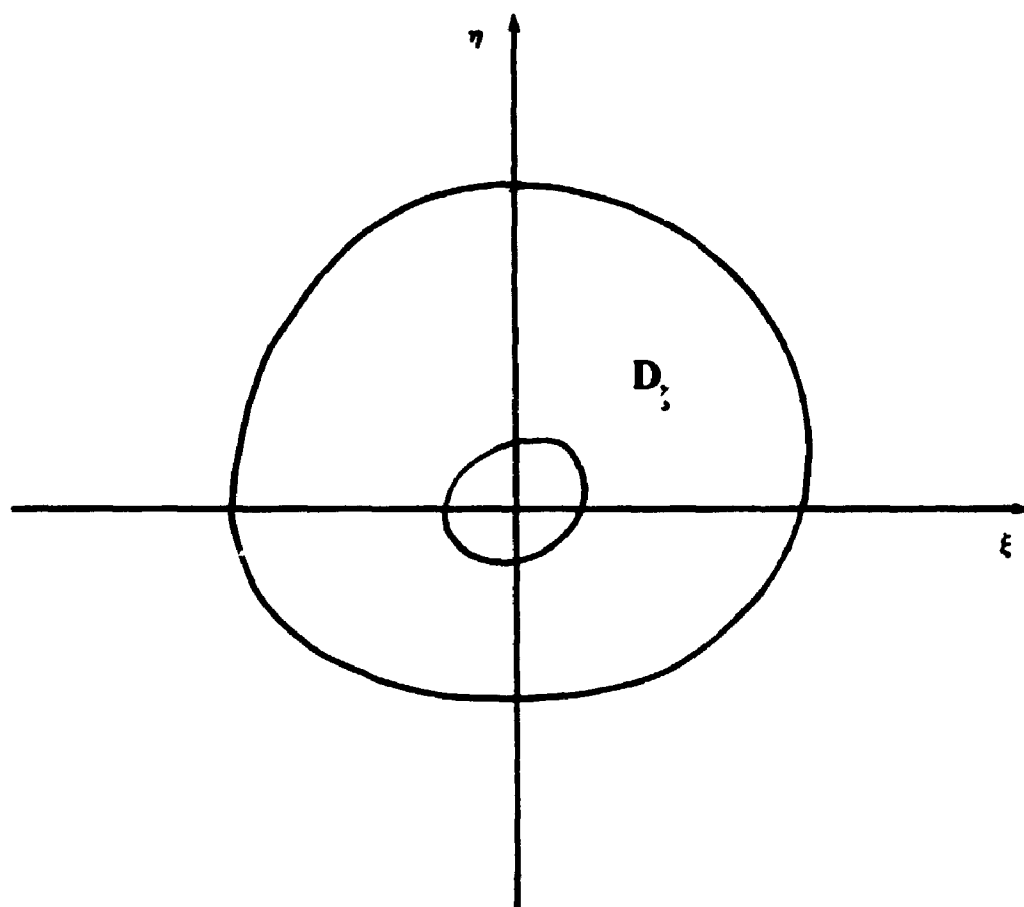


Fig. 2.2.1 Transformed Plane (  $\zeta = \text{EXP}(-iz)$  )

If  $\zeta = re^{i\theta}$ , then  $r = e^y$  and  $\theta = -x$ . The curves  $h$  and  $g$  map into the inner and outer boundaries of the annular domain  $D_\zeta$  shown in figure 2.2.1. The boundary value problem becomes

$$\nabla^2 \Phi(\xi, \eta) = 0 \quad , \quad (\xi, \eta) \in D_\zeta \quad (2.2.1)$$

$$B\Phi(\xi, \eta) = b(\xi, \eta) \quad , \quad (\xi, \eta) \in \partial D_\zeta \quad (2.2.2)$$

where

$$\phi(x, y) = \Phi(\xi(x, y), \eta(x, y)) \quad .$$

In the  $\zeta$ -plane, the original periodicity is reflected by the fact that  $\Phi$  at  $\theta = \theta_0$  assumes the same value as  $\Phi$  at  $\theta = \theta_0 + 2n\pi$ .

We seek an approximate solution  $\Phi_n$  to the boundary value problem (2.2.1) - (2.2.2).

In keeping with the notation of section 1.2, the approximate solution is sought from the linear space of functions  $X(\overline{D}_\zeta)$ , which is at our disposal. Specifically, the trial functions  $\Phi_n$  are elements of the  $n$ -dimensional subspace  $M \subset X(\overline{D}_\zeta)$  spanned by the linear combinations of  $n$  known basis functions  $u_1(\xi, \eta), u_2(\xi, \eta), \dots, u_n(\xi, \eta), (\xi, \eta) \in \overline{D}_\zeta$ . It is assumed that the  $u_j(\xi, \eta)$  are linearly independent. This form of the trial solution is responsible for characterizing the approximation as linear. If the field equation is homogeneous, as is the case in (2.2.1), then the linear boundary approximation method is usually formulated by choosing the basis functions  $u_j(\xi, \eta)$  from a complete set of particular solutions to this equation. We are seeking the linear combination  $\Phi_n \in M$  which minimizes the residual

$$E(\Phi_n) = \|B\Phi_n - b\|_{\omega_\zeta}$$

where  $b \in X(\partial D_\zeta)$  and  $\|\cdot\|_{\omega_\zeta}$  is a suitable norm for  $X(\partial D_\zeta)$ .

Formally, this is a problem of best approximation (b.a.). That is, for a given integer  $n$  and  $f \in X(\partial D_\zeta)$ , the best approximation  $\Phi_n^*$  satisfies

$$\|B\Phi_n^* - b\|_{\omega_\zeta} \leq \|B\Phi_n - b\|_{\omega_\zeta}$$

for all  $\Phi_n \in M$ .

The linear spaces  $X(\Omega)$ ,  $\Omega = \overline{D}_\zeta$ ,  $D_\zeta$ ,  $\partial D_\zeta$  may be quite general. A natural choice for the approximation of functions is the  $L_p(\Omega)$  space. That is, the linear space of  $p$ -integrable functions defined on  $\Omega$  and having the property

$$\int_{\Omega} |f(s)|^p ds < \infty$$

where  $f \in L_p(\Omega)$ . In section 2.5, we take  $X(\overline{D}_\zeta)$  to be the space of functions harmonic on  $D_\zeta$  and continuous on  $\overline{D}_\zeta$ .

The above formulation reduces the problem of exactly solving a given boundary value problem, to the determination of a best approximation to boundary data, in a normed linear space.

Many pertinent questions immediately arise concerning the practical implementation of these ideas. The following considerations figure prominently in any practical calculation of a best approximation:

- (i) the choice of norm
- (ii) the existence of a best approximation
- (iii) the uniqueness of a best approximation
- (iv) the degree of convergence for increasing value of  $n$
- (v) the characterization of a best approximation.

In the remainder of this chapter we address each of these topics in turn, as they pertain to the linear approximation problem. Once a norm has been decided upon (section 2.3), the question of existence of a best approximation is paramount. The uniqueness of the approximation is not a major concern from the practical point of view. All that matters

is that at least one best approximation can be found. Nevertheless, we shall see (section 2.4) that both existence and uniqueness follow quite simply from the well established theory of best approximations in which we have set our problem.

The characterization of best approximations is concerned with the finding of certain properties the best approximation might satisfy. We are concerned with this feature only insofar as such characterizations might suggest numerical algorithms to use (see section 2.7).

Finally, we would hope that by increasing the number of terms in the approximation, we might increase the accuracy of our results. In other words, can we establish convergence in the norm with increasing  $n$ , and if so, how fast is the solution converging? As we shall see, establishing the convergence in norm is an easy matter (once we have demonstrated the completeness of our basis functions (section 2.5)); but determining a useful measure of the degree of approximation is another matter. Still, in section 2.6 we present some rough error estimates for particular problems from our class.

The final section of this chapter (2.8) will be devoted to the nonlinear approximation problem.

### 2.3 The Choice of Norm

The best approximation problem outlined in the last section, requires that we choose a norm for the space  $L_p$ . It is usual to norm  $L_p(\Omega)$  by one of the following  $p$ -norms:

$$\|f\|_p = \left( \int_{\Omega} |f(s)|^p ds \right)^{\frac{1}{p}}, \quad 1 \leq p < \infty \quad (2.3.1)$$

$$\|f\|_{\infty} = \max_{\Omega} |f(s)| \quad (2.3.2)$$

(The vector notation has been suppressed in (2.3.1). In fact, for two-dimensional problems with  $\Omega = \partial D$ ,  $ds$  may be taken as an element of arclength and the vector notation is unnecessary.)

When it comes to numerical computations, it is usually the case that the expressions (2.3.1), (2.3.2) must be replaced by their discrete analogues:

$$\|A\|_p = \left( \sum_{i=1}^m |f(s_i)|^p \right)^{\frac{1}{p}}, \quad 1 \leq p < \infty \quad (2.3.3)$$

$$\|A\|_\infty = \max_{1 \leq i \leq m} |f(s_i)| \quad (2.3.4)$$

As the number of data points  $m$  is increased, it is expected that the discrete norms tend to the continuous case. More will be said on the connection between the discrete and continuous cases later (section 2.7).

We have chosen the case  $p = 2$ . This choice has a number of advantages. The best approximation problem, find the minimum of  $E(\phi_n)$ , is linear in the unknown coefficients  $b_j$  if the condition  $B\phi = b$  is linear in  $\phi$  and its normal derivative. In this case, (2.3.3) (with  $p = 2$ ) leads to the direct solution of a system of linear equations. This type of best approximation problem is usually called a linear least squares approximation.

The choice of approximating functions can often simplify the resulting linear system. If the approximating functions are orthogonal, the system of equations is diagonal - the simplest case. For the problems under consideration in this work, the resulting system of linear equations is full, having no special band structure. There is a price to be paid for choosing to solve a full matrix. It is both computationally less efficient than any method which exploits a band structure and the effects of an ill-conditioned full matrix can be disastrous. Each of these questions is addressed more fully in section 2.7. But first, let us establish the existence of a solution to our least squares problem.

## 2.4 Existence of a Best Approximation

The theory of best approximations on normed linear spaces is well developed and contains some very general results on the existence of a b.a.. For the case of linear approximation, the existence follows from a standard argument. The relevant theorem on this matter is the following:

**Theorem 2.4.1** Let  $M$  be a finite dimensional subspace of a normed linear space  $X$ . Then, there exists a best approximation  $\Phi^* \in M$  to any  $b \in X$  (see Davis (1975), for example).

In the case of the Dirichlet boundary value problem corresponding to (2.2.1) - (2.2.2), we have the b.a. problem: for given  $n$ , find  $\Phi_n^*$  such that

$$\|\Phi_n^* - b\|_{\omega_\zeta} \leq \|\Phi_n - b\|_{\omega_\zeta}$$

for all  $\Phi_n \in M$ . That is, we must solve for  $\Phi_n^*$  such that

$$\|\Phi_n^* - b\|_{\omega_\zeta} = \inf_{\Phi_n \in M} \|\Phi_n - b\|_{\omega_\zeta}.$$

For the purposes of this section, assume the normed linear space  $X$  to be the space  $L_p(\partial D_\zeta)$  defined in section 2.2 and with corresponding norm given in section 2.3.  $M \in L_p(\partial D_\zeta)$  is the  $n$ -dimensional linear space

$$M = \left\{ \Phi_n = \sum_{j=1}^n b_j \mu_j(\xi, \eta) \quad , \quad b_j \in R \quad , \quad (\xi, \eta) \in \partial D_\zeta \right\}.$$

Theorem 2.4.1 is directly applicable to this situation and hence a b.a. does exist.

It is even possible to conclude that the b.a.  $\Phi_n^*$  is unique, if the norm is the  $L_2$  norm.

In fact, if  $1 < p < \infty$ , then the  $L_p$  norms possess the property that they are strictly convex; and strict convexity is sufficient to establish uniqueness of a b.a. (see Clarkson (1936)).

## 2.5 Completeness of the Set of Basis Functions

In section 2.2 it was mentioned that the basis functions  $u_j(\underline{x})$  were to be chosen from a complete set of particular solutions to the Laplace equation. This is essential if we are to establish that the sequence of best approximations  $\{\phi_n^*\}$  converges in some sense to the solution of the given boundary value problem. In this section, we make precise the notions of completeness and convergence and then go on to construct a complete set of basis functions for the class of problems under consideration. More accurately, the development which follows pertains to the approximation of particular functions defined on a two-dimensional closed region (namely functions harmonic on a domain and continuous on  $\overline{D}$ ). As such, the results obtained may be taken as applicable to the Dirichlet boundary value problem from our class.

To begin with, a number of definitions are in order.

**Definition 2.5.1:** The harmonic function  $p(x, y)$  formed by taking linear combinations of the real and imaginary parts of  $z^k$  is called a harmonic polynomial. The harmonic polynomial of degree  $n$  has the form

$$p_n(x, y) = \sum_{\substack{i, j=0 \\ i+j \leq n}}^n a_{ij} x^i y^j$$

where the coefficients  $a_{ij}$  are real. In what follows,  $p_n(x, y)$  may be simply referred to as the harmonic polynomial in the variable  $z$ .

**Definition 2.5.2:** A function  $f(x, y)$  in a normed linear space  $X(\Omega)$  is said to be approximated arbitrarily closely by a sequence of functions  $u_1(x, y), u_2(x, y), \dots$ , in  $X(\Omega)$  if for given  $\varepsilon > 0$ , there exists an integer  $n$  and constants  $h_1, h_2, \dots, h_n$  such that

$$\left\| f(x, y) - \sum_{j=1}^n h_j u_j(x, y) \right\| < \varepsilon$$

If every  $f \in X$  can be so approximated then the sequence is said to be closed or complete<sup>1</sup>.

For example, we have occasion to use the maximum or  $p$ -infinity norm and shall say that a function having the above property is uniformly approximable by the functions  $\{u_j(x, y)\}$ . It is well to bear in mind that, if a function  $f(x, y)$  is uniformly approximable, then it is possible to construct a sequence  $\phi_n(x, y)$  (which is a linear combination of the  $u_j(x, y)$ ,  $j = 1, 2, \dots, n$ ) which converges uniformly to  $f(x, y)$  on  $\Omega$ .

If the sequence  $\{u_j(x, y)\}$  is complete and  $\{\phi_n^*\}$  is a sequence of best approximations to  $\phi \in X$ , then the sequence  $\{\phi_n^*\}$  converges in norm to  $\phi$ ,

$$\lim_{n \rightarrow \infty} \|\phi - \phi_n^*\| = 0.$$

We now proceed to construct a complete set of functions for the class of elliptic boundary value problems at hand. Once again, it is instructive to work in the transformed  $\zeta$ -plane ( $\zeta = e^{-iz}$ ). We consider two cases separately. In the first case (figure 2.5.1)  $D_\zeta$  is a simply connected domain covering the origin. This case corresponds, in the  $z$ -plane, to the boundary value problem (2.1.5) - (2.1.7) where the lower boundary  $y = h(x)$  is absent (ie a semi-infinite  $z$ -domain). In the second case to be examined  $D_\zeta$  is a doubly connected annular domain surrounding the origin (figure 2.5.2). This corresponds in the  $z$ -plane to the usual domain already discussed.

For the simply connected case, we have the following result on approximation of harmonic functions. Let  $D_\zeta$  be a simply connected domain whose boundary is a Jordan curve. Let  $f(\xi, \eta)$  be an arbitrary function harmonic on the closed region  $\overline{D_\zeta}$ . Runge (1885) has established the uniform approximation of a sequence of harmonic polynomials to arbitrary harmonic functions, on closed subsets of the domain in question. (Or at least this

---

<sup>1</sup> Properly defined a sequence  $\{u_n\}$  is closed in a linear space  $X$  if it satisfies definition 2.5.2 and is complete if  $L(u_n) = 0$  implies  $L = 0$ , where  $L$  is a linear functional from the conjugate space of  $X$ . In normed linear spaces, the two definitions are equivalent.



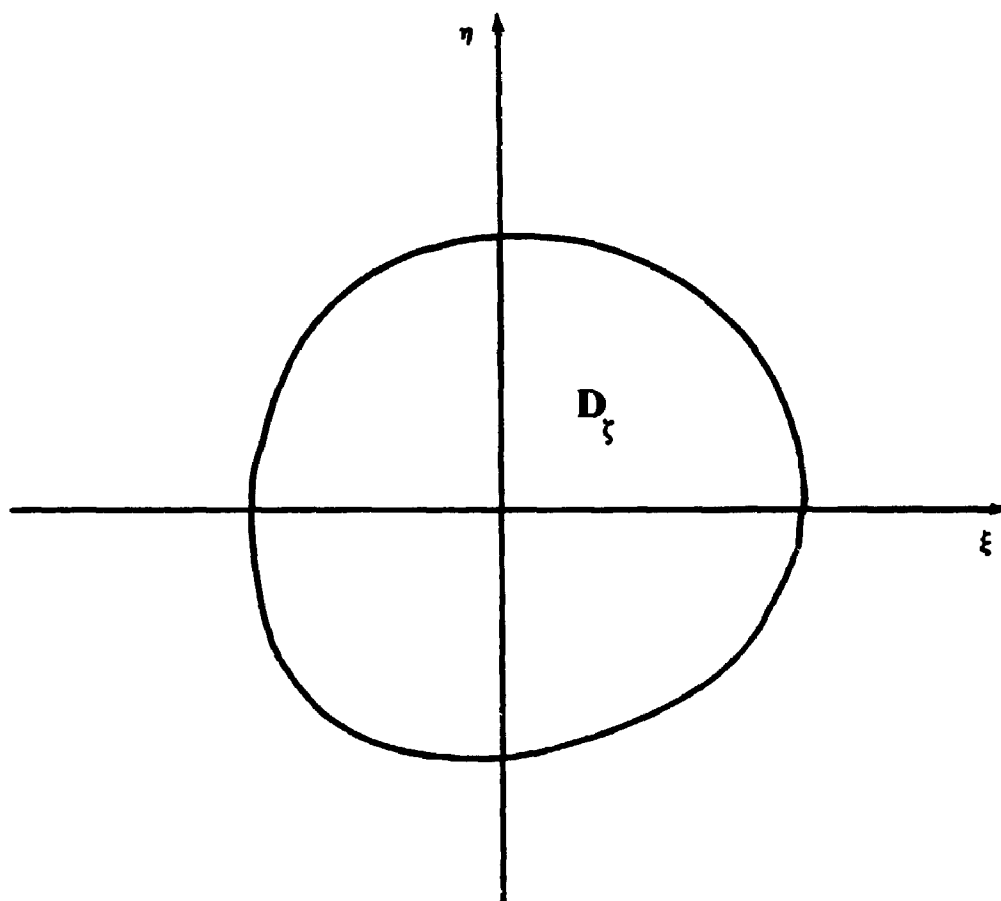


Fig. 2.5.1 Simply Connected Case

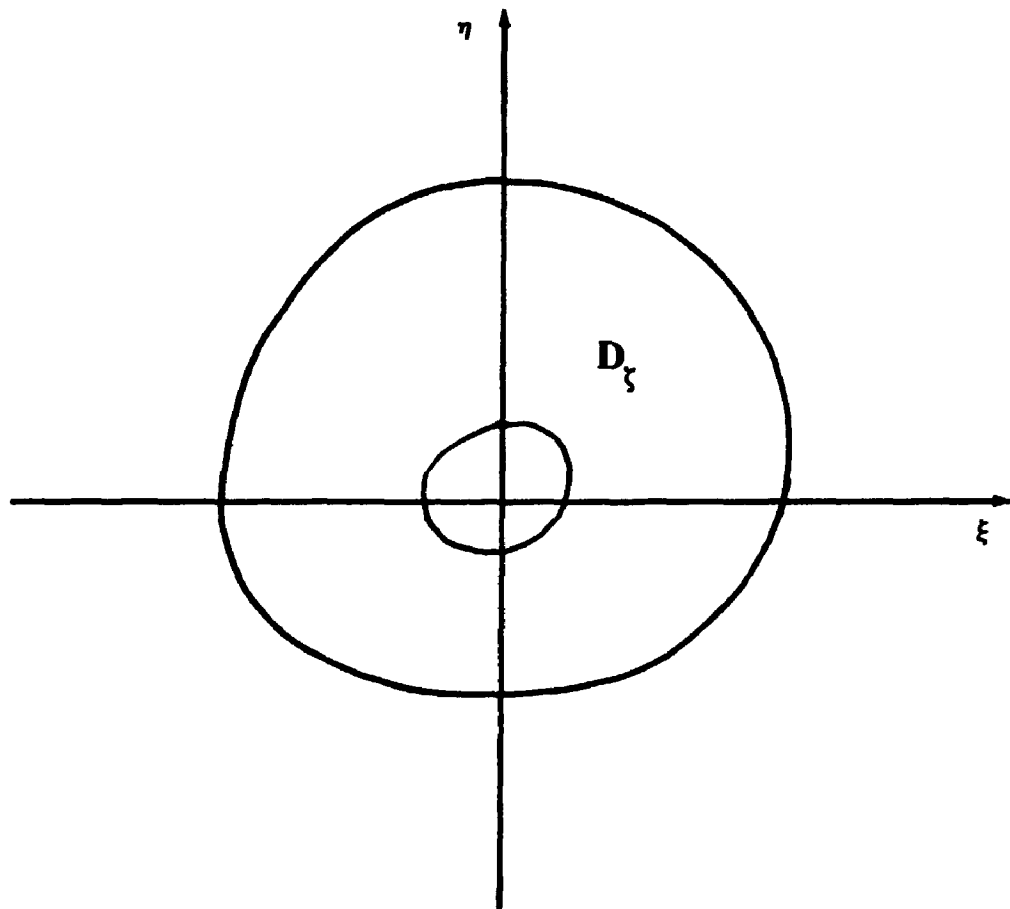


Fig. 2.5.2 Doubly Connected Case

result can be inferred from his results on uniform approximation of analytic functions by complex polynomials.) Walsh (1928(a)) extended this result to allow for uniform approximation, on the corresponding closed region  $\overline{D_\zeta}$ , by harmonic polynomials to functions harmonic in  $D_\zeta$  and continuous in  $\overline{D_\zeta}$ .

Walsh (1929) has established the following generalization of Runge's work to the case of multiply-connected domains.

**Theorem 2.5.1** Let  $D_\zeta$  be a finitely connected domain bounded by the non-intersecting Jordan curves  $C_0, C_1, \dots, C_p$ , and such that  $C_1, \dots, C_p$  lie interior to  $C_0$ . Let  $f(x, y)$  be an arbitrary function, continuous on the closed region  $\overline{D_\zeta}$  and harmonic on the domain  $D_\zeta$ . Given  $\epsilon > 0$ , there exist harmonic polynomials  $p(\xi, \eta)$ ,  $q(\xi, \eta)$  in  $\zeta$  and  $(\zeta - \zeta_j)^{-1}, j = 1, \dots, p$ , respectively, and real constants  $A_1, \dots, A_p$ , such that

$$|f(\xi, \eta) - \left\{ p(\xi, \eta) + q(\xi, \eta) + \sum_{j=1}^p A_j \log |\zeta - \zeta_j| \right\}| < \epsilon$$

for all  $(\xi, \eta) \in \overline{D_\zeta}$ . The points  $\zeta_j$  are arbitrarily chosen to reside one in each of the holes surrounded by  $C_1, \dots, C_p$ . (The theorem includes the previous result on simply connected domains as a subcase.)

Thus, if  $X(\overline{D_\zeta})$  is the linear space of functions which are harmonic on  $D_\zeta$  and continuous on  $\overline{D_\zeta}$ , and where the norm is the maximum norm, then the theorem implies that the set

$$\{\Re, \Im(\zeta^k), (\zeta - \zeta_j)^{-k}, \log |\zeta - \zeta_j|, j = 1, \dots, p, k = 0, 1, 2, \dots\} \quad (2.5.1)$$

is complete for the space  $X(\overline{D_\zeta})$ .

In our case, the domain  $D_\zeta$  is doubly connected (ie  $p=1$ ) and surrounds the origin, so that we may take the set

$$\{\log |\zeta|, \Re, \Im \zeta^k, k = 0, \pm 1, \pm 2, \dots\} \quad (2.5.2)$$

to be complete in the maximum norm.

Of course, the associated Dirichlet problem is not soluble on an arbitrary doubly connected region; but if the boundary  $\partial D_\zeta$  is sufficiently smooth such that a solution  $\Phi$  exists, then Theorem 2.5.1 guarantees the existence of a series

$$\Phi_n = A \log |\zeta| + \sum_{k=-n}^n (a_k \Re + b_k \Im) \zeta^k$$

which converges in the maximum norm to  $\Phi$ . The  $A, a_k, b_k$  depend on  $n$ . That is, for every  $\epsilon > 0$ , there exists an integer  $N$  and real coefficients  $a_k, b_k, A$  such that for  $n > N$  we have

$$\|\Phi(\xi, \eta) - A \log |\zeta| - \sum_{k=-n}^n (a_k \Re + b_k \Im) \zeta^k\|_\infty < \epsilon$$

or more compactly, writing  $c_k = a_k - ib_k$ ,

$$\|\Phi(\xi, \eta) - A \log |\zeta| - \Re \sum_{k=-n}^n c_k \zeta^k\|_\infty < \epsilon.$$

By definition this must be true of the sequence of best approximations  $\Phi_n^m$  derived from the best linear approximation to boundary data. What is more, a similar result must hold for the sequence of best approximations  $\Phi_n^*$  in the  $L_2$  norm. For we have the inequality

$$\begin{aligned} \int_{\partial D_\zeta} |\Phi - \Phi_n^*|^2 ds &\leq \int_{\partial D_\zeta} |\Phi - \Phi_n^m|^2 ds \\ &\leq \text{const.} \left\{ \max_{\partial D_\zeta} |\Phi - \Phi_n^m| \right\}^2. \end{aligned}$$

Note that if the maximum norm is used, the convergence is actually uniform.

Now, the above results pertain to the  $\zeta$ -plane. Under the inverse mapping  $z = i \log \zeta$  to the  $z$ -plane we have that the set of functions

$$\{y, \Re, \Im e^{-ikt}, k = 0, \pm 1, \pm 2, \dots\} \quad (2.5.3)$$

is complete for the space  $X(\bar{D})$  of functions harmonic on  $D$  and continuous on  $\bar{D}$  and extendable in a continuous fashion to be  $2\pi$  periodic on the strip

$$h(x) \leq y \leq g(x) \quad , \quad -\infty < x < \infty .$$

(Recall that  $D$  corresponds to just one period of this strip.)

Thus, we have a complete set of basis functions (2.5.3) (in the maximum or least squares norm). Furthermore, this establishes the convergence in norm of the sequence of best approximations, obtained by using (2.5.3), to the solution of the corresponding Dirichlet boundary value problem.

For a discussion concerning the origin of the logarithmic terms in the complete set (2.5.1), see the appendix A.1.

Finally, it is interesting to note that our boundary approximation scheme can be interpreted as rational approximation in the complex variable of an appropriate complex plane (plus a logarithmic term). For some of the problems to be considered, where the domain of interest is semi-infinite, the approximation corresponds to one of polynomial approximation in the appropriate complex plane.

## 2.6 Degree of Convergence

We have established the existence of a best approximation and in the last section constructed a complete set of basis functions which guarantee convergence in the norm of a sequence of best approximations. It is another matter still to determine the actual degree of convergence. In general, this is a difficult task and we provide here rough estimates for particular problems from our class.

### The Simply Connected Domain

Consider the potential problem

$$\nabla^2 \phi(\mathbf{x}) = 0 \quad \mathbf{x} \in D \quad (2.6.1)$$

$$\phi(\mathbf{x}) = f(\mathbf{x}) \quad \mathbf{x} \in \partial D \quad (2.6.2)$$

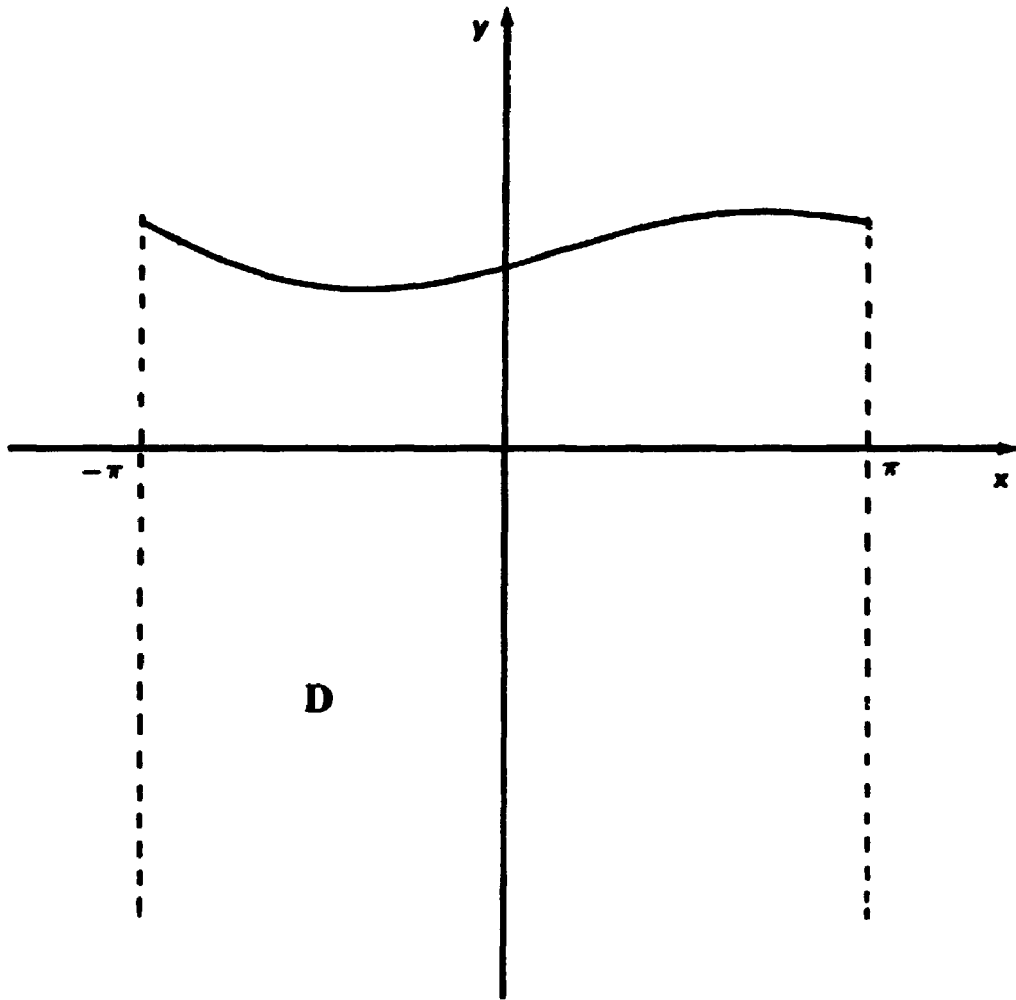


Fig. 2.6.1 Semi-infinite Domain - Physical Plane

where  $D$  is the semi-infinite domain shown in figure 2.6.1 and it is understood that  $\phi$  is bounded at minus infinity. An example is provided by the potential problem associated with the Hele-Shaw flow treated in chapter 4. Assume that  $\partial D$  is an analytic Jordan arc and that  $\phi^{(\rho)}(s) \in \text{Lip } \alpha$ . That is,  $\phi^{(\rho)}$  satisfies a Lipschitz condition of order  $\alpha$  with respect to arclength. More precisely, if  $s$  is the arclength parameter, then  $\phi(x, y)$  is harmonic in  $D$ , continuous in  $\bar{D}$  and

$$\left| \frac{\partial^\rho \phi(s_1)}{\partial s^\rho} - \frac{\partial^\rho \phi(s_2)}{\partial s^\rho} \right| \leq M |s_1 - s_2|^\alpha$$

for all  $s_1, s_2 \in \partial D$ ,  $M$  a constant and  $0 < \alpha \leq 1$ .

Under the transformation  $\zeta = e^{-it}$ ,  $-\pi \leq x \leq \pi$ , the domain  $D$  is mapped into the simply connected domain  $D_\zeta$  surrounding the origin (see figure 2.5.2). Assuming that the Lipschitz condition is invariant under such a mapping, we again have a Dirichlet potential problem for  $\Phi(\xi, \eta)$  in the  $\zeta$ -plane, with  $\Phi^{(\rho)}(t) \in \text{Lip } \alpha$ , where  $t$  is the arclength parameter in the transformed plane.

Now, Walsh, Sewell and Elliott (1949) have established that a sequence of harmonic polynomials  $\Phi_n$  exists such that

$$|\Phi - \Phi_n| \leq \frac{K}{n^{\rho+\alpha}}$$

for all  $(\xi, \eta) \in \bar{D}_\zeta$ , where  $K$  is a positive constant and

$$\Phi_n = \sum_{k=0}^n (a_k \mathcal{R} + b_k \mathcal{I}) \zeta^k.$$

Under the inverse mapping, then, we have the corresponding result in the  $z$ -plane. There exists a harmonic function

$$\phi_n = \sum_{k=0}^n (a_k \mathcal{R} + b_k \mathcal{I}) e^{-kt}$$

such that

$$|\phi - \phi_n| \leq \frac{K}{n^{p+\alpha}} \quad (2.6.3)$$

for all  $(x, y) \in \bar{D}$ . This provides an error estimate for the best approximation problem in the maximum norm to the boundary value problem (2.6.1)-(2.6.2). As in section 2.5, we can establish a similar estimate for the best approximation  $\phi^*$  to boundary data  $f$  in the  $L_2$  norm, namely

$$\|\phi - \phi^*\|_2 \leq \frac{K}{n^{p+\alpha}} \quad x \in \partial D \quad (2.6.4)$$

It remains only to establish the invariance of the Lipschitz condition under a conformal mapping. Let  $s = g(t)$  be an analytic function relating the arclength parameters in the  $z$ - and  $\zeta$ -planes. Then, if  $\phi^{(p)}(s)$  exists and is bounded, so  $\Phi^{(p)}(t)$  exists and is bounded. That is,

$$\Phi^{(p)}(t) = a_1(s)\phi'(s) + a_2(s)\phi''(s) + \dots + a_p(s)\phi^{(p)}(s)$$

where  $a_1(s), a_2(s), \dots, a_p(s)$  are multiples of the derivatives of  $s$  with respect to  $t$ . Now, if  $\phi^{(p)}(s) \in \text{Lip } \alpha$ , then  $\phi^{(p-1)}(s), \phi^{(p-2)}(s), \dots, \phi(s) \in \text{Lip } \alpha$ . This follows from the mean value theorem

$$\left| \frac{\phi^{(p-1)}(s_1) - \phi^{(p-1)}(s_2)}{s_1 - s_2} \right| = |\phi^{(p)}(s_0)| \quad , \quad s_1 \leq s_0 \leq s_2$$

$$\leq M .$$

Therefore

$$|\phi^{(p-1)}(s_1) - \phi^{(p-1)}(s_2)| \leq M |s_1 - s_2|$$

and we have  $\phi^{(p-1)}(s) \in \text{Lip } 1$ . But  $\text{Lip } 1 \subseteq \text{Lip } \alpha$  if  $0 < \alpha \leq 1$ . Therefore,  $\phi^{(p-1)}(s) \in \text{Lip } \alpha$  and so on for the lower derivatives. By the same token,  $a_j(s) \in \text{Lip } \alpha$  for all  $j = 1, 2, \dots, p$ . That each product



$$a_k(s)\phi^{(k)}(s) \in \text{Lip } \alpha$$

follows from the inequality

$$|ab - cd| \leq |a| \cdot |b - d| + |d| \cdot |a - c|.$$

We have the desired result,

$$\Phi^{(r)}(t) \in \text{Lip } \alpha.$$

The estimate (2.6.4) clearly indicates how strongly the degree of convergence is dependent on the smoothness of the boundary data. (2.6.4) by itself is a proof of convergence of the sequence of best approximations  $\{\phi_n^*\}$  for the Dirichlet problem (2.6.1) - (2.6.2).

### The Annular Region

We have the following result (due to Walsh(1928(b))) on the degree of approximation of an analytic function by means of rational functions. Let  $\bar{D}_r$  be the closed annular region bounded by the two Jordan curves  $C_0, C_1$  with  $C_1$  interior to  $C_0$ , and  $C_1$  surrounding the origin. Let  $w = \psi_0(\zeta)$  denote the function which conformally maps the exterior of  $C_0$  onto the exterior of the unit circle and  $w = \psi_1(\zeta)$  denote the function which conformally maps the interior of  $C_1$  onto the interior of the unit circle. Let  $\Gamma_0$  be the curve

$$|\psi_0(\zeta)| = R > 1$$

and  $\Gamma_1$  be the curve

$$|\psi_1(\zeta)| = \frac{1}{R}.$$

Suppose  $g(\zeta)$  is given, single-valued and analytic in the closed region  $\Gamma$  which is interior to  $\Gamma_0$  and exterior to  $\Gamma_1$ . Then, there exists a rational function  $r_n(\zeta)$  of degree  $2n$  such that

$$|g(\zeta) - r_n(\zeta)| \leq \frac{M}{R^n},$$

for all  $\zeta \in \overline{D}_\zeta$ .  $r_n(\zeta)$  has its only pole at the origin and may be taken to be a polynomial in  $\zeta$  plus a polynomial in  $1/\zeta$ , both of degree  $n$ . That is, we may write

$$r_n(\zeta) = \Re \sum_{k=-n}^n c_k \zeta^k,$$

where  $c_k$  is complex.

Now, this result can be used to derive an error estimate for the approximation of harmonic functions on  $\overline{D}_\zeta$ . If  $f(\zeta)$  is the analytic function formed from the single-valued harmonic function  $\Phi(\xi, \eta)$  and its conjugate, then we can express

$$f(\zeta) = g(\zeta) + A \log(\zeta) \quad (2.6.5)$$

for some real constant  $A$  and where  $g(\zeta)$  is a single-valued analytic function (for a proof, see the appendix A.1). Since  $\log|\zeta|$  is bounded in the closed region  $\Gamma$ , we can find constants  $A_1, M_1$ , and  $M_2$  and a rational function  $r_n(\zeta)$  such that

$$|\Re g(\zeta) - \Re r_n(\zeta)| \leq \frac{M_1}{R^n}$$

and

$$|A \log|\zeta| - A_1 \log|\zeta|| \leq \frac{M_2}{R^n}$$

for all  $\zeta \in \overline{D}_\zeta$ . Thus, we have

$$\begin{aligned} |\Re f(\zeta) - A_1 \log|\zeta| - \Re r_n(\zeta)| &= |\Re g(\zeta) + A \log|\zeta| - A_1 \log|\zeta| - \Re r_n(\zeta)| \\ &\leq |\Re g(\zeta) - \Re r_n(\zeta)| + |A \log|\zeta| - A_1 \log|\zeta|| \\ &\leq \frac{M}{R^n} \end{aligned}$$

for all  $\zeta \in \overline{D}_\zeta$ . That is, we have

$$|\Phi(\xi, \eta) - A_1 \log |\zeta| - r_n(\xi, \eta)| \leq \frac{M}{R^n}, \quad (\xi, \eta) \in \bar{D}_\zeta \quad (2.6.6)$$

$r_n(\xi, \eta)$  is the rational function formed from the sum of a harmonic polynomial of degree  $n$  in  $\zeta$  and a harmonic polynomial of degree  $n$  in  $1/\zeta$ .

The larger the region for which  $\Phi(\xi, \eta)$  has a harmonic extension, the greater the rate of convergence. This result is applicable to the Dirichlet problem on the annular region  $\bar{D}_\zeta$  where the solution  $\Phi(\xi, \eta)$  has a harmonic extension to the larger region  $\Gamma$ .

Of course, under the inverse transformation to the  $z$ -plane, we have the degree of convergence estimate

$$\left| \phi(x, y) - A_1 y - \Re \sum_{k=-n}^n c_k \exp -ikz \right| \leq \frac{M}{R^n}.$$

The above result (2.6.6), is included in Walsh (1929), but the nature of the rational function  $r_n(\xi, \eta)$  is not clearly delineated. We prefer the above simple proof which makes use of the existing theory on analytic functions (once the relation (2.6.5) is established).

## 2.7 Characterization of, and Algorithms for, Computing the Best Approximation

Best approximations in different norms can be characterized by certain conditions they must satisfy. Once discovered, these conditions often suggest numerical algorithms to be used for determination of the b.a..

For the class of problems at hand, the best approximation problem outlined in section 2.2 leads to minimization of the error residual

$$E(\phi_n) = \|B\phi_n - f\|_{\omega}$$

where the norm is the  $L_2$  norm

$$\|B\phi_n\|_{\partial D} = \left( \int_{\partial D} (B\phi_n - f)^2 ds \right)^{\frac{1}{2}} \quad (2.6.1)$$

and the boundary  $\partial D$  is made up of the curves  $y = g(x)$  and  $y = h(x)$  (see fig. 2.1.1).

Define the  $L_2$  norm of a function  $f(s)$  on the finite point set  $S_m = \{s_i, i = 1, \dots, m\}$

by

$$\|f\|_2 = \left( \sum_{i=1}^m |f(s_i)|^2 \right)^{\frac{1}{2}} \quad (2.6.2)$$

Then for given  $n$ , the problem of b.a. becomes one of finding

$$\begin{aligned} \min_{b_j} \|B\phi_n - f\|_2 \\ = \min \left\{ \sum_{i=1}^m (B\phi_n(s_i) - f(s_i))^2 \right\}^{\frac{1}{2}}, \quad m \geq n \end{aligned} \quad (2.6.3)$$

where the points  $s_i$  are coordinate points along the  $\partial D$ .

Now, it is the continuous approximation, minimization of (2.6.1), that we wish to solve; but for computational purposes, we must solve the discrete approximation problem (2.6.3) instead. Nevertheless, it is hoped that for  $m$  sufficiently large, the discrete problem is a close approximation to the continuous case. In fact, Rice (1964) has shown that for suitably chosen point sets  $S_m$ , the b.a. on a finite point set converges to the b.a. on a continuous interval in the limit as  $m \rightarrow \infty$ . Davis and Rabinowitz (1961) believe that a value of  $m$  that is two or three times that of  $n$  (the number of basis functions) is sufficient to reflect convergence of the discrete to the continuous cases. It should be noted, moreover, that Davis and Rabinowitz did find examples of poor approximation of their trial solution to points not in the point set  $S_m$  when  $m$  was taken equal to  $n$ . This problem was overcome when  $m$  was two or three times that of  $n$ . (We will have more to say on this important point shortly.)

The best approximation problem (2.6.3) can be expressed in matrix notation

$$\min_{\mathbf{b} \in R^n} \|\mathbf{r}\|_2 = \min_{\mathbf{b} \in R^n} \|A\mathbf{b} - \mathbf{f}\|_2 \quad (2.6.4)$$

where  $A$  is the matrix of coefficients  $Bu_j(s_i)$  in the expression

$$B\phi_n(s_i) = \sum_{j=1}^n b_j Bu_j(s_i) \quad .$$

$\mathbf{b}$  is the  $n$ -dimensional vector of unknowns  $b_j$ ,  $j = 1, \dots, n$ ,  $\mathbf{f}$  is an  $m$ -dimensional vector of known data values  $f_i = f(s_i)$ ,  $i = 1, \dots, m$ . It is, of course, an equivalent process to minimize the quantity  $\|\mathbf{r}\|_2^2$ .

In matrix form, a simple characterization is available, one which immediately suggests a numerical approach. The vector  $\mathbf{b} \in R^n$  which minimizes  $\|A\mathbf{b} - \mathbf{f}\|_2^2$  does so if and only if

$$A^T \mathbf{r} = \mathbf{0}$$

That is,

$$A^T A \mathbf{b} = A^T \mathbf{f} \quad (2.6.5)$$

The equations (2.6.5) (ie one equation for each  $i = 1, \dots, m$ ) are referred to as the normal equations.

Unfortunately, the solution of the normal equations is not the recommended approach. If the matrix  $A$  is ill conditioned, then the product  $A^T A$  is only much more so. In fact, if the condition number of  $A$  is  $\text{cond}(A)$ , then that of  $A^T A$  is  $(\text{cond}(A))^2$  (see Golub (1983)). A badly conditioned matrix can drastically affect the approximate solution of a system of linear equations.

An alternative algorithm (see Golub (1983), for example) for the solution of (2.6.4) which avoids the formation of  $A^T A$  and thus reduces the effects of ill conditioning, is outlined in Appendix A.2 and has been programed in Fortran 5 and implemented in the examples of chapters 3 and 4.

It should be noted that one could attempt to solve the overdetermined system  $A \mathbf{b} = \mathbf{f}$  directly, using Gaussian elimination, but this would not necessarily correspond to a best approximation problem. In fact, there may be no solution if  $\mathbf{f}$  is not in the column space of  $A$ . However, the discrete least squares problem (ie the best approximation problem in the 2-norm) always has a solution.

Whatever the choice of algorithm, the numerical solution of a system of linear equations is subject to roundoff errors. It is always advisable to pre-condition the matrix  $A$ . This can be done by scaling the rows or columns. Row scaling is equivalent to solving a weighted least squares problem where positive weights are added to the definition of the least squares norm. The choice of weights can be very problem dependent and they alter the least squares solution. Column scaling, on the other hand, can be used without altering the solution and is easily applied. If the column vectors  $\mathbf{a}_j$  of  $A$  are scaled by  $\|\mathbf{a}_j\|_2^{-1}$ , then van der Sluis (1969) has shown that the condition number of  $A$  is within a factor of  $n^{\frac{1}{2}}$  of the minimum that is possible with this type of scaling. This form of pre-conditioning is used in our numerical experiments.

Let us consider the case  $m=n$  once again. In this case, the b.a. problem is solved directly by finding the  $\mathbf{b} \in R^n$  such that

$$A \mathbf{b} = \mathbf{f}.$$

This corresponds to interpolation (or collocation) to boundary data by the function  $\phi_n$ . Under the conformal mapping  $\zeta = e^{-z}$ , this is just interpolation to boundary data by harmonic polynomials (plus a log term).

Now, although interpolation may be the simplest and most natural approximation to implement on a finite point set, it is by no means the best. In fact, it can lead to entirely erroneous results. This is often most graphically illustrated if derivatives of the approximate series solution are required, or values of the approximate solution are computed at non-collocation points.

What is more, convergence of the series of approximating functions  $\phi_n$  as  $n \rightarrow \infty$  can not be expected for all sequences of sets of interpolation points. (The classic examples of Runge can be alluded to for the one dimensional interpolation by polynomials.)

Curtiss (1960) has examined the convergence question for complex and harmonic polynomial interpolation to boundary values on the unit circle. It is possible to find sequences of interpolation points on the unit circle that guarantee convergence of the polynomials. He concludes that, for a general smooth boundary  $\Gamma$ , it is most probable that sequences of sets of interpolation points do exist for which accompanying polynomials converge to the solution of a Dirichlet boundary value problem; but to determine these sequences would require knowledge of the proper placement of the interpolation points on the corresponding unit circle. That is, the unit circle could be obtained under a conformal map, but in general such a map would not be known and in the event that it was, an alternative method of solution would undoubtedly be employed.

Since the present method is essentially approximation by harmonic polynomials in the appropriate complex plane, it is expected that these same conclusions apply to interpolation in our case.

In short, the simplest approach, that of interpolation to boundary values, should be used with some caution.

On the other hand, Fenton and Rienecker (1982) solve a water wave problem using interpolation to harmonic functions and do not report any difficulties with differentiating the series that they obtain, even though such derivatives are used in subsequent calculations.

Levin (1980) has also noted the difficulties encountered with collocation approximations of harmonic functions and has suggested a means of computing the best interpolation set. Unfortunately, the technique requires knowledge of the Green's function for the given domain, and cannot therefore be of assistance in the present situation.

In the following chapters we will present some results which dramatically reveal the gross errors that can be incurred with the naive use of interpolation. We shall see that even a slightly overdetermined ( $m > n$ ) system leads to perfectly acceptable behaviour.

## 2.8 Best Approximation by Nonlinear Functions

### Introduction

In this section, a nonlinear approximation scheme is examined. This time, the approximating functions take the form

$$\phi_n(x) = b_0 + \sum_{j=1}^n b_j \gamma(b_{j+n}; x) \quad (2.8.1)$$

where  $\gamma$  is a nonlinear function of the parameters  $b_{j+n}$ . For (2.8.1) to represent the trial solution in the boundary approximation method, the functions  $\gamma$  must be particular solutions of the Laplace equation. (And of course, we impose the restriction that they be  $2\pi$ -periodic in the variable  $x$ .)

Unlike the linear case, a general theory of best approximation by nonlinear functions is lacking. The questions of existence, uniqueness and degree of convergence of a sequence of best approximations are not readily answered. In this section, the trial functions (2.8.1) are considered to be a logical next step to the preceding work on linear approximation.



As such, it is hoped that the discussion herein provides motivation for (if not rigorous justification for) the choice of functions  $\gamma$  to be used. The numerical results presented in chapter 4 speak for themselves.

The particular choice of functions  $\gamma$  is suggested by the following example. As usual, it is instructive to work in the complex  $\zeta$ -plane ( $\zeta = \exp(-iz)$ ). Let  $D_\zeta$  be a simply connected domain, covering the origin and having boundary  $\partial D_\zeta$ . Consider the Dirichlet problem

$$\nabla^2 \Phi(\xi, \eta) = 0 \quad , \quad (\xi, \eta) \in D_\zeta \quad (2.8.2)$$

$$\Phi(\xi, \eta) = b(\xi, \eta) \quad , \quad (\xi, \eta) \in \partial D_\zeta \quad (2.8.3)$$

This corresponds to a potential problem on a periodic semi-infinite domain in the  $(x, y)$ -plane. A worked example is provided by the time-dependent Hele-Shaw problem given in chapter 4.

Recall that a linear approximation of the solution to (2.8.2),(2.8.3) assumes the form

$$\Phi_n(\xi, \eta) = \Re \sum_{j=0}^n c_j \zeta^j.$$

That is, the analytic function  $f(\zeta)$  constructed from  $\Phi$  and its harmonic conjugate, takes as an approximation, a complex polynomial in  $\zeta$ :

$$f_n = \sum_{j=0}^n c_j \zeta^j.$$

Runge (1885) has established the following: Let  $\partial D_\zeta$  be a Jordan curve and let  $f(\zeta)$  be analytic on  $\overline{D}_\zeta$ . Given  $\varepsilon > 0$ , there exists a rational function

$$R(\zeta) = A_0 + \sum_{j=1}^n \frac{A_j}{\zeta - \zeta_j} \quad (2.8.4)$$

whose poles lie exterior to  $\overline{D}_\zeta$  and for which

$$|f(\zeta) - R(\zeta)| < \varepsilon$$

for all  $\zeta \in \overline{D}_\zeta$ .

This is a special case of the more general rational approximation

$$R_{pq}(\zeta) = \frac{\sum_{j=0}^p a_j \zeta^j}{\sum_{j=0}^q b_j \zeta^j}$$

In one dimension, the merits of rational approximation are well known. For many problems the results are superior to those obtained by polynomial approximation, both in terms of accuracy and efficiency. For analytic functions of a complex variable, the rational function  $R_{pq}$  constructed from the power series of a known analytic function is usually referred to as a Pade approximation, and often exhibits striking convergence properties.

The success of Pade approximation may in part be due to its interpretation as a means of accelerating the convergence of the given power series. For example, Shanks (1955) has established the following. He defines a sequence of operators  $e_k(A_n)$  to be applied to the sequence of partial sums of given power series,

$$A_n = \sum_{i=0}^n a_i \zeta^i.$$

(For the definition of  $e_k(A_n)$  for general  $k$ , see Shanks (1955). Note, however, that for  $k = 1$ ,

$$e_1(A_n) = \frac{(A_{n+1}A_{n-1} - A_n^2)}{(A_{n+1} + A_{n-1} - 2A_n)}$$

which is the well known Aitken acceleration. The field of numerical analysis is replete with examples where this expression is used with considerable success as a means of accelerating a sequence's convergence.)

Now, Shanks goes on to show that

$$e_k(A_n) \equiv R_{kn}$$

for  $k \geq n, n = 1, 2, \dots$ . Of course, it is the real part of  $R_{p,q}$  that would be used as an approximation, since we are given information about  $\Phi$  only, in the form of a boundary condition. In addition, a particular choice of degree of the rational function would have to be made.

The expression (2.8.4) is the partial fraction decomposition of a rational function of two polynomials of degree  $n$ , having simple poles only. It is the real part of this expression which is taken as the nonlinear approximation for the harmonic function  $\Phi$ . That is, let

$$\Phi_n = \Re \left\{ A_0 + \sum_{j=1}^n \frac{A_j}{(\zeta - \zeta_j)} \right\} \quad (2.8.5)$$

If the original problem is symmetric about  $x = 0$  then the poles  $\zeta_j$  lie along the real axis in the  $\zeta$ -plane and the constants  $A_j$  are also real. Furthermore, if we restrict the singularities to lie along  $\theta = 0$  only, then we can choose  $\zeta_j = \exp(b_{j+n})$  and the expression for  $\phi_n$  has a particularly simple form in the  $(x, y)$ -plane. We have, after some simplification

$$\phi_n = b_0 + \sum_{j=1}^n b_j \left( \frac{\sinh(b_{j+n} - y)}{\cosh(b_{j+n} - y) - \cos x} \right) \quad (2.8.6)$$

For  $D$  defined by  $-\pi < x < \pi, -\infty < y < g(x)$ , the only singularities in  $\phi_n$  are along the  $y$ -axis at the points  $y = b_{j+n}, j = 1, \dots, n$ . We can restrict the singularities to lie outside  $D$  if  $b_{j+n} > g(0)$  for each  $j = 1, \dots, n$ .

#### Existence of a best approximation

Hobby and Rice (1967) have studied the problem of best approximation by nonlinear functions of the form (2.8.1). In the context of the problem just discussed, let the approximating functions take the form (2.8.6) with

$$\chi(b; x) = \frac{\sinh(b - y)}{\cosh(b - y) - \cos x}.$$

A  $\gamma$ -polynomial of order  $n$  is defined to be any element of  $L_p(\partial D)$  of the form (2.8.1) with  $(x, y) \in \partial D$  and  $\alpha \leq b_{1+n} < b_{2+n} \dots < b_{2n} \leq \beta$ . Denote by  $P_{\gamma, n}$ , the set of all  $\gamma$ -polynomials of order  $n$ . Then, given  $f \in L_p(\partial D)$ , the best approximation problem becomes: find  $\phi^* \in P_{\gamma, n}$  such that

$$\|\phi^* - f\|_p = \inf_{\phi_n \in P_{\gamma, n}} \|\phi_n - f\|_p.$$

Now, the existence of a  $\gamma$ -polynomial of best approximation does not immediately follow, as in the linear case, for the set  $P_{\gamma, n}$  is not closed. For example, let  $\gamma_1, \gamma_2 \in P_{\gamma, n}$ ,  $\gamma_1 \neq \gamma_2$  where

$$\gamma_i = \frac{\sinh(b_i - y)}{\cosh(b_i - y) - \cos x}, \quad i = 1, 2.$$

Clearly

$$\frac{\gamma_1 - \gamma_2}{b_1 - b_2} \in P_{\gamma, n}$$

but

$$\frac{\partial \gamma_1}{\partial b_1} = \lim_{b_2 \rightarrow b_1} \frac{\gamma_1 - \gamma_2}{b_1 - b_2} = \frac{1 - \cosh(b_1 - y) \cos x}{(\cosh(b_1 - y) - \cos x)^2}$$

which is an unlikely candidate for  $P_{\gamma, n}$ . That is, there exists a sequence of  $\gamma$ -polynomials whose limit is not contained in  $P_{\gamma, n}$  and therefore  $P_{\gamma, n}$  is not closed.

This difficulty is circumvented by enlarging the set of  $\gamma$ -polynomials in an appropriate way to include derivatives of the polynomials up to order  $n-1$  (see Hobby and Rice (1967)).

The set extended in this way admits of a best approximation, but the usefulness of the result remains in question. For example, the set of singularities  $b_{j+n}$ ,  $j = 1, \dots, n$  would not normally be restricted to the closed set  $[\alpha, \beta]$ , but would in fact be chosen from the set  $(g(0), \infty)$ .

The above discussion should at least provide some indication of the complexities of a nonlinear approximation problem.

### Computation of a best approximation

Most of the work on  $\gamma$ -polynomials has been confined to a discussion of the characterization of the best approximation in the maximum norm. None of this is applicable to the case at hand, as we have chosen to use the  $L_2$  norm. Nevertheless, some excellent algorithms exist for the solution of the discrete least squares problem. In the examples of the following chapters, it is assumed that a best approximation exists and we attempt to find one using a general nonlinear solver.

The function to be minimized in the least squares sense is  $f(\mathbf{b})$  where

$$f_i(\mathbf{b}) = \phi_n(x_i) - b_0 - \sum_{j=1}^n b_j \gamma(b_{j+n}; x_i)$$

for  $i = 1, \dots, m$ ,  $m \geq 2n + 1$  and  $x_i \in \partial D$ . That is, we wish to find  $\mathbf{b} \in R^{2n+1}$  which minimizes

$$F(\mathbf{b}) = \|\mathbf{f}(\mathbf{b})\|^2 = \mathbf{f}^T \mathbf{f}$$

The minimum of  $F(\mathbf{b})$  is usually found by an iterative procedure. Assume  $\mathbf{b}^k$  is an approximation for the minimum at the  $k^{\text{th}}$  step and find a correction vector  $\mathbf{h}^k$  so that

$$\mathbf{b}^{k+1} = \mathbf{b}^k + \mathbf{h}^k.$$

Expand  $F(\mathbf{b})$  in a Taylor series about  $\mathbf{b}$  (dropping the superscript  $k$ ),

$$F(\mathbf{b} + \mathbf{h}) = e^{\mathbf{h}^T \cdot \nabla} F(\mathbf{b})$$

and retain terms up to second order in the components of  $\mathbf{h}$ . Now,  $\mathbf{h}$  should be chosen to minimize the quadratic expression  $F(\mathbf{b} + \mathbf{h})$  and a necessary condition for this is that

$$\frac{\partial F}{\partial h_j} = 0 \quad , \quad j = 1, 2, \dots, 2n + 1.$$

Hence, the current correction  $\mathbf{h}$  is the solution to the linear system

$$\left( J^T J + \sum_{i=1}^m f_i H^i \right) \mathbf{h} = -J^T \mathbf{f} \quad (2.8.7)$$

where  $J$  is the Jacobian matrix of first derivatives and  $H^i$  is the Hessian matrix of second derivatives,

$$J_{ij} = \frac{\partial f_i}{\partial b_j}$$

$$H_{ij}^k = \frac{\partial^2 f_k}{\partial b_i \partial b_j}$$

Now, most nonlinear solvers are characterized by the way in which the second derivative terms of (2.8.7) are treated. If these terms are omitted altogether, the method is called a Gauss algorithm. It is a first order method only. If the function  $f(\mathbf{b})$  is given explicitly, as it is in our case, the second derivative terms may be computed directly at each step of the iterative process. This can be a time consuming business, but the convergence is now second order.

A method which works very well is the Levenberg-Marquardt algorithm (Levenberg (1944) and Marquardt (1963)). It involves solving for a new direction vector  $\mathbf{h}$  from the linear system

$$(J^T J + \lambda I) \mathbf{h} = -J^T \mathbf{f} \quad (2.8.8)$$

where  $\lambda \geq 0$  is an arbitrary parameter. The inclusion of  $\lambda$  introduces a bias in  $\mathbf{h}$  toward the steepest descent vector  $-2J^T \mathbf{f}$  of the function being minimized,  $F(\mathbf{b})$ . The only difficulty with this method involves a suitable choice of the parameter  $\lambda$ . Apart from this, equations (2.8.8) are seen to be just the normal equations for the residual

$$\mathbf{r} = \begin{bmatrix} J \\ \sqrt{\lambda} I \end{bmatrix} \mathbf{h} + \begin{bmatrix} \mathbf{f} \\ \mathbf{0} \end{bmatrix}.$$

This means that the nonlinear solver reduces, at each iteration, to the solution of a system of linear equations which can be solved in a stable manner using the Householder algorithm outlined in the appendix A.2.

In the computations involving the nonlinear approximation method which appear in the later chapters a Levenberg-Marquardt algorithm is used. In fact, routines were written in Fortran 5 to solve the Gauss method, the more accurate method with the second derivatives retained, and the Levenberg-Marquardt scheme. When tested on a sample potential problem, the Levenberg-Marquardt algorithm proved to be the more robust, provided a "good" choice of  $\lambda$  was used. In view of this difficulty with the choice of  $\lambda$ , a packaged routine was eventually adopted. The MINPACK (Argonne National Laboratory, 1980) routine automatically adjusts the parameter  $\lambda$  as the iterations proceed. The numerical results for the nonlinear approximations presented in chapters 3 and 4 are generated using this MINPACK routine. It is required of the user to supply a subroutine which computes both  $J$  and the residuals  $f_i$ .

The proposed nonlinear scheme, to the best of our knowledge, has not been tested on a potential problem in the manner outlined and certainly not incorporated into a calculation of a moving boundary problem. However, Menikoff and Zemach (1983) do make use of a partial fraction representation similar to (2.8.4) in connection with their calculation of a Rayleigh-Taylor instability. In that paper, they solve the accompanying potential problem via conformal mapping to a canonical domain (the upper half plane) and a Green's function approach. They do not use the partial fraction representation in the time-dependent calculation. Rather, it is used only as a model for approximating interior values of the potential, stream function, velocities etc. once these quantities along the free boundary have already been found by alternative means.

## **CHAPTER 3**

### **Electrochemical Machining Problems**

#### **3.1 Introduction**

Electrochemical machining (ECM) is the name given to the process of machining one metal part by another using electrolysis. The cathode is a shaped tool and the anode is the workpiece from which metal is being eroded. The two electrodes are separated by a small gap which is filled with an electrolyte and an electric potential is applied across the gap (see figure 3.1.1). The electrolyte is pumped through the gap to assist in removal of the eroded material. In addition, the workpiece (anode) is fed toward the tool at a constant rate to maintain a small gap size. Eventually a steady-state condition is reached whereafter the two electrodes maintain a fixed position relative to each other.

There are actually two different processes associated with just the steady-state. The direct free boundary problem involves the determination of the steady anode shape for a given cathode or tool shape. The indirect or inverse free boundary problem involves the determination of the tool shape required to produce a specified steady anode shape. The inverse problem is ill-posed, and in general, the more difficult of the two steady-state problems to solve numerically.

The mathematical description of the unsteady process is a moving boundary problem from the class considered in chapter 1. It is one of the simplest MBP from our class. The



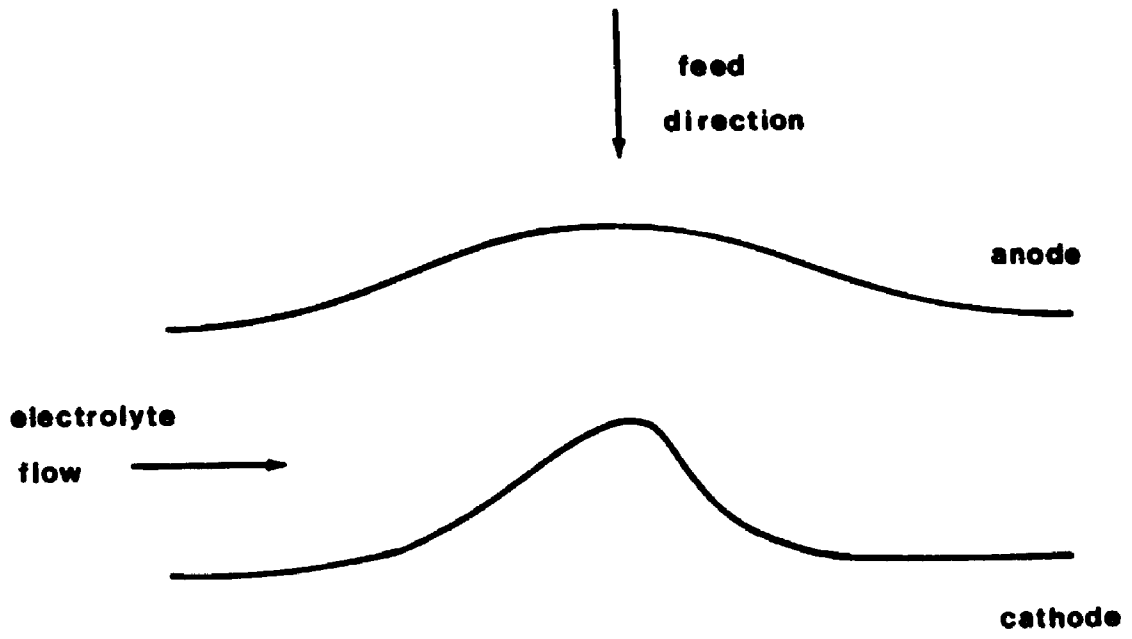


Fig. 3.1.1 Definition: Sketch

potential is constant along the boundaries and the problem is well posed. In addition, the steady state conditions can be ascertained a priori, given suitable initial conditions. This provides an excellent example for illustrating the full time dependent numerical scheme.

In this section, we present the precise mathematical description for two-dimensional ECM and in the following section, previous work on the problem is reviewed. Section 3.3 addresses the potential problem related to the steady-state problem. Since it is possible to derive an analytic solution, this provides a quantitative means of testing the linear approximation method discussed in chapter 2, before incorporating it into a time dependent calculation. An example of the full time dependent ECM problem is treated numerically in section 3.4. The numerical solution of the evolution equation is outlined there.

The potential for application of boundary approximation techniques to the difficult inverse problems cannot be overlooked. In section 3.5, we briefly examine the numerical solution of the inverse ECM problem.

### Mathematical Description of ECM

The complete description of the machining problem must derive from Maxwell's equations governing the electromagnetic field. McGeough and Rasmussen (1974) provide the details. Briefly, the curl of the electric field is taken to be zero so that we may define a scalar electric potential  $\phi$  such that

$$\mathbf{E} = \nabla\phi$$

which together with

$$\nabla \cdot \mathbf{E} = 0$$

gives

$$\nabla^2\phi = 0 \tag{3.1.1}$$

between the two electrodes. The potential is prescribed on the two electrodes,

$$\phi = V \quad \text{on the anode} \tag{3.1.2a}$$

$$\phi = 0 \quad \text{on the cathode} \tag{3.1.2b}$$

In addition, an equation relating the dissolution rate of the anode and the electric field is necessary. Faraday's law and Ohms law together yield

$$\frac{d\mathbf{x}}{dt} = M\nabla\phi - \alpha\mathbf{j} \quad (3.1.3)$$

where  $\mathbf{x}$  is a point on the anode surface and  $M$  is a constant referred to as the electrochemical machining constant. The term  $\alpha\mathbf{j}$  is due to the fact that the anode is fed toward the cathode at a constant rate  $\alpha$ .

Finally, we assume periodicity in the  $x$ -direction, the period being  $2L$ . The description can be nondimensionalized as follows. Let

$$\begin{aligned} x^* &= \frac{x\pi}{L} & \phi^* &= \frac{\phi}{V} & \alpha^* &= \frac{\alpha L}{\pi M V} \\ y^* &= \frac{y\pi}{L} & t^* &= \frac{\pi^2 M V t}{L^2} \end{aligned}$$

The equations (3.1.1), (3.1.2) and (3.1.3) become (dropping the starred notation):

$$\nabla^2\phi = 0 \quad (3.1.4)$$

$$\phi = 1 \quad \text{on the anode} \quad (3.1.5)$$

$$\phi = 0 \quad \text{on the cathode} \quad (3.1.6)$$

$$\frac{d\mathbf{x}}{dt} = \nabla\phi - \alpha\mathbf{j} \quad \text{on the anode} \quad (3.1.7)$$

together with an assumed  $2\pi$ -periodicity in the  $x$ -direction.

The equation (3.1.7) yields immediately the two equations

$$\frac{dx}{dt} = \phi_x \quad (3.1.8)$$

$$\frac{dy}{dt} = \phi_y - \alpha \quad (3.1.9)$$

representing the velocity of the  $x, y$  coordinates of a point on the anode surface.

Alternatively, we may express the evolution equation in the form:

$$g_t = \phi_y - g_x \phi_x - \alpha \quad (3.1.10)$$

obtained from the total derivative of  $y = g(x, t)$  with respect to  $t$  and using equations (3.1.8) and (3.1.9). (The function  $g$  represents the anode profile.)

Yet another form of the free surface equation is possible. From  $\phi = 1$  on  $y = g(x)$  and on taking differentials,

$$\begin{aligned} -\phi_t &= \phi_x \frac{dx}{dt} + \phi_y \frac{dy}{dt} \\ -\phi_t &= \phi_x^2 + \phi_y^2 - \alpha \phi_y. \end{aligned}$$

This can be expressed succinctly in terms of the complex potential,  $w = \phi + i\psi$ . That is, we may write

$$\Re\left(\frac{\partial w}{\partial t}\right) = -\left|\frac{\partial w}{\partial z}\right|^2 - \alpha \Im\left(\frac{\partial w}{\partial z}\right) \quad (3.1.11)$$

### 3.2 Review

The early work on ECM was applied to the direct time independent problem and almost always pertained to two-dimensional problems. Collett, Hewson-Browne and Windle (1970) obtained some semi-analytic solutions in series form for restricted cathode tool shapes. Hewson-Browne (1971) extended the method to cover a wider class of two-dimensional tool shapes. A complex variable technique was applied by Nilson and Tsuei (1975) and again a similar approach was considered by Hougard (1977). Both of these latter techniques involved an iterative numerical scheme for computing the required anode shape by working in the plane of the complex potential as opposed to the physical plane. The transformed domain is rectangular and consequently many of the difficulties associated with the free boundary are avoided.

Sloan (1986) also makes use of a coordinate transformation. The physical domain is mapped onto a square and the transformed equations discretized by finite differences. The result is a system of nonlinear equations which are treated by a Newton like iteration.

The full time dependent ECM problem, as formulated by Mcgeough and Rasmussen (1974) has since been treated by a variety of numerical techniques. Some two-dimensional annular machining problems were computed by Christiansen and Rasmussen (1976) using an integral equations approach. A Kantorovich method was used by Forsyth and Rasmussen (1979) to compute some two dimensional problems. Each of these methods is front-tracking. That is, the shape of the workpiece is followed from time step to time step via the evolution equation for the free surface. The potential field for a new domain is computed at each step of the process.

A finite element approximation of a variational inequality formulation is provided in Elliott (1980) (The existence and uniqueness of the solution of the general ECM problem are included in this paper). The approach has considerable appeal since it provides a fixed domain method for the computation of a moving boundary problem. Forsyth and Rasmussen (1980) use a similar derivation of the variational inequality method and compare their numerical results with a conventional front-tracking calculation (Forsyth and Rasmussen (1979)). If a high degree of accuracy is required, the two approaches involve comparable computing times. However, the variational inequality requires a zero feed rate which represents a significant disadvantage in planar machining problems. Furthermore, the variational approach is not readily adapted to electroforming problems, unlike the front-tracking methods. They conclude that the front-tracking approach is the more flexible of the two and consequently is the preferred method.

The two dimensional inverse problem of ECM is solved exactly by Krylov (1968), for some simple geometries. Nilson and Tsuei (1974) present a general method for obtaining exact solutions to the two dimensional case. They go on (Nilson and Tsuei (1976)) to provide a numerical method to be used when the data is available only in discrete

form. Lacey (1985) provides another analytic approach to the inverse problem, which the author then extends to the three dimensional case. Finally, we mention that there are a number of other areas of applied mechanics where inverse problems arise having a similar mathematical description. We note in passing, the applications of porous cooling and radiation melting treated by Goldstein and Siegel (1970) and Siegel (1973).

### 3.3 Steady State ECM

The example which we consider is related to a steady state electrochemical machining problem for which an analytic solution exists.

Let  $y = g(x)$  and  $y = h(x)$  be the anode and cathode surfaces, respectively, of an electrolytic cell. Then the steady state electrochemical machining problem consists of finding the electric potential  $\phi(x, y)$  and either  $g(x)$  or  $h(x)$  such that

$$\begin{aligned} \nabla^2 \phi &= 0 & , & \quad h(x) < y < g(x) \\ \phi &= 0 & , & \quad y = h(x) \\ \phi &= 1 & , & \quad \phi_y - g_x \phi_x = \alpha & , & \quad y = g(x) \end{aligned} \quad (3.3.1)$$

$h(x)$  is given in the direct steady-state problem whereas  $g(x)$  is known in the inverse problem.

Suppose that  $g(x)$  is given. Then, together with the two conditions along  $y = g(x)$  we have a Cauchy problem for  $\phi$ , which can be readily solved by a complex variables approach (see Nielson and Tseui (1974)). The second condition can be expressed more compactly as

$$\psi = -\alpha x$$

where  $\psi$  is the harmonic conjugate to  $\phi$ . This follows from the steady form of the free surface condition, (3.1.11). If we interchange the dependent and independent variables by mapping to the plane of the complex potential,

$$w = \phi + i\psi$$

then both  $x$  and  $y$  are known functions of  $\psi$  along  $\phi = 1$ . By direct analytic continuation in the  $w$ -plane, we have

$$z = \frac{i(w-1)}{\alpha} + ig\left(i\frac{(w-1)}{\alpha}\right)$$

For example, if

$$g(x) = a \cos x + \frac{1}{\alpha}$$

then we have

$$z = i\left[\frac{w}{\alpha} + a \cosh\left(\frac{w-1}{\alpha}\right)\right] \quad (3.3.2)$$

Along  $\phi = 0$ , then, we have

$$h(\psi) = a \cosh\left(\frac{1}{\alpha}\right) \cos\left(\frac{\psi}{\alpha}\right) \quad (3.3.3)$$

$$x(\psi) = -\frac{\psi}{\alpha} + a \sinh\left(\frac{1}{\alpha}\right) \sin\left(\frac{\psi}{\alpha}\right).$$

In this section, we assume that both  $h(x)$  and  $g(x)$  are given as above and then we calculate  $\phi$ . It is an easy matter, then, to check the accuracy of the solution obtained by the boundary approximation method.

The approximate solution may be written in the form

$$\phi_n = b_p + b_{2p}y + \sum_{j=1}^{p-1} [b_j \sinh(jy) + b_{j+p} \cosh(jy)] \cos(jx) \quad (3.3.4)$$

where  $n = 2p$ . Because of the symmetry, we need only consider  $0 \leq x \leq \pi$ . Along  $y = g(x)$  we use the uniform  $x$ -grid

$$0 = x_1 < x_2 < \dots < x_{\frac{n}{2}} = \pi.$$

Along  $y = h(x)$  we use a uniform  $\psi$  grid, which because of equation (3.3.3) leads to a nonuniform  $x$  grid

$$0 = x_{\frac{n}{2}+1} < x_{\frac{n}{2}+2} < \dots < x_m = \pi.$$

With this notation,  $n$  represents the total number of unknowns and  $m$  is the total number of boundary points.

The results are mainly presented in the form of the root mean square error between the calculated solution and the analytic solution as given by (3.3.2); we compare not only  $\phi$  but also  $\phi_x$  and  $\phi_y$ . We define the root mean square error

$$E_\phi = \left[ \frac{1}{m} \sum_{i=1}^m (\phi_n(x_i) - \phi(x_i))^2 \right]^{\frac{1}{2}}$$

where  $x_i$  is a point on  $y = g(x)$  or  $y = h(x)$ , and  $\phi(x_i)$  is the exact boundary value of  $\phi$  at  $x_i$ . Expressions for the root mean square error in  $\phi_x$  and  $\phi_y$  are defined similarly with the computed values of  $\phi_x$  and  $\phi_y$ , determined from term by term differentiation.

We also calculate the maximum error defined by

$$E_{\max} = \max_{x_i} |\phi_n(x_i) - \phi(x_i)|$$

where the points  $x_i$  are the  $m$  boundary points with an additional one hundred points on the boundary at uniform intervals of  $x$ . These additional points are not used in the determination of the series coefficients. They serve only to test the usefulness of  $\phi_n$  as a functional approximant.

An interesting feature of the exact solution (3.3.2) from the computational point of view is the existence of singularities outside of the domain. We have shown in chapter 2 that the rate of convergence should decrease as a singularity moves closer to the boundary of the solution domain, and our numerical calculations confirm this. There are singularities in the solution, corresponding to the zeros of  $\frac{dz}{dw}$ . These are given by



$$\phi = 1 - \alpha \sinh^{-1}\left(\frac{1}{a}\right), \quad \psi = 0$$

$$\phi = 1 + \alpha \sinh^{-1}\left(\frac{1}{a}\right), \quad \psi = \pm\pi.$$

If  $\alpha$  decreases or  $a$  increases, the singularities move closer to  $\phi = 1$ . When  $\alpha = 1/\sinh^{-1}\left(\frac{1}{a}\right)$ ,

the singularity corresponding to  $\psi = 0$  resides on the cathode  $y = h(x)$ . For  $\alpha > 1/\sinh^{-1}\left(\frac{1}{a}\right)$ , the singularities remain outside of the domain of interest and their position in the  $x$ - $y$  plane is given by

$$y = \frac{1}{\alpha} - \sinh^{-1}\left(\frac{1}{a}\right) + \sqrt{1 + a^2}, \quad x = 0$$

$$y = \frac{1}{\alpha} + \sinh^{-1}\left(\frac{1}{a}\right) - \sqrt{1 + a^2}, \quad x = \pm\pi.$$

The linear system was solved using a singular value decomposition to minimize the effects of ill conditioning on the computed solution. However, when the computations were repeated with the Householder algorithm and with  $m \geq 2n$ , no difference was detected in the results to the significant figures shown in the tables.

In table 3.3.1, we use  $\alpha = 2.5$  and  $a = 0.5$ , which places the closest singularity a vertical distance of 3.52 from  $h(x)$ . The results in the table show that for  $n$  fixed ( $n = 50$ ), an increase in  $m$  gives a rapid decrease in the errors for  $\phi_x$  and  $\phi_y$ . However, we see that  $m:n = 2:1$  is sufficient. As expected, the errors in  $\phi_x$  and  $\phi_y$  are larger than the errors in  $\phi$ . The results of the collocation approximation are unacceptable and may well be suffering from the effects of severe ill conditioning.

Table 3.3.1 The Ratio  $m:n$ 

$$\alpha = 2.5, a = 0.5$$

$n$	$m$	$E_{\max}$	$E_{\phi}$	$E_{\psi}$	$E_{\eta}$	Condition Number
50	50	0.283 E+00	0.382 E-08	0.181 E+01	0.787 E+00	6.56 E+11
50	60	0.104 E+00	0.809 E-05	0.624 E+00	0.129 E+00	2.82 E+10
50	80	0.845 E-04	0.279 E-04	0.493 E-03	0.507 E-03	1.63 E+09
50	100	0.839 E-04	0.278 E-04	0.495 E-03	0.503 E-03	1.62 E+09
50	150	0.834 E-04	0.277 E-04	0.495 E-03	0.501 E-03	1.53 E+09

Another important issue is the convergence of the approximation with increasing number of terms in the expansion. For table 3.3.2, we used  $\alpha = 3.0$  and  $a = 0.25$  which places the singularity further away from  $h(x)$  (at a distance of 4.52) than in table 3.3.1. A constant ratio of 2:1 ( $m:n$ ) was maintained and the corresponding results for collocation are given simultaneously. In both cases, the convergence of the method is clearly demonstrated. For all choices of  $n$  the overdetermined system performs better than collocation. As  $n$  grows larger, the condition number of the matrix equation is enormous and the results of collocation in particular may become suspect. The condition number of the overdetermined system is in all cases less than that of the corresponding collocation system.

**Table 3.3.2 Convergence for Increasing n**

$$\alpha = 3.0, a = 0.25$$

n	m	$E_{max}$	$E_{\phi}$	$E_{\psi}$	$E_{\eta}$	Condition Number
10	20	0.248 E-02	0.107 E-02	0.485 E-02	0.612 E-02	1.23 E+01
	10	0.742 E-02	0.232 E-13	0.121 E-01	0.479 E-02	2.42 E+01
20	40	0.172 E-04	0.705 E-05	0.642 E-04	0.675 E-04	2.50 E+02
	20	0.285 E-03	0.496 E-13	0.804 E-03	0.224 E-03	1.72 E+03
30	60	0.167 E-06	0.687 E-07	0.941 E-06	0.966 E-06	6.75 E+03
	30	0.254 E-04	0.902 E-13	0.102 E-03	0.233 E-04	2.15 E+05
40	80	0.198 E-08	0.793 E-09	0.145 E-07	0.148 E-07	1.89 E+05
	40	0.927 E-06	0.957 E-13	0.538 E-05	0.979 E-06	1.08 E+07
50	100	0.254 E-10	0.100 E-10	0.230 E-09	0.233 E-09	5.37 E+06
	50	0.492 E-08	0.100 E-12	0.420 E-07	0.632 E-08	1.30 E+08
60	120	0.343 E-12	0.174 E-12	0.376 E-11	0.383 E-11	1.54 E+08
	60	0.348 E-09	0.137 E-12	0.238 E-08	0.536 E-09	8.16 E+09

In table 3.3.3, we present some results showing the effects of the proximity of the singularity to  $h(x)$ . With  $n = 30$ ,  $m = 60$  and  $a = 0.25$ , we calculated solutions for three different feedrates  $\alpha$ . The results show that the accuracy of the solution deteriorates as the singularity moves closer to  $y = h(x)$ . This is in agreement with our error estimates of chapter 2.

**Table 3.3.3 Effects of Proximity of Singularity**

$$a = 0.25, n = 30, m = 60$$

Feedrate $\alpha$	$\Delta y$	$E_{\max}$	$E_{\phi}$	$E_{\psi}$	$E_{\eta}$	Condition Number
1.0	0.450	0.333 E-03	0.117 E-03	0.139 E-02	0.143 E-02	1.60 E+07
2.0	2.44	0.177 E-05	0.449 E-06	0.603 E-05	0.618 E-05	4.94 E+04
3.0	4.52	0.167 E-06	0.687 E-07	0.941 E-06	0.966 E-06	6.75 E+03

$\Delta y$  represents the vertical distance of the singularity from  $h(0)$

The series coefficients for the two cases  $\alpha = 1.0, 3.0$  are presented in table 3.3.4. As expected, they show a steady decrease in magnitude as the series index increases. The convergence is somewhat stronger in the case  $\alpha = 3.0$ , where the singularities are further from the boundary. For large series index  $j$ , the coefficients of the hyperbolic sine ( $b_j$ ) and the hyperbolic cosine  $b_{j+p}$  are very close in magnitude and opposite in sign. This is quite apparent in the case  $\alpha = 1.0$ . It is, of course, a consequence of the nearness in magnitude of the two hyperbolic functions for large argument. If the approximate form

$$\phi_n = c_p + c_{2p}y + \sum_{j=1}^{p-1} [c_j e^{jy} + c_{j+p} e^{-jy}] \cos(jx)$$

had been used instead of (3.3.4), then

$$c_j = \frac{1}{2}(b_j + b_{j+p}) \quad , \quad c_{j+p} = -\frac{1}{2}(b_j - b_{j+p}).$$

In this case, the coefficients of the growing exponential (if  $y > \lambda$ ) are very small, as is necessary for rapid convergence. In practice, the matrices proved to be better conditioned when the form (3.3.4) was used.

This behaviour in the coefficients does, however, raise an interesting point. Could we manage with fewer terms in the expansion when summing the series? That the larger number of unknowns gives better convergence is clear from the table 3.3.2; but could we truncate the series  $\phi_n$  if it is to be re-evaluated at some  $(x,y)$  point or if the derivatives of the series are to be used? It was indeed found, that for some  $(x,y)$  values this was the case. However, the difficult questions of where to truncate the series and for what values of  $(x,y)$ , made this process impossible to automate.

**Table 3.3.4 The Coefficients (n=30, p=15)**

$\alpha = 1.0$	$a = 0.25$	$\alpha = 3.0$	$a = 0.25$
$b_1 = 0.296E + 00$	$b_{16} = -0.389E + 00$	$b_1 = 0.257E + 00$	$b_{16} = -0.798E + 00$
$b_2 = 0.120E + 00$	$b_{17} = -0.116E + 00$	$b_2 = 0.118E + 00$	$b_{17} = -0.686E - 01$
$b_3 = 0.608E - 01$	$b_{18} = -0.611E - 01$	$b_3 = 0.214E - 01$	$b_{18} = -0.281E - 01$
$b_4 = 0.374E - 01$	$b_{19} = -0.374E - 01$	$b_4 = 0.833E - 02$	$b_{19} = -0.725E - 02$
$b_5 = 0.251E - 01$	$b_{20} = -0.251E - 01$	$b_5 = 0.260E - 02$	$b_{20} = -0.279E - 02$
$b_6 = 0.179E - 01$	$b_{21} = -0.179E - 01$	$b_6 = 0.101E - 02$	$b_{21} = -0.971E - 03$
$b_7 = 0.130E - 01$	$b_{22} = -0.130E - 01$	$b_7 = 0.375E - 03$	$b_{22} = -0.382E - 03$
$b_8 = 0.950E - 02$	$b_{23} = -0.950E - 02$	$b_8 = 0.150E - 03$	$b_{23} = -0.149E - 03$
$b_9 = 0.665E - 02$	$b_{24} = -0.665E - 02$	$b_9 = 0.601E - 04$	$b_{24} = -0.604E - 04$
$b_{10} = 0.425E - 02$	$b_{25} = -0.425E - 02$	$b_{10} = 0.246E - 04$	$b_{25} = -0.245E - 04$
$b_{11} = 0.233E - 02$	$b_{26} = -0.233E - 02$	$b_{11} = 0.981E - 05$	$b_{26} = -0.982E - 05$
$b_{12} = 0.101E - 02$	$b_{27} = -0.101E - 02$	$b_{12} = 0.359E - 05$	$b_{27} = -0.358E - 05$
$b_{13} = 0.301E - 03$	$b_{28} = -0.301E - 03$	$b_{13} = 0.103E - 05$	$b_{28} = -0.103E - 05$
$b_{14} = 0.457E - 04$	$b_{29} = -0.457E - 04$	$b_{14} = 0.168E - 06$	$b_{29} = -0.168E - 06$
$b_{15} = -0.484E - 05$	$b_{30} = 0.100E + 01$	$b_{15} = -0.358E - 08$	$b_{30} = 0.300E + 01$

### 3.4 Unsteady ECM

An example of the unsteady ECM problem is provided in this section. The numerical solution of (3.1.4)-(3.1.7) proceeds in a step-wise fashion as follows. The anode and cathode profiles are assumed given initially and the new anode position after a time  $\Delta t$  is determined by numerically integrating the evolution equation in the form (3.1.10). The

potential problem is then solved by the boundary approximation method at the new time (using (3.3.4)) and the entire process is repeated. Since we have the opportunity to express derivatives with respect to  $x$  and  $y$  at given points explicitly in terms of known functions at those points, it is preferable to approximate the evolution equation by a system of ordinary differential equations. Thus, if the  $x$  variable is discretized

$$-\pi \leq x_1 \leq \dots \leq x_{m-1} \leq x_m \leq \pi \quad ,$$

the equations (3.1.8), (3.1.9) represent a coupled system of ordinary differential equations for each point  $(x_i(t), y_i(t))$  on the free surface. Alternatively, (3.1.10) represents a system of ordinary differential equations (uncoupled) for each coordinate  $g(x_i, t)$  on the free surface. In either case, the evolution equation may be approximated by the system of ordinary differential equations having the form

$$\mathbf{u}_i = \mathbf{f}(\mathbf{u}, t) \quad (3.4.1)$$

We have chosen to solve this system using the following predictor-corrector format:

$$u_i^{(n+1)} = u_i^{(n-1)} + 2\Delta t f_i^{(n)} \quad (3.4.2)$$

$$u_i^{(n+1)} = u_i^{(n)} + \frac{\Delta t}{2} [f_i^{(n)} + f_i^{(n+1)}] \quad (3.4.3)$$

where  $u_i^{(k)}$  represents the solution at  $x_i$  and after  $k$  time steps with  $\Delta t = t^{(k+1)} - t^{(k)}$ . The second equation resembles the Crank-Nicholson method for iteratively solving (3.4.1), but with the starting values predicted according to (3.4.2). (3.4.2) is called the Midpoint formula and (3.4.3) is the Trapezoidal formula. As we shall see in chapter 4, the ordinary differential equations may well be stiff and use of an implicit scheme is therefore recommended. The second order trapezoidal method is well suited for solving such unstable systems.

A variable stepsize  $\Delta t$  is incorporated in an attempt to control the truncation error  $T$ , of the scheme. The truncation error of both predictor and corrector formulae is  $O(\Delta t)^2$

for fixed  $\Delta t$ , and one could simply perform repeated applications of, or iterations on, the equation (3.4.3) until a desired error tolerance is met. This could prove costly, however, as repeated evaluation of the functions  $f_i$  is time consuming. It is preferable to adjust the stepsize in such a fashion that one application each of (3.4.2) and (3.4.3) produces a satisfactory truncation error. The algorithm which accomplishes this is briefly outlined below. The details may be found in Atkinson (1978).

An error tolerance  $\tau$  is specified and the stepsize  $\Delta t$  at each stage is chosen to satisfy the condition

$$\frac{\tau \Delta t}{4} \leq |T| \leq \tau \Delta t \quad (3.4.4)$$

The calculation proceeds with a given stepsize  $\Delta t$ , making one application of the predictor (3.4.2) and one application of the corrector (3.4.3) until the condition (3.4.4) is violated. At this stage, a new stepsize is chosen such that the error tolerance is maintained and the calculation is then restarted. Since the predictor equation requires two starting values (ie. solutions at two prior time steps), we must have a process of restarting the calculation when the stepsize is altered. The single step Euler method

$$u_i^{(n+1)} = u_i^{(n)} + \Delta t f_i^{(n)}$$

is used to restart the calculation. It is of lower order ( $O(\Delta t)$ ) and two iterations in (3.4.3) are necessary to meet the required error tolerance (3.4.4) (again, see Atkinson (1978) for the details).

The first time dependent example to be considered has the steady state solution given by equation (3.3.2). That is, the cathode tool shape is given by



$$x(\psi) = -\frac{\psi}{\alpha} + \alpha \sinh\left(\frac{1}{\alpha}\right) \sin\left(\frac{\psi}{\alpha}\right)$$

$$h(\psi) = a \cosh\left(\frac{a}{\alpha}\right) \cos\left(\frac{\psi}{\alpha}\right).$$

Initially, the anode workpiece is a flat plate represented by  $g(x, 0) = \text{constant}$ . The exact steady profile is given by

$$g = a \cos x + \frac{1}{\alpha}.$$

Three cases, corresponding to different  $\alpha, a$ , and already examined at steady state in the last section, are presented in figures 3.4.1, 3.4.2, and 3.4.3. In each case, we have taken  $n = 30, m = 60$  and the amplitude  $a = 0.25$ . In figure 3.4.1, the feed rate is  $\alpha = 1.0$ . The initial anode profile is  $g(x, 0) = 1.5$ . The lower curve represents the given cathode profile and the upper curves correspond to the anode shapes at selected times. The maximum error  $E_{\max}$  between the computed time dependent profile and the exact steady profile was monitored at each time step. A steady reduction in  $E_{\max}$  was observed. After a time of about  $t=3.2$ , the value of  $E_{\max}$  was  $0.638 \text{ E-}01$  and by  $t=10.0, E_{\max} = 0.892 \text{ E-}04$ . The actual computing cost for the run to  $t=10.0$  was approximately 344 seconds of CPU time.

In figure 3.4.2,  $\alpha = 2.0$  and the initial profile is  $g(x, 0) = 1.0$ . The value of  $E_{\max}$  was observed to be  $0.143 \text{ E-}02$  by a time of 2.0. By  $t = 4.0$ , the difference between computed and exact profiles was only  $0.609 \text{ E-}06$ . The actual computing cost for a run to  $t = 4.0$  was approximately 140 seconds of CPU time.

Finally, the case of feed rate  $\alpha = 3.0$  is shown in figure 3.4.3. With this value of  $\alpha, E_{\max} = 0.363 \text{ E-}06$  already by a time of  $t = 2.0$ . The value is not significantly improved by running the program longer, the limitations in accuracy being governed by the spacial and temporal errors committed.

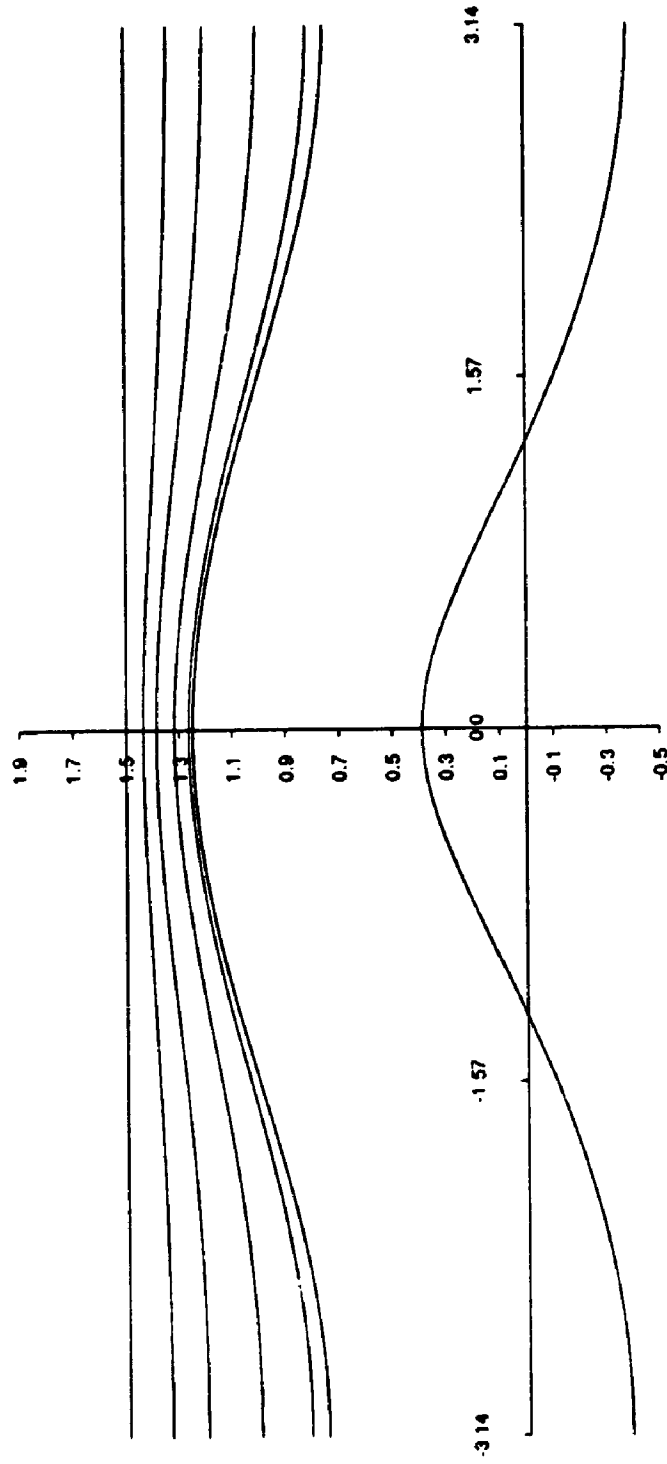


Fig. 3.4.1 Time Evolution of ECM ( $\alpha = 1, a = 0.25$ )

$t = 0.00, 0.40, 0.80, 1.60, 3.20, 10.00$

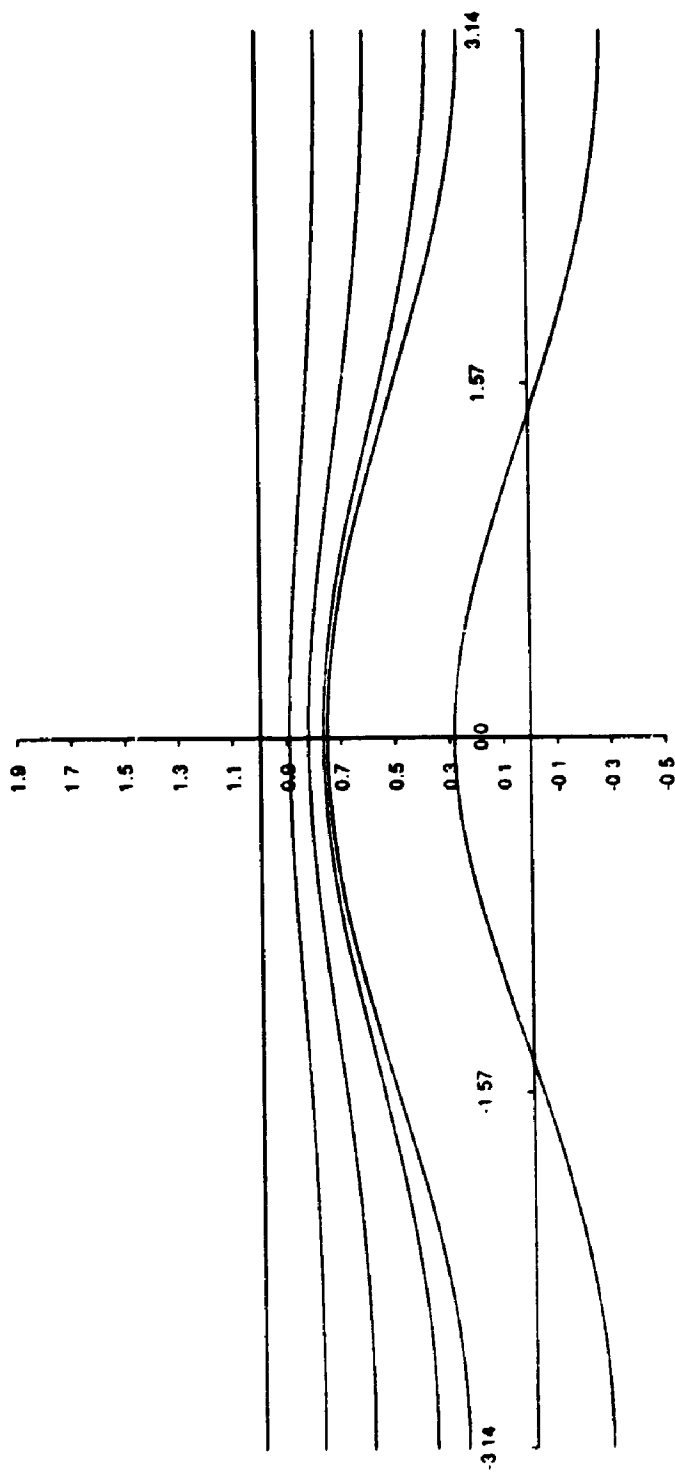


Fig. 3.4.2 Time Evolution of ECM ( $\alpha = 2, a = 0.25$ )

$t = 0.00, 0.20, 0.40, 0.80, 4.00$

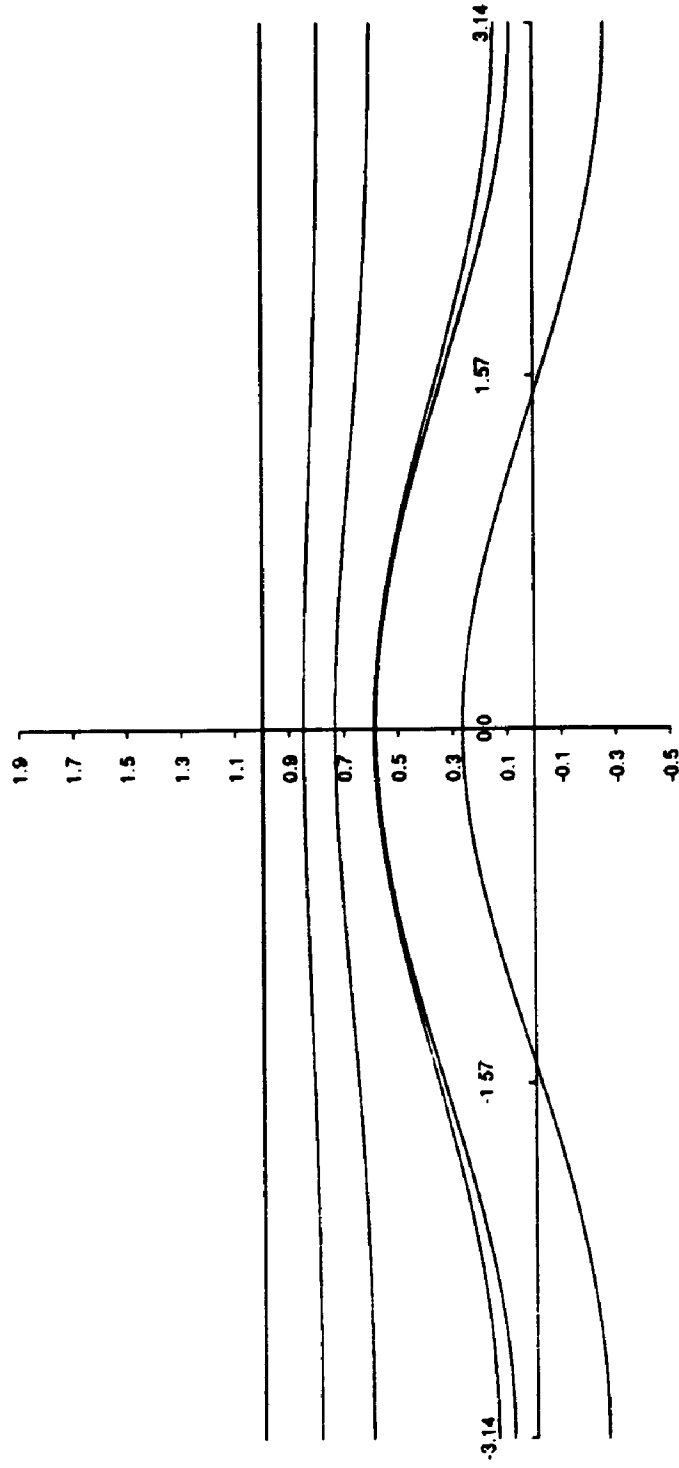


Fig. 3.4.3 Time Evolution of ECM ( $\alpha = 3, a = 0.25$ )

$t = 0.00, 0.10, 0.20, 0.60, 2.00$

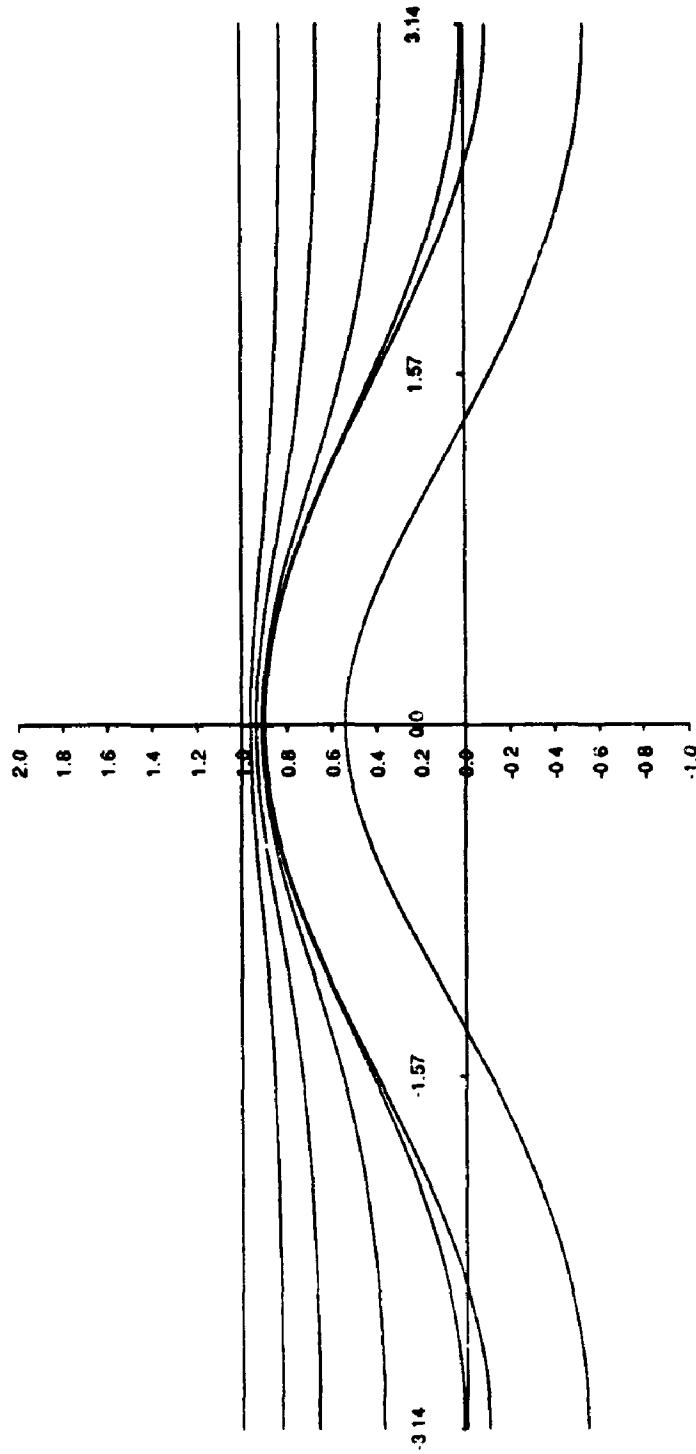


Fig. 3.4.4 Time Evolution of ECM ( $\alpha = 2.5, a = 0.50$ )

$t = 0.00, 0.10, 0.20, 0.40, 0.80, 3.00$

In each of these cases, the error tolerance in the o.d.e. algorithm was  $\tau = 0.01$ . A minimum time step of  $h_{\min} = 0.001$  was imposed. For the two cases,  $\alpha = 1.0, 2.0$ , the maximum time step chosen was  $h_{\max} = 0.05$  and in both cases, the calculation proceeded to the times shown without requiring a reduction in step size. The same values of  $g(x)$  to five decimal places were observed in each case when the maximum time step was reduced by a factor of two. On the other hand, when  $h_{\max}$  was increased by a factor of two, oscillations developed early in the anode profile and persisted through later times in the calculation, clearly showing up the dangers of employing too large a time step.

For the case  $\alpha = 3.0$ , it was necessary to reduce the maximum time step to  $h_{\max} = 0.025$ . This was reasonable, in view of the fact that the larger feed rate produced a more rapid change in anode position.

In figure 3.4.4, the same example is recomputed for a greater value of  $a$ . With  $a = 0.5$  and  $\alpha = 2.5$ , the aspect ratio is close to the allowable limits for this problem and this type of numerical scheme. (We define the aspect ratio as the ratio of the change in  $y$  to the change in  $x$  for the cathode profile.) Twenty-four terms in the series ( $n = 50$ ) and  $m = 100$  boundary points were used. The initial anode shape was taken to be  $g(x, 0) = 1.0$ . By a time of  $t = 3.0$ , the value of  $E_{\max}$  was  $0.343 \text{ E-}05$ . The results are very good, but the computing effort was substantially greater, owing to the elevated cost of solving a large linear system at each time step. The execution time for a complete run to  $t = 4.0$  was 995 seconds of CPU time.

Each of the above examples has been compared with an exact steady state solution. However, when an analytic solution is not available for comparison, another measure of numerical accuracy must be found. Let  $E_M$  and  $E_R$  be the maximum error and the root mean square error, respectively, between the computed value of  $\phi_n$  and the given boundary data. It is a simple matter to calculate either or both of these quantities every time step, and in this way, the accuracy of the potential model can be monitored. Bear in mind, that for a

Dirichlet problem with known boundary data, a maximum principle is in effect. That same principle applies to the boundary approximation and so the maximum error in the computed results is attained along the boundary.

For example, in the calculation with  $a = 0.25$  and  $\alpha = 3.0$ , the value of  $E_R$  after one time step was  $0.682 \text{ E-}07$  and never exceeded  $0.940 \text{ E-}07$  over the entire run. For the case  $a = 0.5$  and  $\alpha = 2.5$ , the initial value of  $E_R$  was  $0.329 \text{ E-}04$  and remained on the same order for the entire calculation. Of course, this is a very stable problem and this sort of consistency should be expected of any respectable numerical routine. However, in an unstable physical problem, simply monitoring the growth of  $E_M$  or  $E_R$  will at least give an indication of when the results of the potential calculation can no longer be trusted.

As a final test of the limits of the scheme, we have taken the cathode tool shape

$$h(x) = \exp\left(-4 \sin^2\left(\frac{x}{2}\right)\right).$$

The aspect ratio is comparable to the last example, but the profile is more sharply peaked near  $x = 0$ .

In figure 3.4.5, the machined profiles are shown for a number of time steps. Without an exact solution for comparison, the accuracy remains in question. Certainly, the computed profiles are smooth and exhibit the correct qualitative behaviour. The root mean square error between the computed solution and the correct boundary data was observed to be on the order of  $10^{-2}$  throughout the entire calculation. The maximum difference between the last profile shown in the figure and the computed profile one time step earlier, is less than 0.001, indicating that steady state has effectively been reached. Running the program again with a larger value of  $n$  would allow for further comparison. However, we are already close to the limits attainable for this problem. With  $n = 40$  for example, the effects of ill conditioning are manifested in small corrugations in the free surface profile.

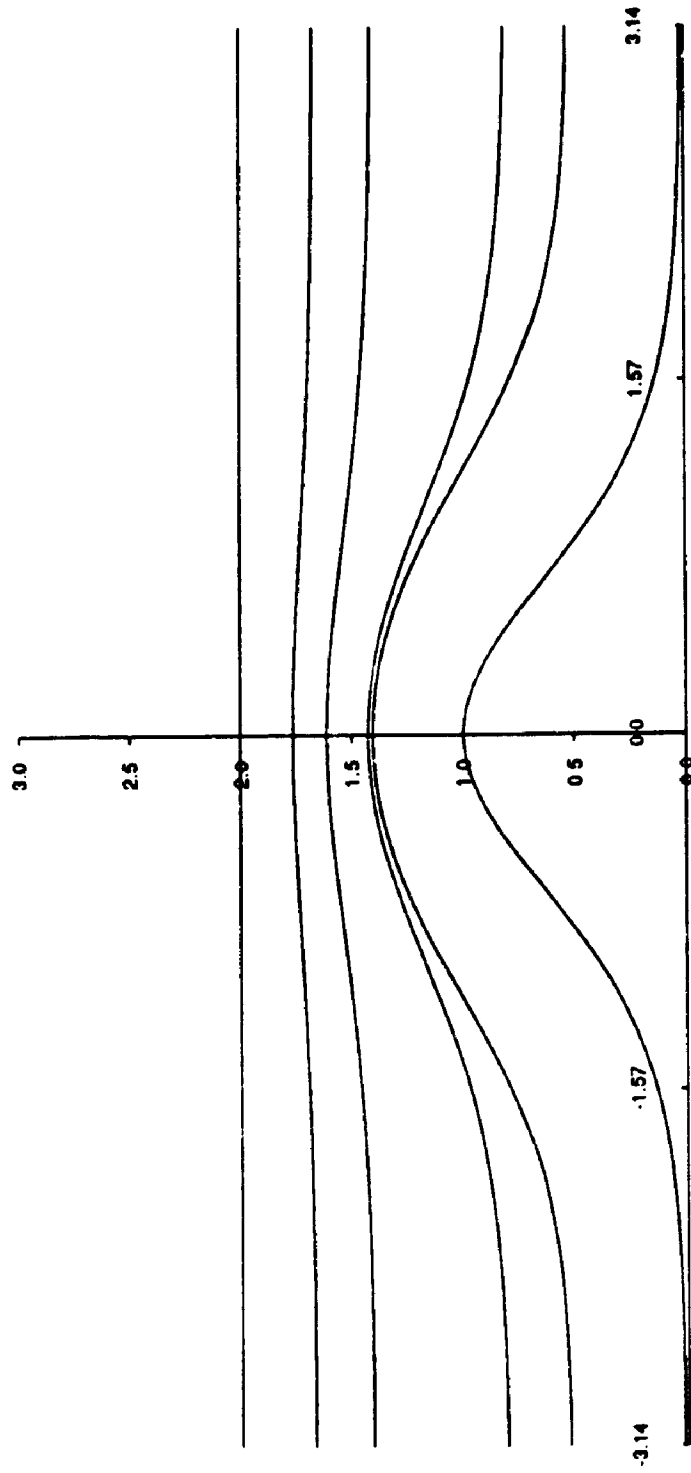


Fig. 3.4.5 Time Evolution of ECM ( $\alpha = 2, h(x) = \text{EXP}(-4 \sin^2(x/2))$ )

$t = 0.00, 0.23, 0.43, 1.03, 4.03$



### 3.5 The Inverse ECM Problem

The inverse problem is perhaps the more significant of the electrochemical machining problems. A specified shape is to be produced by machining and we must determine the shape of the tool required to accomplish this. In this section, we outline two methods we have developed for the inverse problem and report the results of the superior scheme.

The situation is described by equations (3.3.1), where this time, the curve  $y = g(x)$  is specified and  $y = h(x)$  is to be determined. Two boundary conditions are given along the curve  $y = g(x)$ . This is a Cauchy problem and is ill-posed in the sense that small changes in the given data can lead to large changes in the tool shape. What is more, for given  $g(x)$ , there is an infinite family of admissible tool shapes, all of which can be used to machine the desired cathode shape. We will seek the cathode shape corresponding to the zero potential.

This problem has been solved analytically in two dimensions by Nilson and Tsuei (1974) and by Lacey (1985). The latter approach requires knowledge of a special function, the Schwarz function, which is not always available for the given anode geometry. Even at that, it is not always the case that the anode shape can be given in analytic form. We may have only a table of data representing the workpiece profile. In these circumstances, a reliable numerical approach is necessary.

The two conditions,

$$\phi = 1 \quad , \quad \psi = -\alpha x \quad (3.5.1)$$

along  $y = g(x)$  provide given information about the analytic function  $w = \phi + i\psi$  along a curve. In the usual way, then, we could expand  $w$  in series form. The truncated series for the real and imaginary parts of  $w$  take the form

$$\phi_n = h_p + h_{2p} + \sum_{j=1}^{p-1} [h_j \sinh(jy) + h_{j+p} \cosh(jy)] \cos(jx) \quad (3.5.2)$$

$$\psi_p = -b_p x - \sum_{j=1}^{p-1} [b_j \cosh(jy) + b_{j+p} \sinh(jy)] \sin(jx) \quad (3.5.3)$$

The convergence of the series for points  $(x,y)$  away from the boundary is not guaranteed; but we might expect that the series would perform well near the given curve.

The coefficients in the series could be found in the usual way, by discretizing  $x$  and  $y = g(x)$  and applying conditions (3.5.1) simultaneously. Having found the coefficients, we could then solve for the shape of the curve  $y = h(x)$  using the condition

$$\phi(x_i, h(x_i)) = 0 \quad , \quad i = 1, \dots, m \quad (3.5.4)$$

This represents a set of nonlinear equations, one for each  $h_i$ . Each equation in turn could be treated by Newton's method for finding the roots of a nonlinear equation of a single variable.

All this is possible in principle. In practice, we have found the above procedure to be moderately successful at best. For example, we have tested the approach on a variety of cathode shapes of the form

$$g(x) = a \cos x + \frac{1}{\alpha}.$$

For amplitudes less than 0.2, the agreement with the known analytic solution was good. However, for larger amplitudes, the disagreement between computed and exact solutions was substantial and the approach was abandoned.

The above approach to the inverse problem lends itself to applications in three dimensions; but its usefulness is clearly limited. However, an excellent approach, albeit restricted to the two dimensional design problem, is described below.

By working in the plane of the complex potential, it is possible to numerically predict highly distorted cathode tool designs. Recall, the analytic solution could be expressed in the  $w$ -plane in the form

$$z = f(w)$$

where, in particular,

$$x = -\frac{\psi}{\alpha}$$

$$y = g(x)$$

on the anode potential  $\phi = 1$ .  $g(x)$  is periodic in  $x$  and assumed to be symmetric about  $x = 0$ . Consequently,  $g(x)$  is expressible in terms of a Fourier cosine series,

$$g(x) = A_0 + \sum_{k=1}^{\infty} A_k \cos kx \quad (3.5.5)$$

Again, by analytic continuation away from the surface  $\phi = 1$ , we have

$$z = i \left[ \frac{(w-1)}{\alpha} + A_0 + \sum_{k=1}^{\infty} A_k \cosh k \frac{(w-1)}{\alpha} \right] \quad (3.5.6)$$

Given the tabulated function  $g(x)$ , the series (3.5.5) can be truncated and the coefficients solved for by conventional means. Once the coefficients are known, any of the equipotential surfaces corresponding to  $\phi < 1$  can be determined from (3.5.6). It should be noted, however, that the series (3.5.6) may be very slowly convergent, for along  $\phi < 1$ , the series coefficients behave asymptotically like

$$A_k e^{k \left| \frac{\phi-1}{\alpha} \right|}.$$

If a sufficient number of boundary points are available, then fast Fourier transform techniques can be used to efficiently determine the coefficients, even in the case of extremely slow convergence. However, an alternative approximation has proved remarkably successful and does not require a large bank of data, nor does it involve large numbers of coefficients. The approximate form of the solution is taken to be

$$z_n = i \left[ \frac{(w-1)}{\alpha} + b_0 + \sum_{k=1}^n b_k \frac{\cosh \frac{(w-1)}{\alpha}}{1 - b_{k+n} \cosh \frac{(w-1)}{\alpha}} \right] \quad (3.5.7)$$

This form has the required periodicity and along  $\phi = 1$ , reduces to

$$x_n = -\frac{\Psi}{\alpha} \quad (3.5.8)$$

$$y_n = b_0 + \sum_{k=1}^n b_k \frac{\cos \frac{\Psi}{\alpha}}{1 - b_{k+n} \cos \frac{\Psi}{\alpha}}$$

For a given value of  $y = g(x_i)$ ,  $i = 1, \dots, m \geq 2n + 1$ , the parameters  $b_{k+n}$  are determined by minimizing the corresponding least squares residual, using the Levenberg-Marquardt algorithm referred to in section 2.8.

The method is illustrated on the following example, chosen for its sharply spiked profile. Let

$$g(x) = \ln(1 - 2a \cos x + a^2) \quad , \quad 0 < a < 1$$

If this is expanded in a cosine series, we obtain

$$\ln(1 - 2a \cos x + a^2) = \sum_{k=1}^{\infty} \left( -\frac{2a^k}{k} \right) \cos kx .$$

The exact solution of the inverse problem takes the form

$$z = i \left[ \frac{(w-1)}{\alpha} + \ln \left( 1 - 2a \cosh \frac{(w-1)}{\alpha} + a^2 \right) \right] \quad (3.5.9)$$

or in series form

$$z = i \left[ \frac{(w-1)}{\alpha} + \sum_{k=1}^{\infty} \left( \frac{-2a^k}{k} \right) \cosh k \frac{(w-1)}{\alpha} \right] \quad (3.5.10)$$

Note that the series does not converge when  $\phi \leq 1$  unless

$$ae^{|\frac{\phi-1}{\alpha}|} < 1 .$$

In particular, along  $\phi = 0$ , we must have

$$ae^{\frac{1}{\alpha}} < 1 .$$

This is precisely the condition that the branch points of the logarithm remain outside of the domain  $0 \leq \phi \leq 1$ . The branch points are given by those values of  $\phi$  and  $\psi$  for which

$$\cosh\left(\frac{\phi-1}{\alpha}\right) = \frac{1+a^2}{2a}$$

$$\psi = 0.$$

That is,

$$\phi = 1 \pm \alpha \ln\left(\frac{1}{a}\right), \quad \psi = 0.$$

If  $ae^{\frac{1}{\alpha}} < 1$ , the singularities reside outside of the domain of interest. However, for values of  $a$  and  $\alpha$  such that  $ae^{\frac{1}{\alpha}}$  is very nearly 1, the series (3.5.10) is very slowly convergent.

We consider the case  $a = 0.71$ ,  $\alpha = 3.0$ . For this choice of  $a$  and  $\alpha$ , the machined cathode is the sharply spiked profile shown in figure 3.5.1. The parameters of the approximate solution (3.5.8) appear in table 3.5.1. We have chosen  $n = 9$  and  $m = 4$ , so that there are only 19 unknowns. More data points were placed near  $\psi = 0$  in order to resolve the proper anode shape (see figure 3.5.2 for the field pattern). With this choice, the root mean square error between the computed and exact values of  $g(x)$  is  $O(10^{-9})$ . The only singularities in the approximate solution correspond to  $\psi = 0$  and values of the potential beyond the domain  $0 \leq \phi \leq 1$ . A branch point in the exact solution resides at  $\phi = -0.0275$ , very near the zero potential.

The required tool design ( $\phi = 0$ ) is the very steep trough shown in the figures and was computed using the already determined parameters  $b_k$  and equation (3.5.7).

Agreement with the analytic solution is very good for values of  $x \geq 1.5$ . The accuracy in the trough is somewhat poorer owing to the presence of the branch point. The root mean square errors between the computed and exact solutions are shown in table 3.5.2.

To achieve comparable accuracy with the series approximation (3.5.10) would involve on the order of two hundred terms in the series and at least as many boundary points. Also included in the table are the results for several different values of  $n$  to illustrate the convergence of the proposed scheme.

**Table 3.5.1 The Parameters ( $n=9$ )**

Linear terms	Nonlinear terms
$b_0 = 0.408E + 00$	$b_{10} = 0.940E + 00$
$b_1 = -0.112E - 01$	$b_{11} = 0.918E + 00$
$b_2 = -0.314E - 01$	$b_{12} = 0.875E + 00$
$b_3 = -0.638E - 01$	$b_{13} = 0.782E + 00$
$b_4 = -0.113E + 00$	$b_{14} = 0.651E + 00$
$b_5 = -0.171E + 00$	$b_{15} = 0.439E + 00$
$b_6 = -0.203E + 00$	$b_{16} = 0.255E + 00$
$b_7 = -0.175E + 00$	$b_{17} = 0.356E + 00$
$b_8 = -0.205E - 01$	$b_{18} = 0.628E - 01$
$b_9 = -0.155E + 00$	

**Table 3.5.2 Errors in Computed Cathode Coordinates**

n	$E_x$	$E_y$
3	0.242 E+00	0.385 E+00
6	0.620 E-01	0.953 E-01
9	0.182 E-01	0.171 E-01

For larger values of  $n$ , the computation of the nonlinear parameters can be costly, especially if the initial estimate for the parameters is poor. It is generally recommended that the search for a minimum begin with a small value of  $n$  and a large error tolerance. (The error tolerance in the MINPAK routine must be supplied by the user. When the difference between successive iterates meets a specified criterion depending on the tolerance, the routine is terminated.) The error tolerance can then be gradually decreased in stages until further improvement is impossible. At this stage, a larger value  $n$  (and again a crude tolerance) may be chosen with the parameters just calculated being used as part of the initial estimate for the new value of  $n$ . In this way, a calculation involving close to twenty parameters can be performed reasonably efficiently (less than ten minutes of CPU time).

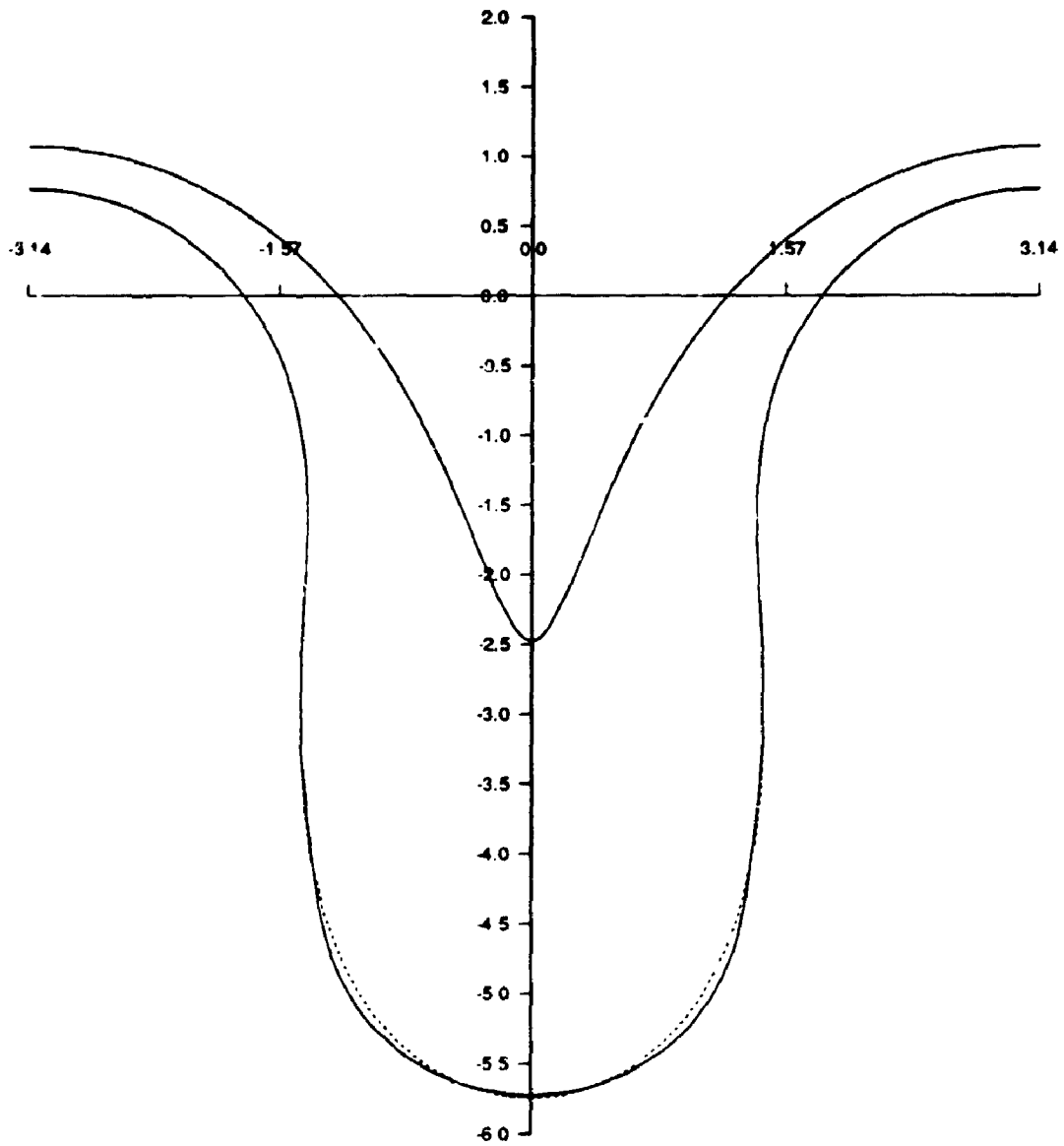


Fig. 3.5.1 Inverse ECM ( $a = 3.0$ ,  $a = 0.71$ )

$$g(x) = \ln(1 - 2a\cos x + a^2)$$



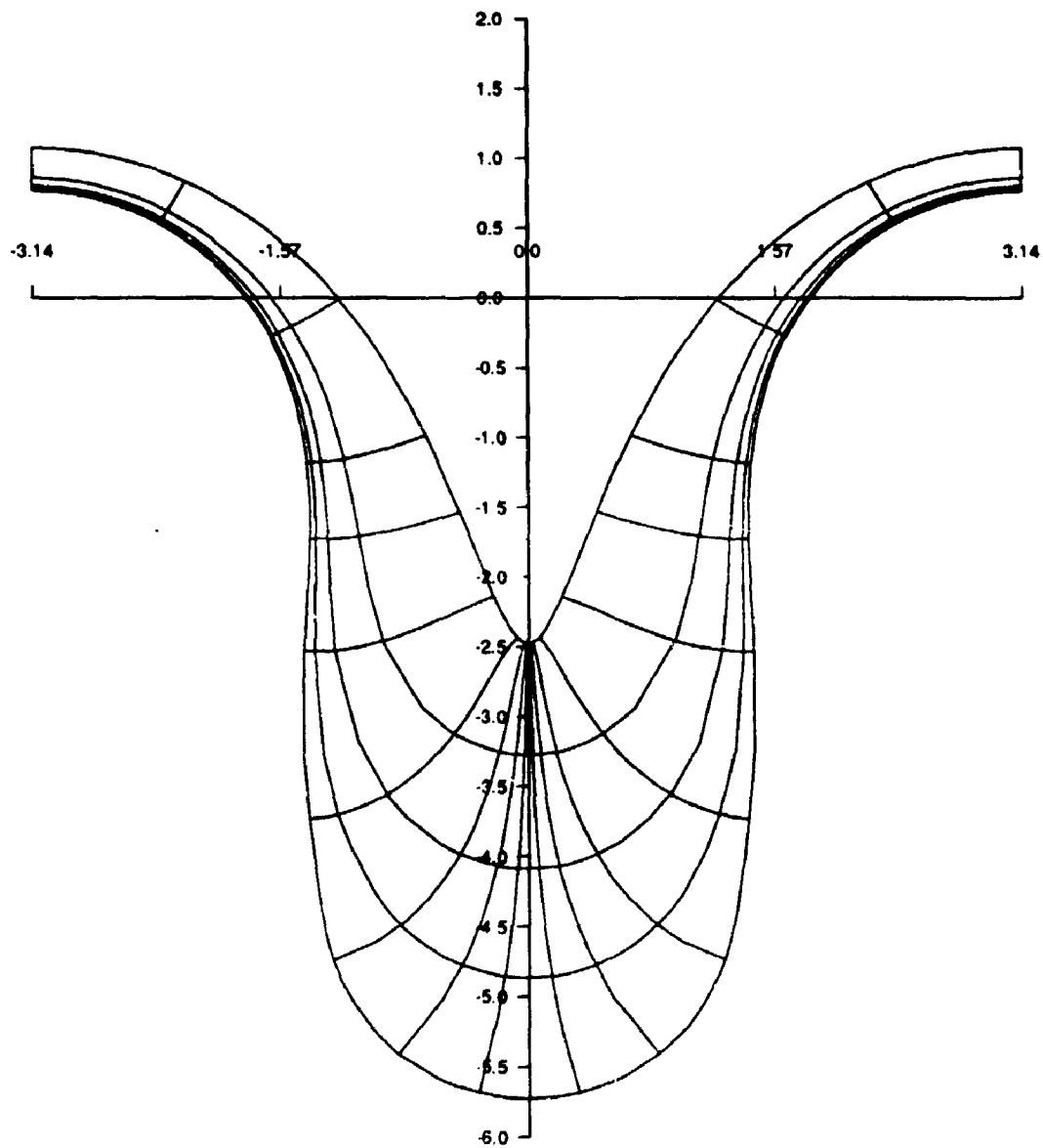


Fig. 3.5.2 Field Lines ( $\alpha = 3.0$ ,  $a = 0.71$ )

$$g(x) = \ln(1 - 2a\cos x + a^2)$$

### 3.6 Discussion

The linear boundary method presented in chapter 2 has been examined as a viable candidate for approximation of the potential problems associated with electrochemical machining. The convergence of the method has been confirmed numerically for a number of steady state configurations. When applied to the unsteady problem, the method is seen to be an efficient and accurate technique for simulating the dissolution of relatively smooth anode shapes. The largest aspect ratios tested are comparable to those examined by Forsyth and Rasmussen (1979) and Sloan (1986), although the former article does include examples having more sharply spiked cathodes, for similar aspect ratio. For these sharply varying shapes, the coordinate transformation of Forsyth and Rasmussen (1979) would be recommended. For smooth profiles, however, the present method is almost certainly more efficient and the programming effort is considerably less.

As the unsteady ECM is physically stable, any numerical instabilities in the time dependent portion should be revealed by reducing the time step. No such instability was observed. The only numerical difficulties arose in the spatial portion of the scheme. There, the matrix systems were often badly conditioned. This is typical of linear boundary methods and it is this feature of the method which limits its applications to relatively simple geometries.

An effective numerical approach to the inverse problem of ECM has been developed. The method, as it stands, is restricted to two dimensional problems, but is capable of handling extremely distorted electrode shapes.

## CHAPTER 4

### Hele-Shaw Flow

#### 4.1 Introduction

In 1898, Hele-Shaw devised an experimental arrangement whereby an incompressible viscous fluid is constrained to move between two close parallel plates. If the plates are made of glass and a dye is introduced into the fluid, then the flow pattern is easily observed. The apparatus is referred to as a Hele-Shaw cell.

Recently, interest has grown in the mathematical description of Hele-Shaw flows for its own sake. Some analytic solutions have been obtained, providing one of the few examples of a nontrivial MBP having a known solution. This permits a quantitative comparison with a numerical simulation of the mathematical model.

Lamb (1945) provides a derivation of the equations governing the flow. If the gap width between the plates is  $b$  and the  $z$ -axis is taken perpendicular to the plates, then the components of velocity of the flow are

$$u = -\frac{1}{2\mu} \frac{\partial p}{\partial x} z(b-z)$$

$$v = -\frac{1}{2\mu} \frac{\partial p}{\partial y} z(b-z)$$

The  $x$  and  $y$  axes lie in the plane of the plates.  $\mu$  is the viscosity of the fluid and  $p$  is the fluid pressure. If the velocities are averaged over the gap-width, we have

$$u_0 = -\frac{b^2}{12\mu} \frac{\partial p}{\partial x}$$

$$v_0 = -\frac{b^2}{12\mu} \frac{\partial p}{\partial y}$$

Then,

$$\mathbf{u} = u_0 \mathbf{e}_1 + v_0 \mathbf{e}_2$$

and the equation of continuity is

$$\nabla \cdot \mathbf{u} = 0.$$

This permits the introduction of a velocity potential

$$\phi = \frac{b^2}{12\mu} p$$

such that

$$\mathbf{u} = -\nabla \phi$$

and

$$\nabla^2 \phi = 0.$$

Thus, the mean velocity in a Hele-Shaw cell can be taken to represent the motion of an ideal fluid. Indeed, it was this observation which lead Hele-Shaw to conduct his experiments. By placing obstacles in the flow field and a dye in the fluid, Hele-Shaw was able to produce a visual representation of the flow pattern around the obstacles. If two immiscible fluids are introduced into the cell, the motion of the interface between the fluids provides a unique opportunity to observe a moving boundary problem.

There is a close analogy between flow in a Hele-Shaw cell and a number of other physical processes. This allows the Hele-Shaw cell to be used as a feasible laboratory model of a variety of phenomena. For example, the motion of an incompressible viscous fluid through an isotropic homogeneous porous medium is governed by the equation of continuity

$$\nabla \cdot \mathbf{u} = 0$$

and Darcy's law

$$\mathbf{u} = -\frac{K}{\mu} \nabla p$$

where  $\mathbf{u}$  is an average velocity,  $p$  is the pressure and  $K$  is the permeability. Clearly, then, the description of flow in a Hele-Shaw cell is identical to the flow in a porous medium of permeability  $K = \frac{b^3}{12}$ . The flow of a viscous fluid through a porous medium has long been a study of both practical and theoretical interest. The practical applications are numerous and include the areas of petroleum engineering, groundwater flows and filtration processes. And once again, the mathematical model of such flows is often a nontrivial moving boundary problem. For example, when one fluid displaces another in a porous medium, the interface between the two fluids is a moving boundary. This occurs in oil recovery processes whereby oil is driven through porous rock by another fluid, usually water. This can lead to the phenomenon of fingering, where the less viscous water rushes into the oil in the form of long fingers of fluid. This of course reduces the effectiveness of the recovery process and is therefore a major concern.

The Hele-Shaw cell can also be used to model the injection moulding of plastics (see Richardson (1972)). It can even be thought of as a model for the basic two-dimensional electrochemical machining problem (see chapter 3). The injection moulding and ECM problems are both examples of what are sometimes referred to as injection problems. These problems are well-posed, and as we have noted, existence of a unique solution has been established.

On the other hand, the suction problem, where fluid is extracted from a Hele-Shaw cell is ill-posed (see Elliott and Ockendon (1982)). Existence and uniqueness have not yet been established in general; but several exact solutions have been constructed in specific cases. It is such a problem as this that we choose to examine in this chapter. In

particular we wish to monitor the interface between two fluids in a Hele-Shaw channel, where suction provides a pressure difference which drives one fluid into the other. Depending on which fluid is the more viscous, the phenomenon of fingering can occur.

In section 4.2 we present a review of the Hele-Shaw problem together with some of the numerical schemes that have been previously implemented. In sections 4.3, 4.4 and 4.5, several examples are presented for which there exist known analytic solutions. We compare the results of linear and nonlinear boundary approximation methods with these known solutions. In one example, the interface between two fluids is followed in time as the instability or fingering process evolves. In another example, the initial free surface profile is slightly perturbed and the ensuing flow is seen to diverge dramatically from that of the first example. This clearly indicates the ill-posedness of the suction problem, and provides a severe test of the numerical technique.

Lastly, in section 4.6, the stability of the time dependent portion of the numerical scheme is briefly examined.

## 4.2 Review

The mathematical description of the instability of the interface between two fluids in a Hele-Shaw cell (or porous medium) was first examined in the classic paper of Saffman and Taylor (1958). Two immiscible fluids occupy a vertically positioned Hele-Shaw cell chosen to lie in the  $(x, y)$ -plane. The initial interface profile is a small sinusoidal displacement

$$y = \epsilon e^{i n x + \sigma t}$$

of wavelength  $2\pi/n$ . The sign of  $\sigma$  determines the relative stability of the interface. Fluid 2 (the lower fluid) is driven vertically upwards with speed  $V > 0$ . Saffman and Taylor used a first order perturbation analysis to show that the interface is unstable for small

initial disturbances of any wavelength if fluid 2 is the less viscous of the two fluids and the effects of gravity and surface tension are neglected. If the effects of gravity are included, the same conclusion holds for sufficiently large  $V$ .

The following condition on  $\sigma$  was derived by Chuoke, van Meurs and van der Poel (1959), when both gravity and surface tension effects are included:

$$\frac{\sigma}{n} (\mu_1 + \mu_2) = (\rho_1 - \rho_2)gK + (\mu_1 - \mu_2)V - n^2\gamma K$$

where  $K$  is the permeability of the medium, here assumed to be the same for both fluids.  $\mu_1, \mu_2$  are the corresponding viscosities,  $\rho_1, \rho_2$  the corresponding densities and  $\gamma$  is the surface tension parameter.  $g$  is the acceleration due to gravity. We have  $\sigma > 0$  if

$$n^2 < \frac{1}{\gamma K} ((\rho_1 - \rho_2)gK + (\mu_1 - \mu_2)).$$

Thus, the flow is still unstable beyond a cutoff value of the wavelength which depends on the parameters of the problem. If the surface tension is small, the interface is unstable to smaller and smaller wavelengths.

Laboratory experiments confirm this instability in the Hele-Shaw problem. (The experimental apparatus is confined to a rectangular channel of fluid sandwiched between two parallel glass plates. The entire apparatus is surrounded by a rigid frame to contain the fluid. In this way, the apparatus can be rotated to a horizontal or vertical position to test the influence of gravity.) Saffman and Taylor (1958) include the results of some experimentation which reveal the eventual growth of a single finger of less viscous fluid emerging from a group of initially smaller unstable fingers. They analyse the shape of the steady finger and find that it should grow to occupy a fraction  $\lambda$  of the channel width. Saffman (1959) actually obtains an exact solution describing the shape of the unsteady finger in terms of the parameter  $\lambda$ . Comparisons with the laboratory experiments seem to suggest a value of  $\lambda = 1/2$ , but there is no satisfactory explanation to account for this

choice. Pitts (1980) repeated the Saffman-Taylor experiments and also found that a single finger would grow to occupy one half of the channel width. He has suggested that, since a film of displaced fluid must remain on the walls of the apparatus, the phenomena should properly be treated as a three-dimensional problem.

In a more recent paper, McLean and Saffman (1981) perform a singular perturbation analysis to test the effects of small surface tension on the steady finger profile. Again, they find an infinite family of finger widths are possible. What is more, Vanden Broeck (1983) has established the existence of countably infinite families of steady-state solutions for each value of the surface tension parameter.

McLean and Saffman (1981) also discuss the stability of the steady finger to small disturbances. Once again, the analysis is linear only. It is found that the fingers themselves are unstable to small disturbances, both with or without the inclusion of surface tension effects. This is in complete discord with experiment, which as we have noted, shows only the development of a long finger reaching a steady profile.

Thus, the theory predicts an infinity of possible solutions, all unstable to small disturbances, whereas experiment has revealed only the development of a unique and stable finger. It is these disagreements between theory and experiment which have prompted a recent interest in the numerical simulation of the Hele-Shaw problem. Perhaps if an accurate numerical scheme were able to follow the development of a finger for a long enough time, we might be able to discern the growth or not of the predicted linear instabilities. It is a challenging undertaking, as the numerical perturbations that are introduced must also be subject to rapid growth.

Meyer (1981) presents a method of lines simulation of the time dependent Hele-Shaw flow in a bubble of fluid. Neglecting the effects of surface tension, he was able to find an exact solution using a complex variables approach, thereby permitting a direct comparison with the numerical results. An interesting feature of the exact solution involves the development of a cusp in the free surface in a finite time. Using similar techniques,



Aitchison and Howison (1985) find exact solutions to several Hele-Shaw flows in an infinite channel. In particular, two solutions are presented, the unsteady solution of Saffman (1959) and one for which the free surface forms a cusp in a finite time. They present comparison of both solutions with a boundary integral simulation. Although efficient, the algorithm is unable to follow the development of the Saffman finger very far in time.

Davidson (1984) includes the effects of surface tension in the numerical solution of a boundary integral equation. The accuracy of the routine is determined by direct comparison with the analytic Saffman solution for the special case of zero surface tension. A saw-tooth instability develops in the computed free surface profile due to the rapid growth of numerical perturbations. This causes the approximate solution to break down well before the steady Saffman profile is reached. The solution is coaxed along by artificially smoothing the instabilities at every time step.

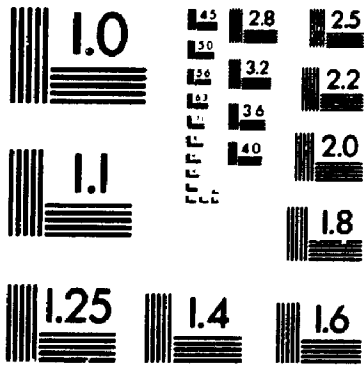
A boundary integral method is used by Degregoria and Schwartz (1986) to simulate channel fingering. The calculation is performed in a transformed plane and for a variety of nonzero values of the surface tension. They are able to produce large steady fingers by applying a sophisticated smoothing technique every few time steps. No comparison is made with an analytic solution for the case of zero surface tension. However, for a very small value of the surface tension parameter, an interesting feature is observed. Instead of a single steady finger developing, the finger tip actually splits in two. The authors attribute this to numerical noise which eventually dominates in the absence of the stabilizing influence of large surface tension.

Finally, an excellent review of viscous fingering in porous media can be found in Homsy (1987). The article also includes some of the more recent experimental observations.

2

of/de

2



### 4.3 Cusping Solution: Linear Approximation

Consider an infinite Hele-Shaw cell containing air and a viscous liquid separated by a sharp interface (see figure 4.3.1). Gravitational effects are ignored. This corresponds to an experimental cell that has been rotated to a horizontal position. The viscosity of the air is considered negligible compared to that of the fluid. The initial shape of the interface is a small sinusoidal disturbance of wavelength  $\lambda$ . Then, if the imposed periodicity is in the  $x$  direction, we can concentrate on the region corresponding to one wavelength

$$-\frac{\lambda}{2} \leq x \leq \frac{\lambda}{2}.$$

The air pressure is a constant, assumed to be zero, while the pressure in the liquid is a function of position at any given time and remains to be determined. The average velocity  $\mathbf{u}$  in the liquid is related to the pressure via the equation

$$\mathbf{u} = -\frac{b^2}{12\mu} \nabla p \quad (4.3.1)$$

$b$  is the gap width between the plates and  $\mu$  is the viscosity of the liquid. Together with the equation of continuity

$$\nabla \cdot \mathbf{u} = 0 \quad (4.3.2)$$

we have

$$\nabla^2 p = 0 \quad (4.3.3)$$

At the interface, ignoring surface tension effects, continuity of pressure implies that the liquid pressure has the constant value zero. As a consequence, the material derivative of the pressure vanishes along the interface giving the kinetic condition

$$\frac{\partial p}{\partial t} = \frac{b^2}{12\mu} \nabla p \cdot \nabla p \quad (4.3.4)$$

Finally, we assume that liquid is being removed at a constant velocity far from the interface.

That is, we have the condition

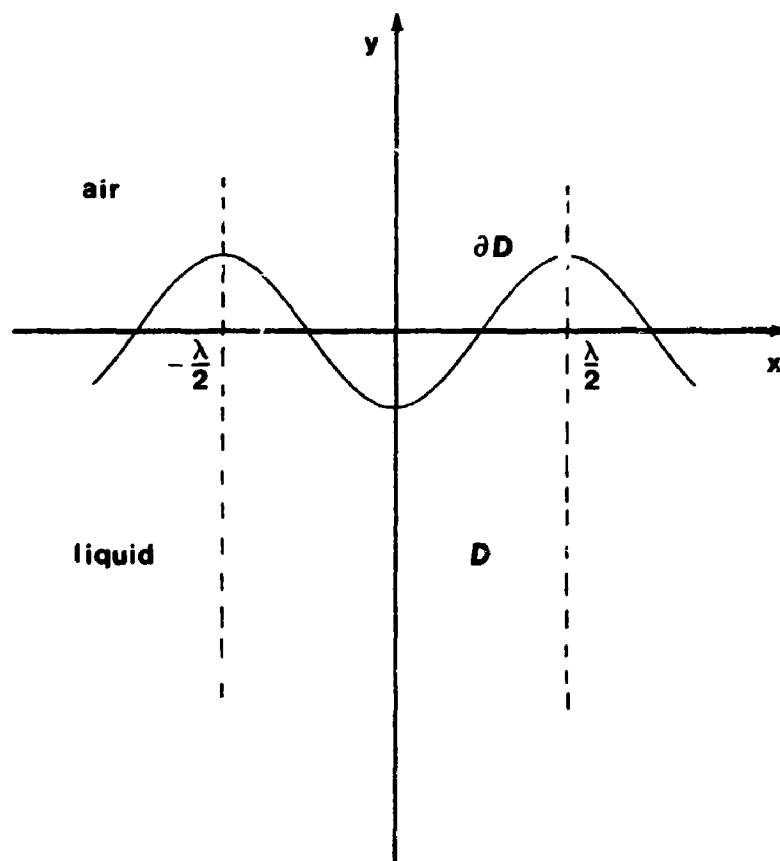


Fig. 4.3.1 Definition Sketch

$$v \rightarrow -V \text{ as } y \rightarrow -\infty \quad (4.3.5)$$

where  $V > 0$  is a constant.

The governing equations can be put into a standard form by introducing the following dimensionless variables.

$$\begin{aligned} \phi = p^* &= \frac{kb^2}{12\mu V} p & x^* &= kx & u^* &= \frac{u}{V} \\ & & y^* &= ky & v^* &= \frac{v}{V} \\ & & t^* &= kVt & k &= \frac{2\pi}{\lambda} \end{aligned}$$

Equations (4.3.6)

In terms of the new variables, then, we have (dropping the star notation)

$$\nabla^2 \phi = 0 \quad , \quad \underline{x} \in D \quad (4.3.7)$$

$$\phi = 0 \quad , \quad \underline{x} \in \partial D \quad (4.3.8)$$

$$\phi_y \rightarrow 1 \text{ as } y \rightarrow -\infty \quad (4.3.9)$$

together with the free surface condition and the initial condition

$$\frac{\partial \phi}{\partial t} = \nabla \phi \cdot \nabla \phi \quad (4.3.10)$$

$$f_0(x, y) = 0 \quad (4.3.11)$$

where (4.3.11) describes the initial position of the free surface.  $D$  is the domain bounded by the interface and the lines  $x = \pm\pi$  and  $\partial D$  denotes the interface itself.

The evolution equation can take several alternative forms. For example, if the free surface can be represented by a function of  $x$  and  $t$ ,

$$y = g(x, t)$$

then differentiation with respect to  $t$  gives

$$g_t = g_x \phi_x - \phi_y \quad (4.3.12)$$

This is the usual Eulerian description; the solution of (4.3.12), for a given value of  $x$ , represents the vertical position of the free surface as a function of time. The actual trajectories of marked particles on the interface are quite different. The values of  $\phi_x, \phi_y$  to be used in the numerical calculation are obtained directly from the series approximation for  $\phi$  from term by term differentiation. In order to avoid a finite difference expression for  $g_x$ , we can differentiate the condition (4.3.8) with respect to  $x$  to obtain the expression

$$g_x = -\frac{\phi_x}{\phi_y}$$

for the derivative of the free surface. The evolution equation becomes

$$g_t = -\frac{\phi_x^2 + \phi_y^2}{\phi_y} \quad (4.3.13)$$

Of course, this form has its limitations, for in a complicated flow problem, the interface may cease to be representable by a single-valued function  $g(x, t)$  of  $x$ . This difficulty can be overcome by using the Lagrangian description of the moving surface

$$\frac{dx}{dt} = -\phi_x \quad (4.3.14)$$

$$\frac{dy}{dt} = -\phi_y$$

The coordinate pair  $(x(t), y(t))$  describes parametrically the path of a specific free surface particle originating at a given coordinate in the  $x$ - $y$  plane. This particular formulation has definite advantages. In the first place, experience has shown that the Lagrangian description of the moving surface permits the numerical simulation of an unstable interface to proceed further in time than the corresponding Eulerian description (see for example, the Rayleigh-Taylor calculation of Menikoff and Zemach (1983)). This may largely be due to the fact that the explicit calculation of  $g_x$  is avoided. Secondly, it is easily

implemented using the present boundary approximation, an advantage that is not readily available to a finite difference calculation, for example. In the calculations which follow, both the Lagrangian and Eulerian descriptions have been used.

Once the initial condition (4.3.11) has been specified the problem is well defined. The two examples chosen in this section and the next admit of an analytic solution. In both cases, the initial profile of the interface has the form

$$y = -\epsilon \cos x + O(\epsilon^2) \quad , \quad 0 \leq \epsilon < 1$$

but their subsequent flow developments are strikingly different. In each example the flow is assumed to be symmetric about the line  $x = 0$ , so that the computations may be restricted to the domain  $0 \leq x \leq \pi$ .

The first example is taken from Aitchison and Howison (1985). The analytic solution may be expressed as

$$z = i(w - b_0(t) - b_1(t)e^w) \quad (4.3.15)$$

where  $w = \phi + i\psi$  is the complex potential. The time dependent coefficients  $b_0, b_1$  are given by

$$b_1 e^{-b_0} = \epsilon$$

$$b_0 - \frac{1}{2} b_1^2 = t - \frac{1}{2} \epsilon^2$$

together with the initial conditions

$$b_0(0) = 0$$

$$b_1(0) = \epsilon \quad , \quad 0 \leq \epsilon < 1.$$

$\epsilon = 0.2$  is used in the following calculations.

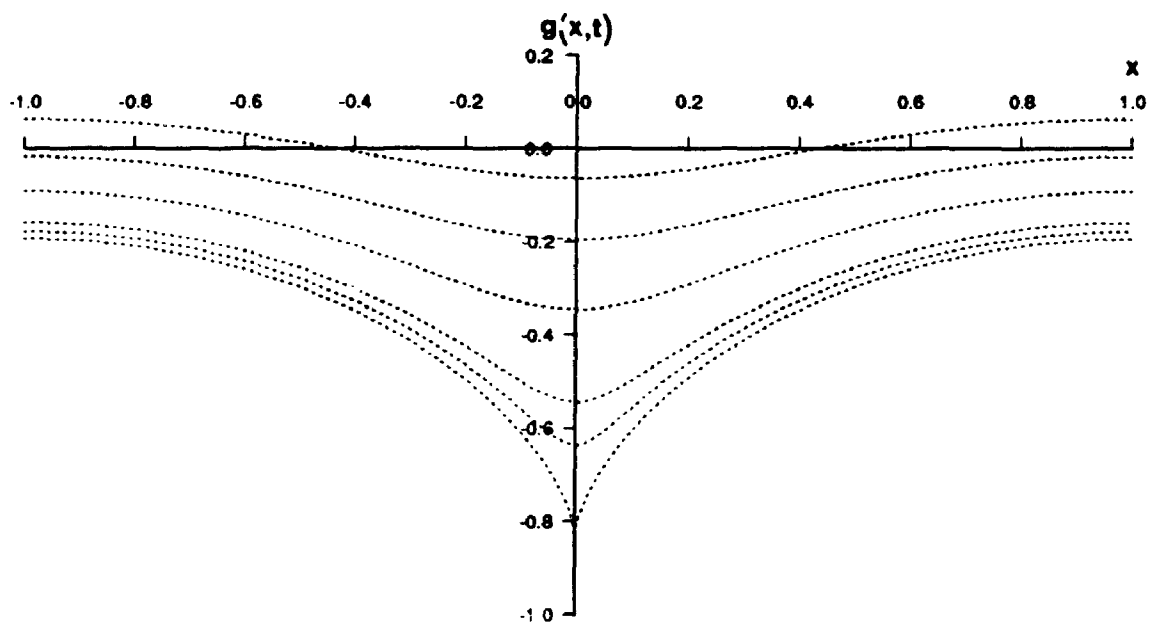


Fig. 4.3.2 Exact Cusping Behaviour



The solution (4.3.15) may be interpreted as a conformal mapping to the plane of the complex potential. This is a slightly different interpretation than that of Aitchison and Howison, but is in keeping with the notation used in our steady state ECM examples. For a discussion of the general procedure whereby analytic solutions to the Hele-Shaw problem may be obtained, see the appendix A.2.

The mapping (4.3.15) ceases to be conformal when a zero of  $dz/dw$  reaches the boundary. Howison, Ockendon and Lacey (1985) have shown that a singularity in  $dw/dz$  resides outside of the domain  $D$  at  $t = 0$ , but that it steadily approaches the boundary as time advances. In a finite time, critical values of  $b_0$  and  $b_1$  can be found for which the free surface has developed a cusp at  $x = 0$ . The vertical flow rate increases dramatically in the vicinity of  $x = 0$  as the cusp time is approached and the slope of the free surface assumes a steep profile (see figure 4.3.2).

It is an extreme test of any numerical scheme to follow the analytic profile as nearly as possible for times close to the cusp time. The linear boundary approximation method has been applied to this example. As in the last chapter, a trial solution is used to determine a best approximation to the boundary data of a potential problem. Then, either (4.3.13) or (4.3.14) is used to advance the solution one time step. The system of ordinary differential equations generated is solved according to the predictor-corrector scheme outlined in the last chapter.

### Test of potential problem

Before presenting the results of the full time dependent flow, we test the validity of the approximation method as a model for the potential portion of the problem. By fixing the time  $t$ , the numerical solution of the potential problem can be directly compared with the known analytic solution at that time. The approximate solution in the linear method takes the form

$$\phi_n = Ay + \Re \sum_{k=n-1}^{n-1} c_k e^{-ik\pi z}$$

The form of the approximation which satisfies the symmetry conditions and the asymptotic boundary condition is

$$\phi_n = y + \sum_{j=0}^{n-1} b_j e^{j\pi y} \cos j\pi x \quad (4.3.16)$$

Note the few changes in notation from that of chapter 2. The upper limit of the summation is  $(n - 1)$ , so that  $\phi_n$  refers to the fact that there are exactly  $n$  unknowns. The  $x$ ,  $y$  and  $t$  variables have all been scaled by an additional factor of  $1/\pi$ . A uniform  $x$ -grid of  $m$  values is used in the calculations. Tables 1 through 4 are designed to test the convergence of the approximation method as  $n$  is increased and to suggest an appropriate ratio of  $m:n$ . In the tables, the root mean square errors of  $\phi$ ,  $\phi_x$ ,  $\phi_y$  and  $g_x$  are monitored along the curve  $y = g(x, t)$ . For example, the root mean square error in  $\phi$  is defined to be

$$E_\phi = \left\{ \frac{1}{m} \sum_{i=1}^m (\phi_n(x_i) - \phi(x_i))^2 \right\}^{\frac{1}{2}}$$

where  $x_i$  is a point on the curve  $y = g(x, t)$ . Expressions for the root mean square error in  $\phi_x$ ,  $\phi_y$  are defined similarly, with the computed value of  $\phi_x$ ,  $\phi_y$  determined from term by term differentiation of the series solution  $\phi_n$ . In the case of  $g_x$ , the computed value was determined from the ratio

$$-\frac{\partial \phi_n}{\partial x} / \frac{\partial \phi_n}{\partial y}.$$

The tables also include the maximum error

$$E_{\max} = \max_x |\phi_n(x_i) - \phi(x_i)|$$

where the  $x_i$  are the  $m$  boundary points plus an additional one hundred points placed along the boundary at uniform intervals of  $x$ . The additional points are not used in the determination of the coefficients, but rather serve as a check on the uniformity or smoothness of the boundary fit.

Table 4.3.1 The Ratio  $m:n$ Cusping Profile - Linear approximation,  $t=0.25$ ,  $n=15$ 

$m$	$E_{\max}$	$E_{\phi}$	$E_{\phi_x}$	$E_{\phi_y}$	$E_{\phi_z}$
15	10.8	0.434 E-12	216.0	42.7	12.5
20	0.458 E-02	0.168 E-02	0.0506	0.0552	0.0572
25	0.448 E-02	0.166 E-02	0.0514	0.0543	0.0602
30	0.442 E-02	0.165 E-02	0.0514	0.0538	0.0596
45	0.432 E-02	0.163 E-02	0.0512	0.0529	0.0591
60	0.441 E-02	0.162 E-02	0.0511	0.0524	0.0588
75	0.449 E-02	0.161 E-02	0.0510	0.0521	0.0586

In Table 4.3.1, we have chosen  $t = 0.25$ ,  $n = 15$  and selected a range of values of  $m \geq n$ . For  $n = m$ , the errors in the derivatives  $\phi_x, \phi_y$  are unacceptable, even though the series  $\phi_n$  itself reproduces the boundary data quite well (as should be the case if ill-conditioning is not a factor). The error in  $\phi$  itself at off-grid points (points along the boundary other than the  $m$  data points) is very great, as indicated by the maximum error. As  $m$  is increased slightly to  $m = 20$ , the errors in the derivatives are already much improved. As  $m$  is increased still further, a modest improvement in the root mean square errors is achieved. A ratio of  $m:n$  of about two or three to one would appear to be more than adequate.

In tables 4.3.2 and 4.3.3, convergence with increasing  $n$  is illustrated for the approximation along the interface at the times  $t = 0.25$  and  $t = 0.3$  respectively. In both cases, the ratio  $m:n$  is three to one. Convergence of the method is clearly established, but there is no doubt that the convergence is slow and only gets worse as time is increased. This is undoubtedly due in large part to the encroachment on the boundary of a singularity in the analytic solution. The error analysis of section 2.6 suggests that the degree of approximation will be poor if the solution cannot be harmonically extended well past the boundary; and in this case, as time increases, the singularity steadily advances toward the boundary. More and more terms would have to be included in the series approximation as the cusp time is approached. However, ill-conditioning rapidly becomes a factor, limiting the size of  $n$ . For example, the condition number for  $n = 60$  in tables 4.3.2 and 4.3.3 is already  $O(10^{13})$ . The accuracy of results for  $n > 60$  would certainly be in question.

**Table 4.3.2 Convergence**

**Cusping Profile - Linear approximation,  $t=0.25$ . The singularity is a vertical distance of 0.067 from the boundary.**

$n$	$m$	$E_{\max}$	$E_{\phi}$	$E_{\phi_x}$	$E_{\phi_y}$	$E_{\psi_x}$
10	30	0.979 E-02	0.404 E-02	0.0893	0.0939	0.105
15	45	0.432 E-02	0.163 E-02	0.0512	0.0529	0.0347
20	60	0.214 E-02	0.752 E-03	0.0303	0.0311	0.0347
25	75	0.111 E-02	0.373 E-03	0.0183	0.0187	0.0209
30	90	0.602 E-03	0.195 E-03	0.0112	0.0114	0.0128
45	135	0.107 E-03	0.327 E-04	0.00272	0.00275	0.00308
60	180	0.296 E-04	0.840 E-05	0.135 E-02	0.139 E-02	0.184 E-02

**Table 4.3.3 Convergence**

Cusping Profile - Linear approximation,  $t=0.30$ . The singularity is a vertical distance of 0.0345 from the boundary.

n	m	$E_{max}$	$E_{\phi}$	$E_{\phi_x}$	$E_{\phi_y}$	$E_{\phi_z}$
15	45	0.242 E-01	0.797 E-02	0.206	0.217	0.295
30	90	0.970 E-02	0.279 E-02	0.123	0.126	0.163
45	135	0.478 E-02	0.127 E-02	0.0783	0.0795	0.101
60	180	0.320 E-02	0.823 E-03	0.0709	0.0679	0.0914

The shape of the boundary itself must play a role in the approximation. As the time to cusp is approached, an inflection point in  $g(x)$  drifts rapidly toward the origin. Now, the inability of approximation schemes to reproduce a sharp inflection point is typical of polynomial approximation in one-dimensional cases where the approximation is not performed piecewise. This difficulty is often lessened in one-dimensional cases if a rational approximation is performed. As we shall see in section 4.3.5, a similar observation may be made for the case of nonlinear approximation in our two-dimensional case.

In table 4.3.4 the efforts of straight forward interpolation or collocation are tabulated. The root mean square residual in the potential is  $O(10^{-13})$  (based on the computed coefficients and resumming the series). The condition number is  $1.24 \times 10^5$  so that ill conditioning is not expected to be a factor. Without question, if the series derivatives are necessary (as they are in the full time dependent approximation), then collocation must be used with some caution.

**Table 4.3.4 Collocation**

**Cusping Profile - Linear approximation,  $t=0.25$ ,  $n=m=15$**

$x$	$\phi_x$	$\phi_{n_x}$	$\phi_y$	$\phi_{n_y}$
0.0	0.0	0.0	0.194 E+01	0.193 E+01
0.143	-0.620 E+00	-0.480 E+00	0.119 E+01	0.120 E+01
0.286	-0.455 E+00	0.163 E+01	0.868 E+00	0.105 E+01
0.5	-0.266 E+00	-0.591 E+02	0.3 E+00	-0.374 E+01
0.786	-0.102 E+00	-0.463 E+03	0.683 E+00	-0.538 E+02
1.0	-0.549 E-14	-0.382 E-09	0.674 E+00	0.958 E+02

#### **Time dependent cusping problem**

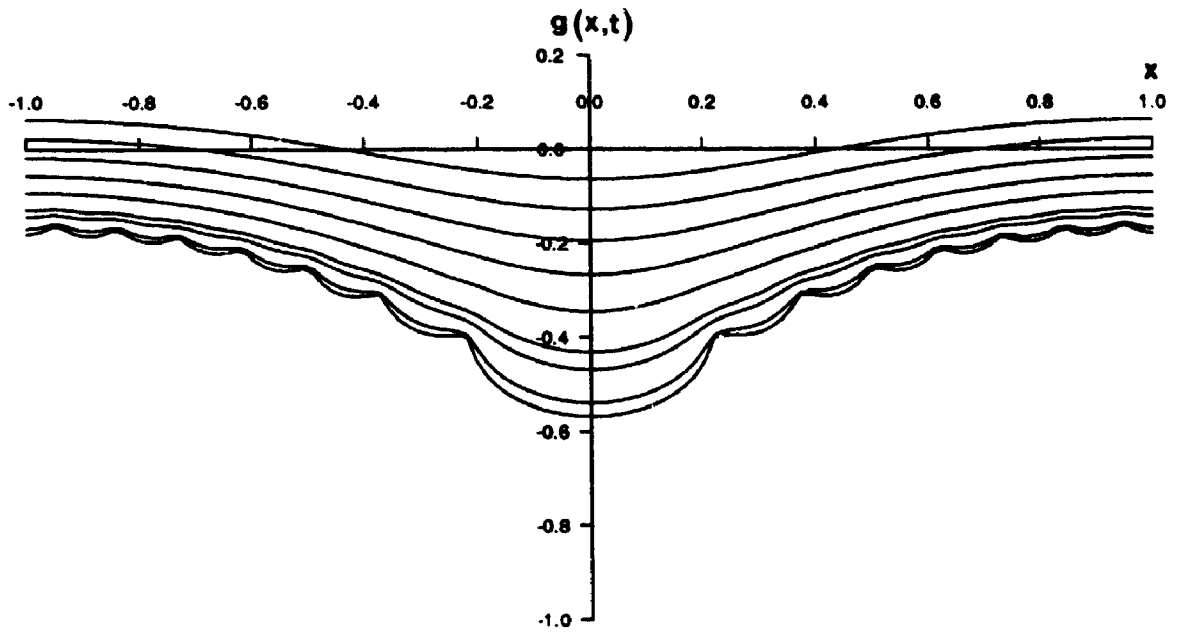
We now turn to the numerical integration of the full time dependent problem. The results are intended to assess the ability of the numerical scheme to monitor the moving boundary and are presented in graphical and tabular form for a variety of parameters. In the first place, the computed free surface profile is plotted for a number of different times and for each of several values of  $n$  and compared directly with the analytic profile at the same times. These plots are generated by joining the boundary points by straight line segments. Symmetry is used to display the entire profile on  $-1 \leq x \leq 1$ . A precise comparison of the computed and analytic values of  $g(0)$  is given in tabular form, clearly indicating the accuracy of the most rapidly varying component on the boundary. Finally,

the growth of the global error with time is compared graphically for each of several values of  $n$ . This is contained in plots of the maximum error between the  $n$  numerical and analytic boundary points.

Computations were performed with both the Eulerian and Lagrangian descriptions of the moving boundary. The Eulerian description was used to produce the interface profiles shown in figure 4.3.3(a). Fifteen unknown coefficients were included in the approximate solution  $\phi_n$ . The approximate and analytic profiles are compared in figure 4.3.3(b). The analytic profile is indicated by a dotted line. The two profiles are indistinguishable at early times, but differ significantly as the cusp time  $t = 0.35945$  is approached. The numerical solution shows a noticeably rounded trough near  $x = 0$ , where a steep cusping behaviour should be developing. More detail on the behaviour of  $g(0)$  is provided in the table 4.3.5.

The time stepping procedure has a step-size control (discussed in chapter 3). If the estimated truncation error of the scheme exceeds a specified tolerance, then the time step-size is decreased. A typical value of the tolerance in these calculations is  $\tau = 0.01$ . The maximum time step was taken as  $h_{\max} = 0.05$ . For the calculation shown in figure 4.3.3 and table 4.3.5, after about  $t = 0.26$  the time step is steadily reduced by the automated routine until a minimum value is reached (unless stated otherwise, the minimum time step was chosen as  $\Delta t = 0.001$ ). It is this stopping criterion which determines the overall run time. (It should be noted however, that if the time step is deliberately fixed in the o.d.e. algorithm, the numerical solution will continue further and may even exceed the cusp time.) For the case shown in figures 4.3.3(a), (b) the minimum time step was reached at approximately  $t = 0.33$ . The last profile shown in the figures is  $t = 0.32$ . The total computing effort was less than ten seconds of CPU time.





**Fig. 4.3.3 (a) Cusping Profiles Using a Linear Eulerian Approximation**

**( $n=15, m=60$ )**

**$t = 0.0, 0.05, 0.15, 0.20, 0.25, 0.27, 0.31, 0.32$**

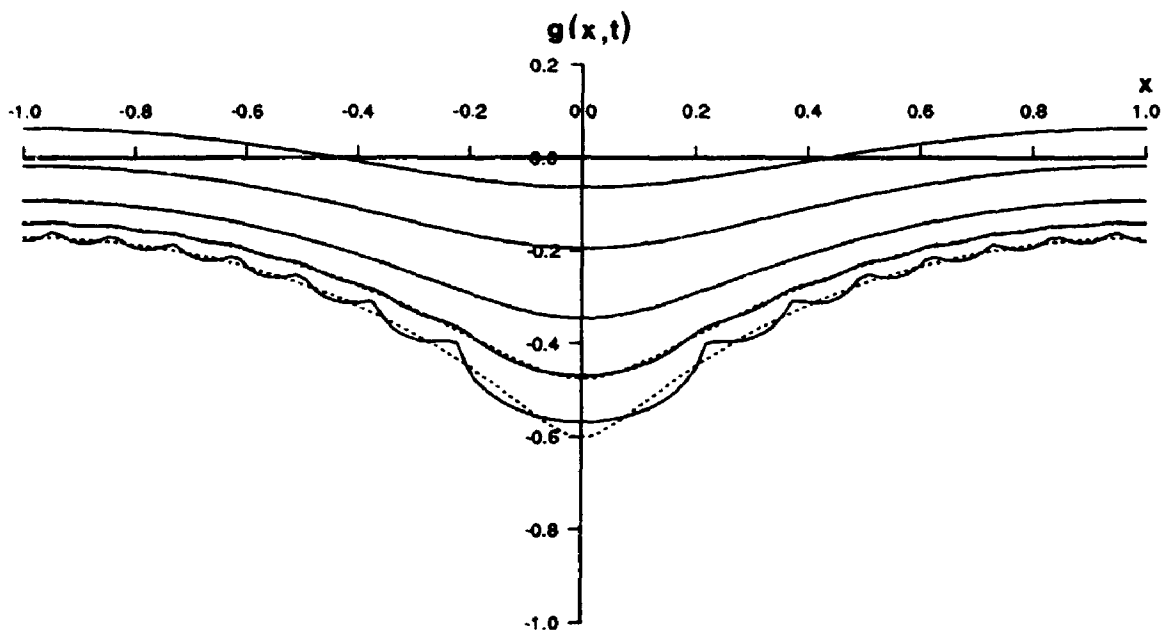


Fig. 4.3.3 (b) Comparison of Exact and Computed Profiles  
( $n=15$ ,  $m=60$ )

$t = 0.0, 0.10, 0.20, 0.27, 0.32$

Notice that corrugations have appeared in the free surface profile by a time of about  $t = 0.25$ . This is a typical feature of many numerical solutions to unstable flow problems. What is particularly interesting in this case is that the instabilities or corrugations manifest themselves quite early.

Now, we have already seen that  $n = 15$  is inadequate to accurately follow the true solution (cf tables 4.3.2, 4.3.3). But, by increasing the value of  $n$ , the numerical scheme only reaches its minimum  $\Delta t$  at an earlier time. For example, with  $n = 30$  the routine reaches the minimum time step at approximately  $t = 0.31$ . The computed profile for a number of times and for the case  $n = 30$  is shown in figure 4.3.4(a) and compared with the analytic solution in figure 4.3.4(b). Both of the cases  $n = 15$  and  $n = 30$  were generated using the same number of boundary points,  $m = 60$ , placed at equal  $x$  intervals along  $0 \leq x \leq 1$ . The case  $n = 30$  exhibits a less rounded trough than  $n = 15$  and the onset of the corrugations is delayed to about  $t = 0.27$ . Nevertheless, once the corrugations have set in, the growth of error is more rapid in the case  $n = 30$ . This is partially responsible for the increased reduction in step size; but even with a fixed time step, the oscillations in the profile are so great that the routine cannot proceed much beyond  $t = 0.33$ . This calculation took about 25 seconds of CPU time.

If the number of series terms is increased even further to  $n = 45$ , the situation is hopeless. The minimum time step is reached at approximately  $t = 0.26$  with greater error than both previous cases. A comparison of the growth of the maximum error with time is given in figure 4.3.5 for  $n = 15, 30, 45$ . Due to the exponential growth of the error, we have chosen to plot  $\log_{10}(E_{\max})$  versus time. Although it is clear that the larger value of  $n$  yields a decreased error at early times, the exponential growth is more marked as  $n$  increases and  $t$  approaches the cusping time. This is believed to be a consequence of the basic instability of the physical model and not so much a feature of the time dependent

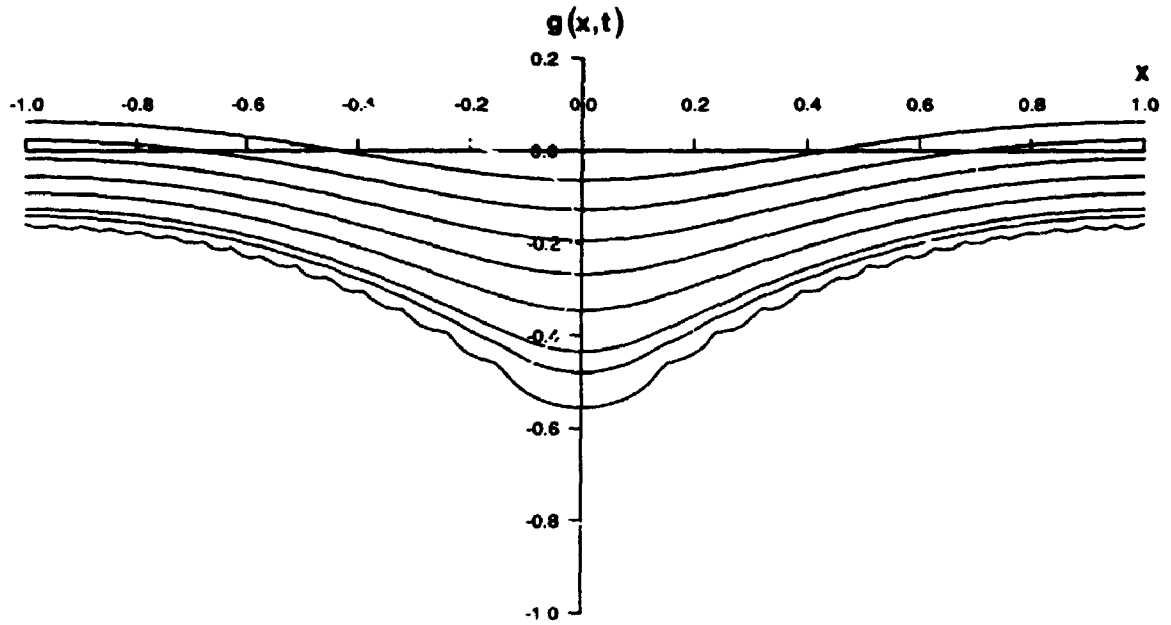


Fig. 4.3.4 (a) Cusping Profiles Using a Linear Eulerian Approximation

( $n=30$ ,  $m=60$ )

$t = 0.0, 0.05, 0.10, 0.15, 0.20, 0.25, 0.27, 0.31$

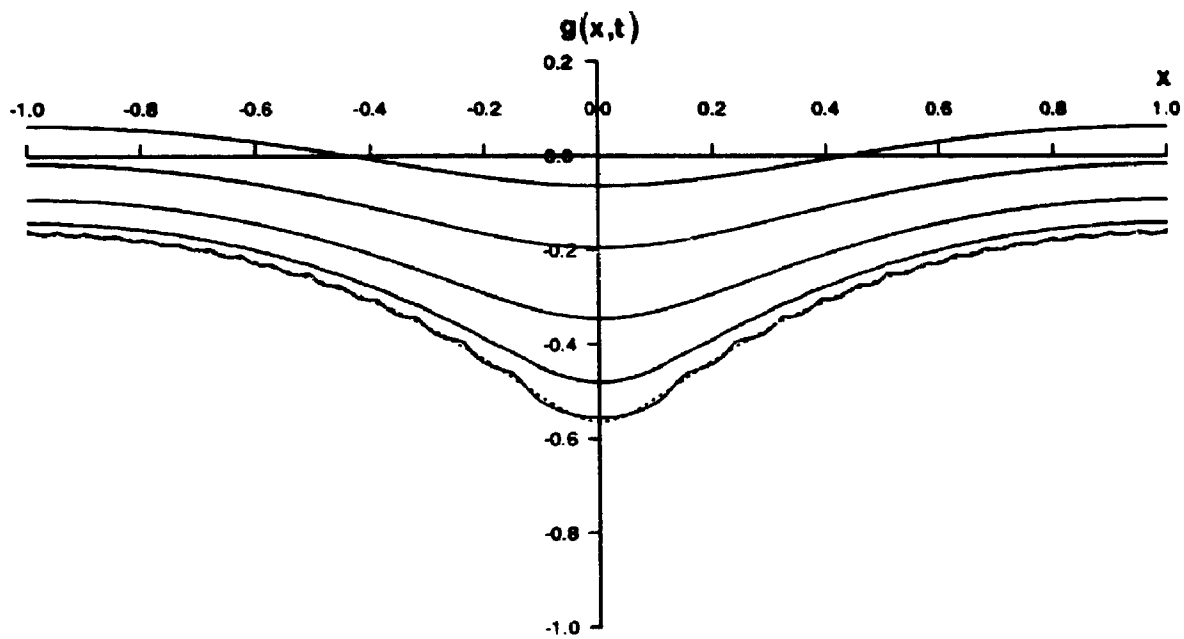


Fig. 4.3.4 (b) Comparison of Exact and Computed Profiles  
( $n=30, m=60$ )

$t = 0.0, 0.10, 0.20, 0.27, 0.31$

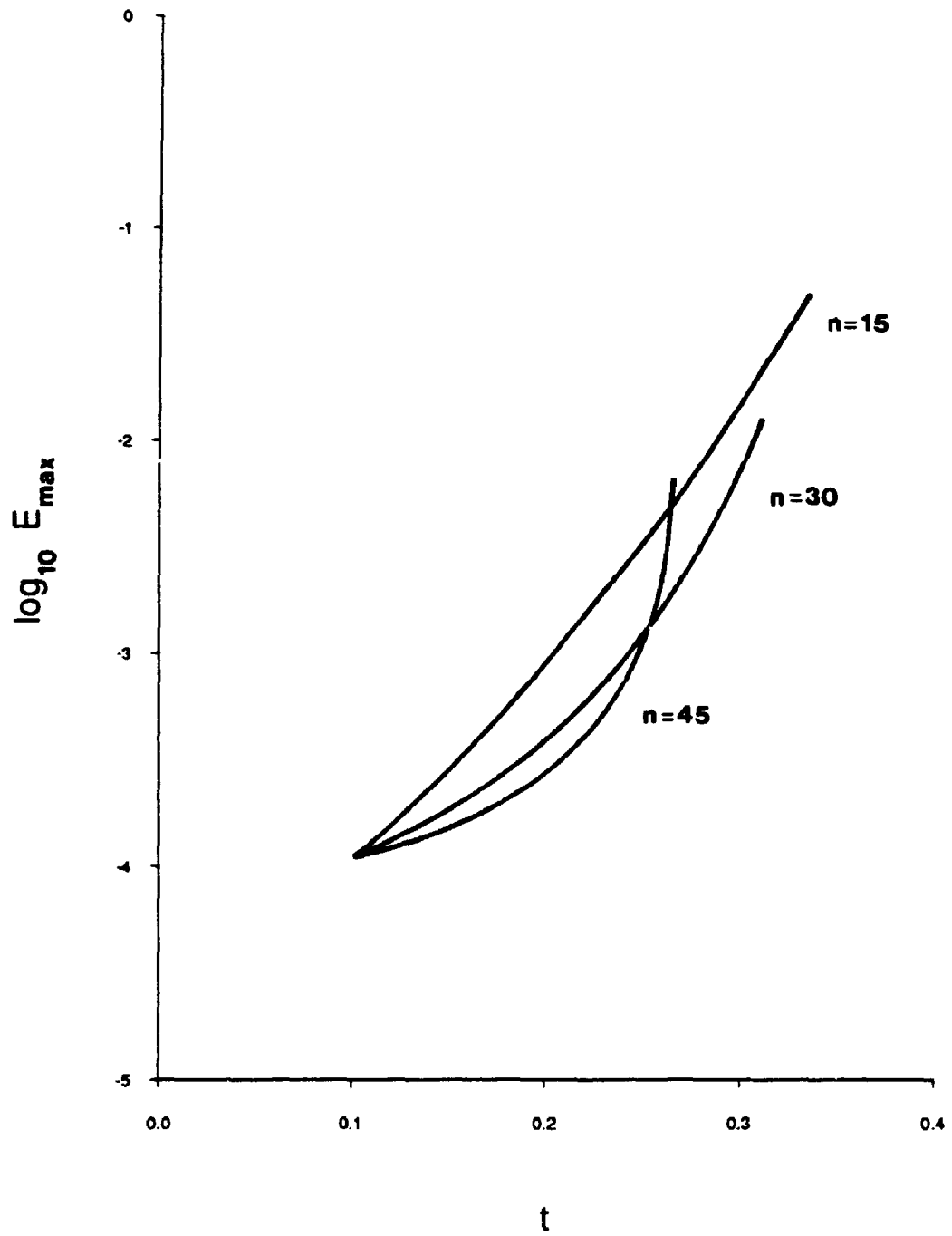


Fig. 4.3.5 Error Growth ( $n=15,30,45$ )  
Cusping Case - Linear Approximation

numerical scheme, since a lesser number of time steps is taken to reach minimum  $\Delta t$  as  $n$  is increased. A direct comparison of computed and exact values of  $g(0)$  is provided in table 4.3.5. After about  $t = 0.30$ , the numerical values begin to deviate from the analytic.

**Table 4.3.5**

**Comparison of Exact and Computed Values of  $g(0,t)$   
Cusping Profile - Linear approximation - Eulerian Description**

t	n=15	n=30	n=45	Exact
0.10	-0.1950	-0.1950	-0.1950	-0.1949
0.20	-0.3452	-0.3458	-0.3458	-0.3454
0.25	-0.4316	-0.4347	-0.4350	-0.4346
0.30	-0.5272	-0.5357	----	-0.5438
0.33	-0.5876	----	----	-0.6313

When the Lagrangian approach is used to advance points on the free surface, similar features as above are observed. After  $t = 0.3$ , the analytic cusping solution spikes sharply in the vicinity of  $x = 0.0$  and yet this is poorly represented by the numerical routine. In fact, the trough is more rounded than the Eulerian case (compare figure 4.3.6(a) where  $n = 30$  with figure 4.3.4(a) where  $n = 30$ ). However, the Lagrangian approach does permit

the routine to proceed further in time before the minimum time step is reached (it continues beyond the actual cusping time to about  $t = 0.37$ ). What is more, the undesirable corrugations do not appear until  $t = 0.35$ , very close to the cusp time.

The case  $n = 30$  ( $m = 60$ ) is shown in figure 4.3.6(a) and compared with the analytic solution in figure 4.3.6(b). The unusual bubble formation in the trough is strikingly different from the true cusping behaviour. The error in the trough by  $t = 0.358$  dominates the error produced by corrugations in the shoulder regions. Unlike the Eulerian case, very little difference was observed between  $n = 15$  and  $n = 30$ . All of the Lagrangian calculations were made with the same error tolerance  $\tau$  (0.01),  $h_{\max}$  (0.01) and  $h_{\min}$  (0.001) as the Eulerian. Once again, it was found that attempts at improvement by increasing to  $n = 45$  terms only resulted in the routine reaching its minimum time step at an earlier value of  $t$  (about  $t = 0.27$ ).

The exponential behaviour of the error growth with time is similar to the Eulerian case and is not repeated here.

It might be thought possible to improve the accuracy of both the Eulerian and Lagrangian descriptions by placing more severe restrictions on the error tolerance  $\tau$  or on the maximum stepsize  $h_{\max}$ . This is not the case. While the Lagrangian approach does permit run times comparable to those of figure 4.3.6 if small changes are made in  $\tau$  or  $h_{\max}$ , both descriptions contribute larger corrugations in the shoulder regions if  $\tau$  or  $h_{\max}$  are reduced. This feature is illustrated in figure 4.3.7, where the Lagrangian case of  $n = 30$  has been recomputed with the same tolerance  $\tau$ , but  $h_{\max}$  reduced to 0.01. The figure shows the computed and analytic profiles at  $t = 0.31$  and the corrugations are noticeably larger than those of the corresponding time shown in either figures 4.3.4(a) or 4.3.6(a). This is believed to be a feature of both the numerical method and the unstable physical problem. By forcing the routine to perform a larger number of time steps, the error accumulated at each step is magnified.



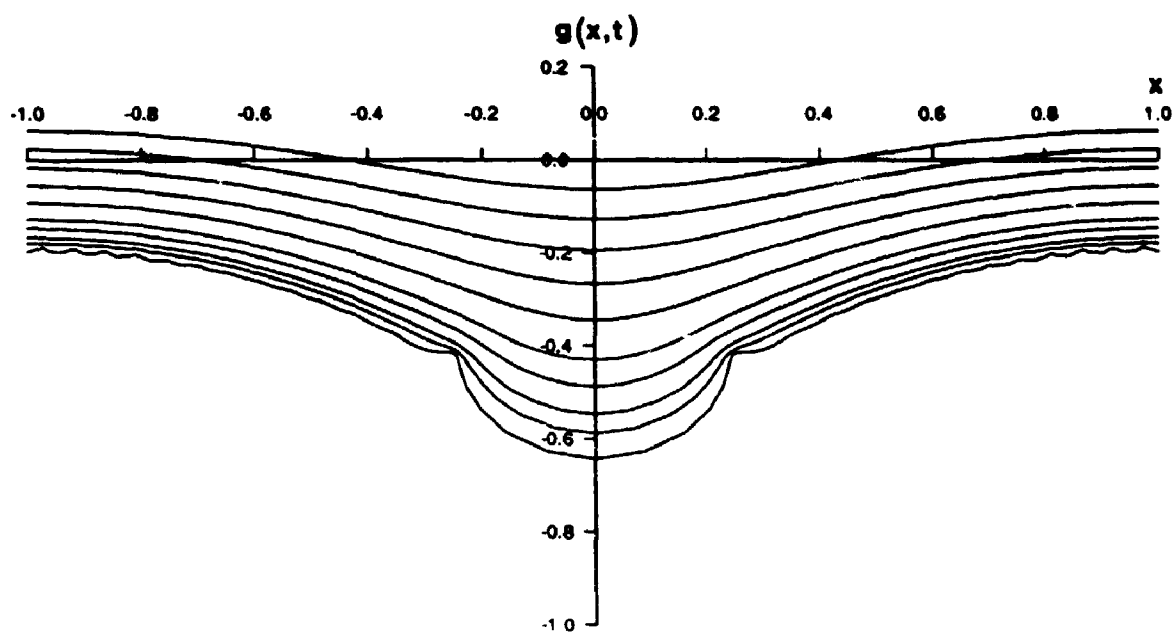


Fig. 4.3.6 (a) Cusping Profiles Using a Linear Lagrangian Approximation

( $n=30$ ,  $m=60$ )

$t = 0.0, 0.05, 0.10, 0.15, 0.20, 0.25, 0.28, 0.31, 0.33, 0.358$

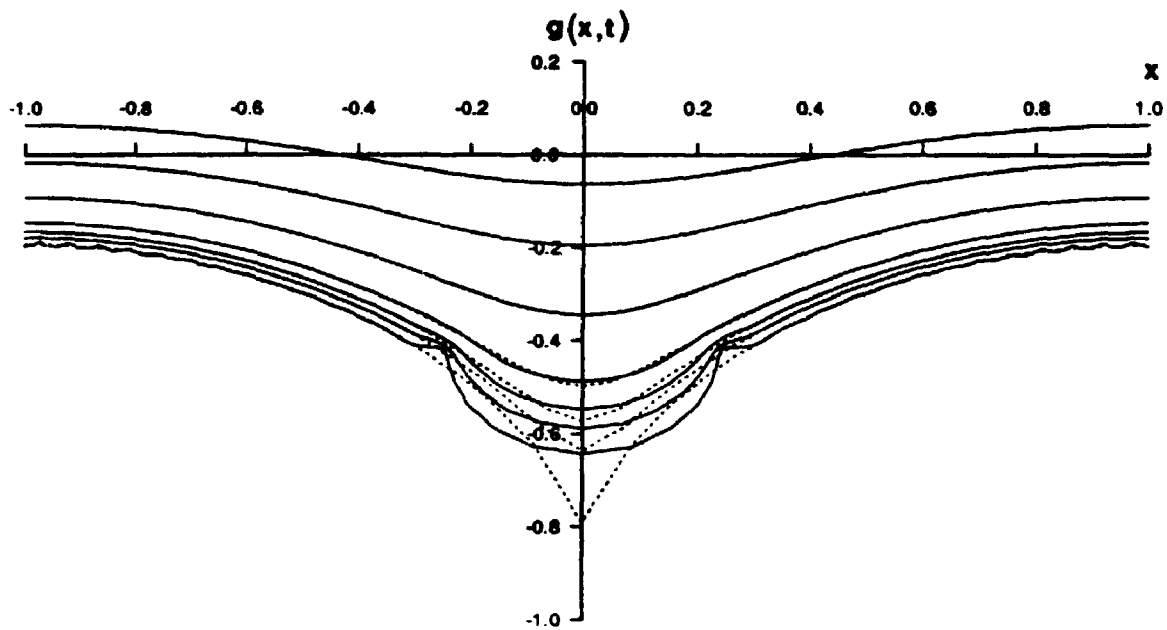


Fig. 4.3.6 (b) Comparison of Exact and Computed Profiles  
( $n=30, m=60$ )

$t = 0.0, 0.10, 0.20, 0.28, 0.31, 0.33, 0.358$

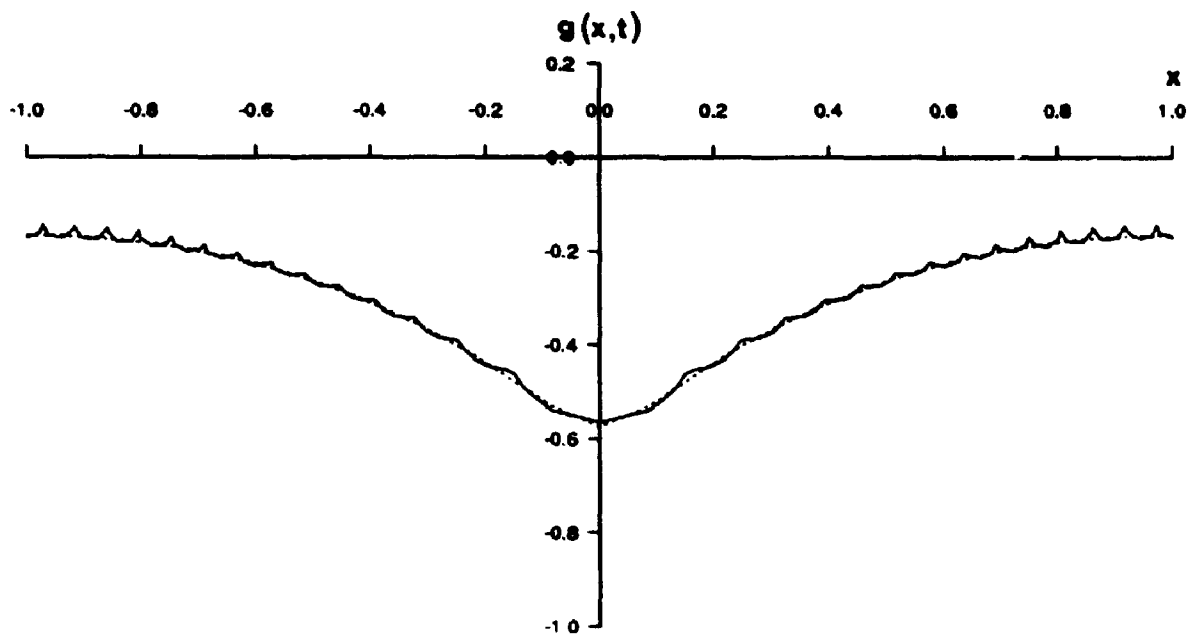


Fig. 4.3.7 Reduced  $h_{\max}$   $t = 0.31$

Presumably, if the perturbations in the free surface profile were kept in check, the routine would have a better chance of emulating the correct cusping behaviour. Now, the magnitude of the corrugations can be artificially damped by smoothing the free surface every few time steps. We have taken our best case (the Eulerian approach with  $n = 30$ ) and smoothed the free surface by fitting it with a cubic spline in the least squares sense. The smoothing was performed every five time steps with the results shown in figure 4.3.8. The profile at the time of about 0.31 can be compared directly with the corresponding profile at  $t = 0.31$  in figure 4.3.4(b). The obvious effect of smoothing is an improved quality of computed free surface; but this is all. The same bubble shaped trough develops as the time advances. At later times the spatial accuracy simply isn't sufficient to model the correct behaviour.

### Conclusions

While it is true that the numerical results do suggest a steepening of the free surface profile in the vicinity of  $x = 0$ , it is clear that the numerical scheme fails to reproduce the true cusping behaviour as the time to cusp is approached. There are several reasons for this. First of all, the presence of the singularity in the analytic solution has a negative influence on the accuracy of the potential approximation. The only recourse for improvement involves increasing the value of  $n$ ; but this only contributes to the basic instability of the flow pattern. That is, increased spatial accuracy leads to an increase in the number of roundoff errors, which may in turn be interpreted as an introduction of smaller wavelength instabilities to the initial perturbation. The linearized stability analysis of Saffman and Taylor (1958) predicts a rapid growth of such errors. Also, the numerical solution of the ordinary differential equations involved requires a steady decrease in the step size as the time to cusp is neared. This is largely due to an increasingly unstable or stiff system of ordinary differential equations. More will be said on the stability of the time dependent scheme in section 4.6.

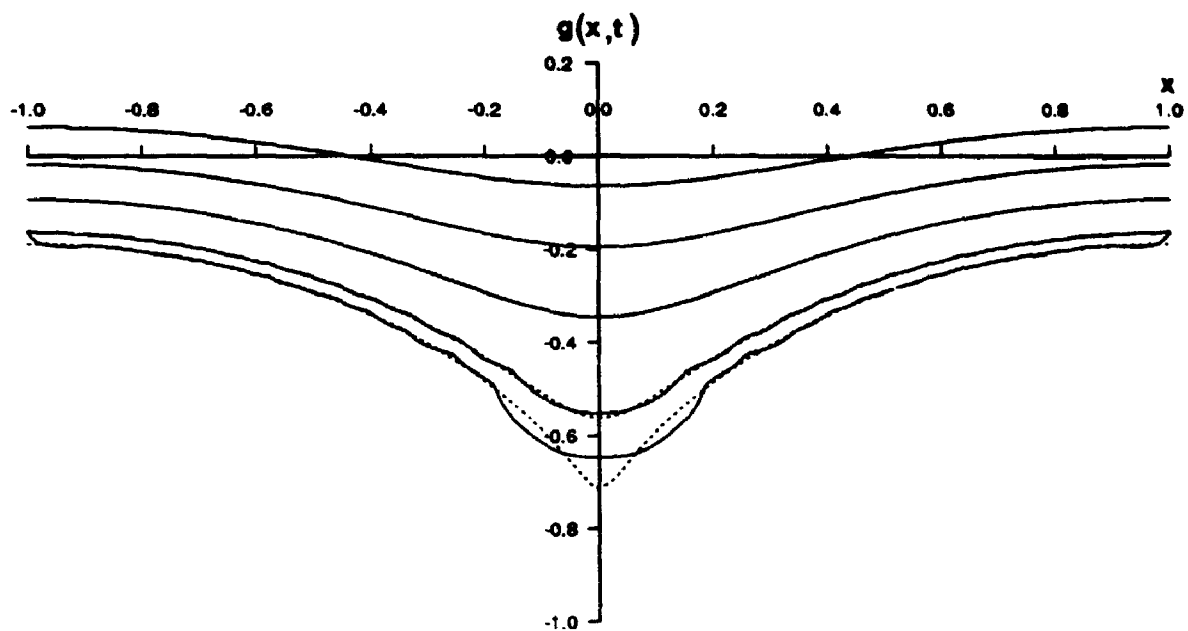


Fig. 4.3.8 Smoothed Cusping Profiles

$t = 0.0, 0.10, 0.20, 0.31, 0.349$

Aitchison and Howison (1985) have tested a boundary integral method on this problem and report much the same findings. That is, they observe a breakdown in the numerical computation as cusp time is neared (for a direct comparison, their times must be divided by  $\pi$ .); and increasing the number of boundary elements only hastens the breakdown. They do however, appear to follow the correct analytic behaviour for longer times (for a direct comparison, their times must be divided by  $\pi$ ). For example, their best numerical calculation involves only sixteen boundary points, yet proceeds to a time of  $t = 0.34$ . They do not present a direct comparison with the analytic solution, but only a graphical illustration of their numerical results.

We point out one advantage that a boundary integral method might possess. It is possible to interpret the discretized form of boundary integral methods as linear approximation schemes of the type discussed here, but with the basis functions being fundamental solutions of the Laplace equation. These solutions take the form

$$u(x, y; \xi, \eta) = \log\left(\frac{1}{r}\right)$$

$$r = \sqrt{(x - \xi)^2 + (y - \eta)^2}$$

where the points  $(\xi, \eta)$  are fixed points on the boundary of the domain. The fundamental solutions are singular solutions and its probable that this is a desirable feature if the solution to the boundary value problem itself has a singularity near the boundary (cf section 4.5).

#### 4.4 Saffman Finger: Linear Approximation

The Hele-Shaw problem is looked at once again in this section, but with a slightly different initial condition. This time, the initial free surface profile is given by

$$y = -\ln(\epsilon \cos x + \sqrt{1 - \epsilon^2 \sin^2 x}) \quad (4.4.1)$$

To first order in  $\epsilon$  we still have

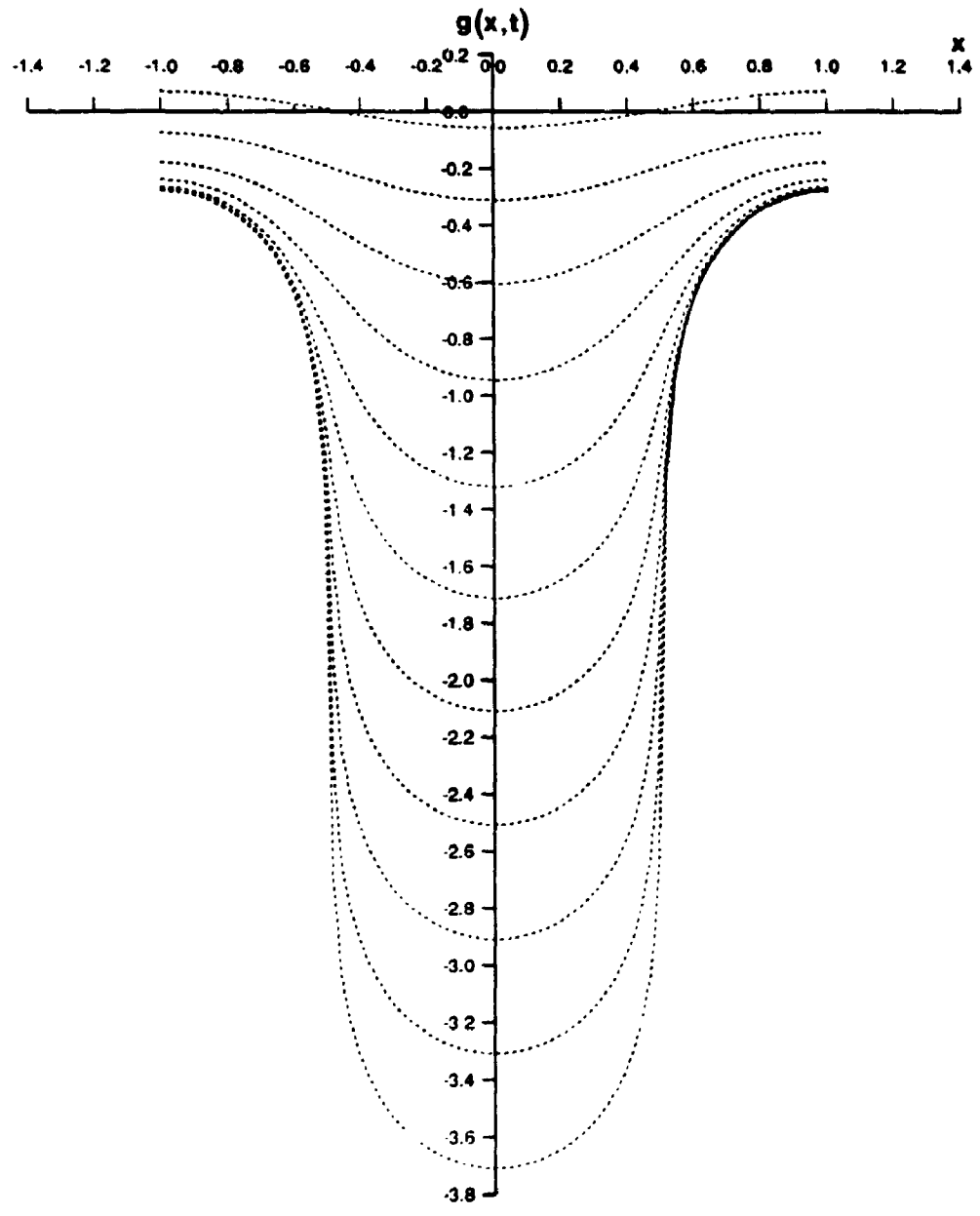


Fig. 4.4.1 Exact Saffman Finger

$$y(x) \approx -\epsilon \cos x.$$

In this example a cusp does not develop, but instead a long finger forms, projecting into the fluid (see figure 4.4.1). Eventually, the tip of the finger assumes a constant profile, filling one half the channel width. The analytic solution takes the form

$$z = i(w - d(t) - \ln(1 + a(t)e^w))$$

where

$$d(t) = t + \frac{1}{2} \ln(1 - \epsilon^2 + \epsilon^2 e^{2t})$$

$$a(t) = \left( \frac{\epsilon^2}{(1 - \epsilon^2)e^{-2t} + \epsilon^2} \right)^{\frac{1}{2}}$$

$\epsilon = 0.2$  is used in the following calculations. This solution was discovered by Saffman (1959) and is commonly referred to as the Saffman finger.

Once again, the example represents a severe test for the numerical method. This time, the free surface profile distorts without limit. And once again, we would expect that the usual attempts to improve the accuracy of the numerical results (ie an increase in the number of series terms and boundary points) only accelerates the growth of errors.

#### **Test of potential problem**

The trial solution is of the form (4.3.16). Before presenting the full time dependent results, the effectiveness of the potential portion of the numerical scheme is examined. As before, we fix the time and solve the potential problem on a known domain (using the analytic free surface profile) and compare with the analytic solution.

Results of the potential problem for given cusping and Saffman profiles at  $t = 0.0$  are compared in table 4.4.1. Even though the same number of unknowns are involved, there is a remarkable discrepancy in the accuracy of the two approximations. Bear in mind



that the two profiles are virtually identical at this time. Ill conditioning is clearly not a factor. The only explanation can be the closeness of the singularity to the domain in the cusping case.

**Table 4.4.1 Comparison of Cusping and Saffman  
Solutions at  $t=0.0$**

	n	m	$E_{max}$	$E_{\phi}$	$E_{\phi_x}$	$E_{\phi_y}$	Condition Number
Cusping	15	45	0.153 E-06	0.754 E-07	0.326 E-05	0.333 E-05	4.87 E+01
Saffman	15	45	0.788 E-12	0.511 E-12	0.233 E-10	0.237 E-10	6.15 E+01

Tables 4.4.2 and 4.4.3 demonstrate convergence of the approximation scheme for the potential problem of the Saffman profile at times 0.25 and 0.5, respectively. Clearly the convergence properties are superior to those of the previous cusping case. Once again, this is largely due to the absence of any external singularity in the Saffman case. For comparison purposes, some collocation results are included in table 4.4.2. The condition numbers are reported to indicate that ill-conditioning is not expected to have influenced the numerical results. Although the collocation results are acceptable, the corresponding overdetermined solution is superior in all cases.

**Table 4.4.2 Convergence**Saffman Profile - Linear approximation,  $t=0.25$ 

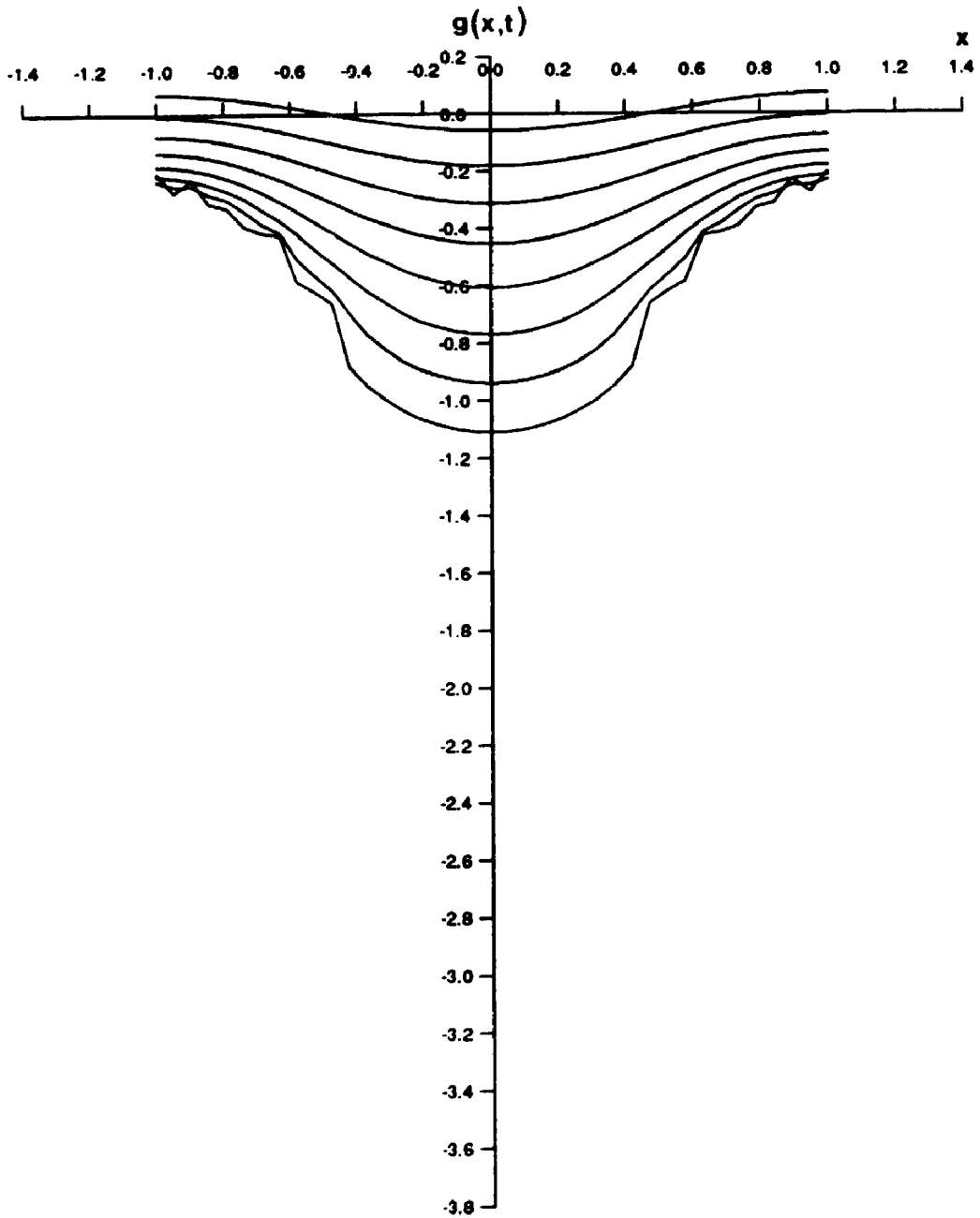
n	m	$E_{\max}$	$E_{\phi}$	$E_{\phi_x}$	$E_{\phi_y}$	Condition Number
10	30	0.550 E-05	0.337 E-05	0.960 E-04	0.985 E-04	1.04 E+02
15	45	0.422 E-07	0.257 E-07	0.110 E-05	0.112 E-05	1.91 E+03
	15	<i>0.150 E-04</i>	<i>0.949 E-14</i>	<i>0.262 E-03</i>	<i>0.725 E-04</i>	<i>6.97 E+04</i>
20	60	0.364 E-09	0.220 E-09	0.126 E-07	0.128 E-07	3.56 E+04
	20	<i>0.153 E-05</i>	<i>0.205 E-13</i>	<i>0.333 E-04</i>	<i>0.778 E-05</i>	<i>5.07 E+06</i>
25	75	0.332 E-11	0.201 E-11	0.144E-09	0.146 E-09	6.73 E+05
	25	<i>0.105 E-06</i>	<i>0.196 E-13</i>	<i>0.282 E-05</i>	<i>0.579 E-06</i>	<i>3.73 E+08</i>
30	90	0.746 E-13	0.196 E-13	0.166 E-11	0.167 E-11	1.28 E+07

**Table 4.4.3 Convergence**  
**Saffman Profile - Linear approximation,  $t=0.50$**

$n$	$m$	$E_{\max}$	$E_{\phi}$	$E_{\phi_x}$	$E_{\phi_y}$
10	30	0.181 E-02	0.105 E-02	0.246 E-01	0.252 E-01
15	45	0.213 E-03	0.122 E-03	0.423 E-02	0.431 E-02
20	60	0.276 E-04	0.156 E-04	0.722 E-03	0.731 E-03
25	75	0.378 E-05	0.214 E-05	0.123 E-03	0.124 E-03
30	90	0.536 E-06	0.303 E-06	0.208 E-04	0.210 E-04

#### **Time dependent Saffman finger**

The Saffman profiles computed using an Eulerian description, are shown in figures 4.4.2(a) and 4.4.3(a) for the cases  $n = 10$  ( $m = 20$ ),  $n = 20$  ( $m = 40$ ), respectively. The constant reductions in stepsize necessary in the cusping calculations were not observed here. In fact, with  $\tau = 0.01$  and  $h_{\max} = 0.01$ , the solution could be computed to about  $t = 0.7$  without the routine decreasing  $\Delta t$ . The results shown in figures 4.4.2 and 4.4.3 are the result of running the program with constant  $\Delta t = 0.01$  until the routine requested a stepsize reduction. For both of the cases  $n = 10, 20$ , the programs were terminated in this way shortly after  $t = 0.7$ . If the stepsize was permitted to decrease after this time, it reached a value of  $\Delta t = 0.001$  within only a few time steps. This was due to the rapid breakup in the free surface near the time of  $t = 0.7$ , as can be seen from the figures.



**Fig. 4.4.2 (a) Saffman Profiles Using a Linear Eulerian Approximation**

**( $n=10, m=20$ )**

**$t = 0.0, 0.10, 0.20, 0.30, 0.40, 0.50, 0.60, 0.70$**

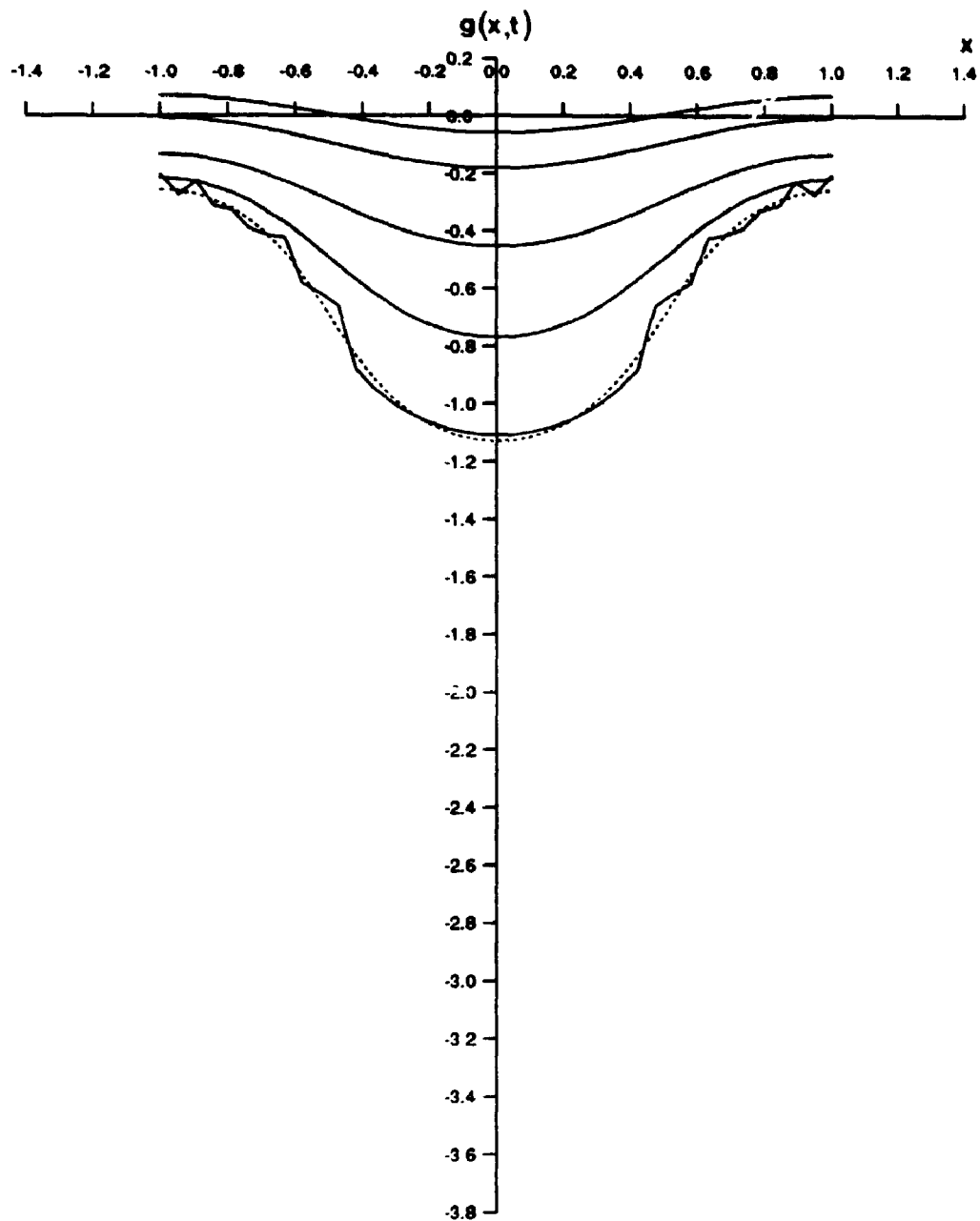


Fig. 4.4.2 (b) Comparison of Exact and Computed Profiles

( $n=10, m=20$ )

$t = 0.0, 0.10, 0.30, 0.50, 0.70$

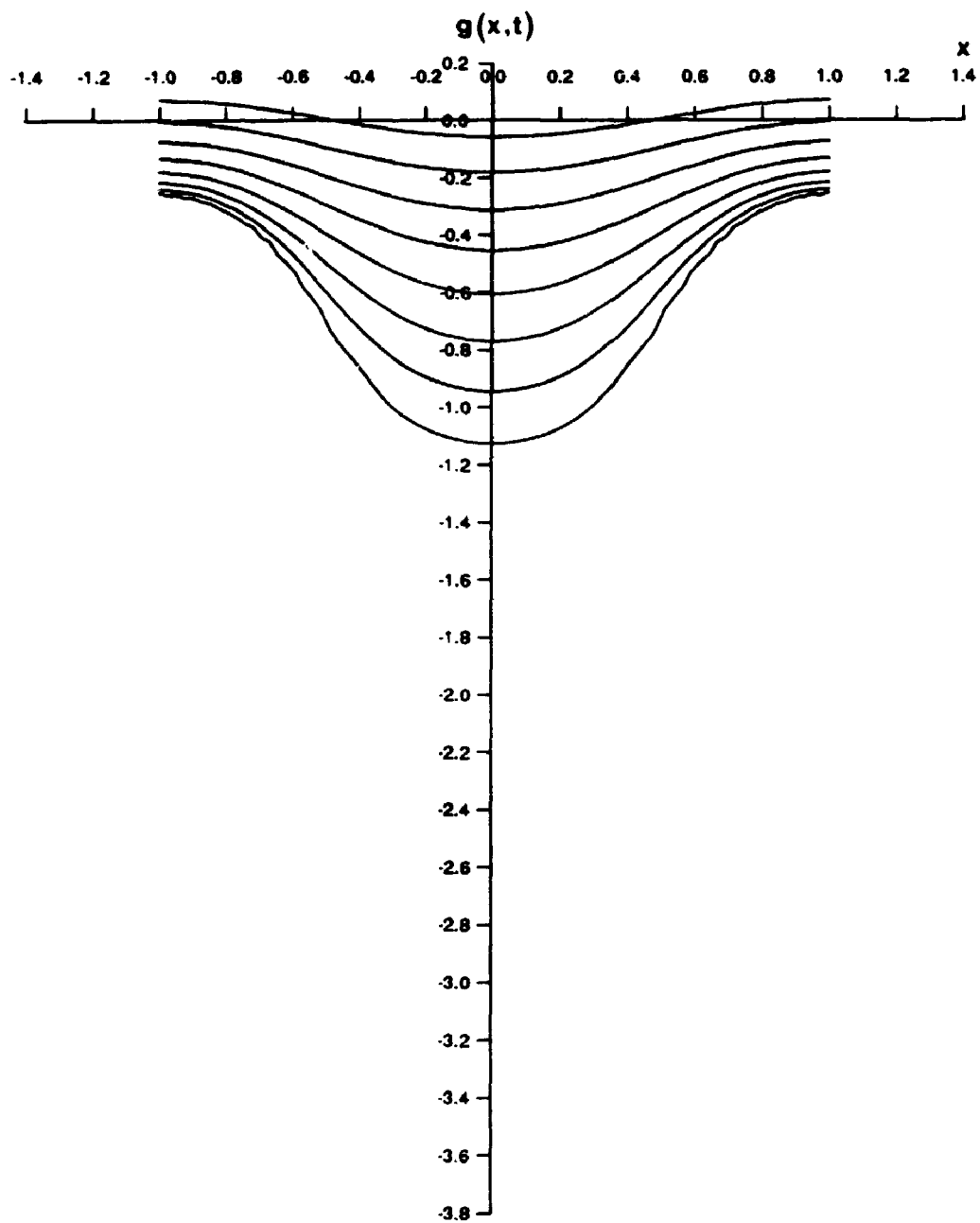


Fig. 4.4.3 (a) Saffman Profiles Using a Linear Eulerian Approximation

( $n=20, m=40$ )

$t = 0.0, 0.10, 0.20, 0.30, 0.40, 0.50, 0.60, 0.70$

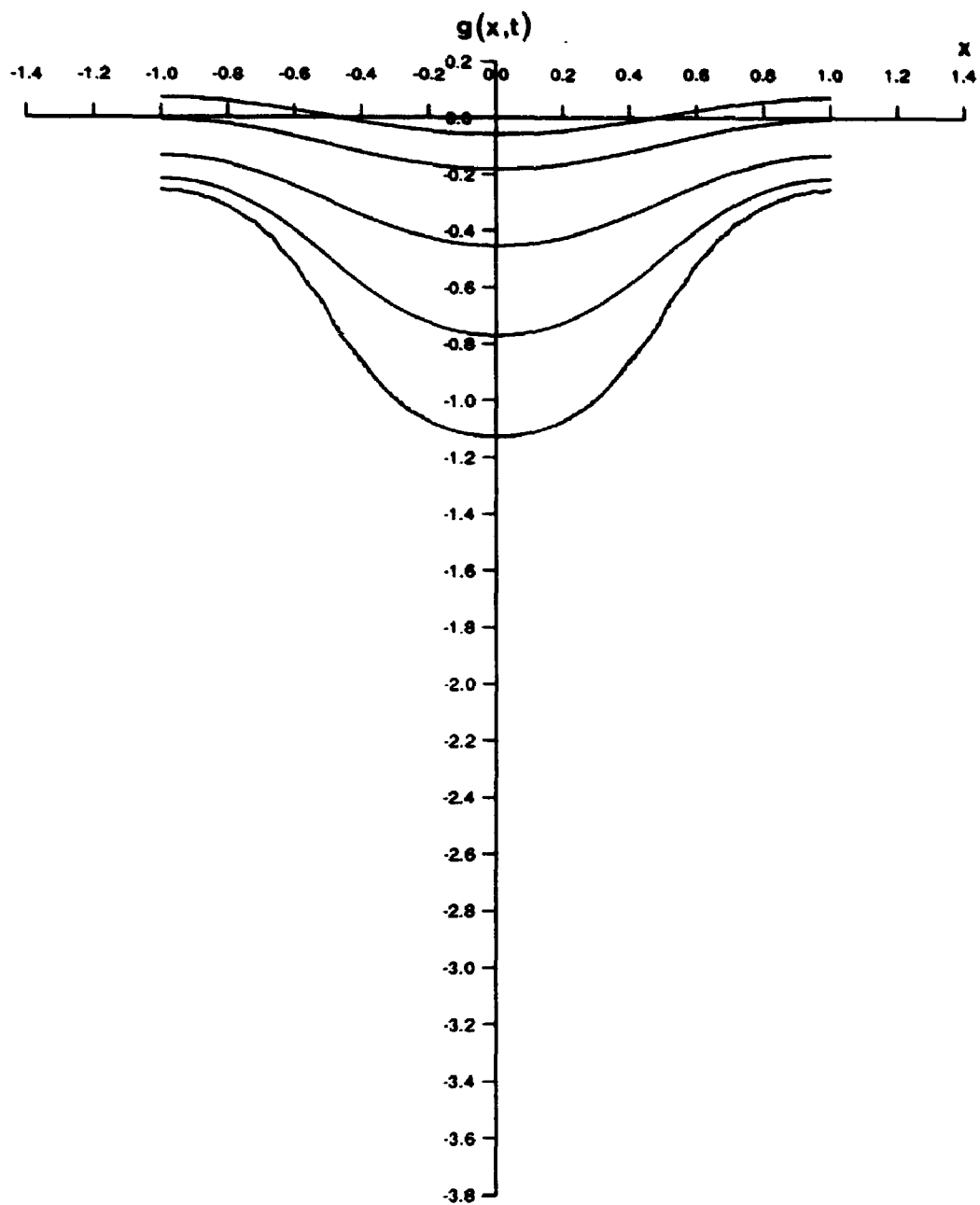


Fig. 4.4.3 (b) Comparison of Exact and Computed Profiles

( $n=20, m=40$ )

$t = 0.0, 0.10, 0.30, 0.50, 0.70$

Comparisons with the analytic profiles are made in figures 4.4.2(b) and 4.4.3(b). Precise values of  $g(0)$  are presented in table 4.4.4. Both of the cases  $n = 10, 20$  perform quite well on this problem, following the exact solution closely until about  $t = 0.6$ . As with the cusping case, corrugations eventually appear in the shoulder region, but at a much later time. The results produced with  $n = 20$  are an improvement over  $n = 10$  and the onset of corrugations has been forestalled a few time steps.

Just as with the cusping case, however, there is a limit to the benefits that can be achieved simply by increasing  $n$ . For example, with  $n = 25$  ( $m = 50$ ), the free surface breaks up sooner, the run being terminated by  $t = 0.5$ ; with  $n = 30$  ( $m = 60$ ), the run was terminated by  $t = 0.3$ . The growth of error with time is plotted in figure 4.4.4 for  $n = 10, 20, 25$ .

Table 4.4.4

Comparison of Exact and Computed Values of  $g(0,t)$

Saffman Profile - Linear approximation - Eulerian description

t	n=10	n=20	n=25	Exact
0.1	-0.181343	-0.181343	-0.181343	-0.181341
0.3	-0.453424	-0.453441	-0.453441	-0.453435
0.5	-0.768663	-0.769947	-0.770530	-0.769968
0.7	-1.10902	-1.12631	----	-1.12958



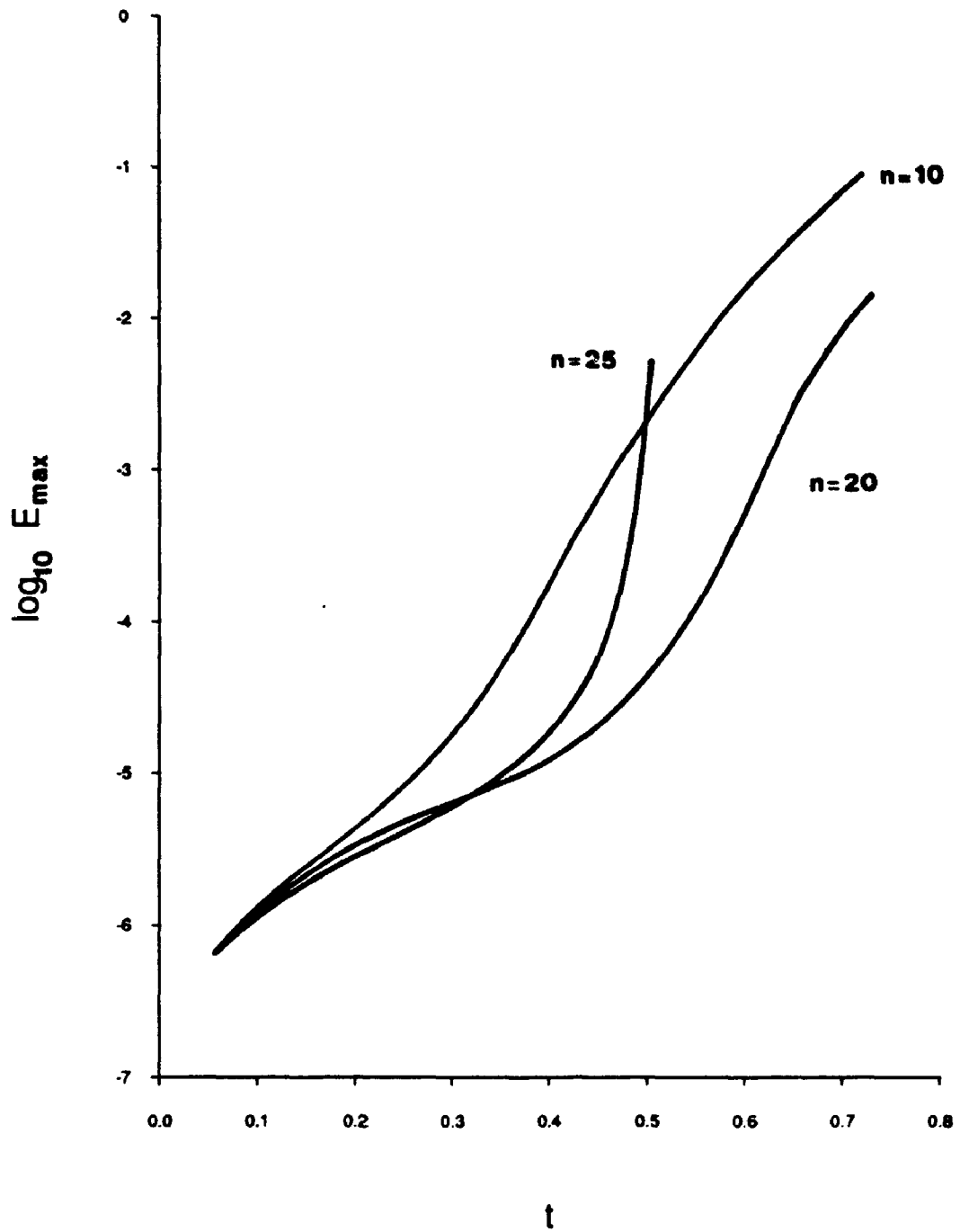


Fig. 4.4.4 Error Growth ( $n=10,20,25$ )  
Saffman Finger - Linear Approximation

The computing costs for these runs varied from less than 10 seconds of CPU time for the case  $n = 10$  to just over one minute for the case  $n = 30$ .

For comparison purposes, the problem was recomputed using a Lagrangian description. For the most part, no appreciable difference in the two approaches was observed. In the Lagrangian description, the points on the boundary drift away from the trough region and upwards along the steepening sides of the Saffman finger. This effect was not very large in the cusping problem, but over the increased time span of the Saffman case, points initially near  $x = 0$  have migrated a substantial distance. In order to produce a respectable profile, a large number of boundary points is required. Alternatively, points could easily be clustered in the vicinity of  $x = 0$  initially, in hope of producing a reasonable profile of the surface at later times. However, when this variable grid was employed, it was found that the simulation did not proceed nearly as far in time as did the Eulerian cases.

### Conclusions

It is possible to accurately follow the Saffman profile for times much greater than those involved in the cusping profile. This must in part be due to the lack of any singularities near the boundary and the more smoothly varying free surface. Of course, the usual difficulties are reported when it comes to attempts at increased accuracy. The basic instability of the problem and the increased ill conditioning of the linear systems being solved, prevents the calculation from proceeding much beyond  $t = 0.7$ . Nevertheless, this is a substantial improvement over the boundary element calculation of Aitchison and Howison (1985) and is almost certainly more efficient. They are able to reach a time of about  $t = 0.4$  only, before the boundary breaks up. Once again, they do not present a direct comparison with the analytic solution, but only give their computed results in graphical form.

#### 4.5 Hele-Shaw Flows: Nonlinear Approximation

The results of the last two sections indicate that the linear approximation method is merely adequate for solving the potential portion of the full moving boundary problem. For example, tables 4.3.3 and 4.4.3 indicate that a reasonable number of terms (less than one hundred) in the series solution is sufficient to accurately solve this portion of the problem on the time scales used. However, the inherent instability of such problems and the poor conditioning of the matrices involved, makes the choice of anything more than fifteen to thirty terms prohibitive. If more terms are included, the errors incurred in the calculation are magnified at an alarming rate.

The results do however suggest, that if a means of accelerating the convergence of the linear method were available, then the solution might be accurately followed for longer times. The nonlinear approximation outlined in section 2.8 may be interpreted as such an acceleration of the linear method. In this section, the nonlinear boundary approximation is applied to the cusping and Saffman solutions.

Following the notation of section 2.8, the trial solution takes the form

$$\phi_n(x) = y + b_0 + \sum_{j=1}^n b_j \gamma(b_{j+n}; x) \quad (4.5.1)$$

where

$$\gamma(b_{j+n}; x) = \frac{\sinh \pi(b_{j+n} - y)}{\cosh \pi(b_{j+n} - y) - \cos \pi x}$$

and we have been careful to include a factor of  $\pi$  as the range of  $x$  values is  $-1 \leq x \leq 1$ .

As usual, the parameters  $b_j$  are solved for in the potential portion by first discretizing  $(x, y)$  and then applying the boundary condition (4.3.8). The nonlinear least squares problem is solved using the Levenberg-Marquardt algorithm outlined in section 2.8.

Preliminary estimates of convergence of the nonlinear scheme, in application to the cusping and Saffman solutions, respectively, are presented in tables 4.5.1 and 4.5.2. These should be compared directly with the tables 4.3.2 and 4.4.3 for the corresponding efforts of the linear approximation. It is clear that the nonlinear approach provides much stronger convergence properties and with a lesser number of parameters. Note that the presence of the singularity in the cusping case is still a factor, this despite the fact that a singularity exists in the approximate solution as well. Thus, even though the Saffman finger is at least as distorted at  $t = 0.5$  as the corresponding cusping profile at  $t = 0.25$ , the Saffman profile yields better convergence. It should also be noted that no appreciable difference was observed in the case of collocation. Nevertheless, in the calculations which follow, we have made it a practice to always solve the least squares problem with  $m$  strictly greater than  $2n + 1$ .

**Table 4.5.1 Convergence - Cusping Profile**

Nonlinear approximation,  $t=0.25$ ,  $m=24$ , number of unknowns =  $2n+1$ .

$n$	$E_{\max}$	$E_{\phi}$	$E_{\phi_x}$	$E_{\phi_y}$
3	0.137 E-03	0.165 E-03	0.460 E-02	0.513 E-02
4	0.153 E-04	0.196 E-04	0.559 E-03	0.898 E-03
5	0.191 E-05	0.230 E-05	0.114 E-03	0.818 E-04

**Table 4.5.1(b) The Nonlinear Parameters**

$n=5, m=24, t=0.25, g(0)=-0.435, g(1)=-0.127$

Coefficients	Singularities
$b_0 = -0.124$	$b_6 = 0.614$
$b_1 = 0.276$	$b_7 = 0.098$
$b_2 = 0.083$	$b_8 = -0.146$
$b_3 = 0.032$	$b_9 = -0.279$
$b_4 = 0.011$	$b_{10} = -0.347$
$b_5 = 0.002$	

**Table 4.5.2 Convergence - Saffman Profile**

Nonlinear approximation,  $t=0.5, m=24,$

number of unknowns =  $2n+1$

n	$E_{\max}$	$E_{\phi}$	$E_{\phi_x}$	$E_{\phi_y}$
3	0.214 E-04	0.202 E-04	0.368 E-03	0.380 E-03
4	0.862 E-06	0.864 E-06	0.203 E-04	0.209 E-04
5	0.347 E-07	0.367 E-07	0.105 E-05	0.108 E-05

**Table 4.5.2(b) The Nonlinear Parameters** $n=5, m=24, t=0.5, g(0)=-0.770, g(1)=-0.217$ 

Coefficients	Singularities
$b_0 = -0.182$	$b_6 = 0.594$
$b_1 = 0.408$	$b_7 = 0.067$
$b_2 = 0.178$	$b_8 = -0.209$
$b_3 = 0.107$	$b_9 = -0.377$
$b_4 = 0.064$	$b_{10} = -0.467$
$b_5 = 0.027$	

The series coefficients and the nonlinear parameters are also recorded in the tables 4.5.1(b) and 4.5.2(b). The linear coefficients exhibit a steady decrease in magnitude with increasing index. The nonlinear parameters are the singularities in the approximate solution and are nicely distributed along the  $y$ -axis outside of the domain  $D$  (see figures 4.5.1(a),(b)). It was found in practice, without exception, that if the initial values of the singularities in the iterative scheme were placed outside of the domain of interest, they remained outside.

### Time Dependent Results

Both the cusping and Saffman calculations were performed using an error tolerance of  $\tau = 0.01$ . In order to minimize the time spent in the nonlinear solver and to forestall possible divergencies, the maximum time stepsize was reduced to  $h_{\max} = 0.005$ . In this way, the coefficients in the series approximation at one time level could be accurately used as first iterates in the nonlinear solver at the next time level. The nonlinear computations proved more robust than the linear in that the smaller time step did not

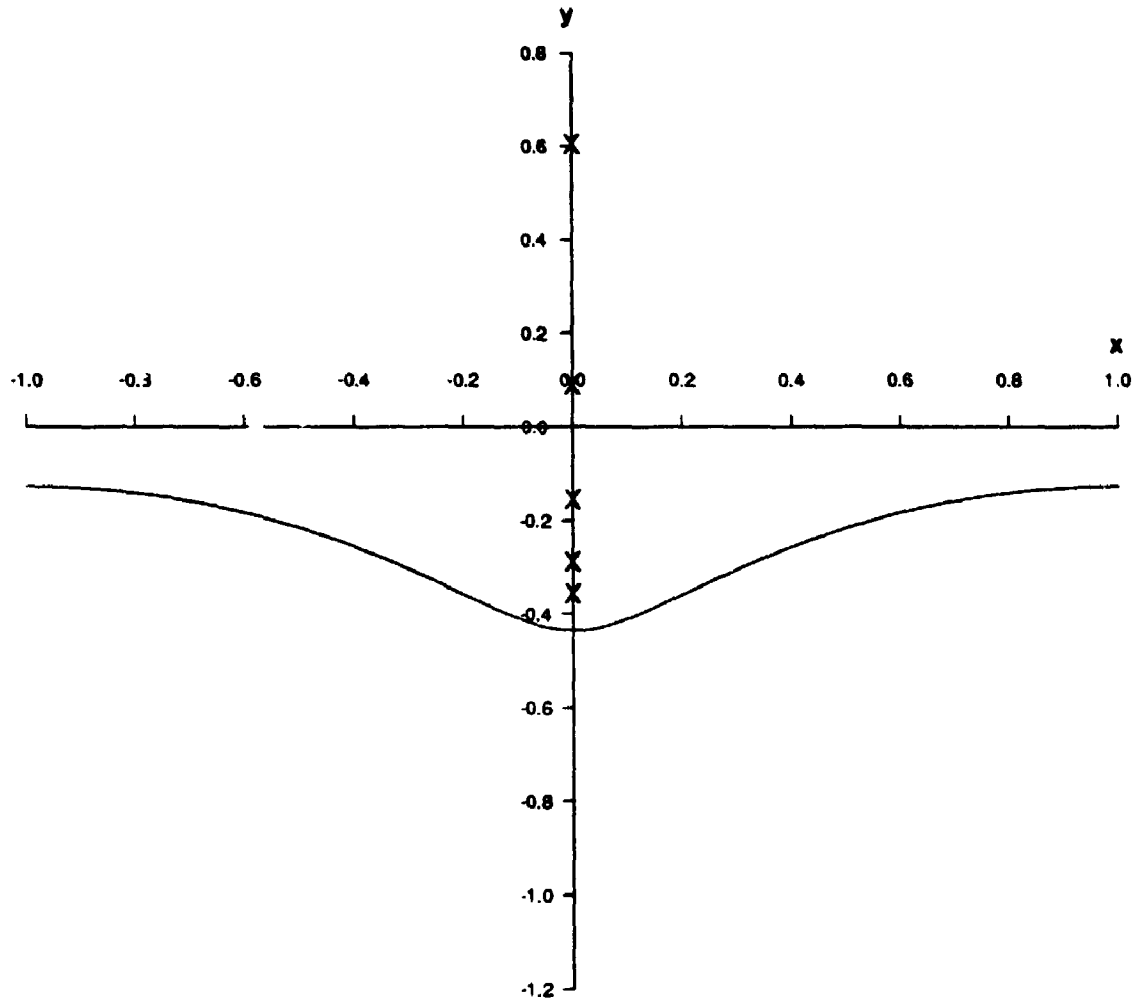


Fig. 4.5.1 (a) Location of the Singularities in the Nonlinear Approximation of the Cusping Problem -  $t = 0.25$

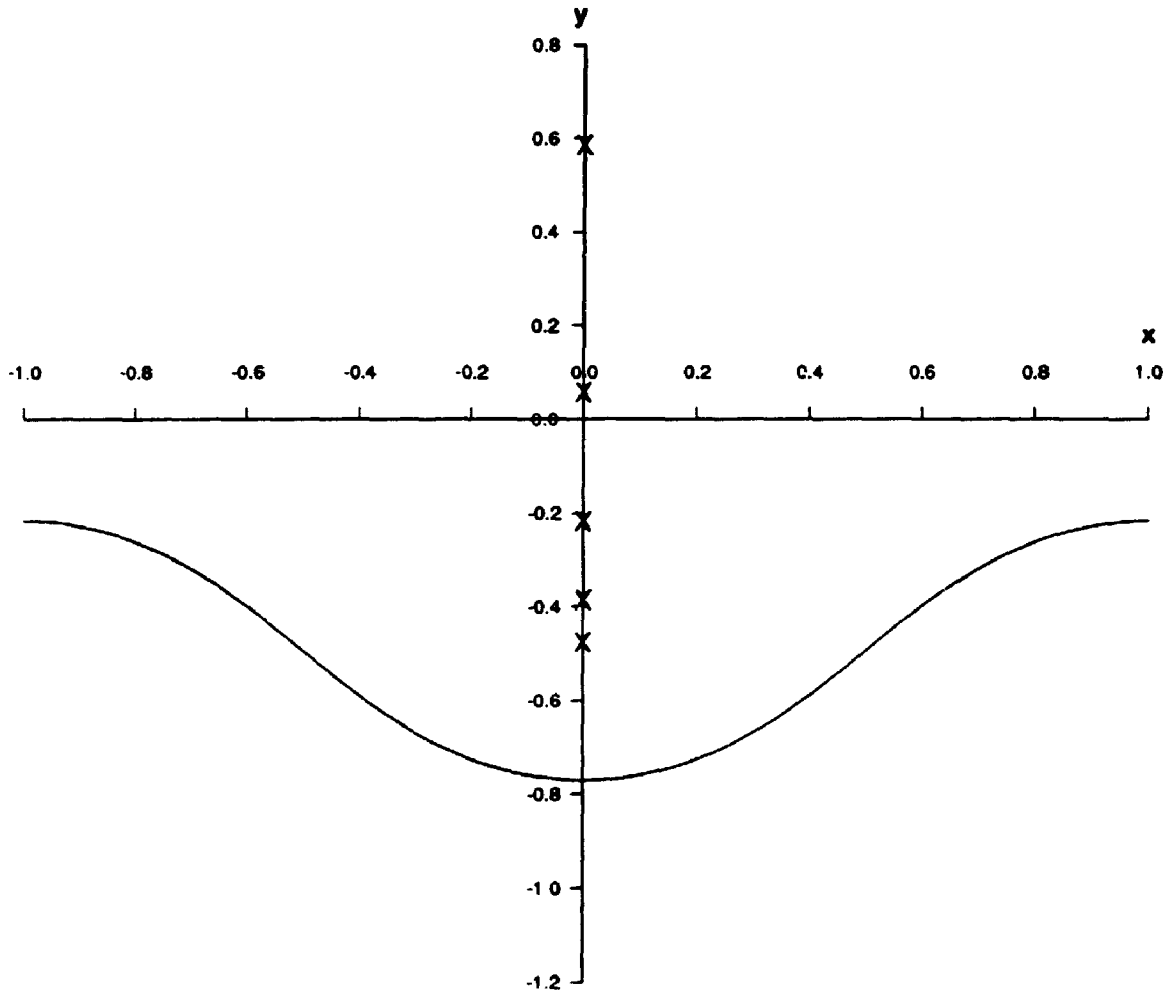


Fig. 4.5.1 (b) Location of the Singularities in the Nonlinear Approximation  
of the Saffman Finger -  $t = 0.50$



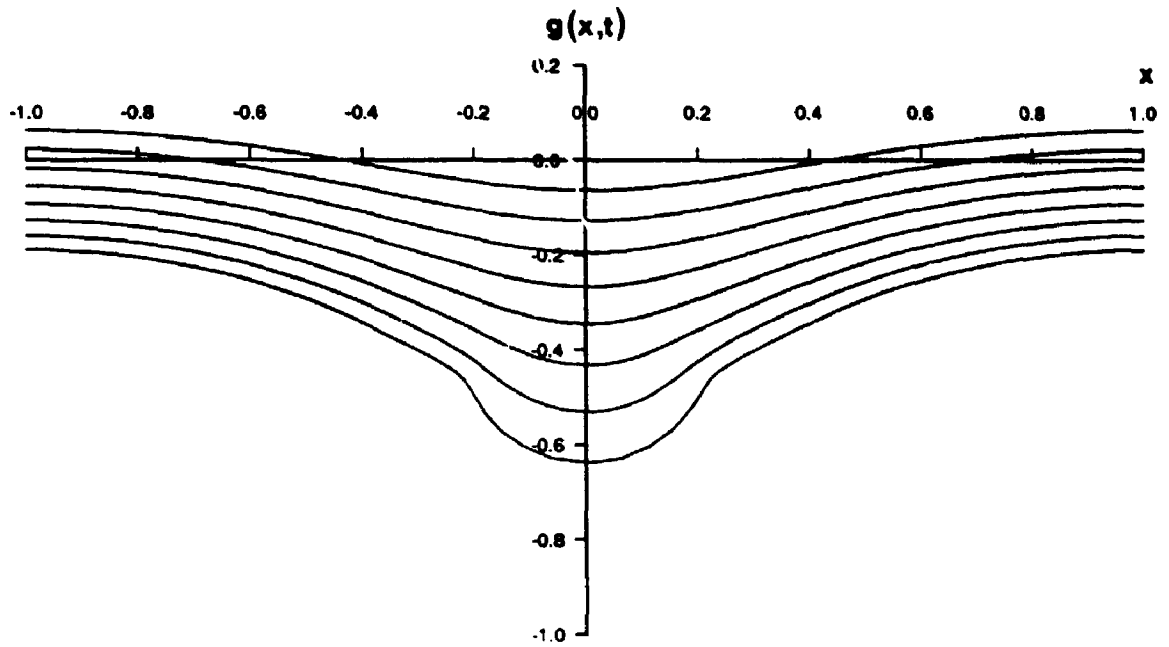


Fig. 4.5.2 (a) Cusping Profiles Using a Nonlinear Lagrangian Approximation

( $n=3, m=29$ )

$t = 0.0, 0.05, 0.10, 0.15, 0.20, 0.25, 0.30, 0.35$

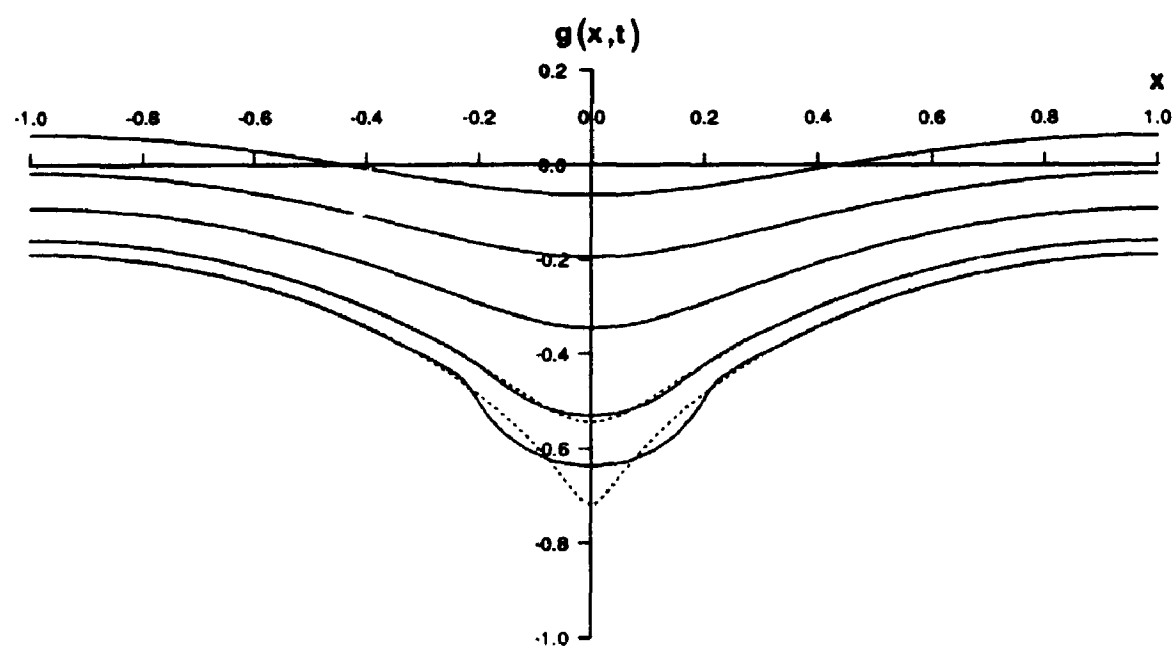


Fig. 4.5.2 (b) Comparison of Exact and Computed Profiles  
( $n=3, m=29$ )  
 $t=0.0, 0.10, 0.20, 0.30, 0.35$

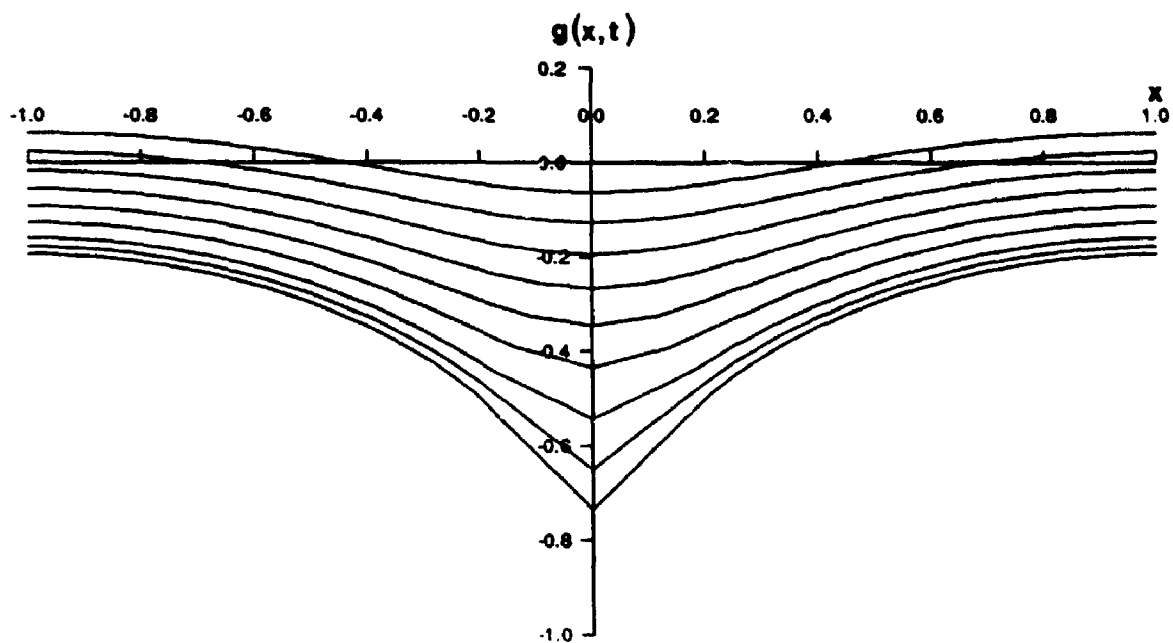


Fig. 4.5.3 (a) Cusping Profiles Using a Nonlinear Lagrangian Approximation

( $n=5$ ,  $m=15$ )

$t = 0.0, 0.05, 0.10, 0.15, 0.20, 0.25, 0.30, 0.33, 0.359$

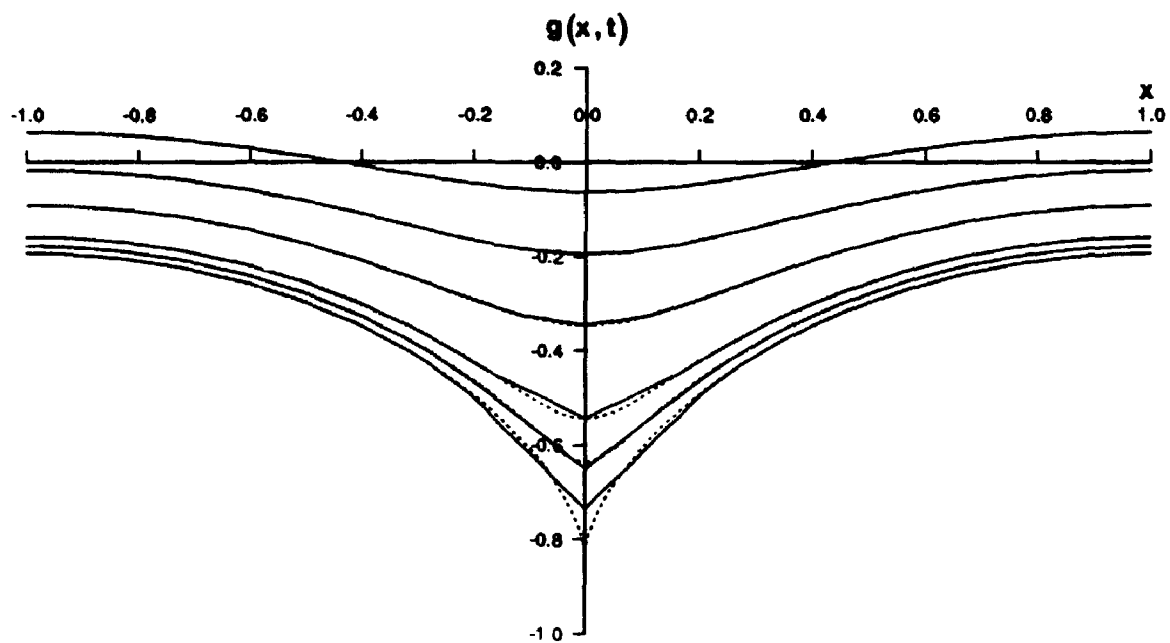


Fig. 4.5.3 (b) Comparison of Exact and Computed Profiles

( $n=5, m=15$ )

$t=0.0, 0.10, 0.20, 0.30, 0.33, 0.359$

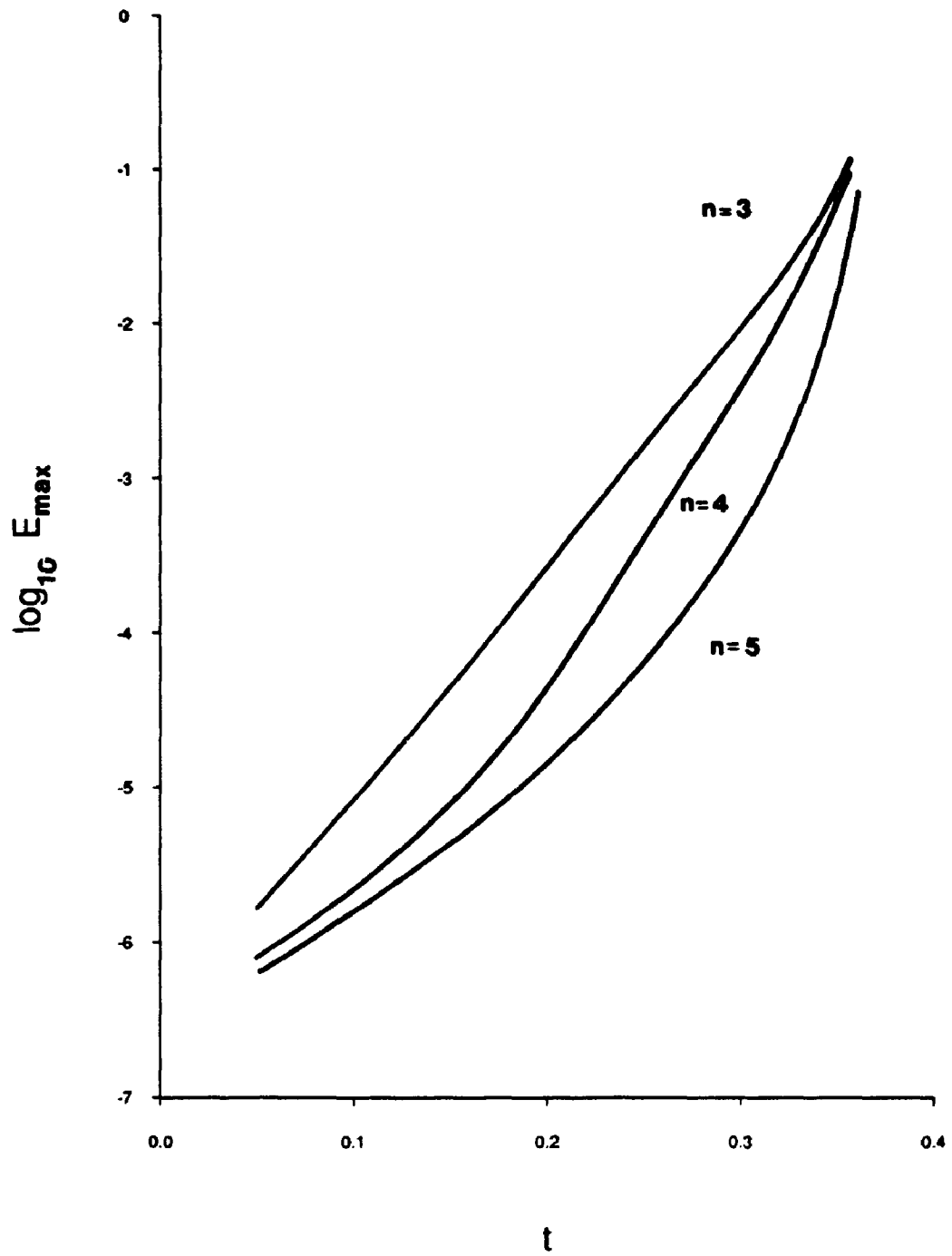


Fig. 4.5.4 Error Growth ( $n=3,4,5$ )  
Cusping Case - Nonlinear Approximation

enhance the onset of numerical instabilities. In fact, with this value of  $h_{\max}$ , no corrugations appeared in the computed cusping profile and developed in the Saffman profile only after very lengthy runs.

The results of a Lagrangian description of the cusping profile for the case  $n = 3$  are presented in figures 4.5.2(a) and (b). Twenty-nine boundary points were used with more points placed near  $x = 0$ . The calculation proceeded beyond the cusp time using the maximum time step. The run was terminated eventually ( $t = 0.365$ ) by an unacceptable growth in the coefficients of the approximating function. Possibly a smaller time step is needed to prevent these divergencies. There is some improvement over the linear approximation method, but there is still a significant deviation from the exact solution.

The value of  $n$  was increased to 4 and 5 with the results of  $n = 5$  shown in figures 4.5.3(a) and (b). The lower bound of  $\Delta t = 0.001$  imposed on the time step was reached at about  $t = 0.34$ , but the calculation was permitted to continue at that fixed time step. It was not found possible to bunch the boundary points near  $x = 0$  for this larger value of  $n$ . Indeed, if points were clustered about  $x = 0$ , the calculation would come to an early end. That is the calculation of the coefficients would proceed normally for a short time and then rapidly they would begin to grow in magnitude each time step until the nonlinear routine could proceed no further. This must in some way be attributed to attempts to accurately follow the sharp inflection point near to  $x = 0$ . As it is, the calculation presented in figure 4.5.3 is devoid of boundary points near  $x = 0$  ( $m = 15$  equally spaced points at  $t = 0$ ) and as such the finer structure of the true profile is not imposed on the calculation. In short, the numerical solution corresponds to a slightly different problem than the analytic solution of Aitchison and Howison. Still, it is remarkable that the calculation exhibits the correct qualitative behaviour and indeed advances the point  $(0, g(0, t))$  most accurately

(see table 4.5.3). A comparison of the error growth with time is shown in figure 4.5.4 for the three cases  $n = 3, 4, 5$ . The computing cost for the longest simulation ( $n = 5$ ) was about seventy seconds of CPU time.

The results of incorporating an Eulerian description into the nonlinear method were much less impressive and have not been detailed here.

**Table 4.5.3**

**Comparison of Exact and Computed Values of  $g(0,t)$**

**Cusping Profile - Nonlinear approximation - Lagrangian description**

t	n=3	n=4	n=5	Exact
0.1	-0.194847	-0.194856	-0.194857	-0.194854
0.2	-0.345052	-0.345362	-0.345429	-0.345403
0.25	-0.432371	-0.433940	-0.434609	-0.434609
0.3	-0.530673	-0.536811	-0.543274	-0.543775
0.33	-0.595028	-0.607521	-0.645615	-0.631304
0.35	-0.637931	-0.654662	-0.708224	-0.718882
0.355	-0.648919	-0.667067	-0.723087	-0.754038
0.359	-0.657728	-0.677061	-0.734415	-0.807569

A Lagrangian treatment of the Saffman finger for the cases  $n = 3, 5$  is presented in figures 4.5.5(a),(b) and 4.5.6(a),(b), respectively. A comparison of the values  $g(0)$  in the

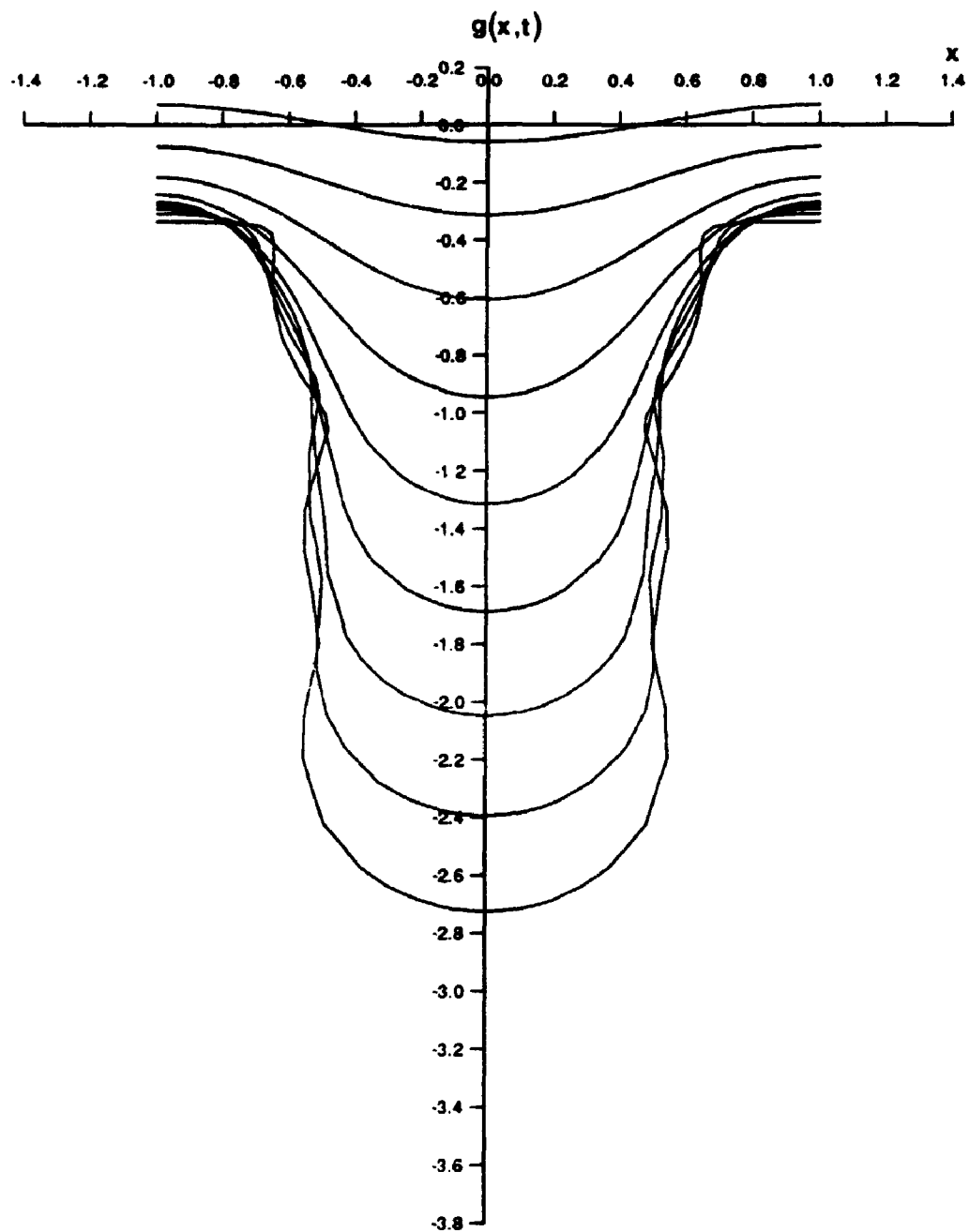


Fig. 4.5.5 (a) Saffman Profiles Using a Nonlinear Lagrangian Approximation

( $n=3$ ,  $m=38$ )

$t = 0.0, 0.20, 0.40, 0.60, 0.80, 1.00, 1.20, 1.40, 1.60$



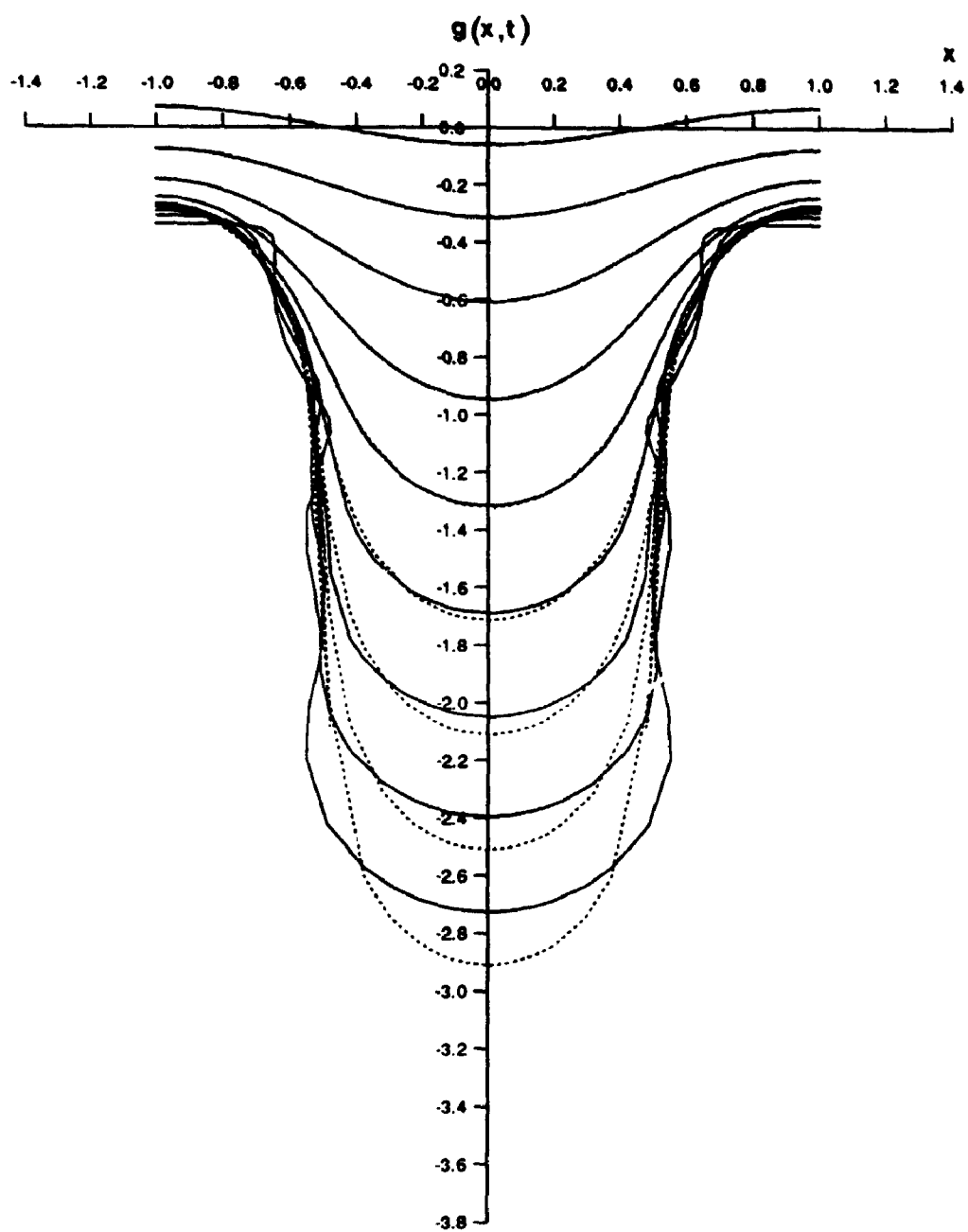


Fig. 4.5.5 (b) Comparison of Exact and Computed Profiles

( $n=3$ ,  $m=38$ )

$t = 0.0, 0.20, 0.40, 0.60, 0.80, 1.00, 1.20, 1.40, 1.60$

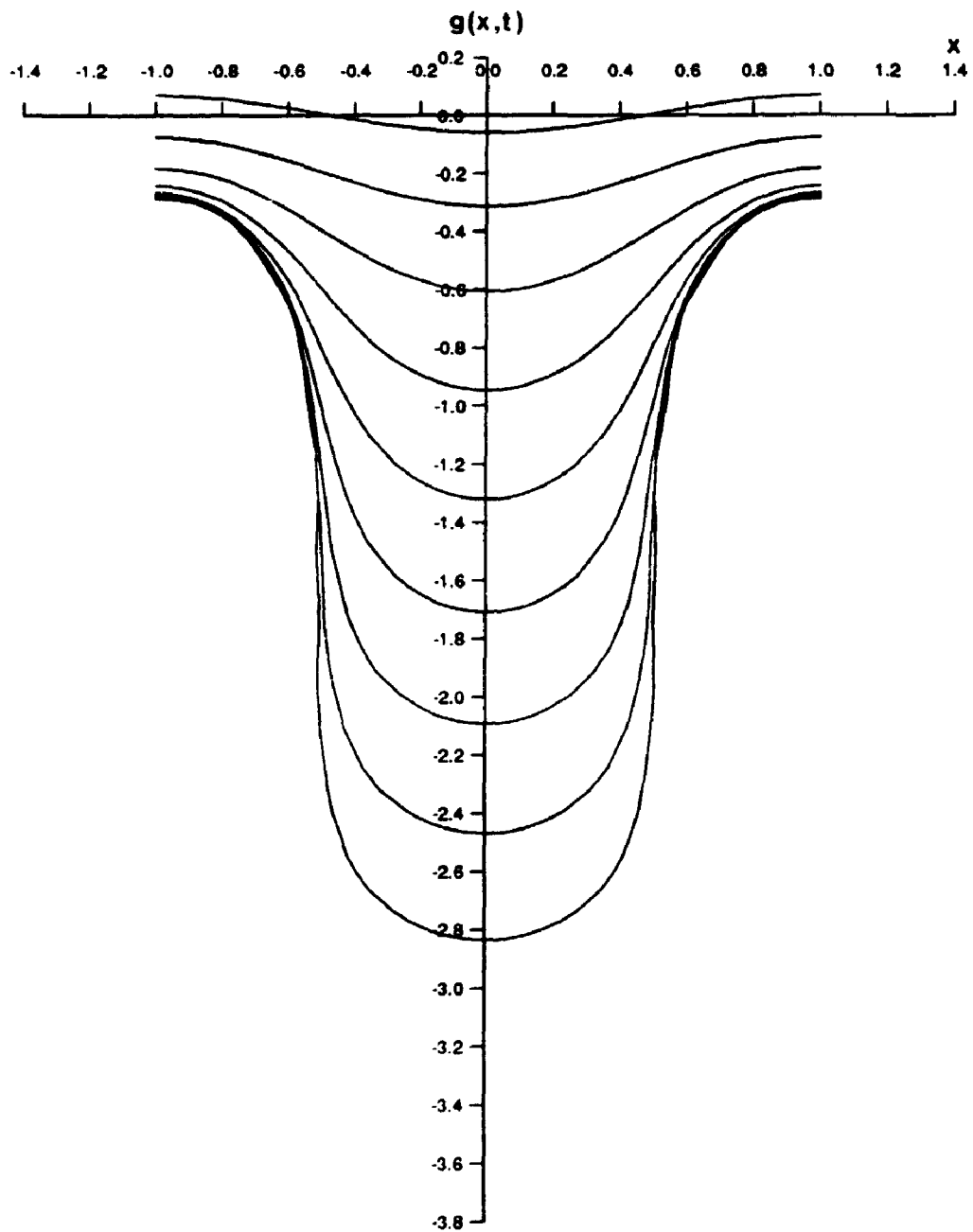


Fig. 4.5.6 (a) Saffman Profiles Using a Nonlinear Lagrangian Approximation

( $n=5$ ,  $m=38$ )

$t = 0.0, 0.20, 0.40, 0.60, 0.80, 1.00, 1.20, 1.40, 1.60$

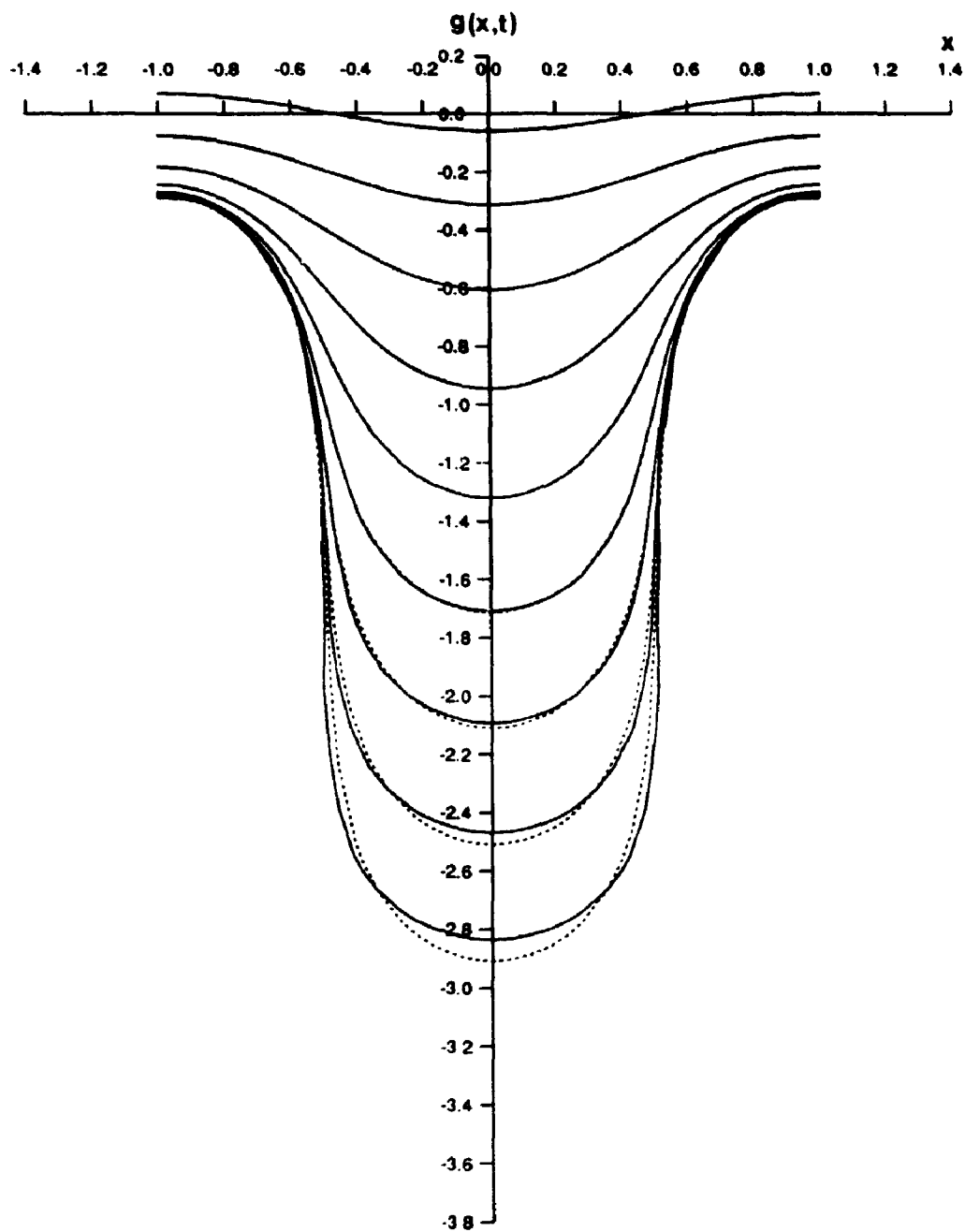


Fig. 4.5.6 (b) Comparison of Exact and Computed Profiles

( $n=5, m=38$ )

$t = 0.0, 0.20, 0.40, 0.60, 0.80, 1.00, 1.20, 1.40, 1.60$

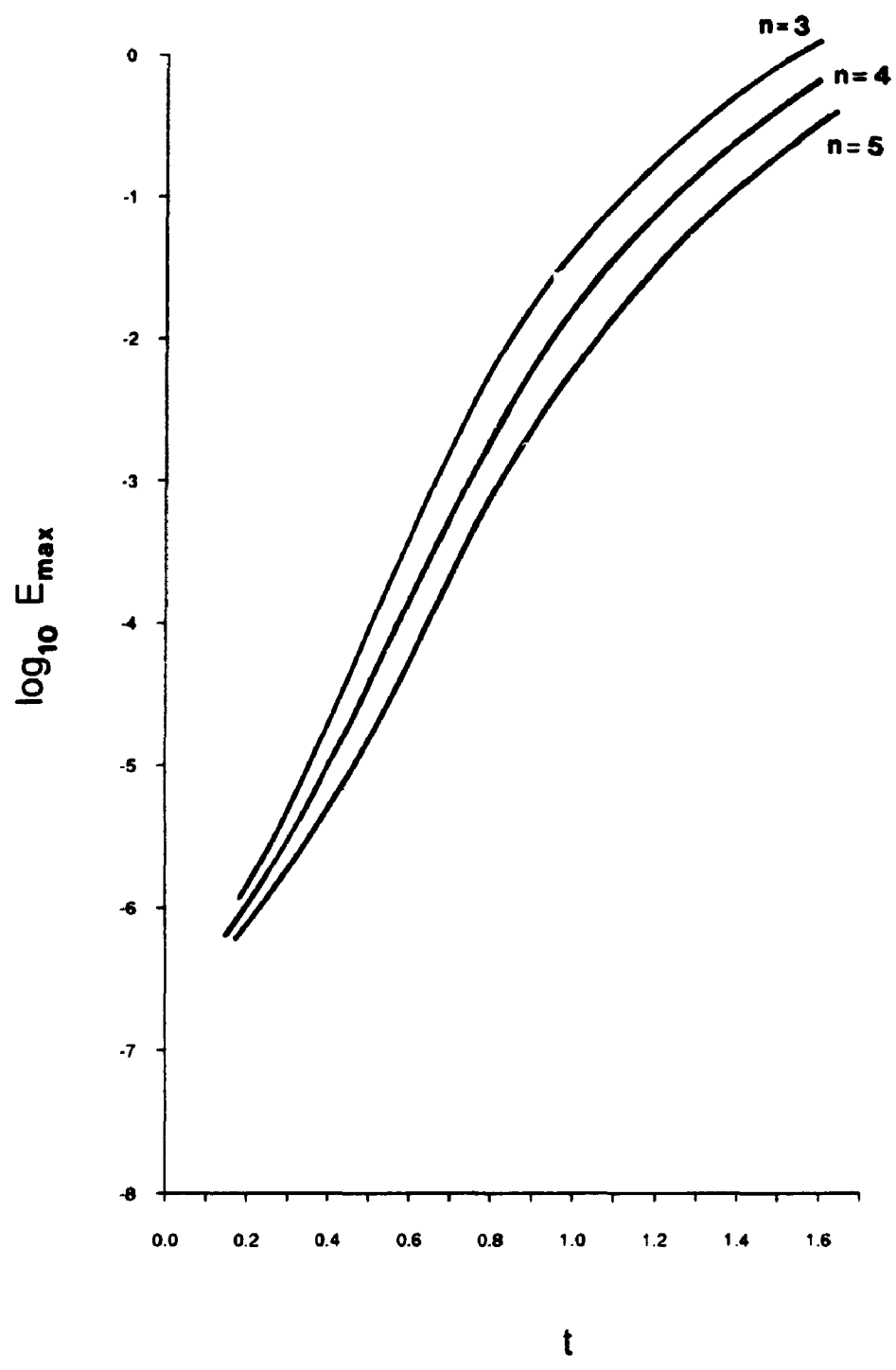


Fig. 4.5.7 Error Growth (n=3,4,5)  
Saffman Finger - Nonlinear Approximation

trough for the cases  $n = 3, 4, 5$  is given in table 4.5.4, and the error growth for these cases is depicted in figure 4.5.7. As with the cusping case, we witness a steady improvement in the accuracy as  $n$  is increased. For  $n = 3$  (figures 4.5.5) we have plotted the profiles at intervals of 0.2 up to a time of  $t = 1.6$ . Irregularities have appeared in vertical walls of the finger by a time of about  $t = 1.2$ ; but still the calculation has accurately followed the analytic profile almost twice as far as the linear calculation of the last section.

The results for  $n = 5$  are in excellent agreement with the exact solution and proceed to a time of over four times that reported by Aitchison and Howison (1985). Corrugations do not develop until a time of about  $t = 1.7$ . The calculation proceeded to the times shown in the figures without decreasing the time step below 0.005. The actual computing cost for this case was only about 200 seconds of CPU time.

Since the boundary points in the vicinity of  $x = 0$  move rapidly in the direction of increasing  $x$  as the moving surface stretches, it was necessary in the case of the Saffman finger to begin the calculation with many boundary points clustered near  $x = 0$ . In fact, for the Saffman profile and the results shown in table 4.5.4 and figure 4.5.6, the calculation was begun with 20 points bunched in the range  $0 \leq x \leq 0.01$ . For example, the boundary point initially located at  $x = 0.0001$ , eventually reached a position of  $x = 0.25$  by the time of  $t = 2.0$ . The migration of the boundary points for this case is depicted in figure 4.5.8. Some experimentation was done with the initial placement of boundary points to produce a respectable distribution at later times. Clearly, much of this could have been avoided by making use of spline interpolation and deftly placing points as the calculation proceeded. Degregoria and Schwartz (1986) use such an elaborate scheme to advance their Lagrangian description in a boundary integral approach to the Hele-Shaw problem. However, this provides artificial smoothing to the free surface and it was decided to avoid

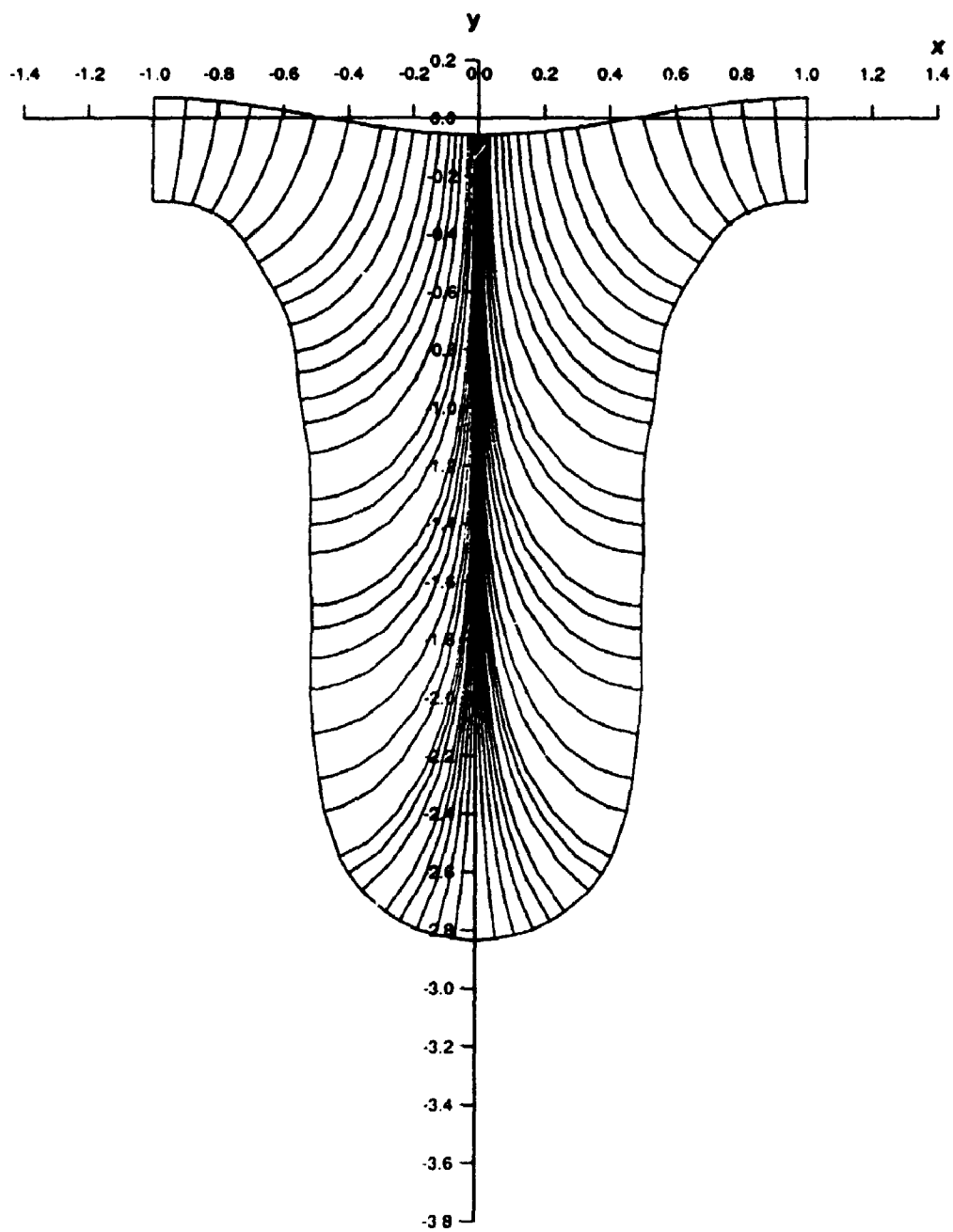


Fig. 4.5.8 The Particle Trajectories Between  $t = 0$  and  $t = 1.6$

this if at all possible and to present the calculation with its natural growth of errors. This rigorously tests the full capabilities and limitations of the nonlinear scheme in direct comparison with the analytic solution.

**Table 4.5.4**  
**Comparison of Exact and Computed Values of  $g(0,t)$**   
**Saffman Profile - Nonlinear approximation - Lagrangian description**

t	n=3	n=4	n=5	Exact
0.2	-0.312507	-0.312506	-0.312506	-0.312504
0.4	-0.605724	-0.605739	-0.605739	-0.605722
0.6	-0.944816	-0.945360	-0.945422	-0.945305
0.8	-1.31407	-1.31892	-1.32022	-1.32019
1.0	-1.68649	-1.70169	-1.70792	-1.71184
1.2	-2.04747	-2.07810	-2.09272	-2.10935
1.4	-2.39449	-2.44287	-2.46909	-2.50862
1.6	-2.72533	-2.79558	-2.83578	-2.90842

Finally, with regard to the local errors incurred with the nonlinear Lagrangian calculation of the Saffman finger, refer to figure 4.5.9. The error between the computed and exact solutions is plotted on the domain  $-1 \leq x \leq 1$  and every 0.2 of a time step up to  $t = 2.0$ . The largest errors are found near  $x = 0.5$  instead of in the shoulder region or the

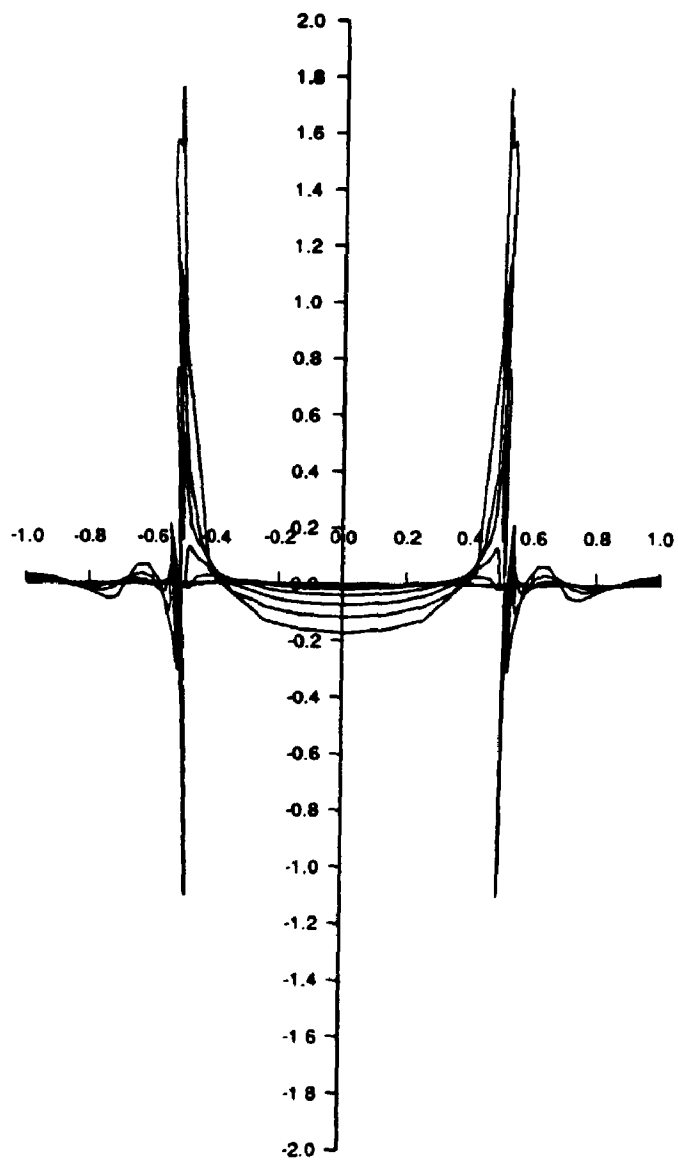


Fig. 4.5.9 Local Error  $t = 0.00$  (0.20) 2.00



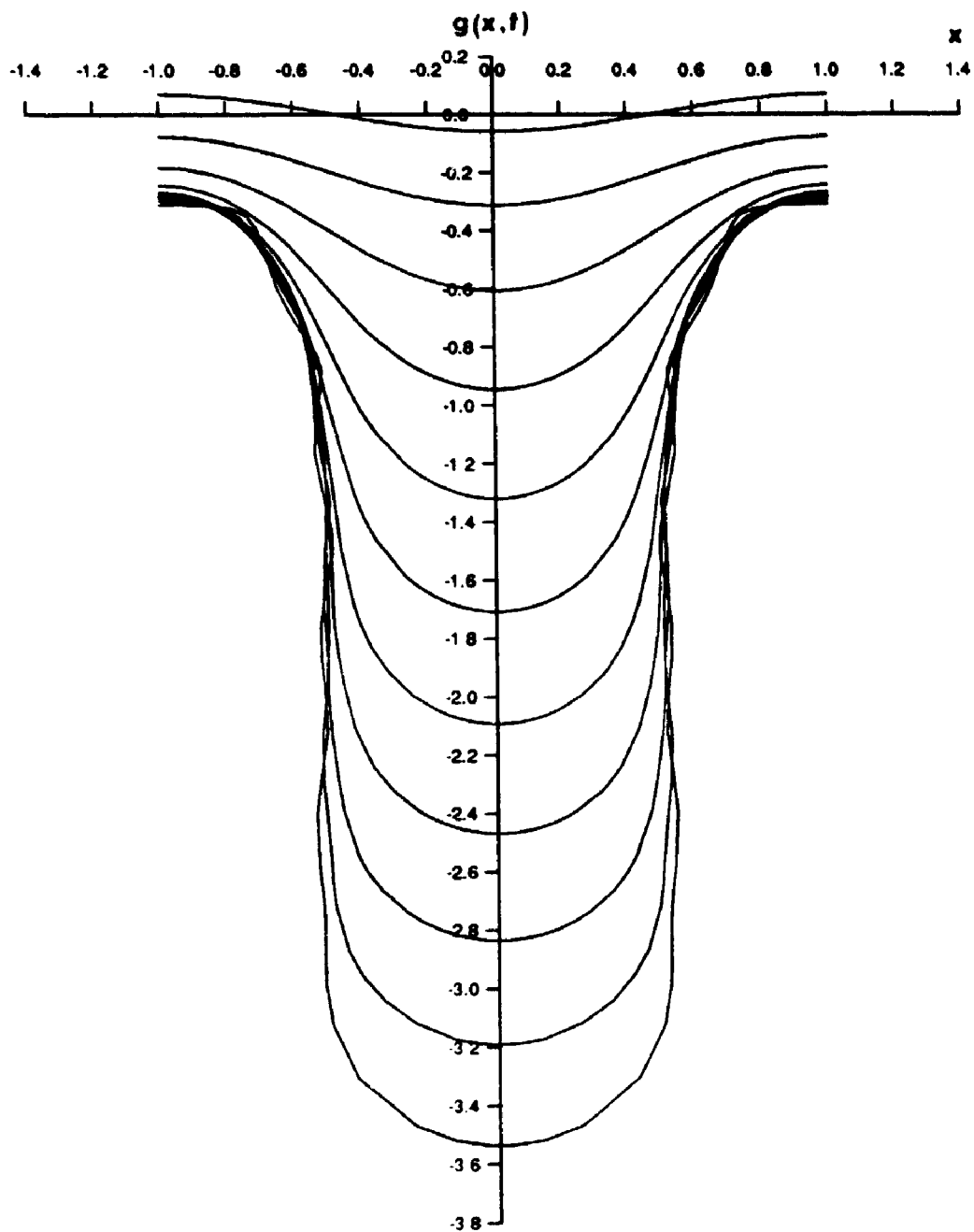


Fig. 4.5.10 (a) Saffman Profiles Using a Nonlinear Lagrangian Approximation

( $n=5$ ,  $m=38$ )

$t = 0.0, 0.20, 0.40, 0.60, 0.80, 1.00, 1.20, 1.40, 1.60, 1.80, 2.00$

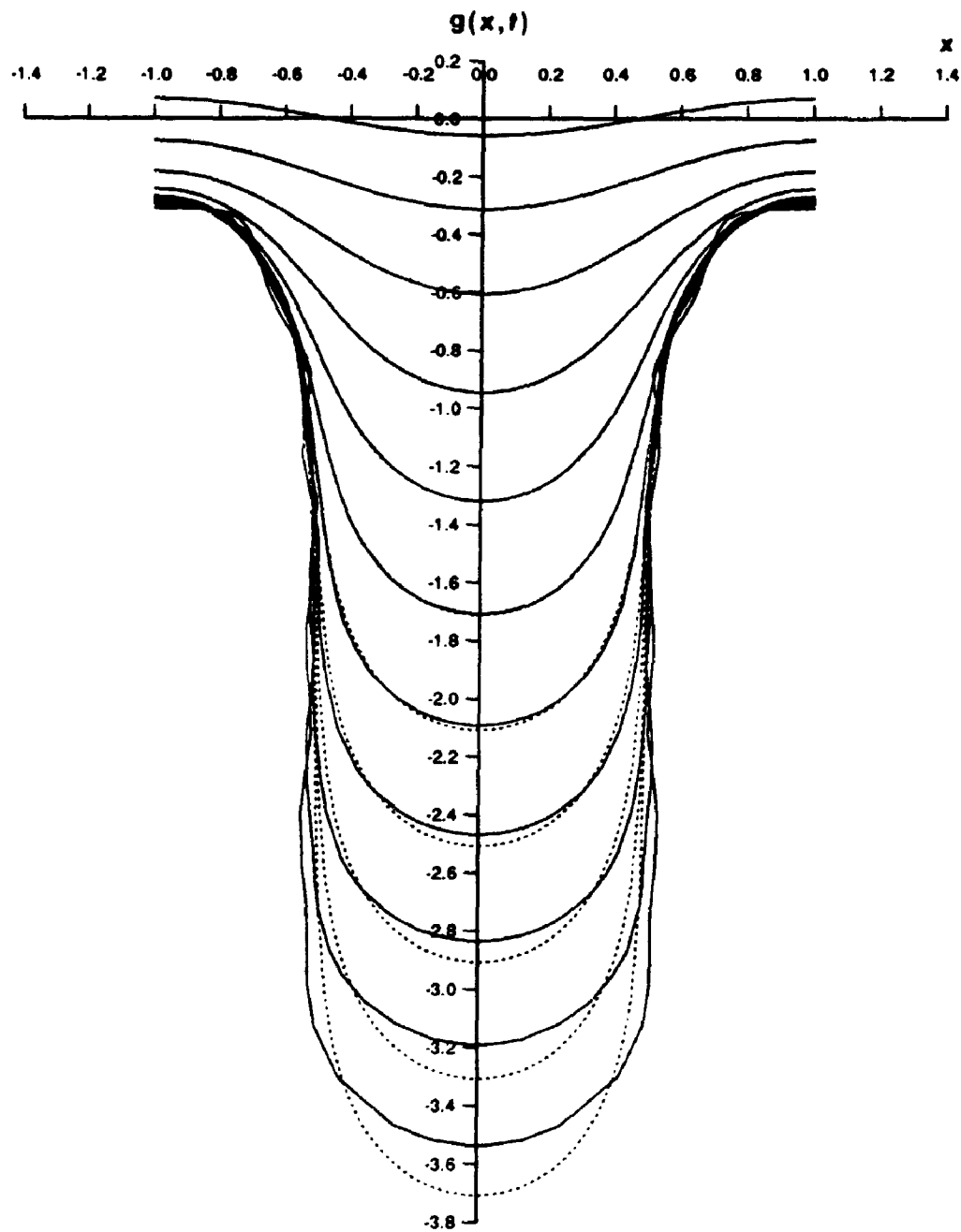


Fig. 4.5.10 (b) Comparison of Exact and Computed Profiles

( $n=5, m=38$ )

$t = 0.0, 0.20, 0.40, 0.60, 0.80, 1.00, 1.20, 1.40, 1.60, 1.80, 2.00$

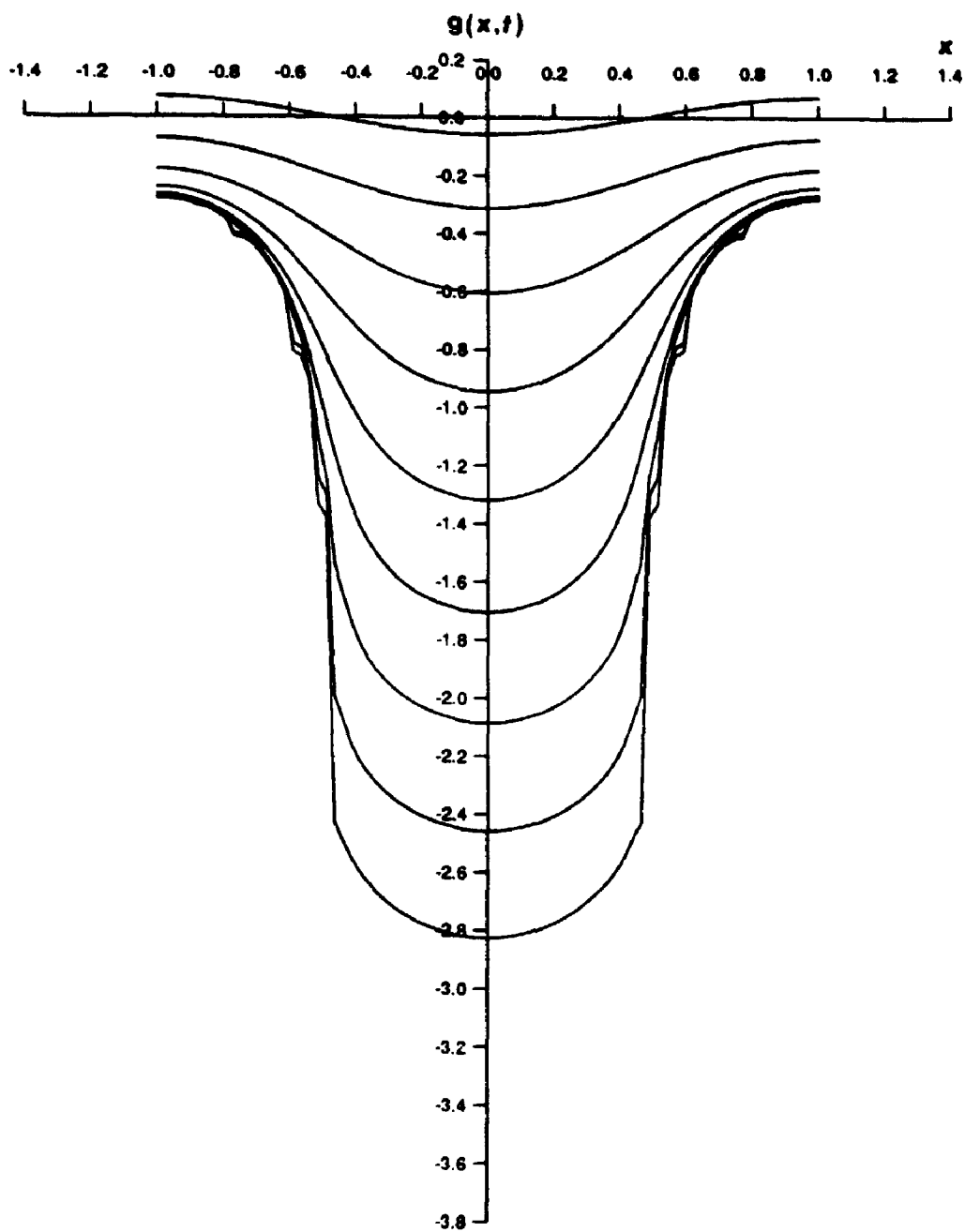


Fig. 4.5.11 Saffman Profiles Using a Nonlinear Eulerian Approximation

( $n=5, m=40$ )

$t = 0.0, 0.20, 0.40, 0.60, 0.80, 1.00, 1.20, 1.40, 1.60$

trough, as was the case with the cusping profile. In figure 4.5.10, the  $n = 5$  calculation is carried to a time of  $t = 2.0$ . The corrugations are clearly evident in the vertical sides of the finger and the trough is beginning to widen more than the exact solution. The widening is probably due to artificial surface tension introduced at later stages of the calculation.

The Eulerian description also worked effectively on this problem, although a breakup in the free surface is evident at earlier times. For example, with  $n = 5$  again and  $m = 40$  equally spaced  $x$ -values, some irregularities show up in the free surface profile by approximately  $t = 1.4$ . This is depicted in figure 4.5.11. The difficulty would appear to lie with the calculation of  $g_x$ . Recall, in the Eulerian description, this derivative was required in the evolution equation and was computed using a ratio of  $\phi_x$  and  $\phi_y$ . The situation was not improved by replacing the calculation of  $g_x$  with a second or even fourth order central difference molecule.

## Conclusions

The nonlinear approximation is vastly superior to the linear one when applied to both the cusping and Saffman problems. In both problems, the nonlinear calculation proceeds farther in time than either of the calculations of Aitchison and Howison (1985). In particular, the Saffman profile has been followed longer than any other calculation that we know of. Although the computing times are greater, none of the calculations are so long as to warrant concern. The growth of error with time for the two methods is compared in figures 4.5.12 and 4.5.13.

In each of the examples presented here, a good first estimate of the nonlinear parameters was found before beginning the time dependent run. The preparation time required to find this first estimate can be comparable to the computing effort of the entire time dependent run. The following procedure was adopted for determining a good estimate to be used at  $t=0.0$ . A guess is made and a crude tolerance is specified in the Levenberg-Marquardt routine. The tolerance is then sharpened, and the minimum is looked for once again, using the values just computed as new estimates of the parameters.

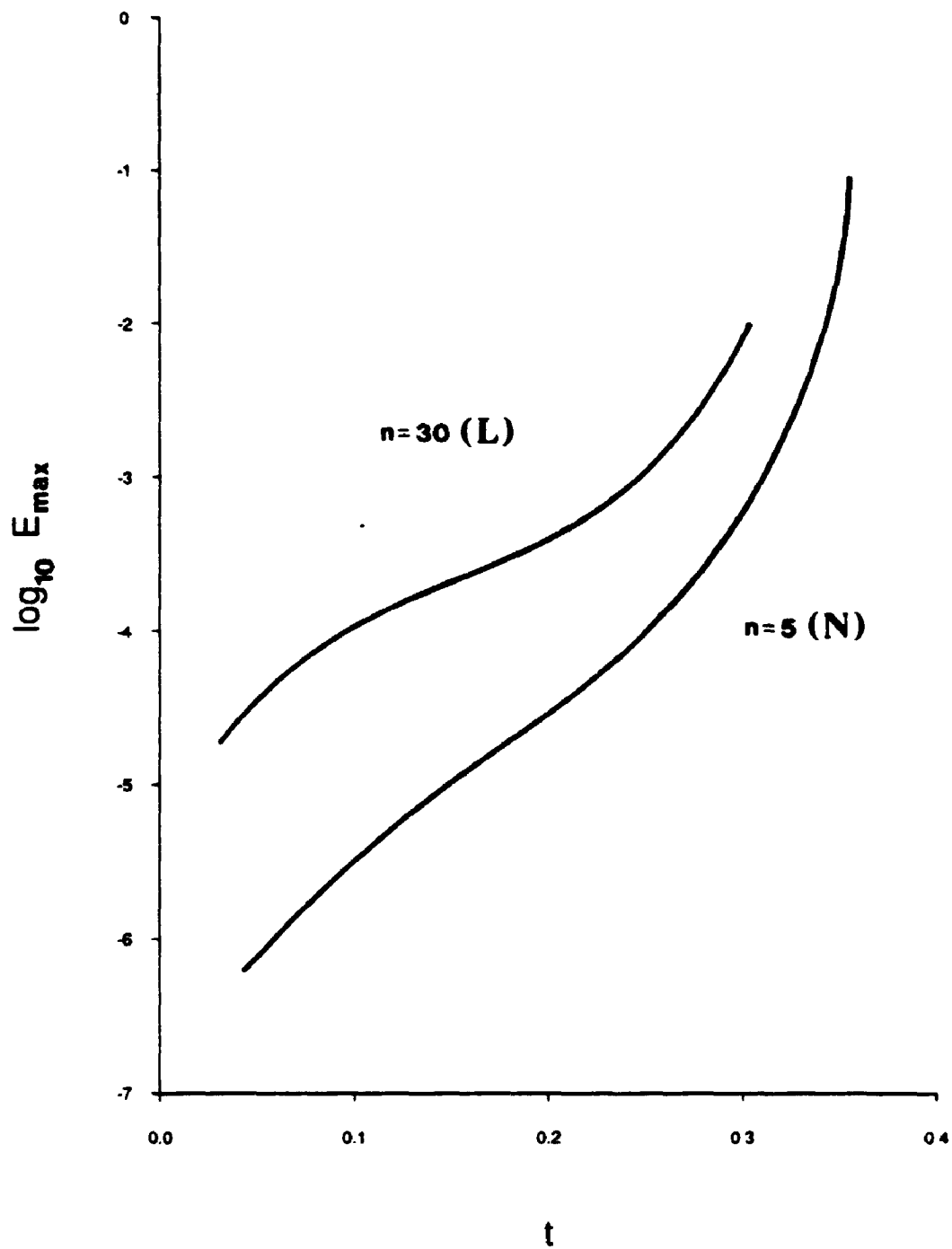


Fig. 4.5.12 Error Growth - Cusping Case  
Linear ( $n=30$ ) vs. Nonlinear ( $n=5$ )

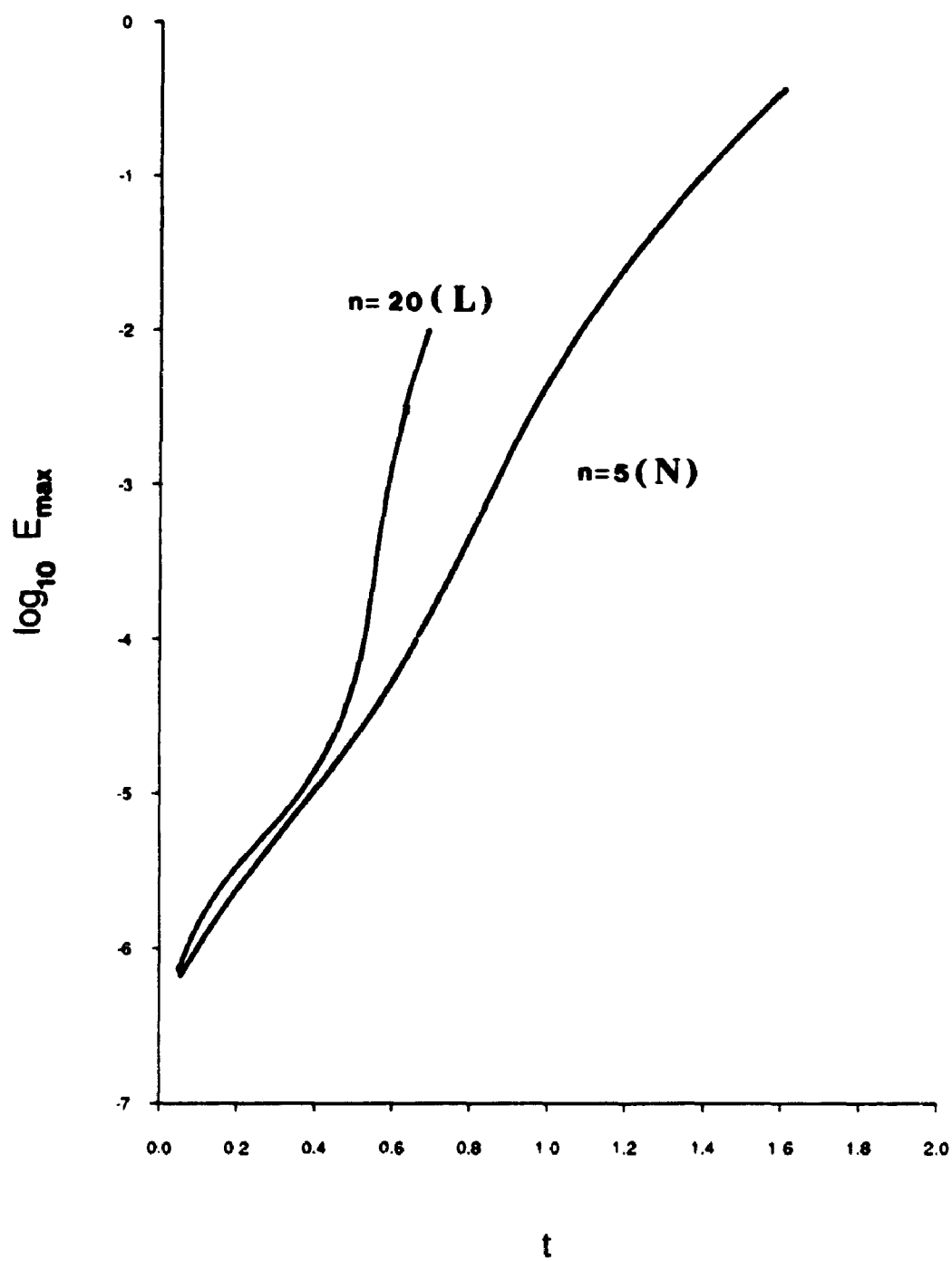


Fig. 4.5.13 Error Growth - Saffman Finger  
Linear ( $n=20$ ) vs Nonlinear ( $n=5$ )

The process was repeated until a near best approximation was found, typically accurate to within  $10^{-10}$  or better. The situation is not as bad as it may appear, since that very first estimate can be made with some measure of confidence. After all, we do know that the nonlinear parameters appear along a straight line only and lie above  $g(0)$ .

Once a good initial estimate has been found for the parameters, the method proceeds quite efficiently, provided the time step is small.

In general, the linear method has proven to work best with an Eulerian description while a Lagrangian approach was preferable for the nonlinear method.

#### 4.6 Stability

The stability of the time dependent scheme is strongly dependent on the system of ordinary differential equations being solved. Consider the Lagrangian description

$$\frac{dx}{dt} = -\phi_x$$

$$\frac{dy}{dt} = -\phi_y$$

or in vector form

$$\frac{d\mathbf{x}}{dt} = \mathbf{F}(\mathbf{x}, t).$$

We can linearize the system, expanding  $\mathbf{F}$  in a Taylor series about the solution at  $t = t^n$ ,

$$\frac{d\mathbf{x}}{dt} = \mathbf{J}\mathbf{x} + \mathbf{c}$$

where  $\mathbf{J}$  is the Jacobian of  $\mathbf{F}(\mathbf{x}, t)$  at  $t = t^n$ , and  $\mathbf{c}$  is a constant vector. If  $\mathbf{J}$  is nonsingular,

then letting

$$\mathbf{u} = \mathbf{x} + \mathbf{J}^{-1}\mathbf{c}$$

gives

$$\frac{d\mathbf{u}}{dt} = J\mathbf{u} \quad (4.6.1)$$

where

$$J = -\begin{pmatrix} \phi_{xx} & \phi_{xy} \\ \phi_{xy} & \phi_{yy} \end{pmatrix}.$$

The solution of (4.6.1) is

$$\mathbf{u} = \begin{pmatrix} a_{11} e^{\lambda_1 t} + a_{12} e^{\lambda_2 t} \\ a_{21} e^{\lambda_1 t} + a_{22} e^{\lambda_2 t} \end{pmatrix}$$

where  $\lambda_1, \lambda_2$  are the eigenvalues of the Jacobian matrix and

$$\lambda_{1,2} = \pm \sqrt{\phi_{xx}^2 + \phi_{xy}^2}.$$

The solution of the linearized system (4.6.1) is dominated by an exponentially growing component. This is as expected, since both the cusping and Saffman profiles exhibit this behaviour. More importantly, from the computational viewpoint, the system (4.6.1) may be stiff. Loosely speaking, a system of ordinary differential equations is stiff or unstable if the magnitude of the largest eigenvalue is large compared to the time scales involved ( $1/T$  if the interval is  $[0, T]$  - see Miranker (1981)). A stiff system of equations can be very difficult to solve numerically. In the present case, we might expect stability problems in the numerical scheme if  $|\lambda|$  should be large or if very long run times are required.

The magnitude of the eigenvalues can be computed from the analytic solution in each of the cusping and Saffman cases. We have

$$\left| \frac{d^2 w}{dz^2} \right|^2 = \phi_{xx}^2 + \phi_{xy}^2.$$

For the cusping case,

$$\lambda^2 = \frac{h_1^2}{(1 - 2h_1 \cos \psi + h_1^2)^3}.$$



Now,  $b_1$  approaches 1 as  $t$  approaches cusping time, so that for  $\psi$  near zero the eigenvalues are quite large. For example, consider the point with  $x$  coordinate  $x=0$  (ie  $\psi = 0$ ). Then,  $\phi_{xy} = 0$  and the solution of the linearized equations is

$$x = c_1 e^{-\phi_{xx} t}$$

$$y = c_2 e^{\phi_{xx} t}$$

where  $\phi_{xx} > 0$  for  $0 < b_1 < 1$ . If  $b_1 = 0.6$ , corresponding to  $t = 0.3$ , then

$$\lambda^2 = \phi_{xx}^2 = 88.$$

On the other hand, for the Saffman case, we have

$$\lambda^2 = a^2(1 + 2a \cos \psi + a^2)$$

and as  $0 < a(t) \leq 1$  for all  $t \geq 0$ , the eigenvalues are never large.

The above linearized analysis suggests that the system of ordinary differential equations associated with the cusping case becomes stiff as times near the cusping time, while the Saffman case is not stiff unless exceptionally long computing times are required. This accounts for the necessary decrease in time step size observed in the case of the cusping profile.

## CHAPTER 5

### Concluding Remarks

Two closely related numerical methods designed for use on a class of MBP have been examined in detail.

The linear method is, in principle, applicable to a greater variety of problems from our class. In practice, it is limited to problems with relatively simple geometries. The nonlinear method, as it stands, is restricted to two dimensional simply connected domains. The obvious expansion, by way of more general rational functions, would extend the method to more complex domains, but the formulation would not be as straight forward. Nevertheless, for the problems to which it pertains, the nonlinear scheme has proven to be an efficient and accurate computing tool. It has performed admirably on some badly posed problems of Hele-Shaw flow.

We have been fortunate to have an analytic solution for comparison. We have found that an increase in  $n$  in the nonlinear method consistently led to an improved approximation in the time dependent calculations, something that was not strictly observed in the linear case. Furthermore, the computed free surface profiles remained quite smooth for long simulations of an unstable flow. These findings should permit the method to be used with some measure of confidence on problems for which an exact solution is unavailable. In this regard, some future directions might include applications of the nonlinear method to the Rayleigh-Taylor instability and to Hele-Shaw problems with surface tension effects included.

Another nonlinear trial function was tested on the inverse ECM problem, with promising results. If a more robust means of computing the nonlinear parameters were available, a larger number of unknowns might be more manageable than with the present technique. For example, the rational function in  $\cos(x)$  might be best expressed as a continued fraction in the variable  $\cos(x)$ . There are very efficient algorithms for dealing with this type of rational approximation. In any event, with the present interest in inverse problems, there is ample scope for future work.

## APPENDIX A.1

### The Logarithmic Term

Many of the results on the approximation of harmonic functions can be readily inferred from already established results on the approximation of analytic functions by complex polynomials. This area has been well explored by a number of mathematicians, notably Runge, Walsh, Keldysch, Mergelyan and Curtiss.

Suppose the function  $f(\zeta)$  to be approximated is analytic on a closed region. Then, the corresponding results on completeness and degree of convergence for the case of harmonic functions can be deduced by taking the real parts of  $f(\zeta)$  and its approximating functions. However, the results for analytic functions are restricted to single-valued functions, even on multiply connected regions.

This represents a limitation if we are to apply the results on analytic functions to general harmonic functions (associated with a boundary value problem) on multiply connected domains. The reason is the following. Although the harmonic function  $\Phi(\xi, \eta)$  is the desired solution to our boundary value problem and, as such, is assumed to be single-valued, there is no reason to believe that the complex function  $f(\zeta)$ , constructed from  $\Phi$  and its harmonic conjugate, is also single-valued.

As an example, we consider the case of a doubly connected domain  $D_\zeta$ , bounded by the two Jordan curves  $C_0$  and  $C_1$  with  $C_1$  interior to  $C_0$ . Assume the origin lies interior to  $C_1$ . Let  $\Phi(\xi, \eta)$  be harmonic on the closed region  $\bar{D}_\zeta$ . Unlike the case for simply connected domains, the integral

$$\int_{(\xi, \eta)}^{(\xi, \eta)} \frac{\partial \Phi}{\partial n} ds \quad (\text{A.1.1})$$

$(n, s$  being the outward normal and tangent) is not independent of path. However, if  $C$  is any simple closed contour surrounding the "hole", then

$$\int_C \frac{\partial \Phi}{\partial n} ds$$

is independent of the contour  $C$ . Let

$$p = \int_C \frac{\partial \Phi}{\partial n} ds.$$

Any two functions  $\Psi, \Psi_1$  of the form (A.1.1) differ at most by an integer multiple of  $p$ . Thus, the harmonic conjugate to  $\Phi(\xi, \eta)$  may be expressed as the multi-valued function

$$\Psi(\xi, \eta) = \Psi_1(\xi, \eta) + mp \quad , \quad m = 0, \pm 1, \pm 2, \dots$$

The multiplicity derives from the number of circuits about the hole. This is precisely the multi-valued nature of the function

$$A \log \zeta \quad \text{if} \quad A = \frac{p}{2\pi}.$$

That is,

$$A \log \zeta = \frac{p}{2\pi} \text{Log} \zeta + imp$$

where  $\text{Log} \zeta$  is the principle branch of  $\log \zeta$ .

Thus, it follows that we can write

$$f(\zeta) = g(\zeta) + A \log \zeta$$

where  $g(\zeta)$  is single-valued. What is more,

$$\begin{aligned} \Phi(\xi, \eta) &= \Re f(\zeta) \\ &= \Re g(\zeta) + A \log |\zeta|. \end{aligned} \quad (\text{A.1.2})$$

A final note is the following. It appears that it is a well known fact that the logarithmic terms must be present in the complete set of harmonic functions for multiply connected domains (see Davis and Rabinowitz (1961), for example). However, its origins and the proof of this fact are more difficult to come by, as is the expression (A.1.2). Axler (1986) mentions this very point and references Walsh (1929) as the only location where a proof was found. Axler refers to the result (that a harmonic function on a multiply connected domain can be expressed as a sum of logarithmic terms and the real part of an analytic function) as the Logarithmic Conjugation Theorem. He provides a proof of this theorem which makes use of the Cauchy Integral Theorem and without recourse to the notion of a multi-valued complex function.

## Appendix A.2

### The Least Squares Algorithm

In what follows, the matrix is  $m \times n$  and assumed to have rank  $n$ . The method proceeds by repeated application of Householder transformations to  $\mathbf{r} = \mathbf{A}\mathbf{b} - \mathbf{f}$ . A Householder transformation applied to a vector  $\mathbf{v} \in R^m$  will produce another vector  $\mathbf{w} \in R^m$  where  $\mathbf{v}$  has been "rotated in  $m$ -space" to coincide with the direction of  $\mathbf{w}$ , the vector length being preserved. The Householder transformation can be represented by an  $m \times m$  matrix  $U$ . In this way the columns of an arbitrary  $m \times n$  matrix  $A$  can be reduced to a simple form.

For example, let  $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n \in R^m$  be the column vectors of matrix  $A$ , so that

$$A = [\mathbf{a}_1 \ \mathbf{a}_2 \ \dots \ \mathbf{a}_n].$$

A Householder transformation  $U_1$  can be found such that

$$\mathbf{b}_1 = U_1 \mathbf{a}_1$$

where

$$\|\mathbf{b}_1\| = \|\mathbf{a}_1\|$$

and

$$\mathbf{b}_1 = (b_{11}, 0, 0, \dots, 0)^T.$$

The vector norm is the usual Euclidean norm

$$\|\mathbf{v}\|^2 = v_1^2 + v_2^2 + \dots + v_m^2.$$

Thus, in the first step toward simplification, we form  $U_1A$ . In a like manner, the column vector  $\mathbf{a}_j$  could be transformed by a Householder matrix  $U_j$ ,

$$\mathbf{b}_j = U_j \mathbf{a}_j$$

so that

$$\|B\mathbf{b}_j\| = \|\mathbf{a}_j\|$$

and

$$\mathbf{b}_j = (b_{1j}, b_{2j}, \dots, b_{jj}, 0, \dots, 0)^T.$$

Thus, at the  $j^{\text{th}}$  step we would form

$$U_j \cdots U_1 A$$

so that in  $n$  steps (or less) we should have reduced  $A$  to upper triangular form. Of course, if at the  $j^{\text{th}}$  step we are not to undo the simplification already achieved in the first  $(j-1)$  columns, we could partition the matrix  $U_j$ ,

$$U_j = \begin{bmatrix} I & 0 \\ 0 & U_j' \end{bmatrix}.$$

$I$  is a  $(j-1) \times (j-1)$  identity matrix and  $U_j'$  is an  $(m-(j-1)) \times (m-(j-1))$  Householder matrix.  $U_j'$  is chosen to reduce the last  $(m-j)$  components of  $\mathbf{a}_j$  to zero, with the  $j^{\text{th}}$  element nonzero. The result is the matrix

$$U_k U_{k-1} \cdots U_1 A = QA$$

where  $k \leq n$  and  $QA$  is upper triangular.

Now, if  $U_j'$  operates on the column vector  $\mathbf{a}_j$  to give

$$\mathbf{b}_j = U_j' \mathbf{a}_j$$

then it can be shown (see Golub (1983)) that  $U_j'$  has the simple form

$$U_j' = I - 2\mathbf{u}\mathbf{u}^T$$



where  $\mathbf{u}$  is a unit vector in the direction of  $\mathbf{b}_j - \mathbf{a}_j$ .

With this form, it is easy to see that each of the matrices  $U_1, U_2, \dots, U_k$  is an orthogonal matrix (ie  $U_j^T U_j = U_j U_j^T = I$ ) and the product

$$Q = U_1 U_2 \cdots U_k$$

is also an orthogonal matrix. Orthogonal matrices have the property that

$$\|Q\mathbf{r}\| = \|\mathbf{r}\|.$$

It is this feature of Householder reductions which permits us to conclude that the above algorithm can be used to solve the best approximation problem

$$\min_{\mathbf{x} \in R^a} \|A\mathbf{b} - \mathbf{r}\|.$$

To see this, we write

$$QA = \begin{bmatrix} T \\ 0 \end{bmatrix} \quad \text{and} \quad Q\mathbf{f} = \begin{bmatrix} \mathbf{c} \\ \mathbf{d} \end{bmatrix}$$

where  $T$  is  $n \times n$  upper triangular,  $\mathbf{c}$  is an  $n \times 1$  column vector and  $\mathbf{d}$  is an  $(m - n) \times 1$  column vector. Then

$$\begin{aligned} \|\mathbf{r}\| &= \|Q\mathbf{r}\| \\ &= \|Q(A\mathbf{b} - \mathbf{r})\| \\ &= \left\| \begin{bmatrix} T \\ 0 \end{bmatrix} \mathbf{b} - \begin{bmatrix} \mathbf{c} \\ \mathbf{d} \end{bmatrix} \right\| \end{aligned}$$

where

$$\begin{bmatrix} T \\ 0 \end{bmatrix} \mathbf{b} - \begin{bmatrix} \mathbf{c} \\ \mathbf{d} \end{bmatrix} = \begin{bmatrix} \mathbf{t} \\ \mathbf{d} \end{bmatrix}.$$

Now,

$$\left\| \begin{bmatrix} \mathbf{t} \\ \mathbf{d} \end{bmatrix} \right\|$$

is minimized when  $\mathbf{t} = \mathbf{0}$ . Hence,  $\|\mathbf{r}\|$  is minimized when  $\mathbf{b} \in R^a$  is found such that

$$T\mathbf{b} = \mathbf{c}.$$

### Appendix A.3

#### Exact Solutions of Some Hele-Shaw Problems

The exact solutions presented for the cusping and Saffman finger profiles of chapter 4 can be found (expressed in different coordinates) in Howison and Aitchison (1985). Similar techniques were used by Meyer (1981) to derive an analytic solution. The method is briefly outlined below. It involves a conformal mapping of the given problem in the x-y physical plane to an auxiliary plane. For our purposes, we consider an interchange of dependent and independent variables so that the mapping is to the plane of the complex potential. (In the following, values of x, y and t must be divided by  $\pi$  to agree with the dimensions of chapter 4.)

For the case of the cusping profile, the mapping

$$z = i[w - b_0(t) - b_1(t)e^w]$$

satisfies the boundary conditions and it remains only to find  $b_0, b_1$  satisfying the free surface condition,

$$\Re\left(\frac{\partial w}{\partial t}\right)_{\phi=0} = \left|\frac{\partial w}{\partial z}\right|_{\phi=0}^2.$$

This may be written in the form

$$\Re\left(\frac{\bar{\partial z}}{\partial w} \frac{\partial z}{\partial t}\right)_{\phi=0} = -1.$$

Application of this condition to the mapping function leads to the coupled ordinary differential equations

$$b_0' - b_1 b_1' = 1$$

$$b_1' - b_0' b_1 = 0.$$

These can be solved, together with the boundary conditions  $b_0(0) = 0$ ,  $b_1(0) = \epsilon$  to yield

$$b_0 - \frac{1}{2} b_1^2 = t - \frac{1}{2} \epsilon^2$$

$$b_1 e^{-b_0} = \epsilon.$$

Now, the cusp develops in a finite time when a zero of  $\frac{dz}{dw}$  reaches the boundary  $\phi = 0$ .

We have

$$\frac{dz}{dw} = i[1 - b_1(t)e^{-w}]$$

so that  $\frac{dz}{dw} = 0$  when  $\psi = 0$  and  $\phi = \ln \frac{1}{b_1}$ . Therefore, the singularity reaches the boundary at

a value of  $t$  for which  $b_1(t) = 1$ . From the expressions above for  $b_0$  and  $b_1$  we find that the cusp develops by a time of

$$t = \frac{1}{2}(\epsilon^2 - 1) - \ln \epsilon.$$

For  $\epsilon = 0.2$  and  $t$  properly scaled by a factor of  $\frac{1}{\epsilon}$ ,

$$t_{\text{cusp}} \approx 0.35951.$$

It is shown in Howison, Ockendon and Lacey (1985) that the cusp is of the form

$$y = x^{\frac{2}{3}}.$$

In a similar manner, the mapping which produces the Saffman profile is given by

$$z = i[\omega - d(t) - \ln(1 + a(t)e^{\omega})] .$$

Substitution into the free surface condition leads to the coupled ordinary differential equations

$$d' = 1 + a^2$$

$$d'a + a' = 2a .$$

Again, these may be integrated, together with the boundary conditions  $a(0) = \varepsilon$ ,  $d(0) = 0$  to give

$$a^2 = \frac{\varepsilon^2 e^{2t}}{1 - \varepsilon^2 + \varepsilon^2 e^{2t}}$$

$$d = t + \frac{1}{2} \ln(1 - \varepsilon^2 + \varepsilon^2 e^{2t}) .$$

## References

- Aitchison, J.M. and Howison, S.D. 1985. Computation of Hele-Shaw Flows with Free Boundaries, *J. Comp. Phys.* **60**, 376-390.
- Ames, W.F. 1965. "Nonlinear Partial Differential Equations in Engineering, vol. 1", Academic Press, New York.
- Atkinson, K.E. 1978. "An Introduction to Numerical Analysis", John Wiley and Sons, In , New York.
- Axler, S. 1986. Harmonic Functions from a Complex Analysis Viewpoint, *Mathematical Monthly* **93**, 247-258.
- Chuoque, R.L., van Meurs, P. and van der Poel, C. 1959. The Instability of Slow, Immiscible, Viscous Liquid-liquid Displacement in Permeable Media, *Trans. AIME* **216**, 188-194.
- Christiansen, S. and Rasmussen, H. 1976. Numerical Solution for Two-dimensional Annular Machining Problems, *J. Inst. Maths. Applics.* **18**, 295-307.
- Clarkson, J.A. 1936. Uniformly Convex Spaces, *Trans. Am. Math. Soc.* **40**, 396-414.
- Collatz, L. 1966. "The Numerical Treatment of Differential Equations", Springer-Verlag, New York.
- Collett, D.E., Hewson-Browne, R.C. and Windle, D.W. 1970. A Complex Variable Approach to Electrochemical Machining Problems, *J. Eng. Math.* **4**, 29-37.
- Curtiss, J.H. 1960. Interpolation with Harmonic and Complex Polynomials to Boundary Values, *J. Math. Mech.* **9**, 167-192.
- Davidson, M.R. 1984. "Computational Techniques and Applications", J. Noye and C Fletcher, Eds., Elsevier Science Publishers B.V. (North-Holland), 317-326.
- Davis, P.J. 1975. "Interpolation and Approximation", Dover Publications, New York.
- Davis, P.J. and Rabinowitz, P. 1961. Advances in Orthonormalizing Computation, appearing in "Advances in Computers, vol. 2", Academic Press, New York, pp 55-133.
- Degregoria, A.J. and Schwartz, L.W. 1986. A Boundary-integral Method for Two-phase Displacement in Hele-Shaw Cells, *J. Fluid Mech.* **164**, 383-400.

- Elliott, C.M. 1980. On a Variational Inequality Formulation of an Electrochemical Machining Moving Boundary Problem and its Approximation by the Finite Element Method, *J. Inst. Maths. Applics.* **25**, 121-131.
- Elliott, C.M. and Ockendon, J.R. 1982. "Weak and Variational Methods for Moving Boundary Problems", Pitman Books Ltd., Boston.
- Fenton, J.D. and Rienecker, M.M. 1982. A Fourier Method for Solving Nonlinear Water-wave Problems: Application to Solitary-wave Interactions, *J. Fluid Mech.* **118**, 411-443.
- Forsyth, P. and Rasmussen, H. 1979. Solution of Time Dependent Electrochemical Machining Problems by a Coordinate Transformation, *J. Inst. Maths Applics.* **24**, 411-424.
- Forsyth, P. and Rasmussen, H. 1980. A Kantorovich Method of Solution of Time-dependent Electrochemical Machining Problems, *J. Comp. Meth. Appl. Mech. Eng.* **23**, 129-141.
- Forsyth, P. and Rasmussen, H. 1980. Comparison of Variational Inequality and Front-tracking Solution of an Electrochemical Machining Problem, *Utilitas Mathematica* **17**, 3-15.
- Goldstein, M.E. and Siegel, R. 1970. Conformal Mapping for Heat Conduction in a Region with an Unknown Boundary, *Int'l. J. Heat Mass Trans.* **13**, 1632-1636.
- Golub, G.H. and Van Loan, C.F. 1983. "Matrix Computations", The John Hopkins University Press, Baltimore, Maryland.
- Gottlieb, D. and Orszag, S.A. 1977. "Numerical Analysis of Spectral Methods: Theory and Applications", CBMS-NSF Reg. Conf. Ser. in Appl. Math. **26**, Philadelphia.
- Hele-Shaw, H.S. 1898. The Flow of Water, *Nature* **58**, 34-36.
- Hewson-Browne, R.C. 1971. Further Applications of Complex Variable Methods to Electrochemical Machining Problems, *J. Eng. Math.* **5**, 233-240.
- Hobby, C.R. and Rice, J.R. 1967. Approximation from a Curve of Functions, *Arch. Rational Mech. Anal.* **24**, 91-106.
- Homsy, G.M. 1987. Viscous Fingering in Porous Media, *Ann. Rev. Fluid Mech.* **19**, 271-311.
- Hougaard, P. 1977. Some Solutions of a Free Boundary Problem Related to Electrochemical Machining, Report # S9, The Danish Centre for Applied Mathematics and Mechanics.
- Howison, S.D., Ockendon, J.R. and Lacey, A.A. 1985. Singularity Development in Moving-boundary Problems, *Quart. J. Mech. Appl. Math.* **38**, 343-360.
- Krylov, A.L. 1968. The Cauchy Problem for the Laplace Equation in the Theory of Electrochemical Metal Machining, *Soviet Physics - Doklady* **13**, 15-17.

- Lacey, A.A. 1984. Design of a Cathode for an Electromachining Process, *IMA J. Appl. Math.* **34**, 259-267.
- Lamb, Horace, Sir 1945. "Hydrodynamics", Dover Publications, New York.
- Levenberg, K. 1944. A Method for the Solution of Certain Nonlinear Problems in Least Squares, *Quart. Appl. Math.* **2**, 164-168.
- Levin, D. 1980. Corrected Collocation Approximations for the Harmonic Dirichlet Problem, *J. Inst. Maths. Applics.* **26**, 65-75.
- Marquardt, D.W. 1963. An Algorithm for Least-squares Estimation of Nonlinear Parameters, *J. Soc. Indust. Appl. Math.* **11**, 431-441.
- McGeough, J.A. and Rasmussen, H. 1974. On the Derivation of the Quasi-steady Model in Electrochemical Machining, *J. Inst. Maths. Applics.* **13**, 13-21.
- McLean, J.W. and Saffman, P.G. 1981. The effect of Surface Tension on the Shape of Fingers in a Hele-Shaw Cell, *J. Fluid Mech.* **102**, 455-469.
- Menikoff, R. and Zemach, C. 1983. Rayleigh-Taylor Instability and the Use of Conformal Maps for Ideal Fluid Flow, *J. Comp. Phys.* **51**, 28-64.
- Meyer, G.H. 1981. Hele-Shaw Flow with a Cusping Free Boundary, *J. Comp. Phys.* **44**, 262-276.
- Miranker, W.L. 1981. "Numerical Methods for Stiff Equations", D. Reidel Publishing Co., Boston.
- Nilson, R.H. and Tsuei, Y.G. 1974. Inverted Cauchy Problem for the Laplace Equation in Engineering Design, *J. Eng. Math.* **8**, 329-337.
- Nilson, R.H. and Tsuei, Y.G. 1975. Free Boundary Problem of ECM by Alternating-field Technique on Inverted Plane, *J. Comp. Meth. Appl. Mech. Eng.* **6**, 265-282.
- Nilson, R.H. and Tsuei, Y.G. 1976. Free Boundary Problem for the Laplace Equation with Application to ECM Tool Design, *J. Appl. Mech.* **98**, 54-58.
- Pitts, E. 1980. Penetration of Fluid into a Hele-Shaw Cell: the Saffman-Taylor Experiment, *J. Fluid Mech.* **97**, 53-64.
- Rice, J.R. 1964. "The Approximation of Functions, vol. 1", Addison-Wesley, Reading, Mass..
- Rienecker, M.M. and Fenton, J.D. 1981. A Fourier Approximation Method for Steady Water Waves, *J. Fluid Mech.* **104**, 119-137.
- Richardson, S. 1972. Hele-Shaw Flows with a Free Boundary Produced by the Injection of Fluid into a Narrow Channel, *J. Fluid Mech.* **56**, 609-618.
- Runge, C. 1885. Zur Theorie der Eidentigen Analytischen Funktionen, *ACTA Mathematica* **6**, 229-245.

- Saffman, P.G. 1959. Exact Solutions for the Growth of Fingers from a Flat Interface Between Two Fluids in a Porous Medium or Hele-Shaw Cell, *Quart. J. Mech. Appl. Math.* **12**, 147-150.
- Saffman, P.G. and Taylor, Geoffrey, Sir 1958. The Penetration of a Fluid into a Porous Medium or Hele-Shaw Cell Containing a More Viscous Liquid, *Proc. R. Soc. London, Ser. A* **245**, 312-329.
- Siegel, R. 1973. Shape of Porous Cooled Region for Surface Heat Flux and Temperature Both Specified, *Intl. J. Heat Mass Trans.* **16**, 1807-1811.
- Shanks, D. 1955. Non-linear Transformations of Divergent and Slowly Divergent Sequences, *J. Math. Phys.* **34**, 1-42.
- Sloan, D.M. 1986. Numerical Solution of a Free Boundary Problem by Continuation, *J. Comp. Appl. Math.* **14**, 279-288.
- Vanden-Broeck, J. 1983. Fingers in a Hele-Shaw Cell with Surface Tension, *Phys. Fluids* **26**, 2033-2034.
- van der Sluis, A. 1969. Condition Numbers and Equilibration of Matrices, *Numer. Mat.* **14**, 14-23.
- Walsh, J.L. 1928(b). *Journal fur Mathematik* **159**, 197-209.
- Walsh, J.L. 1928(b). On the Degree of Approximation to an Analytic Function by Means of Rational Functions, *Trans. Am. Math. Soc.* **30**, 838-847.
- Walsh, J.L. 1929. The Approximation of Harmonic Functions by Harmonic Polynomials and by Harmonic Rational Functions, *Bull. Am. Math. Soc.* **35**, 499-544.
- Walsh, J.L. Sewell, W.E. and Elliott, H.M. 1949. On the Degree of Polynomial Approximation to Harmonic and Analytic Functions, *Trans. Am. Math. Soc.* **67** 381-420.