

GENERAL PROGRAMS FOR  
LEAST pTH AND NEAR MINIMAX APPROXIMATION

GENERAL PROGRAMS FOR  
LEAST pTH AND NEAR MINIMAX APPROXIMATION

by

Jadranka K. Rizoniko - Popović, Dipl. Eng.

A Thesis

Submitted to the Faculty of Graduate Studies  
in Partial Fulfilment of the Requirements  
for the Degree  
Master of Engineering

McMaster University

October 1972

MASTER OF ENGINEERING (1972)  
(Electrical Engineering)

McMASTER UNIVERSITY  
Hamilton, Ontario

TITLE: General Programs for Least pth and Near Minimax Approximation

AUTHOR: Jadranka K. Rizoniko-Popović, Dipl. Eng. (University of Belgrade)

SUPERVISOR: J. W. Bandler, B.Sc. (Eng.), Ph.D. (University of London),  
D.I.C. (Imperial College)

NUMBER OF PAGES: vii, 144.

SCOPE AND CONTENTS:

User-oriented computer programs in FORTRAN IV for discrete least pth approximation with a single specified function, and more generalized discrete least pth approximation with various specifications, which may also be used for nonlinear programming, are presented. Values of p up to  $10^6$  can be used successfully in conjunction with efficient gradient minimization algorithms such as the Fletcher-Powell method and a recent method due to Fletcher. It has already been demonstrated how efficiently extremely near minimax results can be achieved on a discrete set of sample points using this approach and the programs written verify this. The programs may be applied to a wide variety of design problems with a wide range of specifications. They are suitable for electrical network and system design and such problems as filter design.

### ACKNOWLEDGEMENTS

The author wishes to express appreciation to Dr. J. W. Bandler for his guidance and encouragement in the preparation of this Thesis. The author is also sincerely thankful to C. Charalambous and other colleagues for fruitful discussions and help.

Financial support provided by the National Research Council of Canada through grants A7239 and C154 and by the Department of Electrical Engineering at McMaster University in the form of a teaching assistantship is gratefully acknowledged.

I would like to thank Mary Heatherington for typing the manuscript.

## TABLE OF CONTENTS

	<u>Page</u>
<u>CHAPTER I</u> Introduction.....	1
<u>CHAPTER II</u> The Optimization Problem.....	3
<u>CHAPTER III</u> Methods of Minimization.....	6
3.1      Direct search methods.....	6
3.2      Gradient methods of minimization.....	7
3.3      The Fletcher-Powell method.....	11
3.4      The Fletcher method.....	13
<u>CHAPTER IV</u> Some Design Objectives.....	17
4.1      Error criterion.....	17
4.2 $L_p$ and $\ell_p(n)$ norms.....	20
<u>CHAPTER V</u> Practical Least pth Approximation.....	23
5.1      Scaling.....	23
5.2      Quadratic interpolation for a function of a single variable.....	25
5.3      Computer program.....	33
5.4      Examples.....	42
Test 1 .....	43
Test 2 .....	46
Test 3 .....	49

	<u>Page</u>
Test 4 .....	50
<u>CHAPTER VI</u>	
6.1     Generalized Least pth Approximation.....	55
6.2     The objective function.....	55
6.3     Case I - specification violated.....	58
6.4     Case II - specification is satisfied.....	59
6.5     Scaling.....	61
6.6     Discussion.....	63
6.7     Computer program.....	63
6.7     Examples.....	76
Test 1 .....	76
Test 2 .....	77
Test 3 .....	81
Test 4 .....	88
<u>CHAPTER VII</u>	
Conclusions.....	91
APPENDIX I.....	94
APPENDIX II.....	102
APPENDIX III.....	112
APPENDIX IV.....	115
APPENDIX V.....	117
APPENDIX VI.....	128
APPENDIX VII.....	138
REFERENCES.....	141

## LIST OF FIGURES

	<u>Page</u>
Fig. 4-1      Continuous approximation problem.	18
4-2      Discrete approximation problem.	19
5-1      Iterative searching procedure for finding maxima on two subintervals	28
5-2      Flowchart of subroutine NEWSET	31
5-3      The organization of FMCLP	37
5-4      Flowchart of subroutine FUNCT	39
5-5      Flowchart of subroutine FMCLP	40
6-1      Example of a design problem with upper and lower specifications	56
6-2      Example of a design problem when the specification is violated. Case I is applicable.	58
6-3      Example of a design problem when the specification is satisfied. Case II is applicable.	60
6-4      The organization of FMLPO	68
6-5      Flowchart of function subprogram EPSNP	69

	<u>Page</u>
Fig. 6-6 Flowchart of subroutine ERRO	70
6-7 Flowchart of subroutine FUNCG	72
6-8 Flowchart of subroutine FMLPO	74
6-9 Response of the five section transmission-line filter for the unconstrained design problem	87
6-10 Response of the five section transmission-line filter for the constrained design problem	90
7-1 Cascaded five section transmission-line network	140

## CHAPTER I

### Introduction

Two complete user-oriented computer programs in FORTRAN IV are presented which utilize some new ideas on discrete least pth approximation [1] and, more recent ideas on generalized least pth approximation [2]. Least pth approximation with  $p=2$  gives a discrete least squares approximation. With sufficiently large values of  $p$  an optimal solution very close to the optimal minimax solution can be obtained. Values of  $p$  up to  $10^6$  have been successfully employed. Gradient minimization algorithms due to Fletcher and Powell and, more recently, to Fletcher are used. The user has to write all the required specifications, the approximating functions and weighting functions in a straightforward way.

The first program is described in Chapter V with a listing attached in the Appendices and is applicable to design problems with a single specification. Quadratic interpolation, if desired, is employed to bring the discrete approximation solution closer to the solution of the continuous approximation problem. Numerical examples for which the minimax solutions are known were chosen to illustrate the work of the program. The solutions obtained are in excellent agreement with the known ones.

The second program is described in Chapter VI with instructions as to its use and the listing being given in the Appendices. The

program is directly applicable to such problems as meeting or exceeding design specifications on several disjoint closed intervals as in filter design and allows for situations more general than the conventional problem of approximating a single continuous function on a closed interval. There is no restriction on the number of variable parameters, discrete point sets and number of intervals. The examples which demonstrate that the program works, since the methods have already been tested, were chosen in system modelling and multi-section lossless transmission-line network design. Although the program is not written for nonlinear programming we found it is also applicable for problems with parameter constraints.

The programs are run on the CDC 6400 at the McMaster University Computer Center.

## CHAPTER II

### The Optimization Problem

The general optimization problem is to maximize or minimize a function  $f(\underline{a})$ , where

$$(2.1) \quad \underline{a} \triangleq \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ \vdots \\ a_k \end{bmatrix}$$

is a vector of  $k$  independent parameters or design variables. Since the problem of maximizing  $f(\underline{a})$  is equivalent to that of minimizing  $-f(\underline{a})$  only the minimization problem is considered. The function to be minimized is called the objective function and will be denoted by the symbol  $U$ .

In general, a linear or nonlinear set of constraints may have to be satisfied either during the optimization process or by the optimal solution. Each parameter might be constrained by an upper and lower bound

$$(2.2) \quad a_l \leq a_i \leq a_u$$

where

$$(2.3) \quad \underline{a}_\ell \stackrel{\Delta}{=} \begin{bmatrix} a_{\ell 1} \\ a_{\ell 2} \\ \cdot \\ \cdot \\ \cdot \\ a_{\ell k} \end{bmatrix} \quad \text{and} \quad \overline{a}_u \stackrel{\Delta}{=} \begin{bmatrix} a_{u1} \\ a_{u2} \\ \cdot \\ \cdot \\ \cdot \\ a_{uk} \end{bmatrix}$$

are lower and upper bounds, respectively. A possible set of  $h$  implicit constraint functions may be set as

$$(2.4) \quad \xi(\underline{a}) \geq 0$$

where

$$(2.5) \quad \xi \stackrel{\Delta}{=} \begin{bmatrix} c_1(\underline{a}) \\ c_2(\underline{a}) \\ \cdot \\ \cdot \\ \cdot \\ c_h(\underline{a}) \end{bmatrix}$$

and  $0$  is a zero vector. Any vector  $\underline{a}$  which satisfies the constraints is called feasible. A region of feasible points  $\underline{a}$  for which (2.2) and (2.4) are satisfied is a feasible region  $R$ .

The problem of minimizing  $U$  subject to the constraints is to locate a point  $\underline{a}$  such that

$$(2.6) \quad \underline{U} \stackrel{\Delta}{=} f(\underline{a}) = \inf_{\underline{a} \in R} f(\underline{a}) \leq f(\underline{a})$$

for any feasible  $\underline{a}$  in the neighborhood of  $\underline{U}$ . A point  $\underline{a}$  is called a constrained local minimum.

Methods which guarantee convergence to a global minimum for all functions are not available. Since we cannot generally guarantee to find a global minimum we restrict our discussion to local minima. The distinction between local and global minima is not essential from the optimization point of view; it is, however, important when the results of applying optimization techniques have to be interpreted.

Most optimization strategies reduce a constrained problem into an essentially unconstrained one. This can be accomplished by transforming the parameters and leaving the objective function unaltered [ 3 ], [ 4 ], or by modifying the objective function by introducing some kind of penalty [ 5 ] - [ 8 ]. Other methods for handling constraints employ changes in strategy when constraint violations occur, but they are not found to be so efficient and accurate. Further details may be found in some of the general references [ 8 ] - [ 11 ].

## CHAPTER III

### Methods of Minimization

#### 3.1 Direct search methods

Many direct search methods for optimization were proposed between 1960 and 1965. Direct search methods are based on a sequential examination of trial solutions which, by simple comparisons, give an indication for a further searching procedure. These methods require only the evaluation of the function at a given point and thus are applicable to general continuous functions. In general, they do not give a rapid rate of ultimate convergence and hence are inefficient for finding a minimum. Methods of this type are useful in the early stages of optimization and can provide information about a region in which a minimum is located. Many direct search methods for unconstrained optimization are described by Kowalik and Osborne [2] and by Box et al [3].

In the last few years, however, relatively few papers have been published on algorithms that do not require the calculation of any derivatives. The main subjects of these recent papers are estimating derivatives by numerical differences in order to use Davidon's [14] algorithm, lattice approximations, and fitting a quadratic function to calculate function values.

### 3.2 Gradient methods of minimization

In this section it is assumed that appropriate partial derivatives exist and methods are described which utilize partial derivative information to determine the direction of search for a minimum. The partial derivatives may be defined analytically at each point or must be estimated from the differences of values of  $f(\bar{a})$ .

For a continuous objective function  $f(\bar{a})$  with continuous first and second derivatives, we can write, using Taylor's theorem

$$(3.1) \quad f(\bar{a} + \Delta \bar{a}) = f(\bar{a}) + (\nabla f)^T \Delta \bar{a} + \frac{1}{2} \Delta \bar{a}^T G(\bar{a} + \lambda \Delta \bar{a}) \Delta \bar{a}$$

where  $0 \leq \lambda \leq 1$ ,

$$(3.2) \quad \Delta \bar{a} \stackrel{\Delta}{=} \begin{bmatrix} \Delta a_1 \\ \Delta a_2 \\ \vdots \\ \vdots \\ \Delta a_k \end{bmatrix}$$

represents an incremental vector,

$$(3.3) \quad \nabla f \stackrel{\Delta}{=} \begin{bmatrix} \frac{\partial f}{\partial a_1} \\ \frac{\partial f}{\partial a_2} \\ \vdots \\ \vdots \\ \frac{\partial f}{\partial a_k} \end{bmatrix}$$

is the gradient vector containing the first partial derivatives and

$$(3.4) \quad \underline{G} \triangleq G(\underline{x}) \triangleq \begin{bmatrix} \frac{\partial^2 f}{\partial a_1^2} & \frac{\partial^2 f}{\partial a_1 \partial a_2} & \cdots & \frac{\partial^2 f}{\partial a_1 \partial a_k} \\ \frac{\partial^2 f}{\partial a_2 \partial a_1} & \ddots & & \cdot \\ \cdot & \ddots & \ddots & \cdot \\ \cdot & \ddots & \ddots & \frac{\partial^2 f}{\partial a_k^2} \end{bmatrix}$$

is the matrix of second partial derivatives, the Hessian matrix.

For twice continuously differentiable functions  $\underline{G}$  exists and is symmetric, i.e.,  $\underline{G}^T = \underline{G}$ . A symmetric matrix  $\underline{G}$  is positive definite if for all  $\underline{a} \neq 0$

$$(3.5) \quad \underline{a}^T \underline{G} \underline{a} > 0.$$

If  $\underline{G}$  is positive definite for all  $\underline{x}$  then according to (3.1)

$$(3.6) \quad f(\underline{x} + \Delta \underline{x}) > f(\underline{x}) + (\nabla f)^T \Delta \underline{x}$$

for all  $\underline{x}$ . At a local minimum

$$(3.7) \quad \nabla f(\underline{x}) = 0$$

must hold. Assuming the first and second derivatives exist, a point is a minimum if the gradient vector is zero and the Hessian matrix is positive definite at that point.

Since the first partial derivatives of a function vanish at the minimum the second-order terms in the Taylor series expansion dominate in the vicinity of the minimum. It follows that the only methods which will minimize a general function quickly and efficiently are those which work well on a quadratic form and are guaranteed to converge eventually for a general function. All others may be slow, at least in the vicinity of the minimum, and often elsewhere.

The general quadratic function of  $k$  variables can be written

$$(3.8) \quad Q(\underline{a}) = c + \underline{b}^T \underline{a} + \frac{1}{2} \underline{a}^T \underline{A} \underline{a}$$

with gradients

$$(3.9) \quad \underline{g} = \underline{b} + \underline{A} \underline{a}$$

where  $\underline{A}$  is a  $k \times k$  symmetric matrix, and  $c$  and  $\underline{b}$  are constants. If  $\underline{A}$  is positive definite

$$(3.10) \quad \underline{a}^* = -\underline{A}^{-1} \underline{b}$$

and  $\underline{a}^*$  is unique and can be found in a finite number of steps.

We are interested in a quadratic approximation, such as (3.8), to the objective function  $f(\underline{a})$ , because a quadratic is the simplest differentiable function that has a well-defined minimum. In this case in the  $i$ th iteration the obvious quadratic approximation to  $f(\underline{a})$  is

$$(3.11) \quad f(\underline{a}) \approx f(\underline{a}_i) + \underline{g}_i^T \Delta \underline{a}_i + \frac{1}{2} (\Delta \underline{a}_i)^T \underline{G}_i \Delta \underline{a}_i$$

where  $\underline{g}_i$  is the first derivative vector and  $\underline{G}_i$  is the second

derivative matrix of  $f(\bar{x})$  at point  $\bar{x}_i$ . The iterative process known as Newton iterative procedure is defined by the equation

$$(3.12) \quad \bar{x}_{i+1} = \bar{x}_i - (\bar{G}_i)^{-1} \bar{g}_i.$$

The generalized Newton-Raphson method [15] has fast convergence eventually, but frequently fails to converge from a poor initial estimate of the minimum.

A quasi-Newton iterative algorithm that involves linear searches does not require second-order derivatives of the function to be evaluated [16]. The term quasi-Newton means that an attempt is made to simulate equation (3.12), and a matrix  $\bar{H}_i$  is used, which is an estimate of  $(\bar{G}_i)^{-1}$ . Because linear searches are made,  $\bar{x}_{i+1}$  may be defined by the formula

$$(3.13) \quad \bar{x}_{i+1} = \bar{x}_i - \lambda_i \bar{H}_i \bar{g}_i,$$

where  $\lambda_i$  is a multiple which is calculated at each iteration.

The recent papers on minimization algorithms propose modifications to the iteration (3.12), so that it is not necessary for  $\bar{x}_0$  to be close to the solution. Gradient methods that are used at the present time are summarized in a general theory published in reference [17].

The ideas resulting from equation (3.13) have been successfully employed in an iterative procedure in the gradient optimization methods of Fletcher and Powell [18]. Generally acknowledged to be among the most powerful minimization methods currently available,

when first derivatives are analytically defined, is this method and a more recent one due to Fletcher [19]. As they are both used in a program presented in this thesis, a brief discussion of the methods follows.

### 3.3 The Fletcher-Powell method

The method presented by Fletcher and Powell is an iterative method for finding a local minimum of a general function  $f(\tilde{a})$  where the gradient vector  $\tilde{g}$  (3.9) is defined analytically. It is based upon Davidon's variable metric method [4].

In this method the matrix  $\tilde{G}$  from (3.8) is not evaluated directly but instead another matrix  $\tilde{H}$  is introduced and used as an approximation to  $\tilde{G}^{-1}$ . The increment  $\tilde{\delta}$  is taken along the direction of search  $\tilde{s}$

$$(3.14) \quad \tilde{s} = -\tilde{H}\tilde{g}$$

so that

$$(3.15) \quad \tilde{\delta} = \alpha \tilde{s}$$

where  $\alpha$  is such that  $f(\tilde{a} + \lambda \tilde{s})$  is a minimum with respect to  $\lambda$  along  $\tilde{a} + \lambda \tilde{s}$ . The approximating matrix  $\tilde{H}$  is updated at each iteration using the formula

$$(3.16) \quad \tilde{H}_{i+1} = \tilde{H}_i + \frac{\tilde{\delta}_i \tilde{\delta}_i^T}{\tilde{\lambda}_i} - \frac{\tilde{H}_i \tilde{\gamma}_i \tilde{\gamma}_i^T \tilde{H}_i}{\tilde{\lambda}_i \tilde{\lambda}_i \tilde{H}_i \tilde{\lambda}_i}$$

where

$$(3.17) \quad \delta_i = \tilde{x}_{i+1} - \tilde{x}_i = -\alpha_i H_i g_i$$

and

$$(3.18) \quad \gamma_i = \tilde{g}_{i+1} - \tilde{g}_i$$

are the changes in  $\tilde{x}$  and  $\tilde{g}$  during the  $i$ th iteration. If  $H_0$  is positive definite initially then it can be proved inductively that all subsequent  $H_i$  are also positive definite. Since  $H_0$  is initially chosen as the identity matrix, then all  $H_i$  will be positive definite. As a consequence the method converges. When the method is applied to the quadratic function it leads to a minimum in  $k$  iterations and it is further proved that the matrix  $H_k$  converges to the inverse of the Hessian matrix  $G^{-1}$ .

The algorithm needs to find the multiple  $\alpha$  at each iteration. A sufficiently accurate minimum is obtained by evaluating the function and gradient for different values of  $\lambda$  and using cubic interpolation. The linear search usually requires considerable computing effort.

It is practicable to apply this method to find a local minimum of a general function whose first derivatives can be evaluated quickly, even if only a poor initial approximation to a solution is known. Both the proof of convergence and success in practice depend on accurate location of the minimum in the linear searches.

The termination criteria for the Fletcher-Powell iterative procedure are as follows:

1. If the function value has not decreased in the last iteration step, the search for the minimum is terminated provided the gradient vector is already sufficiently small.
2. If the elements of vector  $\alpha$  change by very small amounts, and at least  $k$  iterations are performed, the minimization is terminated.
3. If the number of iterations exceeds an upper bound pre-assigned by the user, further computation is bypassed, and the appropriate message indicates poor convergence.
4. If one of the successive linear searches indicates that no minimum exists, further calculation is bypassed and the message indicates that no minimum exists.

### 3.4 The Fletcher method

A method presented by Fletcher came as a result of further development of the Fletcher-Powell method. He retains the positive definiteness of  $H$  in the new algorithm, as it ensures a reduction in the function value on each iteration, but a multiple  $\alpha$  is not calculated at a point where the function is minimized locally along the direction of search.

The idea of replacing the linear search, which is time consuming, by a simpler way of reducing the function at each iteration

has the consequence of abandoning the property of quadratic termination which guarantees fast ultimate convergence. If linear search is abandoned, something else has to force a sufficiently large decrease in  $f(\hat{x})$  at each iteration to guarantee convergence. The change in function in one iteration would be expected by Taylor's series to be approximately  $\hat{g}^T \hat{\delta}$  when  $\hat{\delta}$  is small, but much less than that in absolute value when the position of the minimum along the line is overestimated. When that happens, the change in  $f(\hat{x})$  relative to  $\hat{g}^T \hat{\delta}$ , i.e.  $\Delta f / \hat{g}^T \hat{\delta}$ , cannot become arbitrarily small.

The corrections are determined by

$$(3.19) \quad \hat{\delta} = -\lambda \tilde{H} \hat{g}$$

and trying the values of  $\lambda = 1, w, w^2, \dots$  for  $0 < w < 1$  will eventually produce a  $\hat{\delta}$  which satisfies the preassigned ratio  $\Delta f / \hat{g}^T \hat{\delta}$ .

Another important feature is to retain a strict positive definiteness of the approximating matrices if these ideas are to be used. Failure may be caused by slow rate of convergence and due to round-off errors in  $\tilde{H}$  which may become singular. So a new formula for updating  $\tilde{H}$  was found

$$(3.20) \quad \tilde{H}_{i+1} = \tilde{H}_i - \frac{\hat{\delta}_i \hat{\gamma}_i \tilde{H}_i}{\hat{\delta}_i^T \hat{\gamma}_i} - \frac{\tilde{H}_i \hat{\gamma}_i \hat{\delta}_i^T}{\hat{\delta}_i^T \hat{\gamma}_i} + \left(1 + \frac{\hat{\gamma}_i^T \tilde{H}_i \hat{\gamma}_i}{\hat{\delta}_i^T \hat{\gamma}_i}\right) \frac{\hat{\delta}_i \hat{\gamma}_i^T}{\hat{\delta}_i^T \hat{\gamma}_i}$$

where  $\hat{\delta}_i$  and  $\hat{\gamma}_i$  are defined in (3.17) and (3.18), respectively.

The algorithm takes into account the above discussion on  $\tilde{H}$  and  $\lambda$ . It is desirable to have some guarantee that the  $\tilde{H}$  matrices tend to  $G^{-1}$  so that the ultimate convergence for quadratic functions can be proved. A suitable property referred as a Property 1 has been defined in [19] which requires that, for quadratic functions, the eigen values of  $\tilde{H}$  must tend monotonically to those of  $G^{-1}$  in a certain sense. The formula (3.20) possesses the required Property 1. It is shown that the formula (3.16) also possesses that property [19]. The two formulae can be combined together and they can be used to minimize quadratic functions and thus permit as close an approximation to  $G^{-1}$  as possible. The test which predicts which formula will be used is the following

$$(3.21) \quad \delta^T \tilde{x} \geq \tilde{x}^T H \tilde{x} .$$

If (3.21) is true the formula (3.20) is used, if not the formula (3.16) is used. In fact when equality holds no preference is important but it was decided to use the formula (3.20) to avoid possible singularity in  $H$ .

In most cases for  $\lambda = 1$  convergence will occur and so the method requires only one evaluation of  $f$  and  $g$  per iteration which results in the efficiency of the algorithm. If  $\lambda = 1$  is not successful in reducing  $f$  sufficiently, then a  $\lambda$  is determined by cubic interpolation and the largest of this value and 0.1 is used. Thus the possibility of obtaining a local minimum along a line is retained when it is necessary, in which case more than one evaluation of  $f$  and

$\hat{g}$  is used in one iteration. The property of quadratic termination has been replaced by a property in which the approximating matrix  $\hat{H}$  has to tend to the inverse Hessian  $\hat{G}^{-1}$  in a certain sense.

The termination criteria for the Fletcher iterative procedure are based on the same ideas used for the termination of the Fletcher-Powell method.

## CHAPTER IV

### Some Design Objectives

#### 4.1 Error criterion

Most electrical network design problems can be formulated as approximation problems. Evaluation of the difference between a specified function and an approximating function leads to an error criterion. Let us define a weighted error function or deviation between a specified function and an approximating function as

$$(4.1) \quad e(\underline{a},x) \stackrel{\Delta}{=} w(x)(F(\underline{a},x) - S(x))$$

where

$S(x)$  is the real or complex specified function,

$F(\underline{a},x)$  is the real or complex approximating function,

$w(x)$  is the positive weighting function,

$\underline{a}$  is a  $k$  dimensional vector representing the adjustable parameters.

The approximation problem when all the functions from (4.1) are continuous and real is illustrated on Fig. 4-1.

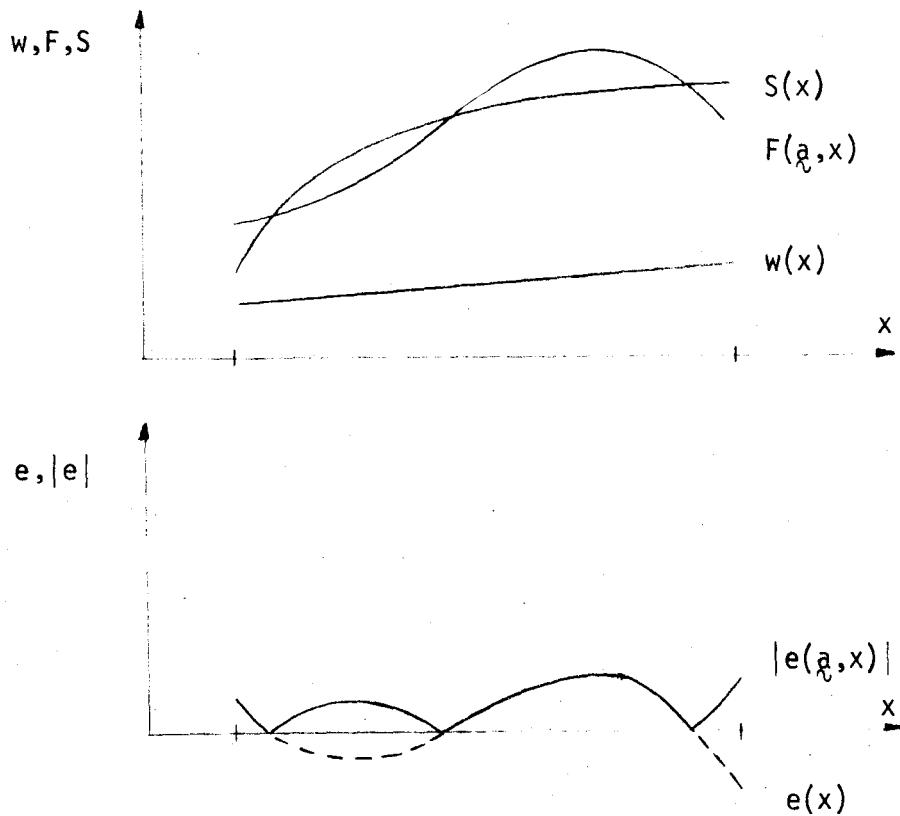


Fig. 4-1 Continuous approximation problem

In general, more can be done for a finite point set of independent variables than for the whole interval. With a discrete set some values of the independent variables can be isolated or approximation can be done on several disjoint closed intervals as is shown in Fig. 4-2. Of course, in a case when a curve has to be fitted to experimental data discrete approximation may be the only

approach possible. It is appropriate then to consider discrete approximation.

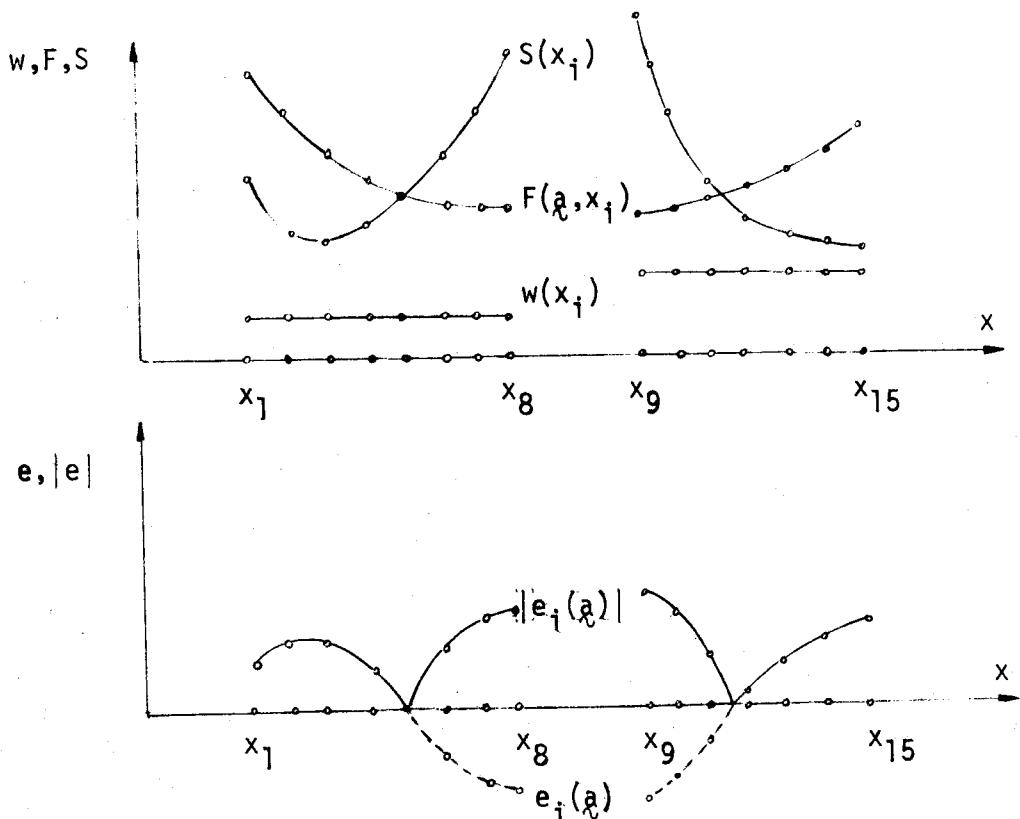


Fig. 4-2 Discrete approximation problem

In many problems of interest a relatively modest number of well-chosen sample points result in solutions to the design problem sufficiently good for all practical purposes.

A discrete version of the error function (4.1) is given by

$$(4.2) \quad e_i(\underline{a}) \stackrel{\Delta}{=} e(\underline{a}, x_i) = w(x_i)(F(\underline{a}, x_i) - S(x_i)) , \quad i \in I$$

where  $I$  is a sample set of  $n$  points.

A purpose of the weighting function  $w(x)$  is to emphasize or deemphasize various parts of the difference between the approximating function and the specification to suit the designer's requirement. Thus the optimum with respect to one weighting function may not be an optimum with respect to another.

#### 4.2 $L_p$ and $\ell_p(n)$ norms

We may define a norm on a closed interval  $[x_a, x_b]$

$$(4.3) \quad \|e\|_p \stackrel{\Delta}{=} \left\{ \int_{x_a}^{x_b} |e(a, x)|^p dx \right\}^{\frac{1}{p}} \quad p \geq 1$$

called the  $L_p$  norm for the continuous case, and a norm

$$(4.4) \quad \|e\|_p^n \stackrel{\Delta}{=} \left\{ \sum_{i \in I} |e_i(a)|^p \right\}^{\frac{1}{p}} \quad p \geq 1$$

where

$$(4.5) \quad e(a) = \begin{bmatrix} e_1(a) \\ e_2(a) \\ \vdots \\ \cdot \\ \vdots \\ e_n(a) \end{bmatrix}$$

called the  $\ell_p(n)$  norm as a discrete approximation to (4.3). The process of minimizing the objective function defined as the  $L_p$

or  $\ell_p(n)$  norm is called least  $p$ th approximation. The solution is termed a best approximation with respect to  $\|e\|_p$  for the continuous or  $\|\tilde{e}\|_p^n$  for the discrete case. A widely known case is for  $p = 2$ , namely least squares approximation. An important limiting case where  $p$  approaches infinity is minimax, Chebyshev or  $L_\infty$  and  $\ell_\infty(n)$  approximation. In this case we may define

$$(4.6) \quad \|e\|_\infty \stackrel{\Delta}{=} \max_{[x_a, x_b]} |e(a, x)|$$

and

$$(4.7) \quad \|\tilde{e}\|_\infty^n \stackrel{\Delta}{=} \max_{i \in I} |e_i(a)|$$

the Chebyshev or uniform norms.

In network and system design, the optimal responses usually turn out to be equal-ripple, although equal-ripple responses do not necessarily yield optimal solutions in the minimax sense.

It remains to point out that the optimization of the objective functions when defined as the uniform norms, by any of the well known methods for differentiable functions is unsatisfactory because such methods attempt to reduce the norm of the gradient function to zero whereas in the minimax problem, the gradient may be discontinuous.

In spite of this some successful algorithms and computer programs with objectives in the form of (4.6) and (4.7) for a certain class of functions have been reported [20] - [22] and a good survey of

practical minimax approximation is presented in [23].

Least pth approximation does not suffer from this difficulty of minimax objectives. Therefore least pth approximation with values of  $p$ ,  $1 < p < \infty$ , has been proposed and used to obtain optimal designs which produce the results close to the minimax response. For the special error criterion of (4.4) where  $p = 2$ , a simplified method has been developed [3], [24] and [25]. Temes and Zai [26] and [27] have extended a least squares method to a least pth method using for minimization in place of  $\|e\|_p^n$  the following objective function

$$(4.8) \quad U = \sum_{i \in I} (e_i(a))^p \text{ for } p \geq 2 \text{ and } p \text{ even.}$$

They developed an algorithm for least pth approximation [27] usable with moderately large values of  $p$ , say of the order of 10. Afterwards, values of  $p$  up to  $10^{12}$  have been successfully employed minimizing the  $\ell_p(n)$  norm in conjunction with efficient gradient minimization algorithms such as the Fletcher-Powell method and the recent method due to Fletcher [1]. Bandler and Charalambous generalized the objectives so that least pth approximation could be used to meet or exceed circuit or system performance specifications [2].

More recent theoretical work has been published on conditions for optimality in least pth approximation with  $p \rightarrow \infty$  [28] from which conditions for a minimax approximation [29] fall out.

## CHAPTER V

### Practical Least pth Approximation

#### 5.1 Scaling

The obvious reason why large values of  $p$  are desirable in the least  $p$ th objective is that the corresponding optimal approximations tend to become minimax approximations. It was pointed out that discrete approximation was more practical, therefore from now on we restrict the discussion to discrete least  $p$ th approximation only.

We consider the minimization of the  $\ell_p(n)$  norm (4.4) which is our objective

$$(5.1) \quad U(a) = \left( \sum_{i \in I} |e_i(a)|^p \right)^{\frac{1}{p}} \quad \text{for } 1 < p < \infty.$$

The scaling of functions and variables in an optimization problem is often crucial to the successful functioning of existing methods. Intuitively, we mean by "well-scaled" that similar changes in the function or variables lead to similar changes in the objective function. In a contour diagram of a well-scaled function of two variables the contour lines would not deviate too far from concentric circles and in this case we would hope that a method of steepest descent [72] would work satisfactorily. When the contours are exactly concentric circles, steepest descent gives the optimum in one step.

To alleviate the ill-conditioning in (5.1) for very large values of  $p$ , a scaling was proposed [1] which modifies all the terms which have to be raised to the power  $p$  such that they are less than or equal to unity. Let

$$(5.2) \quad M(\tilde{a}) \stackrel{\Delta}{=} \max_{i \in I} |e_i(\tilde{a})|$$

and rewrite (5.1) in the form

$$(5.3) \quad U(\tilde{a}) = M(\tilde{a}) \left( \sum_{i \in I} \left| \frac{e_i(\tilde{a})}{M(\tilde{a})} \right|^p \right)^{\frac{1}{p}}.$$

It is always true for all  $i \in I$  that

$$(5.4) \quad \left| \frac{e_i(\tilde{a})}{M(\tilde{a})} \right| \leq 1$$

where there is at least one  $i$  for which equality holds. Large errors will be emphasized by large values of  $p$  and the contribution of the other terms in the summation will become very small.

Assuming that the  $e_i(\tilde{a})$  for all  $i \in I$  are continuous with continuous derivatives, the gradient of (5.3) is given by

$$(5.5) \quad \nabla U(\tilde{a}) = \left( \sum_{i \in I} \left| \frac{e_i(\tilde{a})}{M(\tilde{a})} \right|^p \right)^{\frac{1-p}{p}} \sum_{i \in I} \left| \frac{e_i(\tilde{a})}{M(\tilde{a})} \right|^{p-2} \cdot \text{Re} \left( \frac{e_i^*(\tilde{a})}{M(\tilde{a})} \nabla e_i(\tilde{a}) \right)$$

where

$$(5.6) \quad \nabla \triangleq \begin{bmatrix} \frac{\partial}{\partial a_1} \\ \frac{\partial}{\partial a_2} \\ \vdots \\ \vdots \\ \frac{\partial}{\partial a_k} \end{bmatrix}$$

and \* denotes the complex conjugate.

We note that  $M(a)$  and also  $U(a)$  is positive and  $1 < p < \infty$ .

## 5.2 Quadratic interpolation for a function of a single variable

It is assumed that a sufficient number of sample points have been chosen along the  $x$  axis so that the discrete approximation problem adequately approximates the continuous problem. However, it should be remembered that function evaluations are often by far the most time consuming parts of an optimization process. So the number of sample points should be carefully chosen for the particular problem under consideration. These arguments apply, of course, to any formulation which involves sampling. If the requirement is to concentrate more on minimizing the maximum deviation it is suitable to sample points in the neighbourhood of the maxima of the weighted error function. As one usually cannot know the positions of the maxima in advance, it is common to space the sample points uniformly along the axis. Retaining the

maxima among them and reducing the number of the sample points by removing from the objective function those which do not contribute to the summation considerably may save computation time. Even more can be done if the approximations to the actual maxima replace the sample points in their neighborhood.

It is assumed that, in the neighborhood of the extremum, the function is adequately represented by a quadratic form. This assumption is the same as the ones we made in Chapter III, where we assume that in the neighborhood of the minimum the function to be optimized can be approximated by a quadratic form. The function is evaluated at three points, a quadratic interpolation polynomial is fitted to it, and the maximum of this interpolant is formed. This point replaces one of the initial points. This assumption has been thoroughly tested in practice and found satisfactory.

The basic algorithm involving quadratic interpolation applied to the weighted error function which is a multimodal function may be explained as follows:

Assume that the weighted error function is continuous on the closed interval  $[x_a, x_b]$ . Let

$$(5.7) \quad \{x_i^{(j)}\}, \quad i \in I$$

be the set of  $n$  sample points at the beginning of the  $j$ th iteration. Before the searching procedure for the extrema of the weighted error starts, the end points of the interval are added to the given set (5.7),

if they are not already included. A set  $\{x_i^{(j)}\}$  is used to construct  $n+1$  (or  $n+3$ , if  $x_1 \neq x_a$  and  $x_n \neq x_b$ ) subintervals over  $[x_a, x_b]$ . Each subinterval is divided by a predicted  $n_s$  equidistant grid of points. Let

$$(5.8) \quad \{z_k^{(i,j)}\}, \quad k=1,2,\dots,n_s+1$$

be the set of equidistant grid of points on  $[x_i^{(j)}, x_{i+1}^{(j)}]$  interval. The extrema of the error function are found by the sequential examination of the values of the weighted error function and by comparison with the greatest on the subinterval obtained up to that time. If it happened that both neighbouring points have absolute values of the weighted error less than the current one, then the extremum  $\hat{z}_i^{(i,j)}$  is found by applying the quadratic interpolation

$$(5.9) \quad \hat{z}_k^{(i,j)} = z_k^{(i,j)} + \frac{|e_{k+1}| - |e_{k-1}|}{2|e_k| - |e_{k+1}| - |e_{k-1}|} \cdot \Delta_i$$

where

$$(5.10) \quad e_k \stackrel{\Delta}{=} e_k(z_k^{(i,j)})$$

and

$$(5.11) \quad \Delta_i \stackrel{\Delta}{=} z_{k+1}^{(i,j)} - z_k^{(i,j)} = \frac{x_{i+1}^{(j)} - x_i^{(j)}}{n_s}.$$

This point replaces  $x_i^{(j)}$ , the left end point of the current subinterval.

Immediately after the extreme point is located on the grid, the point on the grid next to the extremum replaces the left end point of the next subinterval. This is done in case there are more than one

extremum on a single subinterval. Thus the other extrema, if they exist, are "removed" to the next subinterval by removing its left end point. However, if it is found that there are no extreme points on the next subinterval, the end point is again set to its previous value.

All the extrema are selected by applying this searching procedure on every subinterval and these points replace the nearby left hand side point obtained up to that time, and the new set  $\{x_i^{(j+1)}\}$ ,  $i \in I$  for the  $(j+1)$ th iteration is thus obtained. This set does not necessarily contain the end points of the interval, but they are included in the searching technique at the next iteration to avoid shrinking the interval.

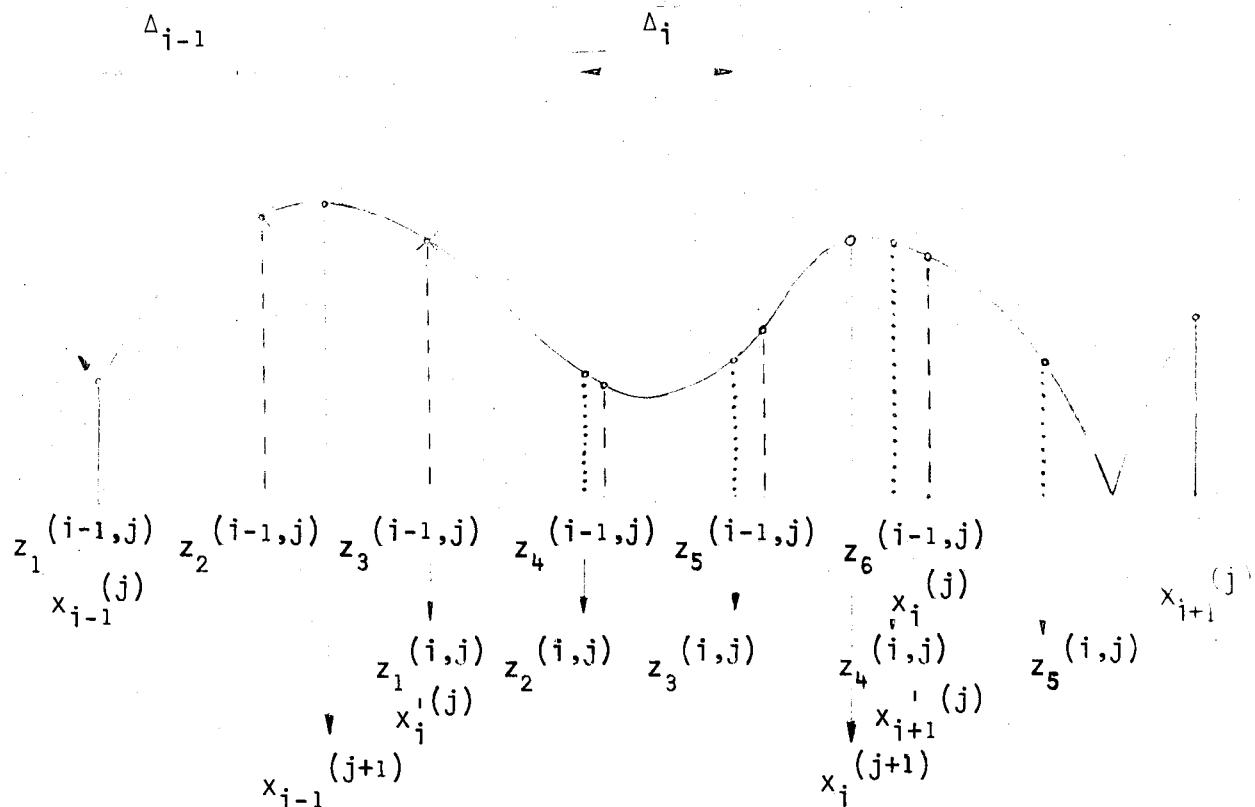


Fig. 5-1. Iterative searching procedure for finding maxima on two subintervals.

In an effort to keep the number of the discrete points in the summation (5.3) constant, at the end of the searching techniques we select  $n$  sample points among  $n+2$  according to the absolute value of the corresponding weighted error function. If the error at one end point of the interval has a considerably large value relative to the other, the other is omitted. But if both end points have considerable errors, two neighbouring points are omitted among four successive internal points where the absolute values of the weighted error function are the smallest.

A scheme of this procedure for  $n_s = 5$  is illustrated in Fig. 5-1.

Selection of the extreme points is significant especially within the first iteration. Once they are found, they do not usually move too far away in the next iterations.

The following is a list of some arguments used for the quadratic interpolation technique in the subroutine NEWSET:

indic(1) may have values 1 or 2.

indic(1)=1 indicates that quadratic interpolation is not applied on a subinterval;

indic(1)=2 indicates that a quadratic interpolation is done, and the left end point of the next subinterval is temporarily removed.

indic(2) may have values 1 or 2.

indic(2)=1 indicates that quadratic interpolation was not

applied on a previous subinterval;

indic(2)=2 indicates that a left end point of the sub-interval was temporarily removed and if indic(1)=2, the left end point will be set to the new extreme point on the subinterval, otherwise it will be fixed to the value it had before the searching technique was applied.

indic(3) may have values 2, 1, -1 and -2.

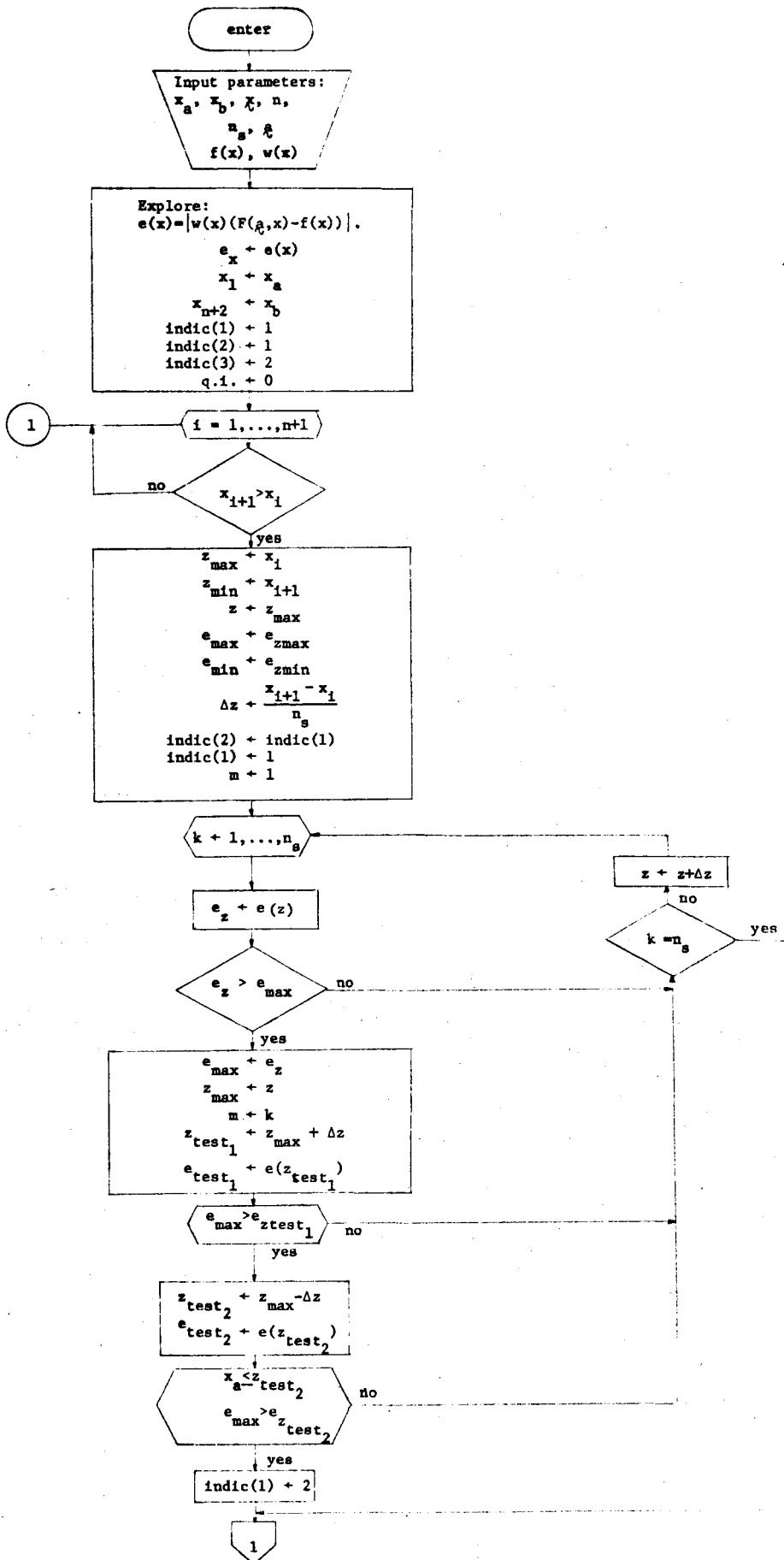
indic(3)=2 indicates that a new set of points does not include the left and the right end points of the interval  $[x_a, x_b]$ ;

indic(3)=1 indicates that  $x_b$  is included in the new set of points;

indic(3)=-1 indicates that  $x_a$  is included in the new set of points;

indic(3)=-2 indicates that both  $x_a$  and  $x_b$  are included in the new set of points.

q.i. indicates the number of quadratic interpolations done on the interval  $[x_a, x_b]$ .



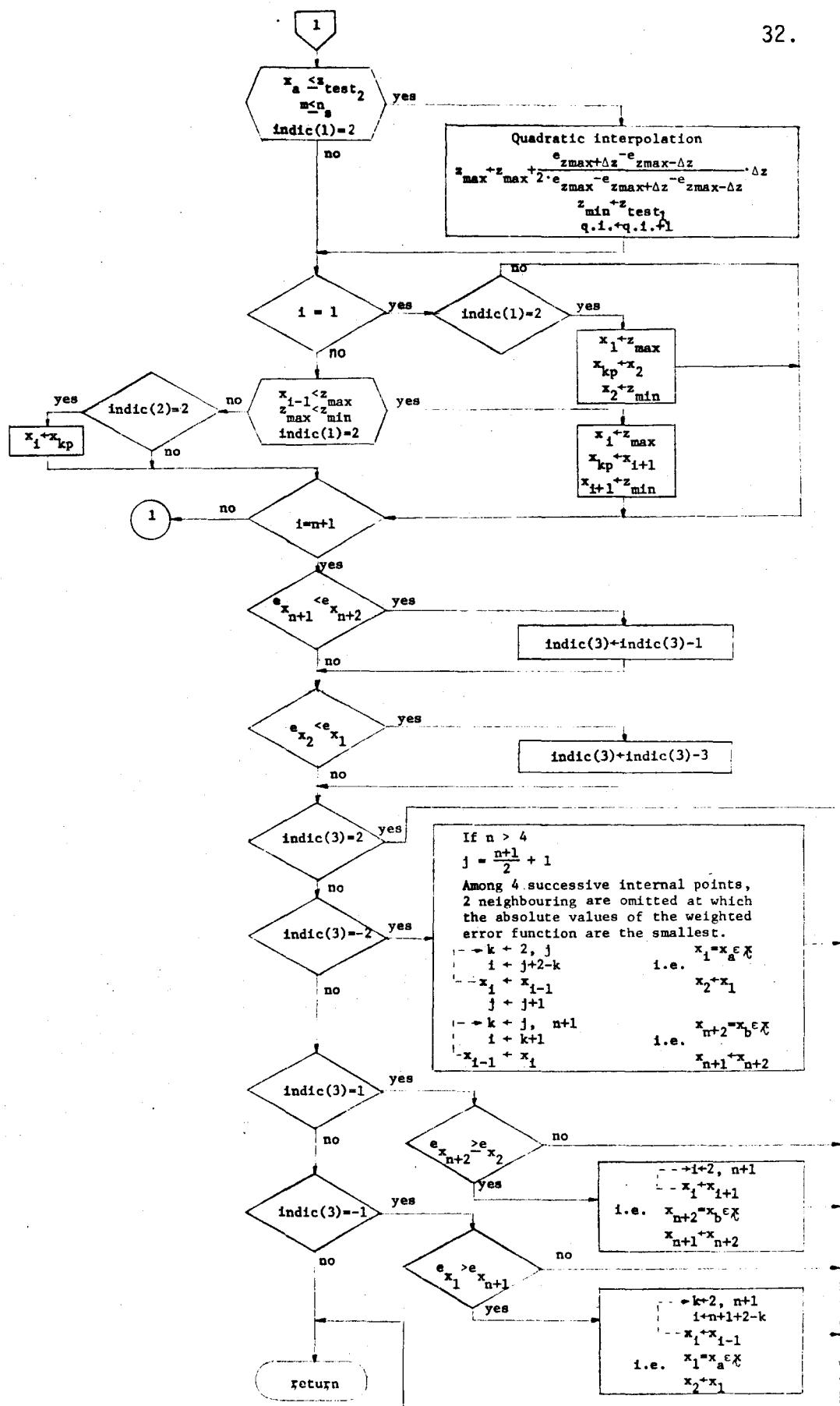


Fig. 5-2 Flowchart of subroutine NEWSET

### 5.3 Computer program

A computer program in FORTRAN IV is written for minimizing a function defined as a least pth objective where all the functions from (4.1) are real. Although the program is used simply by calling the name FMCLP, it contains 14 different subprograms. A list and a brief description of all the subprograms is given below:

FMCLP	Supplies all the data for the function minimization and coordinates the other subprograms in the package;
S	Defines a specified function over an interval;
FAPP	Defines an approximating function over an interval and the gradients w.r.t. variable parameters;
W	Defines a weighting function over an interval;
WERR	Computes a weighted error function at each sample point;
NEWSET	Redefines a sample point set using quadratic interpolation to include all the extreme points;
FUNCT	Finds the absolute maximum weighted error and computes the least pth objective (5.3) and its gradients (5.5);
GRDCHK	Checks the gradients w.r.t. all variable parameters

before the optimization process starts;

FMFPC Minimizes a function using the Fletcher-Powell optimization technique described in the section 3.3;

FMNFC Minimizes a function using the Fletcher optimization technique described in the section 3.4;

INPUT Prints the input data for the optimization process;

FINAL Prints the optimum solution after function minimization;

WRITE1 and WRITE2 Print the intermediate results.

S, W and WERR are function subprograms and the others are subroutine subprograms.

A user of the package FMCLP is supposed to write the following subprograms from the above list: S, FAPP and W where he defines his own design problem in a straightforward way. More detailed instructions about how this should be done to fit the given package are available in Appendix I.

Subroutine NEWSET is written according to the algorithm and the flowchart given in Section 5.2. The other subprograms require more explanation.

The output of the function subprogram WERR has a value of the weighted error at a single point  $x$  for a particular vector  $\underline{a}$ , i.e.

$$(5.12) \quad WERR \stackrel{\Delta}{=} (F(\underline{a}, x) - S(x)) \cdot w(x) .$$

This subprogram is called for each sample point as many times as there are the changes in the vector  $\underline{a}$ .

Subroutine FUNCT keeps the values of the weighted error of each sample point in an n-dimensional array, finds the maximum absolute value among them and computes the objective function and its gradients according to the equations (5.3) and (5.5), respectively. The flowchart of FUNCT is given in Fig. 5-4.

Subroutine GRDCHK checks the gradients of the approximating function, which a user has to supply, indirectly through the gradients of the objective function (5.5), before the optimization procedure starts. The criterion used for checking the gradient is the relative comparison between the analytical and the numerical gradients against a small given number  $\eta$

$$(5.13) \quad \frac{\frac{U(a_i + \Delta a_i) - U(a_i)}{\Delta a_i} - \frac{\partial U}{\partial a_i}}{\frac{U(a_i + \Delta a_i) - U(a_i)}{\Delta a_i}} < \eta , \quad i=1, \dots, k$$

where in particular  $\Delta = 10^{-4}$  and  $\eta = 10^{-1}$ . If it happens that

$$(5.14) \quad \frac{U(a_i + \Delta a_i) - U(a_i)}{\Delta a_i} < 10^{-20} , \quad i=1, \dots, k ,$$

then a small quantity  $10^{-20}$  replaces the denominator in (5.13) to avoid the possible division by zero in that case. The gradients are checked only once, for the starting values of  $x$  and  $\alpha$ . If (5.13) does not hold, the program will terminate and wait for a user's intervention.

FMFPC and FMNFC are subroutines available for function minimization. The user has to specify which method is going to be used, how many times and the stopping criteria for the minimization procedure. There is the possibility of choosing both at one run. For more detailed instructions see Appendix I.

The subroutines INPUT, FINAL, WRITE1 and WRITE2 are called only when certain data or results are desired to be printed out. How this is handled is explained in Appendix I.

All the subprograms are coordinated by FMCLP. Fig. 5-3 illustrates how all these subprograms are tied to FMCLP and is useful in understanding the organization of the whole program. A flowchart of FMCLP, given in Fig. 5-5, shows in more details how this is actually done.

The program terminates when the stopping criteria for the optimization method, Fletcher-Powell or Fletcher method, are satisfied or when the relative change in the objective function in two successive iterations is less than a small prescribed quantity  $\epsilon$ , i.e.,

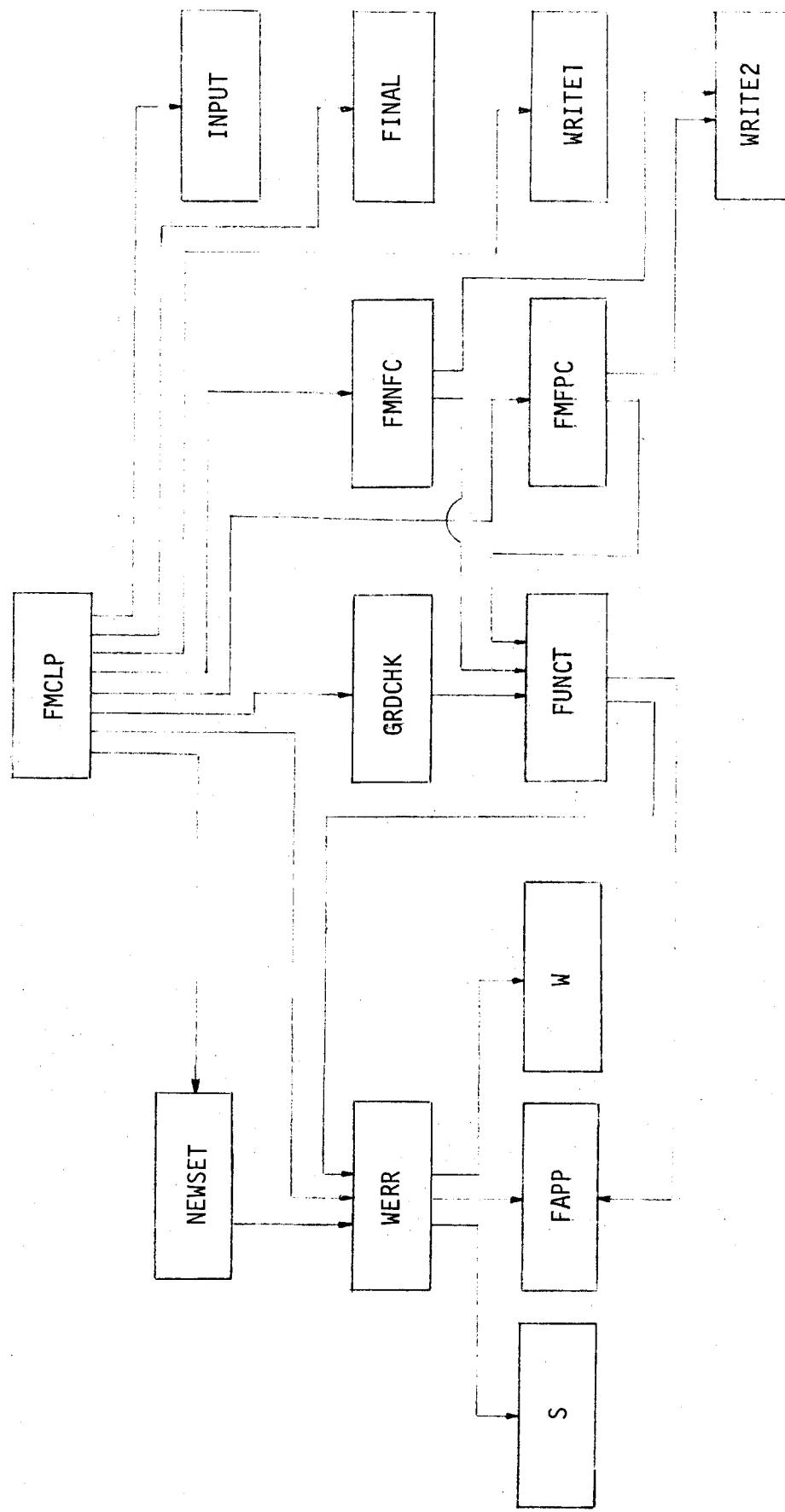


Fig. 5-3 The organization of FMCLP.

$$(5.15) \quad \frac{U(\hat{x}^{j-1}) - U(\hat{x}^j)}{U(\hat{x}^{j-1})} < \epsilon .$$

The listing of the program described in this Chapter is attached in Appendices II and III.

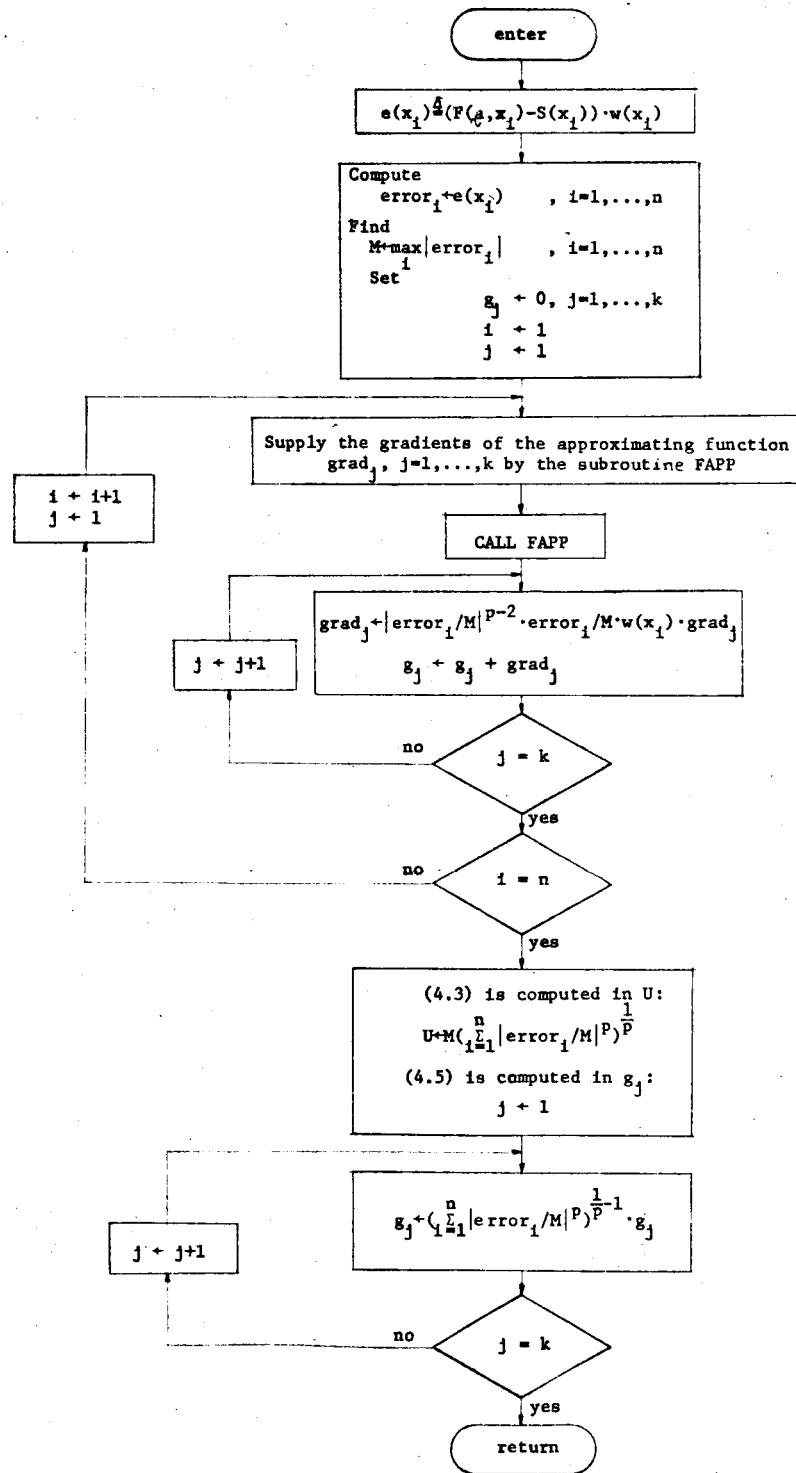
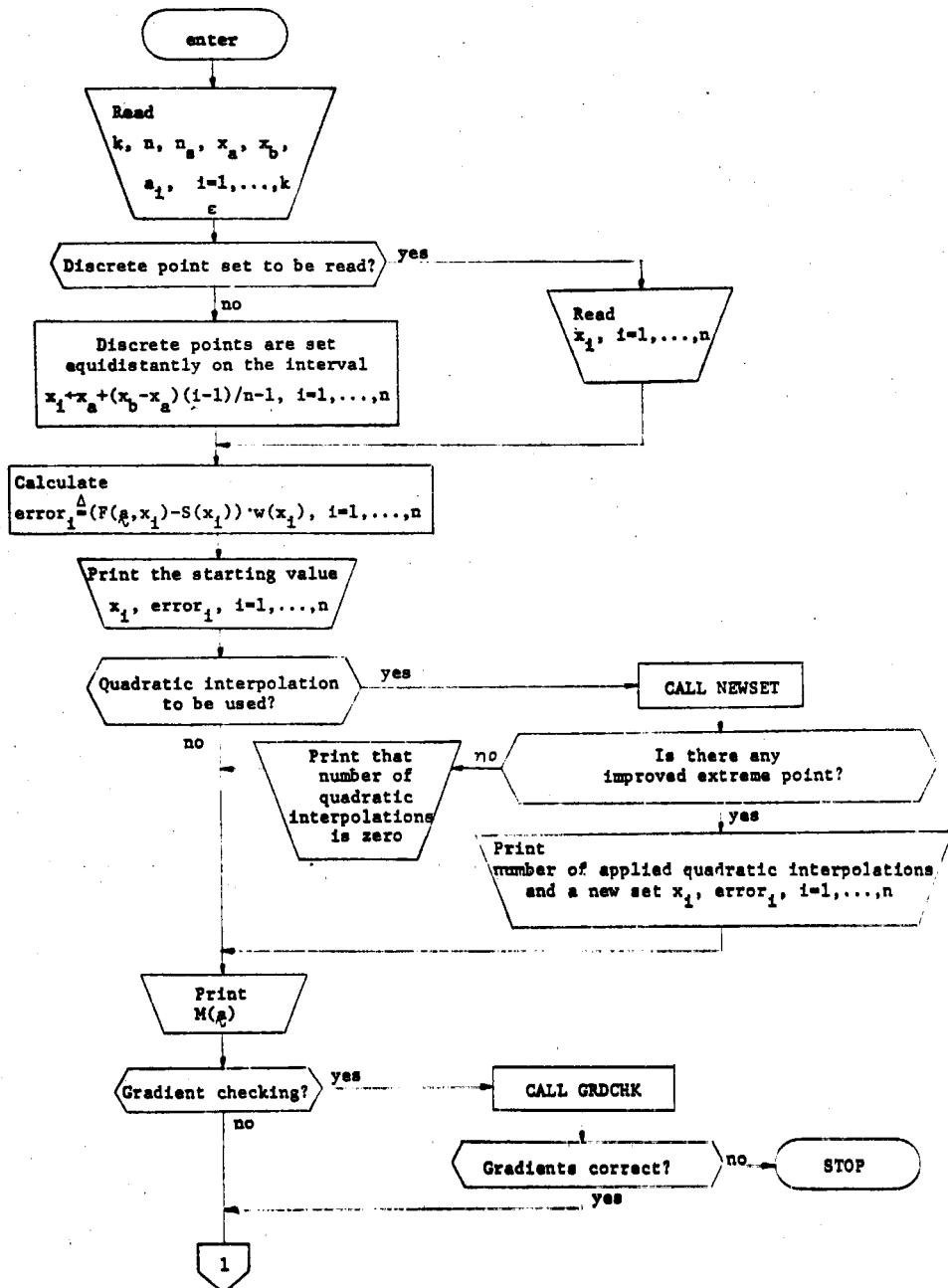


Fig. 5-4 Flowchart of subroutine FUNCT



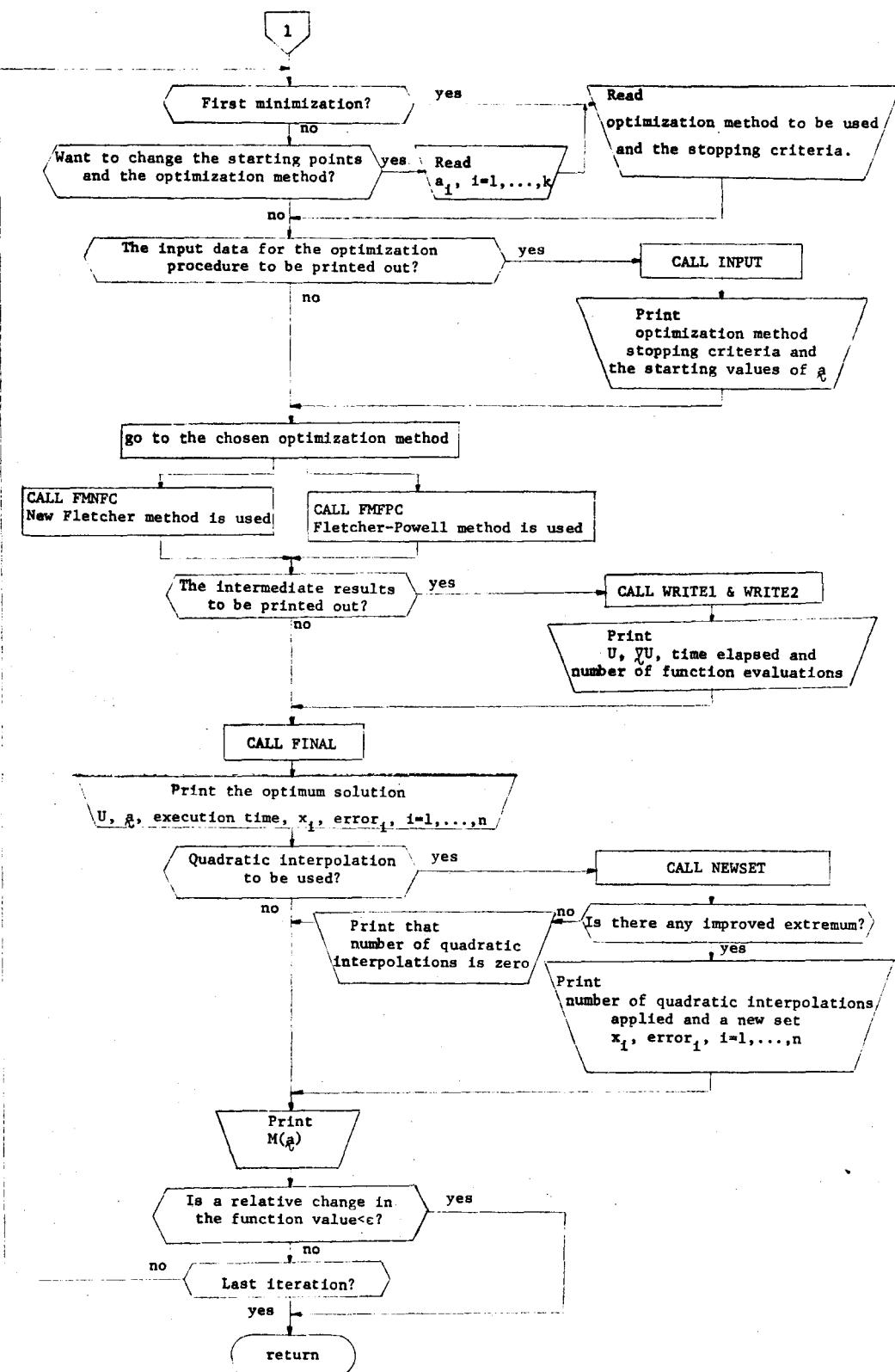


Fig. 5-5 Flowchart of subroutine FMCLP

### 5.4 Examples

As mentioned before, using larger values of  $p$ , the more nearly is the maximum error emphasized and with extremely large values of  $p$  we approach the minimax solution. Two numerical examples, where the rational minimax approximations to the functions are known, were taken for comparison with results obtained using least  $p$ th approximation [22]. The reason why rational approximating functions are chosen in the least  $p$ th objective is that adequate comparison between minimax and least  $p$ th approximation may be done, and there exists a unique rational function

$$(5.16) \quad R_{\ell m}(x) = \frac{P_\ell(x)}{Q_m(x)} = \frac{\sum_{j=0}^{\ell} a_j x^j}{1 + \sum_{j=1}^m b_j x^j}$$

which minimizes the uniform norm. The uniqueness is proved by Chebyshev [30]. In this case the defect has to be included in Chebyshev's theorem. It is known that if the functions form a Chebyshev set, the best approximation to a function over a closed interval is characterized by the maximum error occurring at  $k+1$  points, the sign of error alternating.

The initial approximation and the starting point set in the examples were chosen to be the same as in the minimax approximation problem [22]. The initial approximation was obtained by rational

interpolation where the zeros of the Chebyshev polynomial  $(k+1)$ st order,  $T_{k+1}$  [31], transformed on the interval  $[x_a, x_b]$  are used as supporting points. As a first trial for the point set  $\chi$ , the extrema of the  $(n+1)$ st order Chebyshev polynomial were taken. There is no special reason for choosing these initial points with least pth objective, but it is shown that they provide a good initial guess for the minimax algorithm. A computer program which calculates the starting values of  $\alpha$  and  $\chi$  is given in Appendix IV.

### Test 1

Approximate

$$(5.17) \quad f(x) = e^x - 1.1$$

on the interval  $[-1,1]$  by a rational function (5.16) with  $\ell=1$  and  $m=3$ . Use the weighting function  $w(x)=1$ .

An initial set of 12 points is chosen over an interval  $[-1,1]$  and all the test quantities for the Fletcher-Powell and Fletcher method were set to  $10^{-6}$ . The value of  $p=10^4$  was found to be large enough to bring the least pth close to the minimax solution. The results are presented in Tables 5-1 and 5-2.

From Table 5-1 we may draw the conclusion that good results were obtained with both optimization techniques. Also, it may be noticed that the Fletcher method was more efficient than the

	Starting values	Minimax solution [22]	Least pth solution for p=10 <sup>4</sup>	
			Fletcher-Powell	Fletcher
a <sub>0</sub>	-1.00000x10 <sup>-1</sup>	-9.99928x10 <sup>-2</sup>	-9.99274x10 <sup>-2</sup>	-9.99274x10 <sup>-2</sup>
a <sub>1</sub>	1.04878	1.04887	1.04887	1.04887
b <sub>1</sub>	-4.92007x10 <sup>-1</sup>	-4.91425x10 <sup>-1</sup>	-4.91425x10 <sup>-1</sup>	-4.91425x10 <sup>-1</sup>
b <sub>2</sub>	7.77810x10 <sup>-2</sup>	7.77872x10 <sup>-2</sup>	7.78725x10 <sup>-2</sup>	7.78725x10 <sup>-2</sup>
b <sub>3</sub>	6.04748x10 <sup>-4</sup>	-2.60009x10 <sup>-5</sup>	-2.59776x10 <sup>-5</sup>	-2.60049x10 <sup>-5</sup>
U or M		M=8.68978x10 <sup>-5</sup>	U=8.69170x10 <sup>-5</sup>	U=8.69326x10 <sup>-5</sup>
Number of iterations		4	5	5
Number of function evaluations		-	420	115
Execution time in seconds		-	14	3.5

Table 5-1

Least pth solution for  $p=10^4$

Minimax solution		Fletcher-Powell		Fletcher	
$x$ where the extreme errors occur	extreme errors	working set of $x$ at the optimum	error (a) at the optimum	working set of $x$ at the optimum	error (a) at the optimum
-1.00000	$8.68978 \times 10^{-5}$	-1.00000	$8.68987 \times 10^{-5}$	-1.00000	$8.69123 \times 10^{-5}$
		$-9.59493 \times 10^{-1}$	$3.42715 \times 10^{-5}$	$-9.59493 \times 10^{-1}$	$3.42850 \times 10^{-5}$
$-7.23512 \times 10^{-1}$	$-8.68938 \times 10^{-5}$	$-7.25986 \times 10^{-1}$	$-8.68983 \times 10^{-5}$	$-7.25929 \times 10^{-1}$	$-8.68834 \times 10^{-5}$
		$-7.25825 \times 10^{-1}$	$-8.68983 \times 10^{-5}$	$-7.25829 \times 10^{-1}$	$-8.68833 \times 10^{-5}$
$-1.19232 \times 10^{-1}$	$8.68977 \times 10^{-5}$	$-1.19127 \times 10^{-1}$	$-8.68976 \times 10^{-5}$	$-1.19111 \times 10^{-1}$	$8.69252 \times 10^{-5}$
		$-1.18506 \times 10^{-1}$	$8.68972 \times 10^{-5}$	$-1.18495 \times 10^{-1}$	$8.69248 \times 10^{-5}$
$4.72848 \times 10^{-1}$	$-8.68975 \times 10^{-5}$	$4.73739 \times 10^{-1}$	$-8.69024 \times 10^{-5}$	$4.73912 \times 10^{-1}$	$-8.68853 \times 10^{-5}$
		$4.82411 \times 10^{-1}$	$-8.67694 \times 10^{-5}$	$4.82510 \times 10^{-1}$	$-8.67500 \times 10^{-5}$
$8.65372 \times 10^{-1}$	$8.68972 \times 10^{-5}$	$8.65710 \times 10^{-1}$	$8.69016 \times 10^{-5}$	$8.65783 \times 10^{-1}$	$8.68912 \times 10^{-5}$
		$8.99065 \times 10^{-1}$	$7.93572 \times 10^{-5}$	$8.99070 \times 10^{-1}$	$7.93431 \times 10^{-5}$
		$9.12197 \times 10^{-1}$	$7.15451 \times 10^{-5}$	$9.12197 \times 10^{-1}$	$7.15331 \times 10^{-5}$
1.00000	$-8.68977 \times 10^{-5}$	1.00000	$-8.68953 \times 10^{-5}$	1.00000	$-8.69069 \times 10^{-5}$

Table 5-2

Fletcher-Powell method because it requires fewer function evaluations and is less time-consuming.

### Test 2

Approximate

$$(5.18) \quad f(x) = \frac{\sqrt{(8x-1)^2 + 1} \cdot \tan^{-1}(8x)}{8x}$$

with respect to the weighting function  $w(x)=1$  on the interval  $[-1,1]$  by a rational function with  $\ell=2$  and  $m=2$ .

This example is remarkable because it works near degeneracy defined in the Chebyshev theorem [30], and meets the artificial poles for almost every combination of  $\ell$  and  $m$  in the rational function (5.16) [2]. In true degeneracy of degree one, there is one less extremum of the error curve than is expected with a minimax algorithm, and essentially the numerator and denominator of the rational approximation contain a common linear factor which must be divided out to reduce the rational function to its lowest terms. In the numerical examples it is not at all unusual to find that the numerator has one zero that differs from a denominator zero by some small amount, i.e., to have a rational approximation nearly degenerate of degree one.

The comparative results between the minimax and least pth approximation obtained by the Fletcher optimization technique for  $p=10^4$  are presented in Tables 5-3 and 5-4.

	Starting values	Minimax solution	Least pth solution for $p=10^4$ obtained by the Fletcher method
$a_0$	$1.00000 \times 10^{-2}$	1.41450	1.41448
$a_1$	-3.33600	$-1.06530 \times 10^1$	$-1.06519 \times 10^1$
$a_2$	$4.76782 \times 10^1$	$4.16169 \times 10^1$	$4.16157 \times 10^1$
$b_1$	1.76567	-4.0103	-4.00940
$b_2$	$3.19620 \times 10^1$	$2.82628 \times 10^1$	$2.82620 \times 10^1$
U or M		$M = 2.38113 \times 10^{-2}$	$U = 2.38154 \times 10^{-2}$
Number of iterations		8	13
Number of function evaluations		-	586
Execution time in seconds		-	9.2

Table 5-3

Minimax solution		Least pth solution for $p=10^4$ obtained by the Fletcher method	
x where the extreme errors occur	extreme errors	working set of x at the optimum	errors at the optimum
-1.00000	$-2.38108 \times 10^{-2}$	-1.00000	$-2.38025 \times 10^{-2}$
		$-9.59493 \times 10^{-1}$	$-2.15823 \times 10^{-2}$
		$-4.34766 \times 10^{-1}$	$1.77356 \times 10^{-2}$
$-3.08573 \times 10^{-1}$	$2.38108 \times 10^{-2}$	$-3.15535 \times 10^{-1}$	$2.38026 \times 10^{-2}$
$-6.15510 \times 10^{-2}$	$-2.38109 \times 10^{-2}$	$-6.18947 \times 10^{-2}$	$-2.38007 \times 10^{-2}$
$5.42007 \times 10^{-2}$	$2.38108 \times 10^{-2}$	$5.46342 \times 10^{-2}$	$2.37994 \times 10^{-2}$
		$1.93965 \times 10^{-1}$	$-2.37976 \times 10^{-2}$
$1.96670 \times 10^{-1}$	$-2.38113 \times 10^{-2}$	$1.94011 \times 10^{-1}$	$-2.37976 \times 10^{-2}$
$5.52892 \times 10^{-1}$	$2.38108 \times 10^{-2}$	$5.52292 \times 10^{-1}$	$2.38004 \times 10^{-2}$
		$6.20236 \times 10^{-1}$	$2.27473 \times 10^{-2}$
		$8.09651 \times 10^{-1}$	$1.42303 \times 10^{-2}$
		1.00000	$3.68693 \times 10^{-3}$

Table 5-4

Test 3

There was an attempt by Rice to extend the Chebyshev theorem to include non-linear dependence of  $F(\alpha, x)$  on  $\alpha$  that the same characterization of best approximations holds for all continuous  $S(x)$  [32], but this applies only to a limited class of approximating functions [20]. Unfortunately, there are many useful choices of  $F(\alpha, x)$  which do not fulfil Rice's conditions of best approximations [33].

The third example is another that might be expected to give trouble. Due to Curtis and Powell [33], it is the approximation of  $x^2$  by  $a_1x + a_2e^x$  over  $0 \leq x \leq 2$ . It may be verified that the error function of the approximation

$$(5.19) \quad x^2 \approx 8.465x - 2.0239e^x$$

takes its maximum absolute value at the three points  $x=0$ ,  $x=1.1227$  and  $x=2$ , the error at these points being  $+2.0239$ ,  $-2.0239$  and  $+2.0239$ , respectively. In fact, the best approximation is

$$(5.20) \quad x^2 \approx 0.1842x + 0.4186e^x$$

the maximum absolute error is 0.5382 and this error occurs at just the two points:  $x=0.4064$  and  $x=2$ . Not only do the approximating functions fail to form a Chebyshev set, but also the error curve has only two extrema instead of the three that would normally be anticipated according to the Rice's theorem. The least pth results for this problem are given in Table 5-5 and again show the success of

FMCLP. The estimates of the best approximation agree to four figures with those given by Curtis and Powell.

	Initial approximation	$\begin{bmatrix} 1 \\ 1 \end{bmatrix}$
	Fletcher method:	
	Stopping criteria	$10^{-6}$
n		10
p		$10^5$
a		0.1848
		0.4184
M(a)		0.5382
xM		0.4066
		2.000
U(a)		0.5382
function evaluation		67
execution time in seconds		1.3

Table 5-5

#### Test 4

The computer program is used to solve an electrical engineering design problem. The problem is stated: find a second order model of a fourth-order system with a given transfer function

$$(5.21) \quad G(s) = \frac{s+4}{(s+1)(s^2+4s+8)(s+5)} .$$

The transfer function of the second-order model considered is

$$(5.22) \quad H(s) = \frac{a_3}{(s+a_1)^2 + a_2^2} .$$

Using the inverse Laplace transform the responses for (5.21) and (5.22) are

$$(5.23) \quad s(t) = \frac{3}{20} e^{-t} + \frac{1}{52} e^{-5t} - \frac{1}{65} e^{-2t} (3 \sin 2t + 11 \cos 2t)$$

and

$$(5.24) \quad F(a, t) = \frac{a_3}{a_2} e^{-a_1 t} \cdot \sin a_2 t ,$$

respectively, where

$$\vec{a} = \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix} .$$

The results for different values of  $p$  and different numbers of sampling points  $n$ , but the same product  $n \times n_s$ , over the range  $0 \leq t \leq 10$  are given in Table 5-6, agreeing with those given in [28].

Using the quadratic interpolation the locations of the extreme points were found precisely, and the solutions are closer to the minimax solution. Both solutions, for  $n=10$  and  $n=25$ , are better than for  $n=50$  where the quadratic interpolation was not employed. Moreover, the case  $n=10$  is less time consuming, because this sampling takes the least number of points for the objective function (4.3) in comparison with the other two cases under consideration.

	$n=10, n_s = 5$	$n=25, n_s = 2$	$n=50, n_s = 1$
p=2	$a_1 = 1.27339$ $a_2 = 6.54190 \times 10^{-1}$ $a_3 = 2.17787 \times 10^{-1}$ $M(a) = 2.05859 \times 10^{-2}$ $t_M = 2.51460 \times 10^{-1}$ $U(a) = 4.71528 \times 10^{-3}$ f.e.=41 q.i.=2 1.6 sec	$a_1 = 1.05489$ $a_2 = -7.67814 \times 10^{-1}$ $a_3 = 1.618192 \times 10^{-1}$ $M(a) = 1.32708 \times 10^{-2}$ $t_M = 2.46234 \times 10^{-1}$ $U(a) = 1.21662 \times 10^{-2}$ f.e.=30 q.i.=3 2.3 sec	$a_1 = 1.01687$ $a_2 = 7.8915 \times 10^{-1}$ $a_3 = 1.61435 \times 10^{-1}$ $M(a) = 1.28696 \times 10^{-2}$ $t_M = 2.04081 \times 10^{-1}$ $U(a) = 2.06679 \times 10^{-2}$ f.e.=36 q.i.=0 5.1 sec
p=10	$a_1 = 7.37873 \times 10^{-1}$ $a_2 = 9.2622 \times 10^{-1}$ $a_3 = 1.2862 \times 10^{-1}$ $M(a) = 8.9565 \times 10^{-3}$ $t_M = 1.67727 \times 10^{-1}$ $U(a) = 8.4364 \times 10^{-3}$ f.e.=38 q.i.=3 1.5 sec	$a_1 = 7.46289 \times 10^{-1}$ $a_2 = -9.23825 \times 10^{-1}$ $a_3 = 1.27596 \times 10^{-1}$ $M(a) = 8.76150 \times 10^{-3}$ $t_M = 1.73062 \times 10^{-1}$ $U(a) = 8.2421 \times 10^{-3}$ f.e.=32 q.i.=3 2.5 sec	$a_1 = 7.43325 \times 10^{-1}$ $a_2 = 9.29377 \times 10^{-1}$ $a_3 = 1.2812 \times 10^{-1}$ $M(a) = 8.5446 \times 10^{-3}$ $t_M = 2.04081 \times 10^{-1}$ $U(a) = 9.1834 \times 10^{-3}$ f.e.=31 q.i.=0 4.5 sec

	$n=10, n_s = 5$	$n=25, n_s = 2$	$n=50, n_s = 1$
$p=10^2$	$a_1 = 6.79369 \times 10^{-1}$ $a_2 = 9.55429 \times 10^{-1}$ $a_3 = 1.21997 \times 10^{-1}$ $M(\tilde{a}) = 8.9563 \times 10^{-3}$ $t_M = 1.67727 \times 10^{-1}$ $U(\tilde{a}) = 8.2055 \times 10^{-3}$ f.e.=35 q.i.=1 1.3 sec	$a_1 = 6.8047 \times 10^{-1}$ $a_2 = -9.5471 \times 10^{-1}$ $a_3 = 1.2206 \times 10^{-1}$ $M(\tilde{a}) = 8.1300 \times 10^{-3}$ $t_M(\tilde{a}) = 1.73062 \times 10^{-1}$ $U(\tilde{a}) = 8.1866 \times 10^{-3}$ f.e.=35 q.i.=1 2.5 sec	$a_1 = 6.88905 \times 10^{-1}$ $a_2 = 9.52106 \times 10^{-1}$ $a_3 = 1.2334 \times 10^{-1}$ $M(\tilde{a}) = 7.9450 \times 10^{-3}$ $t_M = 2.04081 \times 10^{-1}$ $U(\tilde{a}) = 8.0045 \times 10^{-3}$ f.e.=35 q.i.=0 5. sec
$p=10^3$	$a_1 = 6.737 \times 10^{-1}$ $a_2 = 9.5590 \times 10^{-1}$ $a_3 = 1.2168 \times 10^{-1}$ $M(\tilde{a}) = 8.1125 \times 10^{-3}$ $t_M = 1.67727 \times 10^{-1}$ $U(\tilde{a}) = 8.1182 \times 10^{-3}$ f.e.=28 q.i.=0 1. sec	$a_1 = 6.77142 \times 10^{-1}$ $a_2 = -9.5556 \times 10^{-1}$ $a_3 = 1.21735 \times 10^{-1}$ $M(\tilde{a}) = 8.0905 \times 10^{-3}$ $t_M = 1.73062 \times 10^{-1}$ $U(\tilde{a}) = 8.0957 \times 10^{-3}$ f.e.=29 q.i.=0 2.3 sec	$a_1 = 6.8510 \times 10^{-1}$ $a_2 = 9.5289 \times 10^{-1}$ $a_3 = 1.2294 \times 10^{-1}$ $M(\tilde{a}) = 7.9009 \times 10^{-3}$ $t_M = 2.04082 \times 10^{-1}$ $U(\tilde{a}) = 7.9068 \times 10^{-3}$ f.e.=26 q.i.=0 4. sec
Total	f.e.=142 for 5.4 sec	f.e.=126 for 9.6 sec	f.e.=128 for 18.6 sec

Table 5-6

Remark: The optimum result for  $a_2$  is true for both positive and negative values from Table 5-6 because this does not effect the approximating function (5.24) where

$$\frac{\sin a_2 t}{a_2} = \frac{\sin(-a_2 t)}{-a_2}$$

## CHAPTER VI

### Generalized Least pth Approximation

#### 6.1 The objective functions

The applicability of least pth approximation may be extended to a wider variety of design problems and a wider range of specifications than considered in Chapter V if the least pth objective function is written in a more general form [ 2 ].

We will evaluate all the functions at a finite discrete set of values of  $x$  taken from one or more closed intervals.

Define real weighted error functions related to the upper and lower specifications, as shown in Fig. 6-1, respectively as

$$e_u(a, x_i) = w_u(x_i)(F(a, x_i) - s_u(x_i)), \quad i \in I_u \quad (6.1)$$

$$e_\ell(a, x_i) = w_\ell(x_i)(F(a, x_i) - s_\ell(x_i)), \quad i \in I_\ell$$

where

- $F(a, x_i)$  is the approximating function;
- $s_u(x_i)$  is an upper specified function;
- $s_\ell(x_i)$  is an lower specified function;
- $w_u(x_i)$  is an upper positive weighting function;
- $w_\ell(x_i)$  is an lower positive weighting function;
- $a$  is an vector containing the  $k$  independent parameters;

$x_i$  is the discretized independent variable;  
 $I_u$  and  $I_\ell$  are appropriate index sets.

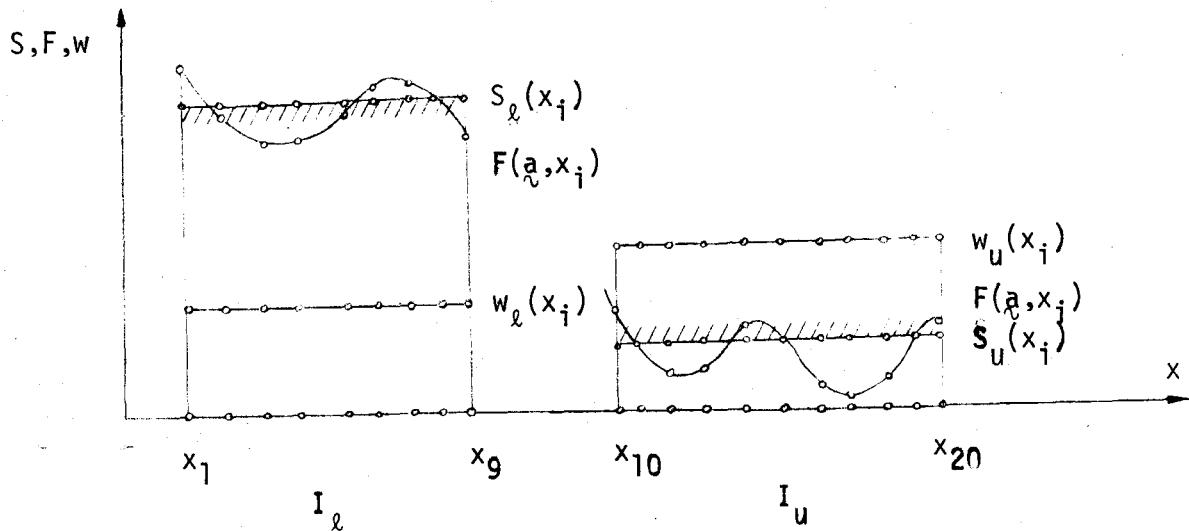


Fig. 6-1 Example of a design problem with  
upper and lower specifications.

A special case of (6.1) when  $S_u = S_\ell$  and  $w_u = w_\ell$  leads to the common form of a real weighted error function defined in (5.1).

We will slightly modify the error functions by introducing a constant  $\xi$  which is used for shifting the level of the originally defined errors functions (6.1) [34] such that

$$e_u'(a, x_i, \xi) \stackrel{\Delta}{=} e_u(a, x_i) - \xi$$

(6.2)

$$e_\ell'(a, x_i, \xi) \stackrel{\Delta}{=} e_\ell(a, x_i) + \xi .$$

The above expressions may be obtained by redefining the

specifications which are going to be used in the error function (6.2) as follows

$$S_u'(x_i, \xi) = S_u(x_i) + \frac{\xi}{w_u(x_i)}$$

(6.3)

$$S_l'(x_i, \xi) = S_l(x_i) - \frac{\xi}{w_l(x_i)}$$

where  $S_u'(x_i, \xi)$  and  $S_l'(x_i, \xi)$  are the artificial upper and lower specified functions, respectively.

Therefore, the modified error functions are in the following forms

$$\begin{aligned} e_u'(a, x_i, \xi) &\stackrel{\Delta}{=} w_u(x_i)(F(a, x_i) - S_u'(x_i, \xi)) \\ (6.4) \quad e_l'(a, x_i, \xi) &\stackrel{\Delta}{=} w_l(x_i)(F(a, x_i) - S_l'(x_i, \xi)). \end{aligned}$$

A certain flexibility in the formulation of the optimization problem is possible when the artificial margin  $\xi$  is introduced. A more detailed discussion will come later.

Let

$$\begin{aligned} e_{u_i}(a, \xi) &\stackrel{\Delta}{=} e_u(a, x_i, \xi), \quad i \in I_u \\ (6.5) \quad e_{l_i}(a, \xi) &\stackrel{\Delta}{=} e_l(a, x_i, \xi), \quad i \in I_l. \end{aligned}$$

It is quite possible that some of the upper and lower error functions tend to  $-\infty$  and  $+\infty$ , respectively, in which case they may be

simply ignored by a proper formulation of the problem.

We will consider two possible cases which may happen in function approximation with upper and lower specifications: the specification may be violated and the specification may be satisfied.

### 6.2 Case I-Specification Violated

When the specification is violated, some of the  $e_{u_i}(x)$  or  $-e_{l_i}(x)$  are positive, as is illustrated in Fig. 6-2. (The illustration is made for  $\xi=0$ , for convenience).

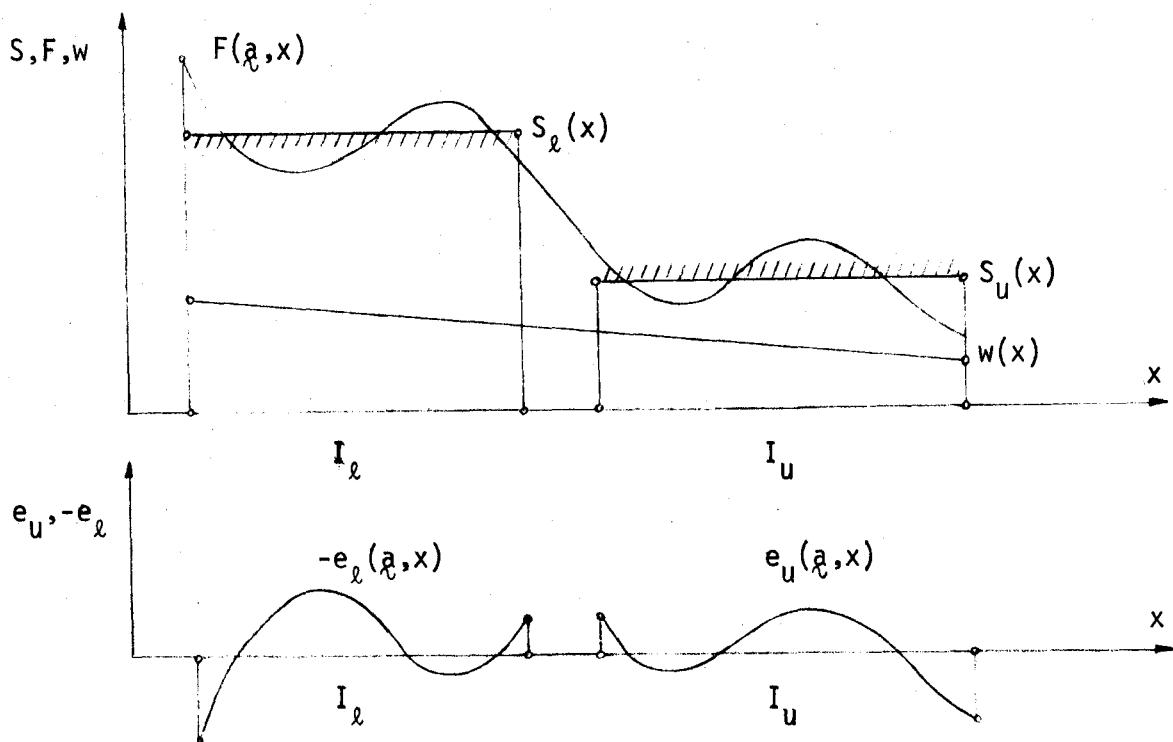


Fig. 6-2 Example of a design problem when the specification is violated.  
Case I is applicable.

We wish to make as small as possible all those  $e_{u_i}(\tilde{a})$  and  $-e_{\ell_i}(\tilde{a})$  which are positive, in an effort to meet the specification, and the following objective function was proposed [34] to be minimized:

$$(6.6) \quad U(\tilde{a}, \xi) = \left( \sum_{i \in J_u(\tilde{a}, \xi)} [e_{u_i}(\tilde{a}, \xi)]^p + \sum_{i \in J_\ell(\tilde{a}, \xi)} [-e_{\ell_i}(\tilde{a}, \xi)]^p \right)^{\frac{1}{p}}$$

where

$$J_u(\tilde{a}, \xi) \stackrel{\Delta}{=} \{i \mid e_{u_i}(\tilde{a}, \xi) \geq 0, \quad i \in I_u\}$$

(6.7)

$$J_\ell(\tilde{a}, \xi) \stackrel{\Delta}{=} \{i \mid -e_{\ell_i}(\tilde{a}, \xi) \geq 0, \quad i \in I_\ell\}$$

and

$$p > 1.$$

If  $J_u$  and  $J_\ell$  are empty then  $U(\tilde{a})=0$  and we have just met or exceeded the artificial specifications. The larger the value of  $p$ , the more nearly would we expect the maximum error to be emphasized, since

$$(6.8) \quad \max_{i,j} [e_{u_i}(\tilde{a}, \xi), -e_{\ell_j}(\tilde{a}, \xi)] = \lim_{p \rightarrow \infty} U(\tilde{a}, \xi), \quad i \in J_u(\tilde{a}, \xi), \quad j \in J_\ell(\tilde{a}, \xi),$$

hence minimization of (6.6) is an effort to meet the specification.

### 6.3 Case II-Specification is satisfied

When the specification is satisfied all the  $-e_{u_i}(\tilde{a}, \xi)$  and  $e_{\ell_i}(\tilde{a}, \xi)$  will be positive as it is shown in Fig. 6-3. Again for

convenience,  $\xi=0$  is chosen,

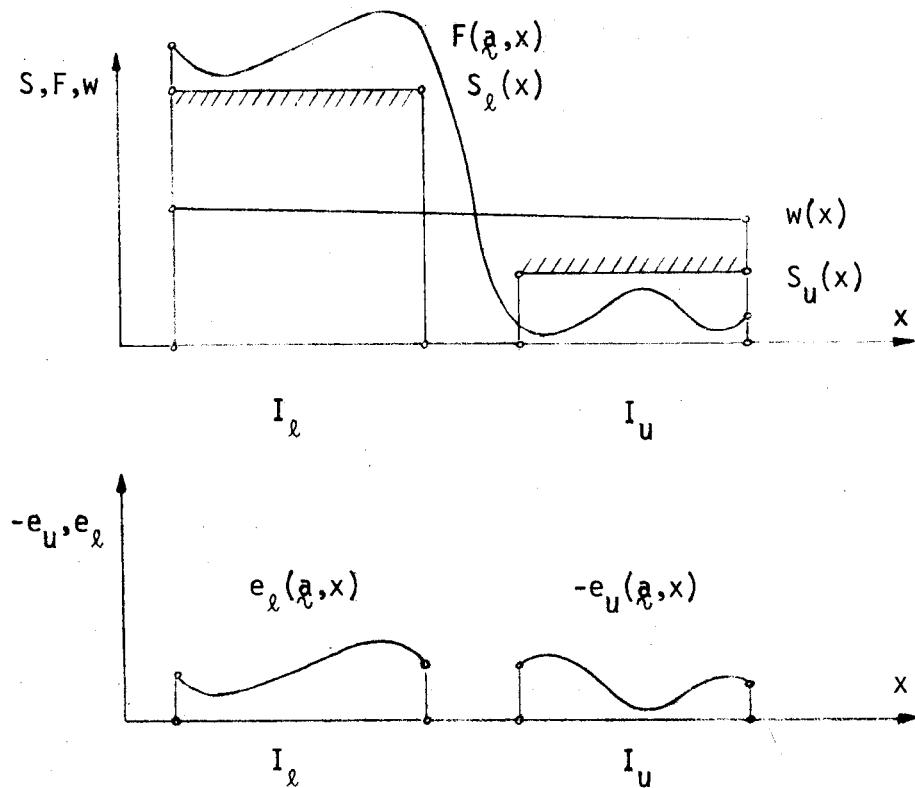


Fig. 6-3 Example of a design problem when the specification is satisfied. Case II is applicable.

Now, in an effort to exceed the specification by the greatest amount, the following objective function was proposed [34] to be minimized

$$(6.9) \quad U(a, \xi) = -\left\{ \sum_i I_u [-e_{u,i}(a, \xi)]^{-p} + \sum_i I_l [e_{l,i}(a, \xi)]^{-p} \right\}^{-\frac{1}{p}}$$

for

$$(6.10) \quad -e_{u_i}'(\tilde{a}, \xi) > 0, \quad i \in I_u$$

$$e_{\ell_j}'(\tilde{a}, \xi) > 0, \quad j \in I_\ell$$

and

$$p \geq 1.$$

Again, the larger the value of  $p$  the more nearly would we expect the minimum error to be emphasized, since

$$(6.11) \quad \max_{i,j} [e_{u_i}'(\tilde{a}, \xi), -e_{\ell_j}'(\tilde{a}, \xi)] = \lim_{p \rightarrow \infty} U(\tilde{a}, \xi), \quad \begin{matrix} i \in I_u \\ j \in I_\ell \end{matrix}$$

#### 6.4 Scaling

If  $p$  is very large we have ill-conditioning in both cases resulting from the numerical evaluation of  $[e_{u_i}']^{\pm p}$  and  $[e_{\ell_j}']^{\pm p}$ . By proper scaling we can alleviate the ill-conditioning and also define only one objective function which is valid for both cases [34].

Let

$$(6.12) \quad M(\tilde{a}, \xi) \triangleq \max_{i,j} [e_{u_i}'(\tilde{a}, \xi), -e_{\ell_j}'(\tilde{a}, \xi)], \quad \begin{matrix} i \in I_u \\ j \in I_\ell \end{matrix}$$

and define an objective function

$$(6.13) \quad U(\tilde{a}, \xi) = M(\tilde{a}, \xi) \left( \sum_{i \in K_u} \left[ \frac{e_{u_i}'(\tilde{a}, \xi)}{M(\tilde{a}, \xi)} \right]^q + \sum_{j \in K_\ell} \left[ \frac{-e_{\ell_j}'(\tilde{a}, \xi)}{M(\tilde{a}, \xi)} \right]^q \right)^{\frac{1}{q}}$$

and its gradients

$$(6.14) \quad \bar{y}U(\tilde{a}, \xi) = \left( \sum_{i \in K_u} \left[ \frac{e_{ui}'(\tilde{a}, \xi)}{M(\tilde{a}, \xi)} \right]^q + \sum_{i \in K_\ell} \left[ \frac{-e_{\ell i}'(\tilde{a}, \xi)}{M(\tilde{a}, \xi)} \right]^q \right)^{\frac{1}{q}} -$$

$$\cdot \left( \sum_{i \in K_u} \left[ \frac{e_{ui}'(\tilde{a}, \xi)}{M(\tilde{a}, \xi)} \right]^{q-1} \bar{y} e_{ui}'(\tilde{a}, \xi) \right.$$

$$\left. - \sum_{i \in K_\ell} \left[ \frac{-e_{\ell i}'(\tilde{a}, \xi)}{M(\tilde{a}, \xi)} \right]^{q-1} \bar{y} e_{\ell i}'(\tilde{a}, \xi) \right)$$

where

$$(6.15) \quad \begin{aligned} K_u &\stackrel{\Delta}{=} \begin{cases} J_u(\tilde{a}, \xi) & \text{if } M(\tilde{a}, \xi) > 0 \\ I_u & \text{if } M(\tilde{a}, \xi) < 0 \end{cases}, \\ K_\ell &\stackrel{\Delta}{=} \begin{cases} J_\ell(\tilde{a}, \xi) & \text{if } M(\tilde{a}, \xi) > 0 \\ I_\ell & \text{if } M(\tilde{a}, \xi) < 0 \end{cases}, \end{aligned}$$

with  $J_u(\tilde{a}, \xi)$  and  $J_\ell(\tilde{a}, \xi)$  defined in (6.7) and  $I_u$  and  $I_\ell$  defined in (6.10), and

$$(6.16) \quad q \stackrel{\Delta}{=} \begin{cases} p > 1 & \text{for } M(\tilde{a}, \xi) > 0, \\ p \geq 1 & \text{for } M(\tilde{a}, \xi) < 0. \end{cases}$$

If  $e_{ui}'(\tilde{a}, \xi)$  and  $e_{\ell i}'(\tilde{a}, \xi)$  for  $i \in I_u$  and  $i \in I_\ell$ , respectively, are continuous with continuous partial derivatives, the proposed objective function is continuous everywhere with continuous partial derivatives. The objective function (6.13) and partial derivatives (6.14) still remain continuous even when, for some  $i$ 's,  $e_{ui}'$  and/or  $e_{\ell i}'$  are discontinuous or continuous with discontinuous derivatives, simply because those points are ignored if  $e_{ui}'/M$  and/or  $-e_{\ell i}'/M$  are negative and  $M > 0$ . This is very suitable for the wide variety of network and system design problems, especially in filter design.

### 6.5 Discussion

The artificial margin  $\xi$  which is constant during optimization does not affect the location of the minimax optimum ( $p \rightarrow \infty$ ). Its important role, however, is evident for a finite value of  $p$ . The value of the parameter  $\xi$  can be chosen so that the  $M$  of (6.12) is always positive or negative during optimization. When  $M$  is positive, only sample points which satisfy the conditions in (6.7) are considered and, therefore, there is a saving in gradient computation. But in this case it may happen that  $M=0$ , when the function is continuous but the derivatives are discontinuous. On the rare occasions when this situation causes a failure of the gradient minimization algorithm, one can change the value of  $\xi$  and restart the optimization process. If the value of  $M$  is chosen to be negative this possible failure is avoided.

### 6.6 Computer program

A computer program in FORTRAN IV has been written, which utilizes the least  $p$ th approach described in previous sections. The program consists of 15 separate subprograms.

FMLPO Supplies all the data for the optimization process and coordinates the other subprograms in the package;

FUNCS Defines upper and lower specified functions;

FCTAPP Defines an approximating function and its gradients

with respect to variable parameters;

- W      Defines upper and lower weighting functions;
- FCT     Calculates artificial upper and lower specified functions according to the equation (6.3);
- EPSNP   Calculates upper and lower weighted error functions;
- ERRO    Selects the weighted error functions of interest for the objective function (6.13) according to the equations given in (6.15);
- FUNCG   Computes the generalized least pth objective function (6.13) and its gradients (6.14) w.r.t. variable parameters;

GRDGHK, FMFPG and FMNFG are subroutine subprograms which have the same role as the subroutines GRDCHK, FMFPC and FMNFC, respectively, in the FMCLP package.

INPUT, FINAL, WRITE1 and WRITE2 are subroutine subprograms which have already been introduced in section 5.3.

FUNCS, W, FCT and EPSNP are function subprograms, and the others are subroutine subprograms.

A user of the package is supposed to write the following subprograms by himself: FUNCS, FCTAPP and W in a straightforward way. More detailed instructions may be found in Appendix V. A user has to arrange the intervals, not necessarily disjoint, such

that each of them has only one specification. For example, if the original design problem has upper and lower specifications for the same values of the independent parameter  $x$ , two intervals with a single specification have to be formed, one with the upper and the other with the lower specification. A two-dimensional array is constructed of the input data, where the first of the two subscripts relates to the type of the specification and the second denotes the appropriate intervals. The elements of the array may have values +1 if there is an upper or -1 if there is a lower specification in a given interval.

Let

$$(6.17) \quad \begin{aligned} e_i''(a, \xi) &\stackrel{\Delta}{=} (+1) \cdot e_u'(a, \xi), \quad i \in I_u \\ e_i''(a, \xi) &\stackrel{\Delta}{=} (-1) \cdot e_l'(a, \xi), \quad i \in I_l \end{aligned}$$

where  $I_u \cap I_l = \emptyset$ .

Then the expressions for the generalized least pth approximation (6.13)-(6.16) may be redefined in a form suitable for programming:

$$(6.18) \quad M(a, \xi) \stackrel{\Delta}{=} \max_i [e_i''(a, \xi)], \quad i \in I \stackrel{\Delta}{=} I_u \cup I_l$$

$$(6.19) \quad U(a, \xi) = M(a, \xi) \left( \sum_{i \in K} \left[ \frac{e_i''(a, \xi)}{M(a, \xi)} \right]^q \right)^{\frac{1}{q}}$$

$$(6.20) \quad \bar{U}(a, \xi) = \left( \sum_{i \in K} \left[ \frac{e_i''(a, \xi)}{M(a, \xi)} \right]^q \right)^{\frac{1}{q}-1} \cdot \left( \sum_{i \in K} \left[ \frac{e_i''(a, \xi)}{M(a, \xi)} \right]^{q-1} \right) \bar{e}_i''(a, \xi)$$

where

$$(6.21) \quad K = \begin{cases} J(a, \xi) & \text{if } M(a, \xi) > 0 \\ I & \text{if } M(a, \xi) < 0 \end{cases}$$

with  $J = J_u^{\Delta}(x, \xi) \cup J_l^{\Delta}(x, \xi)$ , and  $q$  the same as in (6.16).

Function subprogram EPSNP computes the upper or lower weighted error at a single point  $x$  for a current variable vector  $\underline{x}$ . A flowchart of function EPSNP is shown in Fig. 6-5.

The flowcharts of the subroutines ERRO and FUNCG are available in Fig. 6-6 and Fig. 6-7, respectively. The list of the symbols used in the flowcharts are given on the page 67. A scheme presented in Fig. 6-4 illustrates how the subprograms are connected to each other. The flowchart of FMLPO is given in Fig. 6-8.

If the gradients of the approximating function are not supplied correctly, the program will terminate and print out the appropriate message. Also, suitable diagnostic messages are printed out whenever there is any unusual exit.

When the criteria for the optimum have been satisfied, the optimum solution, the discrete point set with the artificial weighted errors in each interval and the execution time in seconds are printed out.

The computer program of the whole package FMLPO is given in Appendices III and VI.

Most of the symbols used in the flowchart are in agreement with the corresponding ones used in the theory. Those which differ or have not been mentioned yet are

$e_i$  corresponds to  $e_i(a, \xi)$ ;

$e_{\max}$  corresponds to  $M(a, \xi)$ ;

$n_I$  is a total number of intervals;

$j_I$  is a current interval;

$k_j$  is a number of discrete point set at  $(j_I)$ th interval;

$x_{1,j}$  is a left end point of  $(j_I)$ th interval

$x_{2,j}$  is a right end point of  $(j_I)$ th interval;

$x_{3,j}$  is a characteristic number of  $(j_I)$ th interval with the information about upper or lower specification;

$m_j$  is an integer which determinates the interval where the selected point belongs;

$\text{grad}_i$  is the gradient of the approximating function  $\nabla F(a, x_i)$ ;

$\epsilon$  is a small quantity used as a stopping criterion of the iterative procedure in FMLPO.

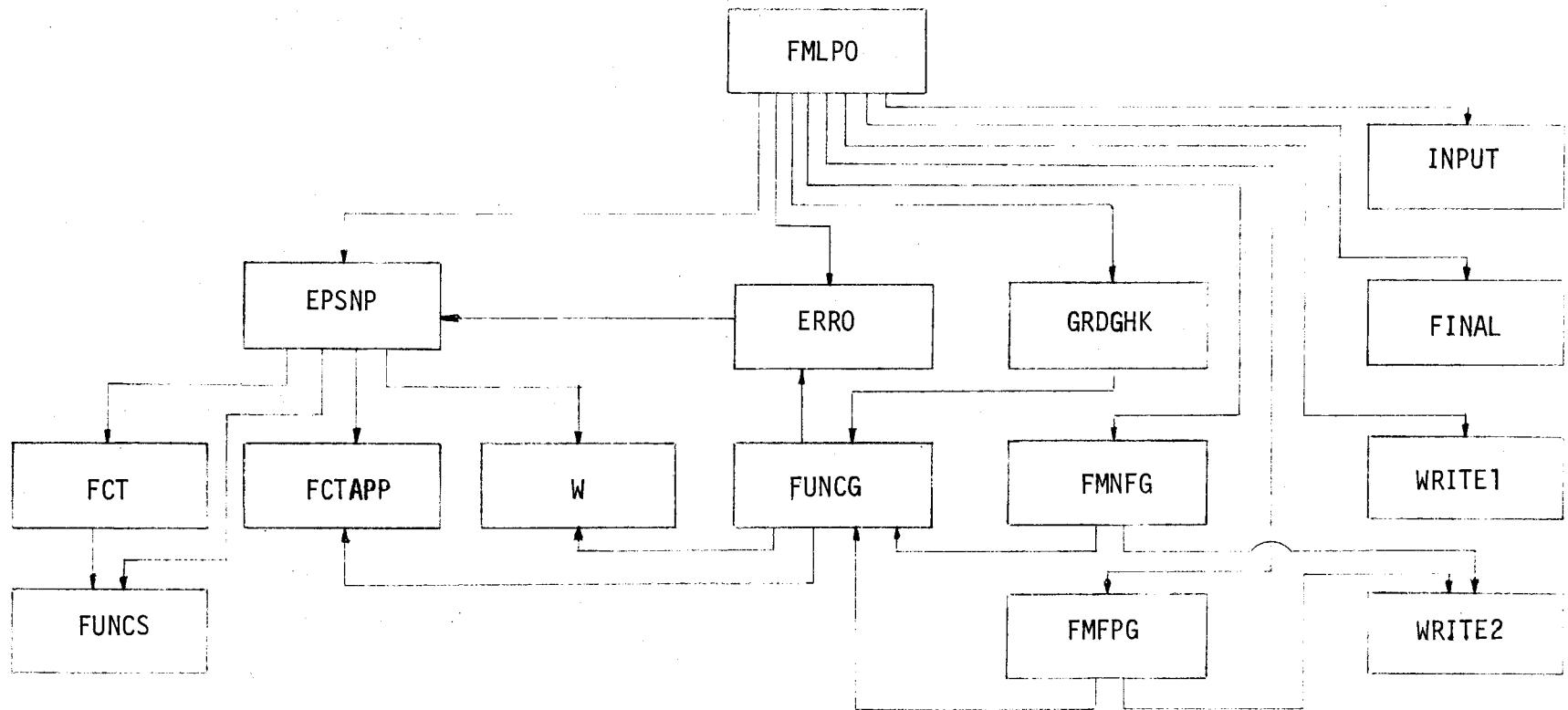


Fig. 6-4 The organization of FMLPO

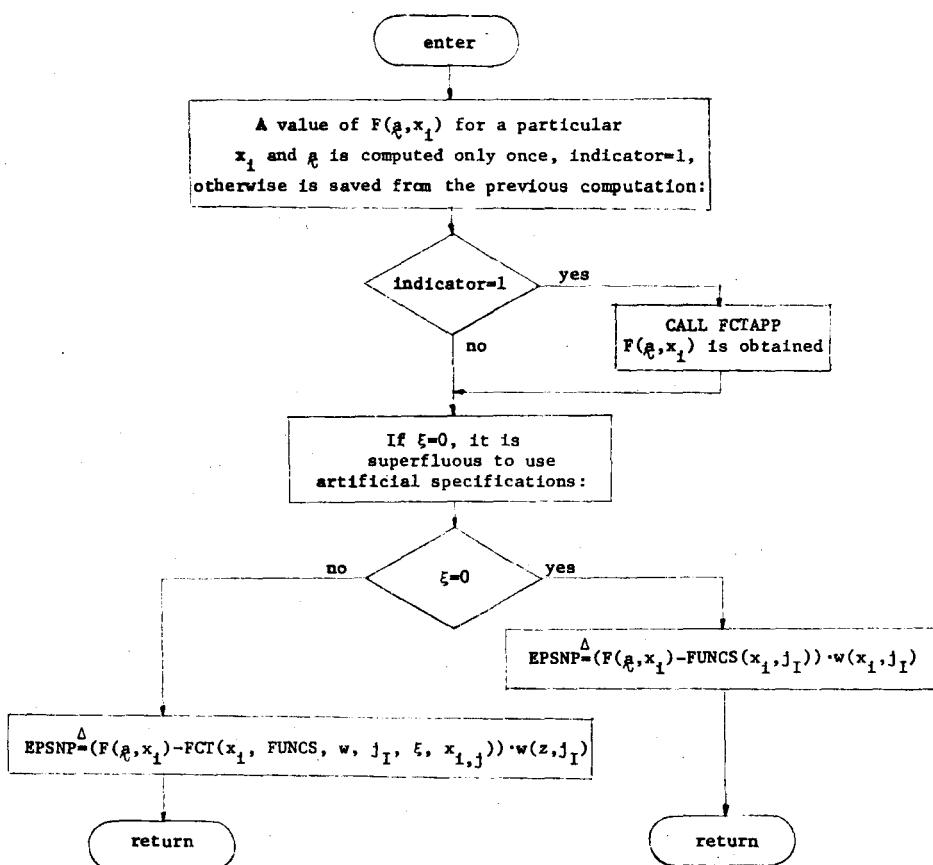
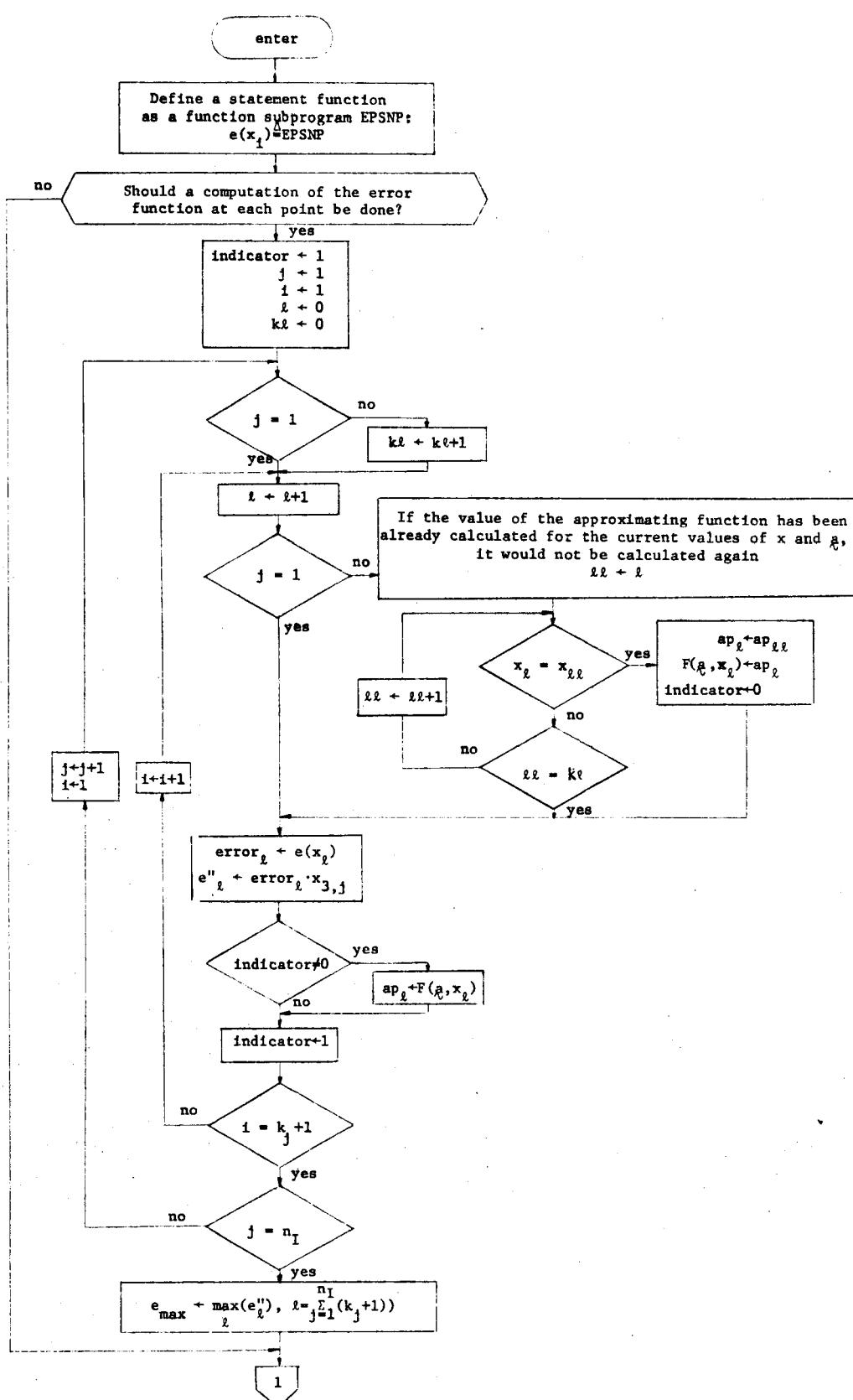


Fig. 6-5 Flowchart of function subprogram EPSNP



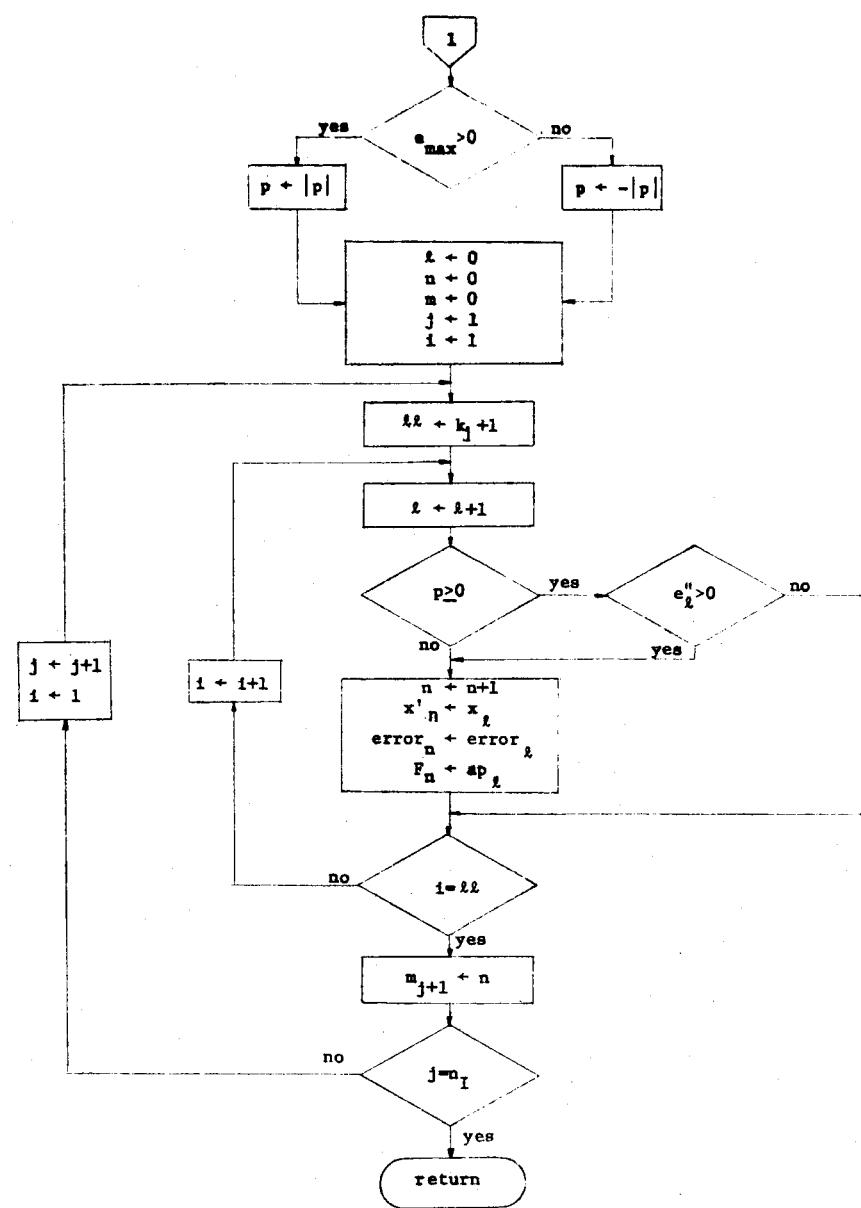
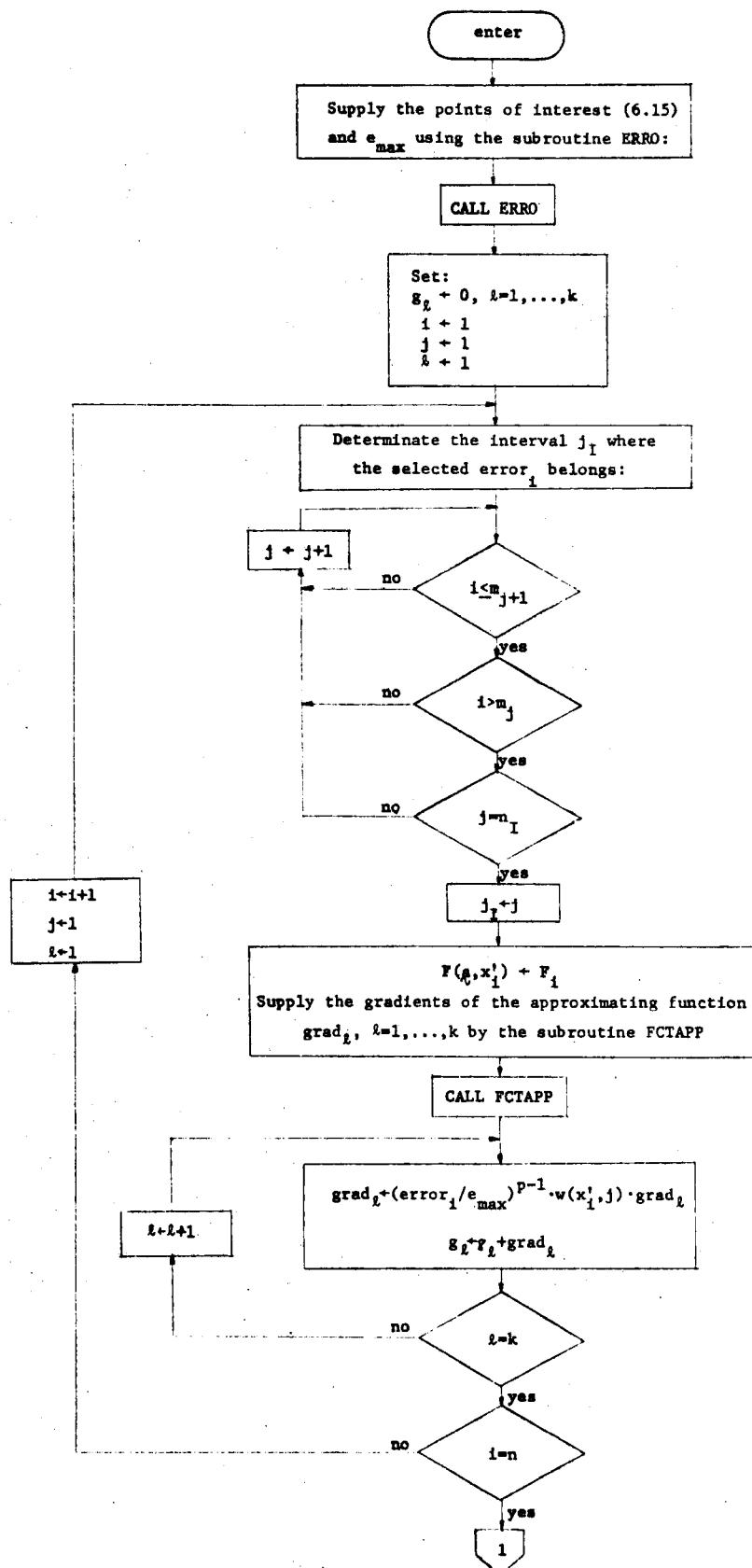


Fig. 6-6 Flowchart of subroutine ERRO



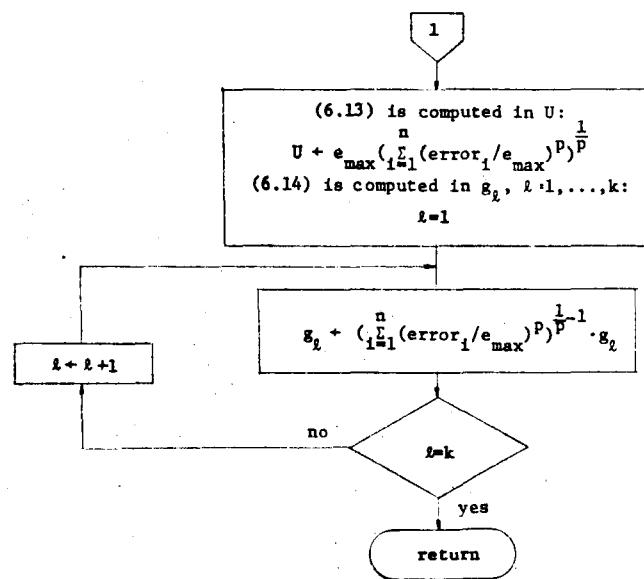
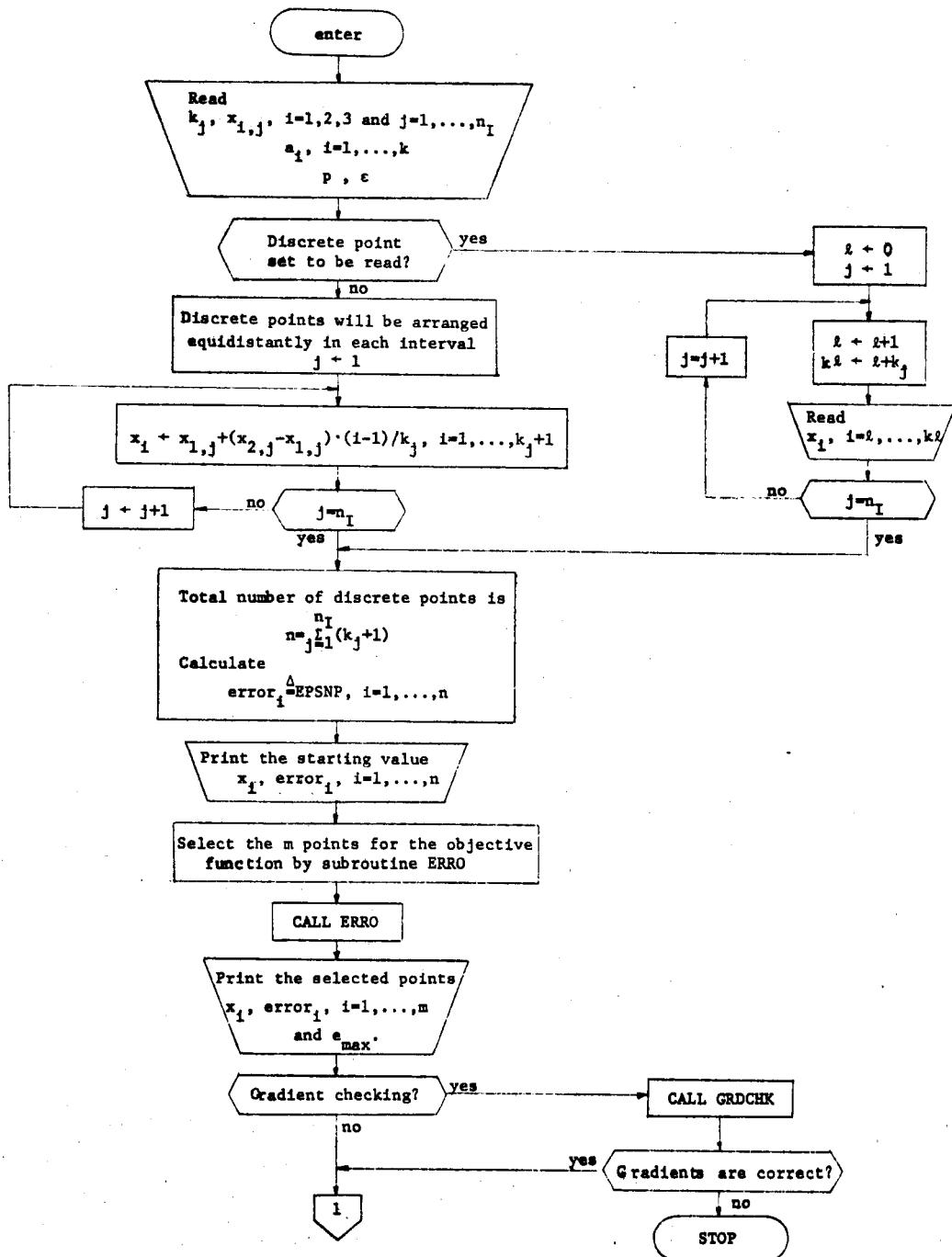


Fig. 6-7 Flowchart of subroutine FUNCg



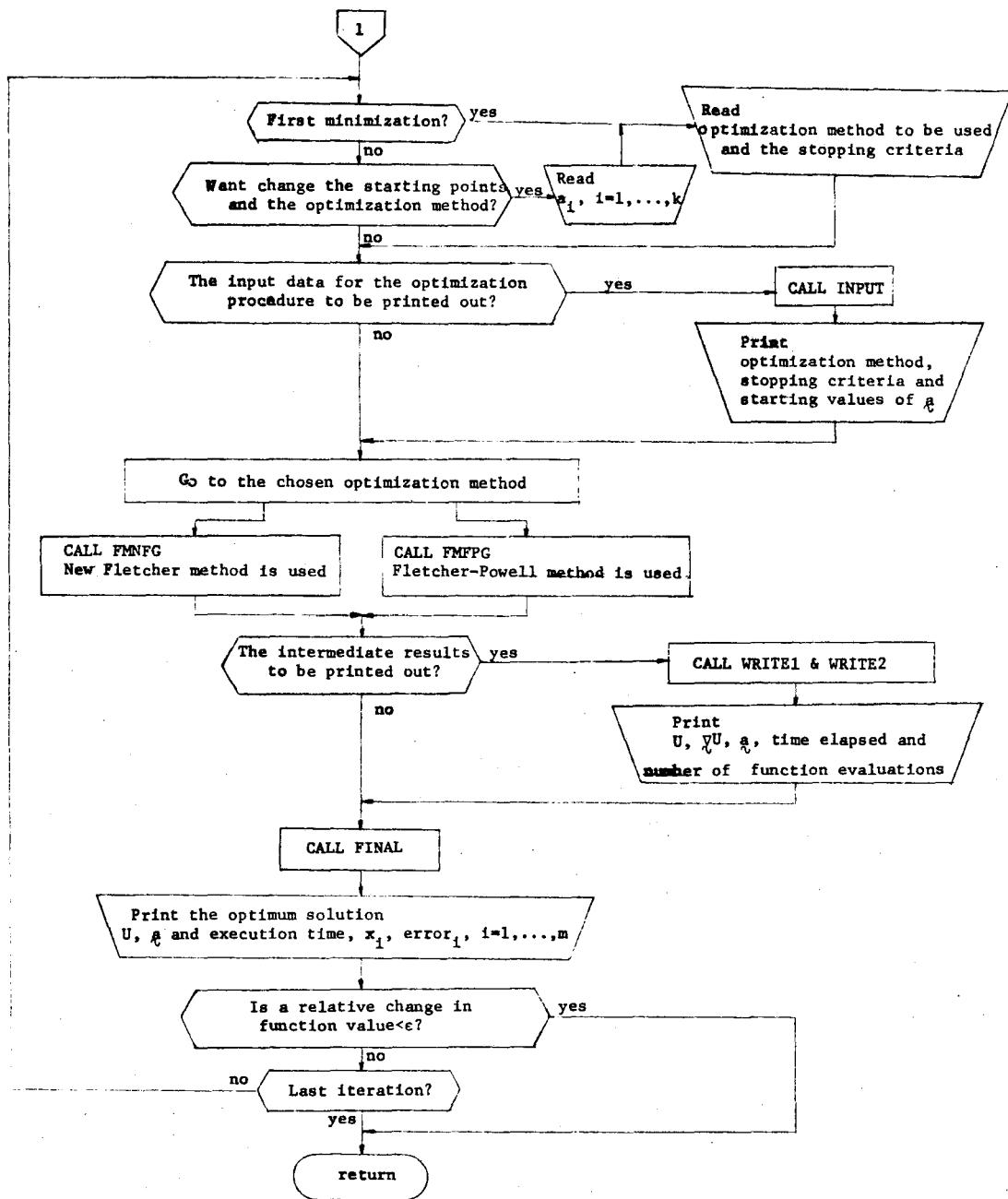


Fig. 6-8 Flowchart of subroutine FMLPO

## 6.7 Examples

Four examples were chosen to illustrate the work of the package. The first example is with the objective function as in case I, the second when it is as in case II, the third the low pass filter design problem without and, fourth, with parameter constraints.

### Test 1

A problem for system modelling, the same as used in Test 4 in Chapter V, is solved again. In this case upper and lower specified functions are chosen to be the same:

$$(6.22) \quad s_u'(t) = s_\ell'(t) = s(t) = \frac{3}{20} e^{-t} + \frac{1}{52} e^{-5t} - \frac{1}{65} e^{-2t} (3\sin 2t + 11\cos 2t)$$

with the artificial margin  $\xi=0$ , over the range  $t \in [0, 10]$ .

When one function approximates the other, it always corresponds to the situation in the case I.

We consider the same interval I twice, once as  $I_u$  and secondly as  $I_\ell$ . 51 uniformly spaced points were chosen over the range  $[0, 10]$  including 0 and 10, and a starting vector

$$\begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} .$$

The Fletcher-Powell and Fletcher optimization techniques were used with the stopping criterion  $\epsilon=10^{-6}$ . The results are

shown in Table 6-2. Also, the list of errors in both intervals and selected errors for further optimization, for  $p=2$  is given in Table 6-1. With  $p=2$ , the following parameter vector is found

$$\begin{bmatrix} 1.01647 \\ 0.78927 \\ 0.161400 \end{bmatrix}$$

The results from Table 6-2 are in excellent agreement with those in Table 5-6 when  $n=50$  and  $n_s=1$ , i.e. the case without quadratic interpolation. Also it is evident that the Fletcher method is more efficient.

### Test 2

The same optimization problem was chosen as in the Test 1, but this time the artificial margin is set to be different than zero, in particular

$$\xi = 2 \times 10^{-2}$$

therefore the artificial upper and lower specifications are distinct. Two possible cases may occur now: the specification may be violated and the specification may be satisfied.

The same initial set of sampling points and the starting value  $\hat{\alpha}$  as in Test 1 was used. From the Table 6-3, we may conclude that the starting objective function belongs to case I. The Fletcher optimization method is used with the same stopping criterion as in the Test 1 and already for  $p=2$  the problem is in case II, also is illustrated

sample points	independent variable in I	error: upper=lower	selected errors: nonnegative upper & nonpositive lower
1	0.0	-8.882 $10^{-16}$	1.288 $10^{-2}$
2	0.2	1.288 $10^{-2}$	6.985 $10^{-3}$
3	0.4	6.985 $10^{-3}$	1.033 $10^{-3}$
4	0.6	-1.173 $10^{-3}$	2.944 $10^{-3}$
5	0.8	-5.906 $10^{-3}$	3.810 $10^{-3}$
6	1.0	-6.635 $10^{-3}$	3.755 $10^{-3}$
7	1.2	-4.701 $10^{-3}$	3.043 $10^{-3}$
8	1.4	-1.736 $10^{-3}$	1.970 $10^{-4}$
9	1.6	1.033 $10^{-3}$	7.867 $10^{-3}$
10	1.8	2.944 $10^{-3}$	7.341 $10^{-7}$
11	2.0	3.810 $10^{-3}$	1.215 $10^{-6}$
12	2.2	3.755 $10^{-3}$	1.263 $10^{-6}$
13	2.4	3.043 $10^{-3}$	1.058 $10^{-6}$
14	2.6	1.970 $10^{-3}$	-8.882 $10^{-16}$
15	2.8	7.868 $10^{-4}$	-1.173 $10^{-3}$
16	3.0	-3.253 $10^{-4}$	-5.906 $10^{-3}$
17	3.2	-1.258 $10^{-3}$	-6.635 $10^{-3}$
18	3.4	-1.963 $10^{-3}$	-4.701 $10^{-3}$
19	3.6	-2.437 $10^{-3}$	-1.736 $10^{-3}$
20	3.8	-2.702 $10^{-3}$	-3.253 $10^{-4}$
21	4.0	-2.795 $10^{-3}$	-1.258 $10^{-3}$
22	4.2	-2.754 $10^{-3}$	-1.963 $10^{-3}$
23	4.4	-2.617 $10^{-3}$	-2.437 $10^{-3}$
24	4.6	-2.419 $10^{-3}$	-2.702 $10^{-3}$
25	4.8	-2.183 $10^{-3}$	-2.795 $10^{-3}$

sample points	independent variable in I	error: upper=lower	selected errors: nonnegative upper & nonpositive lower
26	5.0	-1.933 $10^{-3}$	-2.754 $10^{-3}$
27	5.2	-1.681 $10^{-3}$	-2.617 $\times 10^{-3}$
28	5.4	-1.440 $10^{-3}$	-2.419 $10^{-3}$
29	5.6	-1.215 $10^{-3}$	-2.183 $10^{-3}$
30	5.8	-1.011 $10^{-3}$	-1.933 $10^{-3}$
31	6.0	-8.301 $10^{-4}$	-1.681 $10^{-3}$
32	6.2	-6.723 $10^{-4}$	-1.440 $10^{-3}$
33	6.4	-5.371 $10^{-4}$	-1.215 $10^{-3}$
34	6.6	-4.231 $10^{-4}$	-1.011 $10^{-3}$
35	6.8	-3.284 $10^{-4}$	-8.301 $10^{-4}$
36	7.0	-2.510 $10^{-4}$	-6.723 $10^{-4}$
37	7.2	-1.886 $10^{-4}$	-5.371 $10^{-4}$
38	7.4	-1.391 $10^{-4}$	-4.231 $10^{-4}$
39	7.6	-1.005 $10^{-4}$	-3.284 $10^{-4}$
40	7.8	-7.081 $10^{-5}$	-2.510 $10^{-4}$
41	8.0	-4.848 $10^{-5}$	-1.886 $10^{-4}$
42	8.2	-3.199 $10^{-5}$	-1.391 $10^{-4}$
43	8.4	-2.013 $10^{-5}$	-1.005 $10^{-4}$
44	8.6	-1.182 $10^{-5}$	-7.081 $10^{-5}$
45	8.8	-6.211 $10^{-6}$	-4.848 $10^{-5}$
46	9.0	-2.597 $10^{-6}$	-3.199 $10^{-5}$
47	9.2	-4.265 $10^{-7}$	-2.013 $10^{-5}$
48	9.4	7.341 $10^{-7}$	-1.182 $10^{-5}$
49	9.6	1.215 $10^{-6}$	-6.211 $10^{-6}$
50	9.8	1.263 $10^{-6}$	-2.597 $10^{-6}$
51	10.0	1.058 $10^{-6}$	-4.265 $10^{-7}$

Table 6-1

	$p=2$		$p=10^2$	
	Fletcher-Powell	Fletcher	Fletcher-Powell	Fletcher
$a_1$	1.01647063	1.01647060	$6.8785026 \times 10^{-1}$	$6.8785037 \times 10^{-1}$
$a_2$	$7.8927025 \times 10^{-1}$	$7.8927036 \times 10^{-1}$	$9.5318185 \times 10^{-1}$	$9.5318177 \times 10^{-1}$
$a_3$	$1.61400086 \times 10^{-1}$	$1.61400093 \times 10^{-1}$	$1.2318597 \times 10^{-1}$	$1.2318599 \times 10^{-1}$
$M(a)$	$1.2880047 \times 10^{-2}$	$1.2880048 \times 10^{-2}$	$7.988425 \times 10^{-3}$	$7.988427 \times 10^{-3}$
$t_M$	$2.0 \times 10^{-1}$	$2.0 \times 10^{-1}$	$2.0 \times 10^{-1}$	$2.0 \times 10^{-1}$
$U(a)$	$2.09004705 \times 10^{-2}$	$2.09004705 \times 10^{-2}$	$8.04667205 \times 10^{-3}$	$8.04667205 \times 10^{-3}$
execution time in seconds	7.5	7.2	10	7.4
	$p=10$		$p=10^4$	
	Fletcher-Powell	Fletcher	Fletcher-Powell	Fletcher
$a_1$	$7.4276852 \times 10^{-1}$	$7.4276856 \times 10^{-1}$	$6.844480 \times 10^{-1}$	$6.844475 \times 10^{-1}$
$a_2$	$9.295083 \times 10^{-1}$	$9.295085 \times 10^{-1}$	$9.540875 \times 10^{-1}$	$9.540873 \times 10^{-1}$
$a_3$	$1.2803127 \times 10^{-1}$	$1.2803137 \times 10^{-1}$	$1.2286721 \times 10^{-1}$	$1.2286716 \times 10^{-1}$
$M(a)$	$8.592078 \times 10^{-3}$	$8.59209 \times 10^{-3}$	$7.947445 \times 10^{-3}$	$7.947439 \times 10^{-3}$
$t_M$	$2.0 \times 10^{-1}$	$2.0 \times 10^{-1}$	$2.0 \times 10^{-1}$	$2.0 \times 10^{-1}$
$U(a)$	$9.22275978 \times 10^{-3}$	$9.22275978 \times 10^{-3}$	$7.94802468 \times 10^{-3}$	$7.94802476 \times 10^{-3}$
execution time in seconds	13.7	6.7	8.6	6.8

Table 6-2

in Table 6-3. The optimal solutions for different values of  $p$  are presented in Table 6-4. The artificial margin  $\xi$  does not affect the optimal solution of the variable vector  $\alpha$ . It has the same value as the corresponding vector from the Table 6-2.

The starting vector  $\alpha$  was

$$\begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

and the following vector was found

$$\begin{bmatrix} 0.92369 \\ 0.83516 \\ 0.14720 \end{bmatrix}$$

for  $p=2$ .

### Test 3

The computer program was used in the optimization of a five section cascaded transmission-line low pass filter which has been considered by Carlin [35]. The terminations of the filter are unity, the length of the  $i$ th section  $\ell_i$  and the normalized characteristic impedance of the  $i$ th section  $Z_{0i}$ , such that a maximum insertion loss  $\gamma$  in the passband, from 0 to 1 GHz, is not more than 0.4 dB, while maximizing  $\gamma$  at a point in the stopband. All section lengths were kept fixed at 2.5 cm so that

$\xi = 2 \times 10^{-2}$	Initial values			$p=2$
sample points	upper errors; $lower = upper + 2\xi$	selected errors nonnegative upper nonpositive lower		upper errors; $lower = upper + 2\xi$
1	-2.000x10 <sup>-2</sup>			-2.000x10 <sup>-2</sup>
2	1.293x10 <sup>-1</sup>	1.293x10 <sup>-1</sup>		-8.991x10 <sup>-3</sup>
3	2.057x10 <sup>-1</sup>	2.057x10 <sup>-1</sup>		-1.535x10 <sup>-2</sup>
4	2.380x10 <sup>-1</sup>	2.380x10 <sup>-1</sup>		-2.321x10 <sup>-2</sup>
5	2.429x10 <sup>-1</sup>	2.429x10 <sup>-1</sup>		-2.728x10 <sup>-2</sup>
6	2.304x10 <sup>-1</sup>	2.304x10 <sup>-1</sup>		-2.727x10 <sup>-2</sup>
7	2.070x10 <sup>-1</sup>	2.070x10 <sup>-1</sup>		-2.469x10 <sup>-2</sup>
8	1.772x10 <sup>-1</sup>	1.772x10 <sup>-1</sup>		-2.124x10 <sup>-2</sup>
9	1.445x10 <sup>-1</sup>	1.445x10 <sup>-1</sup>		-1.818x10 <sup>-2</sup>
10	1.115x10 <sup>-1</sup>	1.115x10 <sup>-1</sup>		-1.615x10 <sup>-2</sup>
11	8.009x10 <sup>-2</sup>	8.009x10 <sup>-2</sup>		-1.532x10 <sup>-2</sup>
12	5.178x10 <sup>-2</sup>	5.178x10 <sup>-2</sup>		-1.551x10 <sup>-2</sup>
13	2.741x10 <sup>-2</sup>	2.741x10 <sup>-2</sup>		-1.643x10 <sup>-2</sup>
14	7.360x10 <sup>-3</sup>	7.360x10 <sup>-3</sup>		-1.776x10 <sup>-2</sup>
15	-8.373x10 <sup>-3</sup>			-1.919x10 <sup>-2</sup>
16	-2.007x10 <sup>-2</sup>			-2.054x10 <sup>-2</sup>
17	-2.820x10 <sup>-2</sup>			-2.168x10 <sup>-2</sup>
18	-3.334x10 <sup>-2</sup>			-2.255x10 <sup>-2</sup>
19	-3.608x10 <sup>-2</sup>			-2.314x10 <sup>-2</sup>
20	-3.699x10 <sup>-2</sup>			-2.348x10 <sup>-2</sup>
21	-3.660x10 <sup>-2</sup>			-2.360x10 <sup>-2</sup>
22	-3.533x10 <sup>-2</sup>			-2.356x10 <sup>-2</sup>
23	-3.354x10 <sup>-2</sup>			-2.340x10 <sup>-2</sup>
24	-3.151x10 <sup>-2</sup>			-2.314x10 <sup>-2</sup>
25	-2.944x10 <sup>-2</sup>			-2.284x10 <sup>-2</sup>

CONT'D

sample points	upper errors; lower=upper+2ξ	selected errors nonnegative upper nonpositive lower	upper errors; lower=upper+2ξ
26	$-2.748 \times 10^{-2}$		$-2.251 \times 10^{-2}$
27	$-2.570 \times 10^{-2}$		$-2.218 \times 10^{-2}$
28	$-2.417 \times 10^{-2}$		$-2.186 \times 10^{-2}$
29	$-2.289 \times 10^{-2}$		$-2.155 \times 10^{-2}$
30	$-2.186 \times 10^{-2}$		$-2.128 \times 10^{-2}$
31	$-2.106 \times 10^{-2}$		$-2.103 \times 10^{-2}$
32	$-2.047 \times 10^{-2}$		$-2.082 \times 10^{-2}$
33	$-2.005 \times 10^{-2}$		$-2.063 \times 10^{-2}$
34	$-1.978 \times 10^{-2}$		$-2.048 \times 10^{-2}$
35	$-1.962 \times 10^{-2}$		$-2.035 \times 10^{-2}$
36	$-1.954 \times 10^{-2}$		$-2.025 \times 10^{-2}$
37	$-1.952 \times 10^{-2}$		$-2.017 \times 10^{-2}$
38	$-1.954 \times 10^{-2}$		$-2.011 \times 10^{-2}$
39	$-1.960 \times 10^{-2}$		$-2.006 \times 10^{-2}$
40	$-1.965 \times 10^{-2}$		$-2.003 \times 10^{-2}$
41	$-1.972 \times 10^{-2}$		$-2.000 \times 10^{-2}$
42	$-1.978 \times 10^{-2}$		$-1.999 \times 10^{-2}$
43	$-1.984 \times 10^{-2}$		$-1.998 \times 10^{-2}$
44	$-1.989 \times 10^{-2}$		$-1.9978 \times 10^{-2}$
45	$-1.993 \times 10^{-2}$		$-1.9977 \times 10^{-2}$
46	$-1.997 \times 10^{-2}$		$-1.9977 \times 10^{-2}$
47	$-1.999 \times 10^{-2}$		$-1.9979 \times 10^{-2}$
48	$-2.001 \times 10^{-2}$		$-1.9982 \times 10^{-2}$
49	$-2.002 \times 10^{-2}$		$-1.9986 \times 10^{-2}$
50	$-2.003 \times 10^{-2}$		$-1.9989 \times 10^{-2}$
51	$-2.003 \times 10^{-2}$		$-1.999 \times 10^{-2}$

Table 6-3

Fletcher

	$q=-p=-2$	$q=-p=-10$	$q=-p=-10^2$	$q=-p=-10^4$
$a_1$	$9.2369 \times 10^{-1}$	$7.7204 \times 10^{-1}$	$6.9015 \times 10^{-1}$	$6.8446 \times 10^{-1}$
$a_2$	$8.3515 \times 10^{-1}$	$9.1379 \times 10^{-1}$	$9.5224 \times 10^{-1}$	$9.5408 \times 10^{-1}$
$a_3$	$1.4720 \times 10^{-1}$	$1.3053 \times 10^{-1}$	$1.2338 \times 10^{-1}$	$1.2287 \times 10^{-1}$
$M(\hat{a})$	$-8.9912 \times 10^{-3}$	$-1.1105 \times 10^{-2}$	$-1.1987 \times 10^{-2}$	$-1.205 \times 10^{-2}$
$t_M$	0.2	0.2	0.2	0.2
$U(\hat{a})$	$-1.8963 \times 10^{-3}$	$-9.9831 \times 10^{-3}$	$-1.1899 \times 10^{-2}$	$-1.2051 \times 10^{-2}$
function evaluations	33	21	23	33
execution time in seconds	8.4	5.9	6.5	9.8

Table 6-4

the maximum stopband insertion loss would occur at 3GHz and the normalized characteristic impedances are used as variables.

The insertion loss is given by

$$(6.23) \quad \gamma = 10 \log_{10} \frac{1}{1 - |\rho|^2}$$

where  $\rho$  is the reflection coefficient, therefore the corresponding  $\rho$  to the insertion loss may be considered in the optimization problem because there is a finite upper bound on  $|\rho|$  which is 1. 21 uniformly spaced sample points were used in the passband and a single point at 3GHz. The artificial margin  $\xi$  for one case is set to be zero and for the other 0.02337. The weighting function is set to be 1 everywhere. The starting value of the variable vector  $\alpha$  was

$$\begin{bmatrix} 3.180 \\ 0.443 \\ 4.38 \\ 0.443 \\ 3.180 \end{bmatrix}$$

Results obtained using the Fletcher method for the value of  $p=10^3$  are presented in Table 6-5 and the response is shown in Fig. 6-9. More details about the computation of the reflection coefficient are given in the Appendix VII.

Although, for physical reasons the symmetrical results for the variables are expected, symmetry was not assumed.  $\xi$  does not affect the optimal solution. The reason that it was considered was

## Fletcher method

$p=10^3$	$\xi=0$	$\xi=2.337 \times 10^{-2}$
$a_1$	3.1525	3.1508
$a_2$	$4.4203 \times 10^{-1}$	$4.4165 \times 10^{-1}$
$a_3$	4.4212	4.4194
$a_4$	$4.4159 \times 10^{-1}$	$4.4169 \times 10^{-1}$
$a_5$	3.1526	3.1508
M	$3.9466 \times 10^{-5}$	$-2.3330 \times 10^{-2}$
$x_M$	$3.0700 \times 10^{-1}$	3.0
U	$3.9466 \times 10^{-5}$	$-2.3330 \times 10^{-2}$
Function evaluations	177	79
Execution time in seconds	17	13

Table 6-5

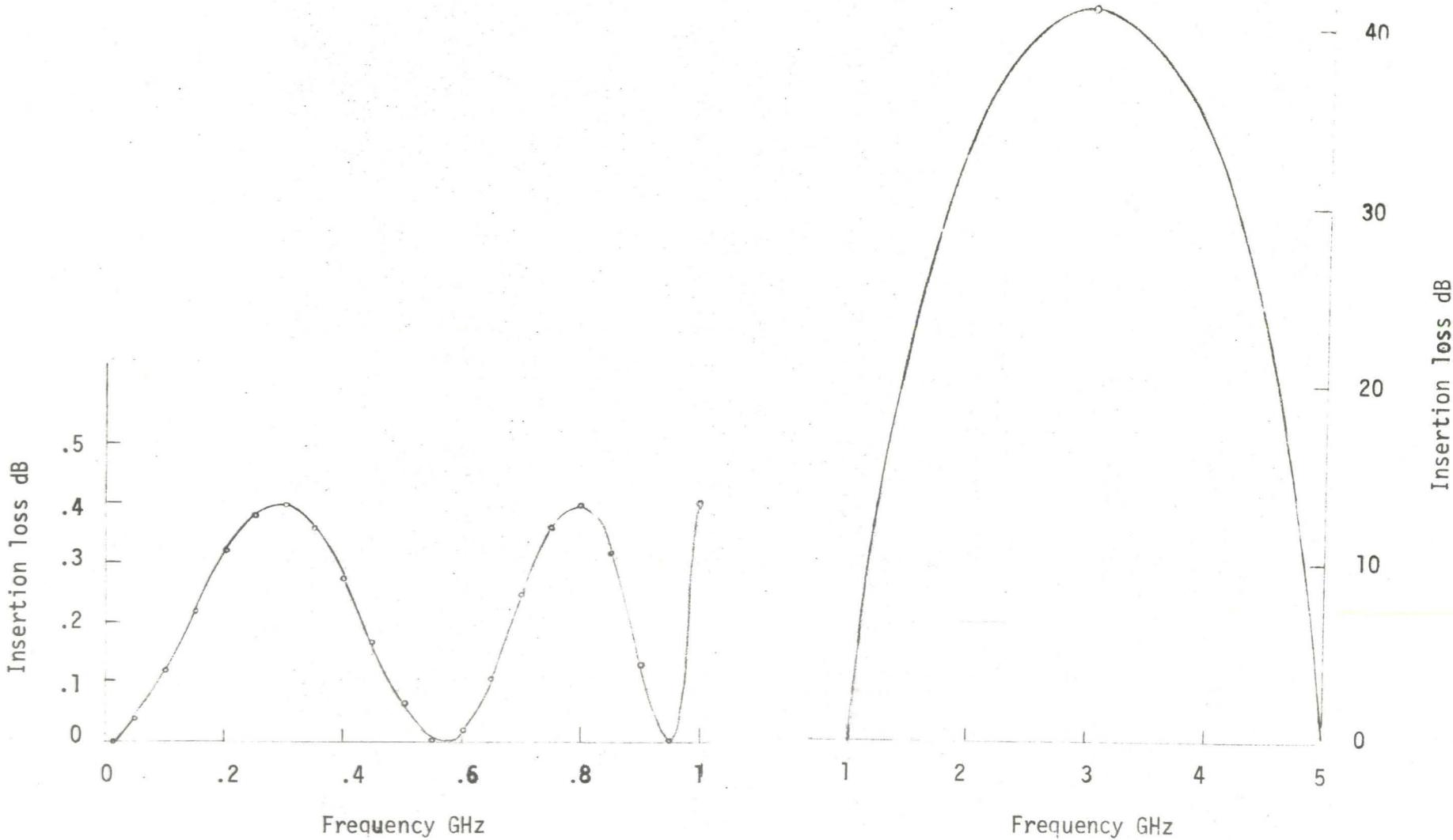


Fig. 6-9 Response of the five section transmission-line filter for the unconstrained design problem.

to bring the objective function into the case when the artificial specification is satisfied.

#### Test 4

Constraints are put on the parameter vector  $\alpha$  such that they are not satisfied at the optimal solution of the unconstrained problem given in Test 3.

Although the package FMLPO is not written for nonlinear programming, the constrained problem also may be considered. The constraints on a parameter may be considered as upper and lower specifications on an approximating function defined as a single variable parameter over the dummy point outside of the working set of points. This dummy point has to be defined as a new interval for each specification, which is not so practical, but it works. The variable parameter vector  $\alpha$  is constrained by

$$(6.24) \quad 0.2 \leq a_i \leq 4.0, \quad i=1,2,\dots,5,$$

i.e., the solutions from Table 6-5 for  $a_3$  are infeasible. For the problem given in Test 3 and considering (6.24), results are presented in Table 6-6. The response is shown in Fig. 6-10.

## Fletcher method

$p=10^3$	$\xi=0$	$\xi=2.337 \times 10^{-2}$
$a_1$	2.9429	2.9422
$a_2$	$4.1069 \times 10^{-1}$	$4.1070 \times 10^{-1}$
$a_3$	4.0	3.9998
$a_4$	$4.1069 \times 10^{-1}$	$4.1070 \times 10^{-1}$
$a_5$	2.9429	2.9422
M	$4.7403 \times 10^{-5}$	$-2.3322 \times 10^{-2}$
$x_M$	$8.02 \times 10^{-1}$	3.0
U	$4.7403 \times 10^{-5}$	$-2.3322 \times 10^{-2}$
Function evaluations	55	157
Execution time in seconds	31	55

Table 6-6

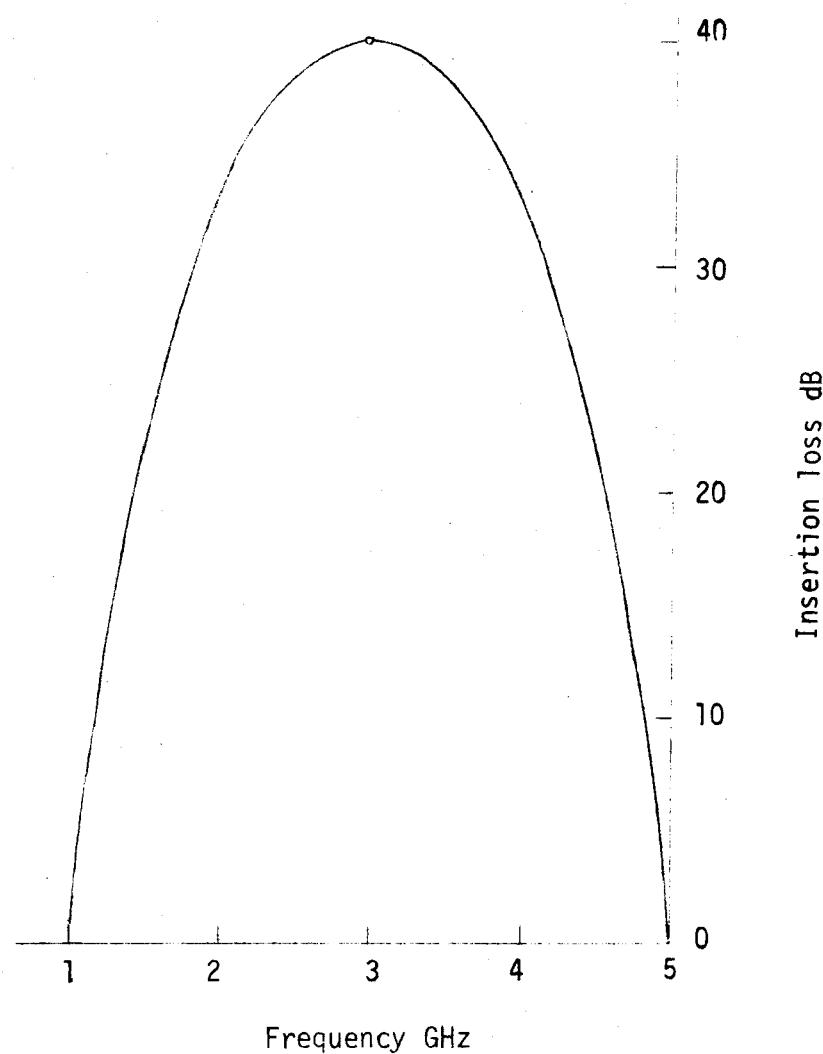
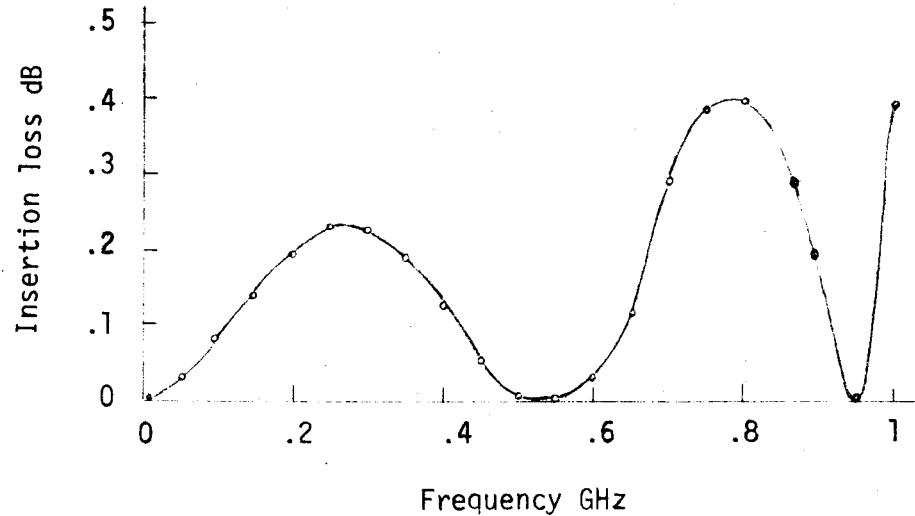


Fig. 6-10 Response of the five section transmission-line filter for the constrained design problem.

## CHAPTER VII

### Conclusions

Two user-oriented computer programs in FORTRAN IV for the discrete nonlinear least pth approximation [1] and generalized least pth approximation [2], respectively, have been presented in this thesis. The instructions on how to use the packages are given, illustrated with simple examples.

The first program treats the case when we have only one specified function. The package can be used for fitting one function or data to another continuous function when the error function happens to be real. If quadratic interpolation is applied, the specified function and the approximating function have to be continuous on a given interval. Employing quadratic interpolation the sampling for the objective function may take fewer points for equally successful solutions in comparison with the case which does not consider quadratic interpolation, and may save some computation time. Using quadratic interpolation the location of the extreme points were found more precisely and the solutions are closer to the minimax solutions for a given value of p. The computer program can be rearranged such that the specified function and the approximating function are both complex, which is useful in electrical engineering design problems.

The second program is an extension of the first program and is applicable to design problems with upper and lower specifications. The program is directly applicable to such problems as meeting or exceeding design specifications on several disjoint closed intervals as in filter design. Although the program is not written for non-linear programming we found that it is also applicable for problems with parameter constraints.

The methods presented abandon the linear programming subproblem which many of the minimax methods do. The advantage over the direct minimax methods is that they use very efficient gradient methods such as Fletcher-Powell [18] and a recent method by Fletcher [19]. From the experimental results, the Fletcher-Powell algorithm was found to be reliable. The method, however, was found to be slow in comparison with the method proposed by Fletcher. The latter method requires fewer function evaluations to reach the optimum and is less time consuming.

The larger the value of  $p$  that is used, the more nearly the minimax solution is obtained, but more function evaluations are required to bring the objective function close to the optimum. For practical purposes smaller values of  $p$  may be used to attain a satisfactory solution, hence the objective function will be minimized faster. We can start with a smaller value of  $p$ , increase it after each complete optimization and terminate when the relative change in the objective function in the successive iterations is less than a prescribed small quantity. This can be a disadvantage if the starting point is close to the minimax optimum, which rarely happens in practice.

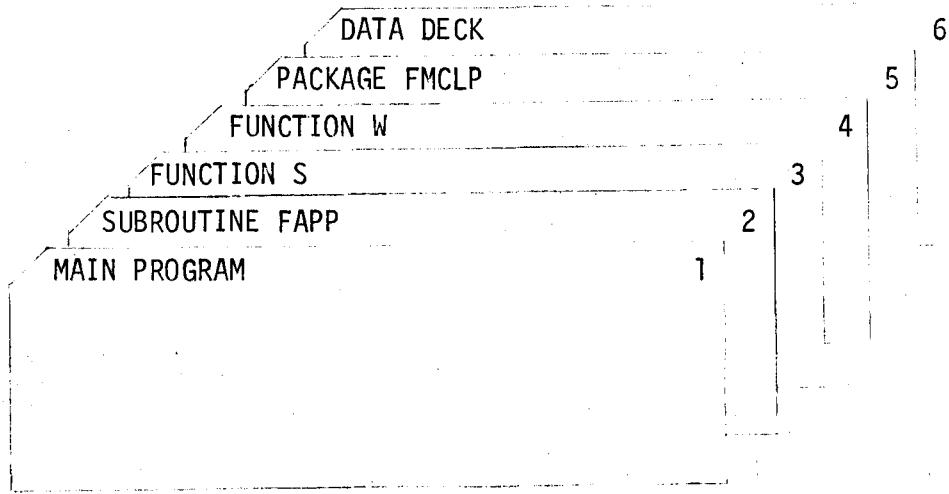
Typically less than 3 minutes of CDC 6400 computer time and a core requirement of about 14K<sub>10</sub> is sufficient to optimize the type of examples given in this thesis to a high degree of accuracy from a good starting point.

## APPENDIX I

Usage of FMCLP ( Function Minimization for Continuous Least pth  
Approximation )

Purpose To minimize the objective function of  $k$  variables as defined  
as a discrete least pth objective with a single speci-  
fication (5.3) using gradient methods.

How to use Set the input deck as follows:



### 1. Main program

Write the main program as indicated below:

- (i) Dimension the following arrays  
A( $K$ ), ASTRT( $K$ ), G( $K$ ), GRAD( $K$ ), Y( $K$ ), PY( $K$ ), DUM1( $K$ ), DUM2( $K$ ), EPS( $K$ ),  
H( $M$ ), X( $N_2$ ), ERROR( $N_2$ ), IPA(ITER)

where

K is the number of variable parameters,

M = K(K+7)/2, N2 = N+2,

N is the discrete point set,

ITER is the maximum number of times the optimization method is used.

(ii) Call the subroutine FMCLP as follows:

CALL FMCLP (A, ASTRT, G, GRAD, Y, PY, DUM1, DUM2, EPS, H, X,  
ERROR, IPA)

(iii) Add STOP and END cards.

## 2. Subroutine FAPP

This subroutine calculates the approximating function and its gradients with respect to variable vector  $\tilde{A}$ .

Write subroutine FAPP as follows:

SUBROUTINE FAPP (X, K, A, APP, GRAD, INDIC)

DIMENSION A(1), GRAD(1)

where X, K, A and INDIC are input and APP and GRAD are output variables.

GO TO (1, 2), INDIC

C Compute the value of the approximating function

1 APP = F( $\tilde{A}$ , X)

RETURN

C Calculate the values of the gradients of the approximating  
function

$$2 \quad \text{GRAD}(1) = \frac{\partial F(\tilde{A}, X)}{\partial A_1}$$

$$\text{GRAD}(2) = \frac{\partial F(A, X)}{\partial A_2}$$

$$\text{GRAD}(K) = \frac{\partial F(A, X)}{\partial A_K}$$

RETURN

END

where INDIC may have the values 1 or 2 which indicate whether the approximating function or all the gradients should be calculated, respectively.

### 3. FUNCTION S

Function S is a subprogram of a single input variable X and defines a specified function

$$S = S(X).$$

Add RETURN and END cards.

### 4. FUNCTION W

Function W is also a subprogram of a single input variable X and defines a weighting function

$$W = w(X).$$

Add RETURN and END cards.

An example which shows how to set these subprograms is the same as Test 3 in the Chapter V. The listings which define this problem follow.

```
SUBROUTINE FAPP(X,N1,A,APP,GRAD,INDIC)
```

```
C  
C          SURROUTINE WHICH CALCULATES APPROXIMATING  
C          FUNCTION AND ITS GRADIENTS WITH RESPECT TO  
C          VARIABLE PARAMETERS  
C
```

```
DIMENSION A(1),GRAD(1)  
GO TO(100,200),INDIC  
100 APP=A(1)*X+A(2)*FXP(X)  
      RETURN  
200 GRAD(1)=X  
      GRAD(2)=FXP(X)  
      RETURN  
END
```

.....

```
FUNCTION S(X)
```

```
C  
C          FUNCTION SUBPROGRAM WHICH DEFINES  
C          SPECIFIED FUNCTION  
C  
S=X**2  
RETURN  
END
```

.....

```
FUNCTION W(X)
```

```
C  
C          FUNCTION SUBPROGRAM WHICH DEFINES  
C          WEIGHTING FUNCTION  
C  
W=1.  
RETURN  
END
```

## 6. Data Deck

Parameters to be supplied as data are defined below:

N	The number of sample points forming the discrete point set.
XA, XB	The left and the right end points of the interval of the independent parameter.
NSUB	The number of subintervals over [XA, XB].
IREAD	Integer which denotes whether or not the discrete set of points in [XA, XB] will be read.  If IREAD=0, the discrete point set will be arranged equidistantly over the interval;  If IREAD=1, the discrete point set will be read from data.
X(I), I=1, N	The discrete point set over the interval.
K	The number of the independent variable parameters  A.
ASTRT(I), I=1, K	Starting values for the K variable parameters.
IGRDCH	Gradients to be checked if IGRDCH=1, it should be set to 0 if gradients are not to be checked.
MET	Optimization method to be called:  if MET=1, Fletcher method will be called;  if MET=2, Fletcher-Powell method will be called.
MAX	Maximum number of permissible iterations.

ITER	Has already been defined in the main program as a length of the working array.
IPA(I), I=1, ITER	Vector containing the values of p for different least pth objectives.
IOPT	Denotes how many times the optimization is repeated with different starting points and/or different optimization techniques.
IPRINT	Intermediate output is printed out every IPRINT iteration; it should be set to 0 if no intermediate output is desired.
IDATA	Input data is printed out if IDATA=1; it should be set to 0 if input data is not to be printed out.
EST	Minimum estimated value of the objective function.
EPS(I), I=1, K	Small test quantities used by the Fletcher method.
EPS1	Small test quantity used by the Fletcher-Powell method.
DIF	Small test quantity used by the subroutine FMCLP.

Setting up the data deck

Conditions	Number of cards	Parameters	Type	Format
-	1	K, N, NSUB, IREAD, IOPT, ITER, IGRDCH	INTEGER	7I10
-	1	XA, XB	REAL	2E16.8
IREAD = 1	As many as required by N	X(I), I = 1, N	REAL	5E16.8
-	1	EST, DIF	REAL	2E16.8
+ -	As many as required by K	ASTRT(I), I = 1, K	REAL	5E16.8
IOPT times MET = 1	1	MET, MAX, IPRINT, IDATA	INTEGER	4I10
	As many as required by K	EPS(I), I=1, K	REAL	5E16.8
MET = 2	1	EPS1	REAL	5E16.8
+ -	As many as required by ITER	IPA(I), I = 1, ITER	INTEGER	8I10

Recommended values for some of the parameters

NSUB = 5

MAX = 100

EPS(I), I = 1, K, each  $10^{-6}$ DIF =  $10^{-4}$ 

EST      A lower bound of the minimum of the objective function may be obtained from physical reasons. If the true minimum is not known, choose EST to be small enough (negative values are allowed). For approximation problems 0 is convenient.

The following list illustrates how to set the data for the example for which listings have been shown. Both optimization methods, Fletcher-Powell and Fletcher, are called.

?	10	10	0	?	4	1
0.0E	00	2.0F	00			
0.0F	00	1.0E	-4			
1.0F	00	1.0F	00			
1	50	1	1			
1.0F	-6	1.0E	-6			
2	10	100	1000			
1.0F	00	1.0E	00			
2	50	1	1			
1.0F	-6					
2	10	100	1000			

Listing of the Package FMCLP

```

SUBROUTINE FMCLP (A,XSTR, G, GRAD, Y, PY, DUM1, DUM2, EPS, H, X, ERROR, IPA) A
  SUBROUTINE WHICH COORDINATES THE OTHER A
  SUBROUTINES IN THE PACKAGE FMCLP A
  EXTERNAL FUNCT A
  EXTERNAL SW A
  DIMENSION A(1), XSTR(1), G(1), Y(1), PY(1), DUM1(1), DUM2(1), EPS A
  1(1), H(1), GRAD(1), X(1), ERROR(1), IPA(1) A
  COMMON /BLK/ K0 A
  LOGICAL CONV, UNITH A
  EPR(Z)=WE3R(Z,G,H,A,GRAD) A
  UNITH=.TRUE. A
  READ (5,40) N1,N2,NPOD,IREAD,IOPT,ITER,IGRDCH A
  READ (5,42) XA,XE A
  XAB=XB-XA A
  NP=N+1 A
  IF (IREAD.EQ.0) IREAD=2 A
  GO TO (3,1), IREAD A
  1 X(2)=XA A
  K=N-1 A
  DELTA=XAB/K A
  DO 2 I=1,K A
  X(I+2)=X(I)+DELTA A
  CONTINUE A
  X(N+1)=XB A
  GO TO 4 A
  3 READ (5,42) (X(I),I=2,NP) A
  4 WRITE (6,50) A
  READ (5,42) EST,CIF A
  READ (5,42) (XSTR(I),I=1,N1) A
  DO 5 I=1,N1 A
  A(I)=XSTR(I) A
  CONTINUE A
  DO 6 I=2,NP A
  5 I=I-1 A
  EPROR(I)=ERR(X(I)) A
  WRITE (6,51) II,X(I),ERROR(I) A
  CONTINUE A
  6 IGRDCH=1 A
  CALL NEWSET (XA,XB,NPOD,S,H,A,GRAD,N,X,ERROR,IOI) A
  WRITE (6,43) IOI A
  IF (IOI.EQ.9,7 A
  7 WRITE (6,45) A
  DO 8 I=2,NP A
  7 I=I-1 A
  EPROR(I)=ERR(X(I)) A
  WRITE (6,51) II,X(I),ERROR(I) A
  CONTINUE A
  8 EPROR(1)=ERR(2) A
  DO 10 I=3,NP A
  8 EPROR(1)=AMAX1 (ERROR(1),ERROR(I)) A
  CONTINUE A
  10 WRITE (6,48) A
  WRITE (6,49) A
  WRITE (6,47) ERROR(1) A
  DATA FOR THE OPTIMALITY A
  FOR THE OPTIMIZATION METHOD USED A
  11 DO 39 K=1,IOPT A
  39 K0=1 A
  12 IF (K-1).LT.12,12,11 A
  11 READ (5,42) (XSTR(I),I=1,N1) A
  12 READ (5,40) MET,MAX,IPRINT,IData A
  IF (MET.EQ.1) READ (5,42) (EPS(I),I=1,N1) A
  IF (MET.EQ.2) READ (5,42) EPS1 A
  READ (5,40) (IPA(I),I=1,ITER) A
  DO 38 KK=1,ITER A
  38
  C   OPTIMIZATION A
  IP=IPA(KK) A
  IF (KK.GT.1) FF=F A
  IF (KR.EQ.0) GO TO 13 A
  DO 14 I=1,N1 A
  A(I)=XSTR(I) A
  CONTINUE A
  14 IF (IGRDCH.NE.1) GO TO 15 A
  CALL GDPCHK (N1,A,G,PY,Y,GRAD,N,X,ERROR,IP,DUM1) A
  15 IF (KR.EQ.0) GO TO 16 A
  IF (IData.EQ.0) GO TO 16 A
  M=2 A
  CALL INPUT (MET,M,MAX,N1,IPRINT,IData,EPS1,EST,EPS,XSTR) A
  16 IF (MET.EQ.0) ME1=4 A
  INDEX=0 A
  GO TO (17,22,31,27), MET A
  17 IF (IPRINT.EQ.1) GO TO 18 A
  CALL WRITE1 (1) A
  18 CALL SECOND (T1) A
  IF (KP.NE.0) GO TO 20 A
  DO 19 I=1,N1 A
  A(I)=DUM1(I) A
  CONTINUE A
  19 CALL FMNFC (N1,A,F,G,H,UNITH,EST,EPS,MAX,IPRINT,IEXIT,GRAD,N,X,ERR A
  20 10P,IP) A
  20

```

```

DO 21 I=1,N1          A  97
DUM1(I)=A(I)          A  98
CONTINUE               A  99
CALL SECOND (T2)       A 100
CALL FINAL (A,F,N1)    A 101
T=T2-T1                A 102
WPITE (6,41) T         A 103
GO TO 22               A 104
IF (IPRINT.EQ.0) GO TO 23   A 105
CALL WRITE1 (2)         A 106
CALL SECOND (T1)       A 107
IF (IKR.NE.0) GO TO 25   A 108
DO 24 I=1,N1           A 109
A(I)=DUM2(I)           A 110
CONTINUE               A 111
25 CALL FMFPC (FUNCT,N1,A,F,G,EST,EPS1,MAX,IER,H,IPRINT,GRAD,N,X,ERR0 A 112
      DO 26 I=1,N1           A 113
      DUM2(I)=A(I)           A 114
26 CONTINUE               A 115
CALL SECOND (T2)       A 116
CALL FINAL (A,F,N1)    A 117
T=T2-T1                A 118
WPITE (6,41) T         A 119
INDEX=INDEX+1           A 120
IF (M.EQ.1) GO TO 28   A 121
GO TO 30               A 122
DO 29 I=1,N1           A 123
A(I)=XSTR1(I)           A 124
CONTINUE               A 125
30 CONTINUE               A 126
C                         A 127
31 KP=0                  A 128
WPITE (6,44) IP        A 129
WRITE (6,45)             A 130
DO 32 I=2,NP             A 131
II=I-1                 A 132
WRITE (6,51) II,X(I),ERROR(I) A 133
CONTINUE               A 134
CALL NEWSET (XA,XB,NPOD,S,W,A,GRAD,N,X,ERROR,IOI) A 135
WPITE (6,43) IOI          A 136
IF (IOI).EQ.35,35,33     A 137
33 WRITE (6,45)             A 138
DO 34 I=2,NP             A 139
II=I-1                 A 140
EPROR2(I)=ERR(X(I))    A 141
WPITE (6,51) II,X(I),ERROR(I) A 142
CONTINUE               A 143
34 EPROR(I)=ERROR(2)     A 144
DO 35 I=3,NP             A 145
EPROR(I)=AMAX1(ERROR(1),ERROR(I)) A 146
35 CONTINUE               A 147
36 WPITE (6,48)             A 148
WPITE (6,49)             A 149
WPITE (6,50) ERROR(1)    A 150
IGROCH=IGROCH+2          A 151
IF (KK-1).EQ.38,38,37     A 152
37 FTST=ABS ((FF-F)/FF)   A 153
IF (FTST.LT.DIF) GO TO 39   A 154
38 CONTINUE               A 155
39 CONTINUE               A 156
RETURN                  A 157
C                         A 158
40 FORMAT (8I10)           A 159
41 FORMAT (1H0,//25X,*EXECUTION TIME IN SECONDS=*,F10.5) A 160
42 FORMAT (5E16.8)           A 161
43 FORMAT (1H1,14X,*NUMBER OF Q.I.=*,I3) A 162
44 FORMAT (1H1,13X,*P =*,I7) A 163
45 FORMAT (*8X,*NEW SET OF INDEPENDENT VARIABLE*,15X,*ERRORS*) A 164
46 FORMAT (*13X,*INDEPENDENT VARIABLE*,15X,*ERRORS*) A 165
47 FORMAT (*13X,E21.12) A 166
48 FORMAT (*15X,*ABSOLUTE VALUE*) A 167
49 FORMAT (*15X,*OF MAXIMUM ERROR*) A 168
50 FORMAT (*1H1,8X,*INITIAL SET OF INDEPENDENT VARIABLE*,9X,*ERRORS*) A 169
51 FORMAT (19,1X,E23.12,E23.12) A 170
END                      A 171
                                         A 172
*****
```

```

FUNCTION WERR (Z,S,W,B,GRAD)
  FUNCTION SUBPROGRAM WHICH CALCULATES
  WEIGHTED ERRCR FUNCTION
EXTERNAL S,W
DIMENSION B(1), GRAD(1)
CALL FAPP (Z,N1,B,APP,GRAD,1)
WERR=(APP-S(Z))&K(Z)
RETURN
END
*****
```

1  
10  
11

```

SUBROUTINE NEWSET (XA,XB,NPOD,S,W,B,GRAD,N,X,ERROR,IQI)
  SUBROUTINE WHICH CALCULATES THE NEW SET OF THE INDEPENDENT
  VARIABLES WHICH INCLUDE ALL THE EXTREMA OF THE WEIGHTED
  ERROR FUNCTION

EXTERNAL S,W
DIMENSION B(1), GRAD(1), X(1), ERROR(1)
EPSR(Z)=ERR(Z,S,W,B,GRAD)
EPSN(Z)=ABS(WERR(Z,S,W,B,GRAD))
IER=0
NN=N+1
NNN=N+2
X(1)=XA
X(NNN)=XB
I=2
IQ=1
IND=1
IQI=0
DO 29 I=1,NN
IF (X(I)-X(I+1)) 1,29,29
1 ZMIN=X(I+1)
ZMAX=X(I)
IQ=IND
IND=1
Z=ZMAX
EMAX=EPSN(Z)
EMIN=EPSN(ZMIN)
DELTA=(ZMIN-ZMAX)/NPOD
NPOD1=NPOD+1
MM=1
MIN=NPOD1
DO 7 K=1,NPOD
ETREN=EPSN(Z)
IF (ETREN-EMAX) 6,6,2
2 EMAX=ETREN
MM=K
ZMAX=Z
ZTEST1=ZMAX+DELTA
ETEST1=EPSN(ZTEST1)
IF (EMAX-ETEST1) 6,6,3
3 ZTEST2=ZMAX-DELTA
IF (XA-(ZMAX-DELTA)) 4,4,6
4 ETEST2=EPSN(ZTEST2)
TF (EMAX-ETEST2) 6,6,5
5 IND=2
GO TO 8
6 Z=Z+DELTA
7 CONTINUE
8 IF (XA-(ZMAX-DELTA)) 9,9,13
9 IF (MM=NPOD) 10,10,13
10 GO TO (13,11), IND
11 Q1=EPSN(ZMAX-DELTA)
C2=EPSN(ZMAX+DELTA)
RA=(EMAX*2.-01-Q2)*2.
IF (RA) 12,13,12
12 ZMAX=ZMAX*(Q2-Q1)*DELTA/RA
IQI=IQI+1
ZMIN=ZTEST1
13 IF (I-1) 14,14,17
14 IF (ZMAX-ZMIN) 15,29,29
15 GO TO (29,16), IND
16 X(1)=ZMAX
XKP=X(2)
X(2)=ZMIN
GO TO 29
17 IF (X(I-1)-ZMAX) 18,21,21
18 IF (ZMAX-ZMIN) 19,21,21
19 GO TO (23,20), IND
20 X(I)=ZMAX
21 IF (X(I)-ZMIN) 22,25,25
22 XKP=X(I)
X(I+1)=ZMIN
GO TO 23
23 GO TO (25,24), IC
24 X(I)=XKP
25 IF (I>NN) 29,26,29
26 IF (X(I)-X(I+1)) 27,29,29
27 IF (EPSN(X(I))-EPSN(X(I+1))) 28,29,29
28 II=II-1
29 CONTINUE
30 IF (EPSN(X(2))-EPSN(X(1))) 30,31,31
31 II=II-3
32 IF (II-2) 32,51,32
33 IF (II+2) 42,33,42
JJ=NN/2+1
IF (N-2) 39,39,34

```

1  
 2  
 3  
 4  
 5  
 6  
 7  
 8  
 9  
 10  
 11  
 12  
 13  
 14  
 15  
 16  
 17  
 18  
 19  
 20  
 21  
 22  
 23  
 24  
 25  
 26  
 27  
 28  
 29  
 30  
 31  
 32  
 33  
 34  
 35  
 36  
 37  
 38  
 39  
 40  
 41  
 42  
 43  
 44  
 45  
 46  
 47  
 48  
 49  
 50  
 51  
 52  
 53  
 54  
 55  
 56  
 57  
 58  
 59  
 60  
 61  
 62  
 63  
 64  
 65  
 66  
 67  
 68  
 69  
 70  
 71  
 72  
 73  
 74  
 75  
 76  
 77  
 78  
 79  
 80  
 81  
 82  
 83  
 84  
 85  
 86  
 87  
 88  
 89  
 90  
 91  
 92  
 93

```

74 IF (EPSN(X(JJ-1))-EPSN(X(JJ+1))) 35,36,36
35 JJ=JJ-1
36 GO TO 37
37 JJ=JJ+1
38 IF (EPSN(X(JJ-1))-EPSN(X(JJ+1))) 38,39,39
39 DO 40 K=2,JJ
40 I=JJ+2-K
X(I)=X(I-1)
CONTINUE
41 JJJ=JJ+1
DO 41 K=JJ,NN
I=K+1
X(I-1)=X(I)
CONTINUE
C
42 IF (II-1) 46,43,46
43 IF (EPSN(X(NN))-EPSN(X(2))) 50,44,44
44 DO 45 I=2,NN
X(I)=X(I-1)
CONTINUE
45 IF (II+1) 51,47,50
46 IF (EPSN(X(1))-EPSN(X(NN))) 50,50,48
47 C
48 DO 49 K=2,NN
I=NN+2-K
X(I)=X(I-1)
CONTINUE
CONTINUE
RETURN
END

```

```

94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125

```

\*\*\*\*\*

SUBROUTINE FUNCT (N1,B,OBJ,G,GRAD,N,X,ERROR,IP)

SUBROUTINE WHICH SELECTS THE MAXIMUM ERROR  
AND COMPUTES THE OBJECTIVE FUNCTION AND ITS  
GRADIENTS W.R.T. THE VARIABLE PARAMETERS  
IN THE LEAST P-TH SENSE

```

EXTERNAL S,W
DIMENSION B(1), G(1), GRAD(1), X(1), ERROR(1)
EPR(Z)=ERR(Z,S,k,B,GRAD)
OBJP=j.
GRADP=f.
DO 1 K=1,N1
G(K)=0.
CONTINUE
NN=N+1
ERROR(1)=0.
DO 2 I=2,NN
ERROR(I)=ERR(X(I))
CONTINUE
DO 4 I=2,NN
IF (ABS(ERROR(I))-ERROR(1)) 4,4,3
ERROR(I)=AES(ERROR(I))
CONTINUE
DO 6 I=2,NN
Z=X(I)
DEC=ERROR(I)/ERRCR(1)
DEL=AES(DEC)
OBJI=DEL**TP
GRADI=OBJI**(TP-2)*DEC
OBJP=OBJP+OBJI
CALL FAPP (Z,N1,E,APP,GRAD,2)
DO 5 K=1,N1
GRAD(K)=GRADI*W(Z)*GRAD(K)
G(K)=G(K)+GRAD(K)
CONTINUE
CONTINUE
PP=1./TP
OBJ=ERROR(1)*(OB,P**PR)
GRP=OBJP**(PR-1.)
DO 7 K=1,N1
G(K)=GRP*G(K)
CONTINUE
RETURN
END

```

```

1
112
113
114
115
116
117
118
119
120
121
122
123
124
125

```

\*\*\*\*\*

SUBROUTINE GRDCHK (N,X,G,PY,Y,GRAD,np,XP,ERROR,IP,DUM1)

SUBROUTINE WHICH CHECKS THE GRADIENTS  
W.R.T. ALL VARIABLE PARAMETERS

```

DIMENSION X(1), G(1), PY(1), Y(1), GRAD(1), XP(1), ERROR(1), DUM1(11)

```

```

1
112
113
114
115
116
117
118
119
120
121
122
123
124
125

```

```

CALL FUNCT (N,X,F,G,GRAD,NP,XP,ERROR,IP)
DO 1 I=1,N
  DELX=1.E-4*X(I)
  X(I)=X(I)+DELX
  CALL FUNCT (N,X,FNEW,PY,GRAD,NP,XP,ERROR,IP)
  Y(I)=(FNEW-F)/DELX
  DUM1(I)=Y(I)
  X(I)=X(I)-DELX
CONTINUE
DO 2 I=1,N
  IF (ABS(Y(I)).LT.1.E-20) DUM1(I)=1.E-20
  PY(I)=ABS((Y(I)-G(I))/DUM1(I))*100.
CONTINUE
WRITE (6,6)
WRITE (6,7)
WRITE (6,8) (I,X(I),I=1,N)
WRITE (6,9)
DO 3 I=1,N
  WRITE (6,10) G(I),Y(I),PY(I)
CONTINUE
DO 4 I=1,N
  IF (PY(I).GT.10.) GO TO 5
CONTINUE
WRITE (6,11)
RETURN
WRITE (6,12)
CALL EXIT
FORMAT (1H1)
FORMAT (1H0,5X,*GRADIENTS CHECKING*,/6X,18(*-*),//,6X,*GRADIENTS
1 HAVE BEEN CHECKED AT THE FOLLOWING POINT*)
FORMAT (10X,*X(*,I2,*)=,E16.8)
FORMAT (//,1H0,5X,*ANALYTICAL GRADIENTS*,5X,*NUMERICAL GRADIENTS*
1,5X,*PERCENTAGE ERROR*,/)
FORMAT (1H0,5X,3(E16.8,9X))
FORMAT (1H0,//,6X,*GRADIENTS ARE 0. K.*)
FORMAT (1H0,//,6X,*YOUR PROGRAM HAS BEEN TERMINATED BECAUSE GRADIE
1NTS ARE INCORRECT*,/6X,*PLEASE CHECK IT AGAIN*)
END

```

```

1      SUBROUTINE FMNFC (N,X,F,G,H,UNITH,FEST,EPS,MAXFN,IPRINT,IEEXIT,GRAD
2      1,NP,XP,ERROR,IP)
3
4      PURPOSE
5      TO FIND A LOCAL MINIMUM OF A FUNCTION OF SEVERAL VARIABLES
6      ASSUMING THAT ITS GRADIENTS CAN BE CALCULATED EXPLICITLY
7      BY THE METHOD OF FLETCHER
8
9      THE METHOD IS DESCRIBED IN THE FOLLOWING ARTICLE
10     R. FLETCHER, A NEW APPROACH TO VARIABLE METRIC ALGORITHMS,
11     COMP. JOURNAL, VOL.13, 1970, PP.317-322.
12
13     DIMENSION X(1), G(1), H(1), EPS(1), GRAD(1), XP(1), ERROR(1)
14     LOGICAL CONV,UNITH
15     COMMON /BLK1/ K0
16     CALL SECOND (T3)
17     K0=U
18     CALL FUNCT (N,X,F,G,GRAD,NP,XP,ERROR,IP)
19     IF (F.LT.FEST) GC TO 23
20     NFN=1
21     ITN=0
22     STEP=1.
23     IDX=N
24     ING=N+N
25     IH=IDG+N
26     IF (,NOT.UNITH) GO TO 2
27     IJ=IH+1
28     DO 1 I=1,N
29     DO 1 J=I,N
30     H(IJ)=0.
31     IF (I.EQ.J) H(IJ)=1.0
32     IJ=IJ+1
33     CONV=.TRUE.
34     GDX=0.
35     DO 6 I=1,N
36     Z=0.
37     IJ=IH+I
38     IF (I.EQ.1) GO TO 4
39     IJ=I-I
40     DO 3 J=I,I
41     Z=Z-H(IJ)*G(J)
42     IJ=IJ+N-J
43     CONTINUE
44     DO 5 J=I,N
45     Z=Z-H(IJ)*G(J)
46     IJ=IJ+1
47     CONTINUE
48     IF (ABS(Z).GT.EPS(I)) CONV=.FALSE.
49     H(IDX+I)=Z
50     GDX=GDX+G(I)*Z
51     CONTINUE

```

```

IF (IPRINT.EQ.0) GO TO 7
IF (MOD(ITN,IPRINT).NE.0) GO TO 7
CALL SECOND (T4)
TIME=T4-T3
CALL WRITE2 (X,N,G,F,NFNS,ITN,TIME)
IEXIT=1
IF (CONV) GO TO 24
IEXIT=2
IF (GDX.GE.0.) GC TO 24
Z=1
IF (ITN.LT.N.AND.UNITH) Z=STEP
W=2.* (FEST-F)/GDX
IF (W.LT.Z) Z=W
STEP=Z
GDX=GDX*Z
DO 9 I=1,N
H(IDX+I)=H(I)+H(IDX+I)*Z
X(I)=X(I)+H(IDX+I)
CONTINUE
CALL FUNCT (N,X,FP,HGRAD,NP,XP,ERRCR,IP)
IF (FP.LT.FEST) GO TO 23
NFNS=NFNS+1
IEXIT=3
IF (ITN.EQ.MAXFN) GO TO 24
GDX=0.
DO 10 I=1,N
H(IDG+I)=H(I)-G(I)
GDX=GDX+H(I)*H(IDX+I)
CONTINUE
DGDGX=GDX-X-GDX
IF (F.GT.FP-.1001*GDX) GO TO 12
IEXIT=4
IF ((GDX.LT.0.. AND. ITN.GT.N) GO TO 24
Z=3.* (F-Z)*GDX+GDX
W=SORT(1.-GDX/Z*GDX+Z)*ABS(Z)
Z=1.-(GDX+W-Z)/(DGDX*Z.*W)
IF (Z.LT.0.1) Z=0.1
DO 11 I=1,N
X(I)=X(I)-H(IDX+I)
CONTINUE
GO TO 14
F=FP
DO 13 I=1,N
G(I)=H(I)
CONTINUE
IF ((GDX.GT.0.) GO TO 15
GDX=GDX
Z=4.
STEP=Z*STEP
GO TO 9
IF (GDX.LT.0.5*GDX) STEP=2.*STEP
DGHDG=0.
DO 19 I=1,N
Z=J.
IJ=I+I
IF (I.EQ.1) GO TO 17
IJ=I-1
DO 16 J=1,I
Z=Z+H(IJ)*H(IDG+J)
IJ=IJ+N-J
CONTINUE
DO 18 J=I,N
Z=Z+H(IJ)*H(IDG+J)
IJ=IJ+1
CONTINUE
DGHDG=DGHDG+Z*H(IDG+I)
F(I)=Z
CONTINUE
IF (DGHDG.LT.0.0) DGHDG=DGDGX*0.01
IF (GDX.LT.DGHDG) GO TO 21
W=1.2*J+DGHDG/DGDGX
DO 20 I=1,N
H(IDX+I)=W*H(IDX+I)-H(I)
CONTINUE
DGDGX=DGDGX+DGHDG
DGHDG=DGDGX
IJ=IH
DO 22 I=1,N
W=H(IDX+I)/DGDGX
Z=H(I)/DGHDG
DO 22 J=1,N
IJ=IJ+1
H(IJ)=H(IJ)+W*H(IDX+J)-Z*H(J)
ITN=ITN+1
GO TO 2
IEXIT=5
IF (IEXIT.EQ.1) K0=1
IF (IPRINT.EQ.0) RETURN
GO TO (25,26,27,26,28), IEXIT
WPITE (6,30) IEXIT
GO TO 29
WPITE (6,31) IEXIT
GO TO 29
WPITE (6,32) IEXIT
GO TO 29

```

```

28      WRITE (6,33) IEXIT
29      CONTINUE
30      RETURN
C      FORMAT (/,1HO,*IEXIT=*,I2,*CRITERION FOR OPTIMUM (CHANGE IN VECTOR
1 X .LT.EPS) HAS BEEN SATISFIED*)
31      FORMAT (/,1HO,*IEXIT=*,I2,*THERE OF THE FOLLOWING THINGS HAS HAPPENED*
1 ENDED*,/,*X*,*1X*,*EPS CHSEN IS TOO SMALL*,/,*X*,*2. GRADIENTS ARE NOT
2 CORRECT*,/,*X*,*3. MATRIX H GOES SINGULAR*)
32      FORMAT (/,1HO,*IEXIT=*,I2,*MAXIMUM NUMBER OF ALLOWABLE ITERATION H
1 HAS BEEN EXCEEDED*)
33      FORMAT (/,1HO,*IEXIT=*,I2,*FUNCTION VALUE LESS THAN MINIMUM ESTIMA-
END

```

\*\*\*\*\*

```

SUBROUTINE FMFPC (FUNCT,N,X,F,G,EST,EPS,LIMIT,IER,H,IPRINT,GRAD,np,
1,XP,ERROR,IP)
C
C PURPOSE
C   TO FIND A LOCAL MINIMUM OF A FUNCTION OF SEVERAL VARIABLES
C   ASSUMING THAT ITS GRADIENTS CAN BE CALCULATED EXPLICITLY
C   BY THE METHOD OF FLETCHER AND POWELL
C
C   THE METHOD IS DESCRIBED IN THE FOLLOWING ARTICLE
C   R. FLETCHER AND M.J.D. POWELL, A RAPIDLY CONVERGENT
C   DESCENT METHOD FOR MINIMIZATION, COMP. JOURNAL,
C   VOL.6, 1963, PP.163-168.
C
C COMMON /BLK/ K0
C
C   DIMENSIONED CMMY VARIABLES
C   DIMENSION H(1), X(1), G(1), GRAD(1), XP(1), ERROR(1)
C
C   COMPUTE FUNCTION VALUE AND GRADIENT VECTOR FOR INITIAL ARGUMENT
C   K0=0
C   CALL SECOND (T3)
C   CALL FUNCT (N,X,F,G,GRAD,np,XP,ERROR,IP)
C   KOUNT=0
C   NUMF=1
C   CALL SECOND (T4)
C   TIME=T4-T3
C   IF (IPRINT.EQ.0) GO TO 1
C   CALL WRITE2 (X,N,G,F,NUMF,KOUNT,TIME)
C   CONTINUE
C
C   RESET ITERATION COUNTER AND GENERATE IDENTITY MATRIX
C   IFR=0
C   KK=0
C   N2=N+N
C   N3=N2+N
C   N1=N3+1
C   K=N31
C   DO 5 J=1,N
C   H(K)=1.
C   NJ=N-J
C   IF (NJ) 6,6,3
C   DO 4 L=1,NJ
C   KL=K+L
C   H(KL)=0.
C   CONTINUE
C   K=KL+1
C   CONTINUE
C
C   START ITERATION LOOP
C   IF (KOUNT.EQ.0) GO TO 7
C   IF (KK.NE.IPRINT) GO TO 7
C   KK=0
C   CALL SECOND (T4)
C   TIME=T4-T3
C   CALL WRITE2 (X,N,G,F,NUMF,KOUNT,TIME)
C   CONTINUE
C   KOUNT=KOUNT+1
C   KK=KK+1
C
C   SAVE FUNCTION VALUE, ARGUMENT VECTOR AND GRADIENT VECTOR
C   OLF=F
C   DO 11 J=1,N
C   K=N+J
C   H(K)=G(J)
C   K=K+N
C   H(K)=X(J)
C
C   DETERMINE DIRECTION VECTOR H
C   K=J+N3
C   T=0.
C   DO 10 L=1,N
C   T=T-G(L)*H(K)
C   IF (L-J) 8,9,9
C   K=K+N-L
C   GO TO 10
C   K=K+1

```

```

10    CONTINUE
11    H(J)=T
11    CONTINUE
C      CHECK WHETHER FUNCTION WILL DECREASE STEPPING ALONG H.
DO 12 J=1,N
HNRM=HNRM+ABS(H(J))
GMRM=GMRM+ABS(G(J))
DY=DY+H(J)*G(J)
CONTINUE
12
C      REPEAT SEARCH IN DIRECTION OF STEEPEST DESCENT IF DIRECTIONAL
C      DERIVATIVE APPEARS TO BE POSITIVE OR ZERO.
IF (DY) 13,57,57
C      REPEAT SEARCH IN DIRECTION OF STEEPEST DESCENT IF DIRECTIONAL
C      VECTOR H IS SMALL COMPARED TO GRADIENT VECTOR G.
IF (HNRM/GMRM-EPS) 57,57,14
C      SEARCH MINIMUM ALONG DIRECTION H
SEARCH ALONG F FOR POSITIVE DIRECTIONAL DERIVATIVE
FY=F
ALFA=2.*(EST-F)/CY
AMBDA=1.
C      USE ESTIMATE FOR STEPSIZE ONLY IF IT IS POSITIVE AND LESS THAN
1. OTHERWISE TAKE 1. AS STEPSIZE
IF (ALFA) 17,17,15
IF (ALFA-AMBDA) 16,17,17
AMBDA=ALFA
ALFA=0.
C      SAVE FUNCTION AND DERIVATIVE VALUES FOR OLD ARGUMENT
FX=FY
DX=DY
C      STEP ARGUMENT ALONG H
DO 19 I=1,N
X(I)=X(I)+AMBDA*F(I)
CONTINUE
19
C      COMPUTE FUNCTION VALUE AND GRADIENT FOR NEW ARGUMENT
CALL FUNCT(N,X,F,G,GRAD,np,XP,ERROR,IP)
NUMF=NUMF+1
FY=F
C      COMPUTE DIRECTIONAL DERIVATIVE DY FOR NEW ARGUMENT. TERMINATE
C      SEARCH, IF DY IS POSITIVE. IF DY IS ZERO THE MINIMUM IS FOUND
DY=0.
DO 20 I=1,N
DY=DY+G(I)*H(I)
CONTINUE
IF (DY) 21,41,24
C      TERMINATE SEARCH ALSO IF THE FUNCTION VALUE INDICATES THAT
C      A MINIMUM HAS BEEN PASSED
IF (FY-FX) 22,24,24
C      REPEAT SEARCH AND DOUBLE STEPSIZE FOR FURTHER SEARCHES
AMBDA=AMBDA*ALFA
ALFA=AMBDA
C      END OF SEARCH LOOP
C      TERMINATE IF THE CHANGE IN ARGUMENT GETS VERY LARGE
IF (HNRM*AMBDA-1.E10) 18,18,23
C      LINEAR SEARCH TECHNIQUE INDICATES THAT NO MINIMUM EXISTS
ITER=2
GO TO 62
C      INTERPOLATE CUBICALLY IN THE INTERVAL DEFINED BY THE SEARCH
C      ABOVE AND COMPUTE THE ARGUMENT X FOR WHICH THE INTERPOLATION
C      POLYNOMIAL IS MINIMIZED
T=0.
IF (AMBDA) 26,41,25
Z=3.*(FX-FY)/AMBDA+DX+DY
ALFA=AMAX1(ABS(Z),ABS(DX),ABS(DY))
DALFA=Z/ALFA
DALFA=DALFA*DALFA-DX/ALFA*DY/ALFA
IF (DALFA) 57,27,27
W=ALFA*SOPT(DALFA)
ALFA=DY-DX+HW
IF (ALFA) 28,29,28
ALFA=(DY-Z+W)/ALFA
GO TO 30
ALFA=(Z+DY-W)/(Z+DX+Z+DY)
ALFA=ALFA*AMBDA
DO 31 I=1,N
X(I)=X(I)+(T-ALFA)*H(I)

```

```

31      CONTINUE
C      TERMINATE, IF THE VALUE OF THE ACTUAL FUNCTION AT X IS LESS
C      THAN THE FUNCTION VALUES AT THE INTERVAL ENDS, OTHERWISE REDUCE
C      THE INTERVAL BY CHOOSING ONE END-POINT EQUAL TO X AND REPEAT
C      THE INTERPOLATION, WHICH END-POINT IS CHOSEN DEPENDS ON THE
C      VALUE OF THE FUNCTION AND ITS GRADIENT AT X
C
C      NUMF=NUMF+1
C      CALL FUNCT(N,X,F,G,GRAD,np,XP,error,IP)
C      IF (F-FX) 32,32,33
C      IF (F-FY) 41,41,33
32      DALFA=0.
33      DO 34 I=1,N
34      DALFA=DALFA+G(I)*H(I)
CONTINUE
35      IF (DALFA) 35,38,38
36      IF (F-FX) 37,36,38
37      IF (DX-DALFA) 37,41,37
FX=F
DX=DALFA
T=ALFA
AMBDAA=ALFA
GO TO 25
38      IF (FY-F) 40,39,40
39      IF (DY-DALFA) 40,41,40
FY=F
DY=DALFA
AMBDAA=AMBDAA-ALFA
GO TO 24
C
C      TERMINATE, IF FUNCTION HAS NOT DECREASED DURING LAST ITERATION
41      IF (OLDF-F+EPS) 57,42,42
C
C      COMPUTE DIFFERENCE VECTORS OF ARGUMENT AND GRADIENT FROM
C      TWO CONSECUTIVE ITERATIONS
42      DO 43 J=1,N
K=N+J
L(K)=G(J)-H(K)
K=N+K
L(K)=X(J)-H(K)
CONTINUE
C
C      TEST LENGTH OF ARGUMENT DIFFERENCE VECTOR AND DIRECTION VECTOR
C      IF AT LEAST N ITERATIONS HAVE BEEN EXECUTED. TERMINATE, IF
C      BOTH ARE LESS THAN EPS
IER=3
43      IF (KOUNT-N) 47,44,44
T=0.
Z=0.
44      DO 45 J=1,N
K=N+J
W=H(K)
K=K+N
T=T+ABS(H(K))
Z=Z+W*H(K)
CONTINUE
45      IF (HNE1-EPS) 46,46,47
46      IF (T-EPS) 62,52,47
C
C      TERMINATE, IF NUMBER OF ITERATIONS WOULD EXCEED LIMIT
47      IF (KOUNT-LIMIT) 48,55,55
C
C      PREPARE UPDATING OF MATRIX H
48      ALFA=0.
DO 52 J=1,N
K=J+N3
W=0.
DO 51 L=1,N
KL=N+L
W=W+H(KL)*H(K)
IF (L-J) 49,50,50
49      K=K+N-1
GO TO 51
50      K=K+1
CONTINUE
K=N+J
ALFA=ALFA+W*H(K)
W(J)=W
CONTINUE
C
C      REPEAT SEARCH IN DIRECTION OF STEEPEST DESCENT IF RESULTS
C      ARE NOT SATISFACTORY
52      IF (Z*ALFA) 53,2,53
C
C      UPDATE MATRIX H
53      K=N31
DO 54 L=1,N
KL=N2+L
DO 54 NJ=L,N
N2=J
H(K)=H(K)+H(KL)*H(NJ)/Z-H(L)*H(J)/ALFA
54      K=K+1
GO TO 6
C      END OF ITERATION LOOP

```

```

55      NO CONVERGENCE AFTER LIMIT ITERATIONS
56      ICR21
57      IF (KK.NE.IPRINT) GO TO 56
58      CALL WRITE2 (X,N,G,F,NUMF,KOUNT)
59      CONTINUE
60      GO TO 62
C
61      RESTORE OLD VALUES OF FUNCTION AND ARGUMENTS
62      DO 58 J=1,N
63      K=N2+J
64      X(J)=H(K)
65      CONTINUE
66      CALL FUNCT (N,X,F,G,GRAD,np,XP,ERROR,IP)
67      NUMF=NUMF+1
C
68      REPEAT SEARCH IN DIRECTION OF STEEPEST DESCENT IF DERIVATIVE
69      FAILS TO BE SUFFICIENTLY SMALL
70      IF (GNRM-EPS) 61,61,59
C
71      TEST FOR REPEATED FAILURE OF ITERATION
72      IF (IER) 62,60,60
73      IER=1
74      GO TO 2
75      IER=2
76      II=IEP+2
77      IF ((II.EQ.2) KO=1
78      IF ((IPRINT.EQ.3)) RETURN
79      GO TO (63,64,65,66), II
80      WRITE (6,68) IER
81      GO TO 67
82      WRITE (6,69) IER
83      GO TO 67
84      WRITE (6,70) IER
85      GO TO 67
86      WRITE (6,71) IER
87      RETURN
C
88      FORMAT (1HO,*ITER=*,I2,* ERROR IN GRADIENTS CALCULATIONS*)
89      FORMAT (1HO,*ITER=*,I2,* CRITERION FOR OPTIMUM HAS BEEN SATISFIED*)
90      FORMAT (1HO,*ITER=*,I2,* MAXIMUM NUMBER OF ALLOWABLE ITERATIONS HAS
91      BEEN EXCEEDED*)
92      FORMAT (1HO,*ITER=*,I2,* CHANGE IN ARGUMENTS GETS TOO LARGE, LINEAR
93      SEARCH INDICATES THAT NO MINIMUM EXISTS*)
94      END

```

HSP200L //// END OF LIST ////

### APPENDIX III

#### Listings of the common subroutines for FMCLP and FMLPO packages

Add the subroutines INPUT, FINAL, WRITE1 and WRITE2 to both packages, FMCLP and FMLPO. They are similar to ones presented in the package for function minimization [36] programmed by V. K. Jha.

Listings of the subroutines INPUT, FINAL, WRITE1 and WRITE2 follow.

```

SUBROUTINE INPUT (MET,M,MAX,N,IPRINT,IData,EPS1,EST,EPS,XSTRT)
    PRINTS THE INPUT DATA
    FOR THE OPTIMIZATION PROCESS
DIMENSION XSTRT(1), EPS(1)
WRITE (6,5)
IF (MET.NE.1.AND.MET.NE.2) GO TO 4
INDEX=0
GO TO (1,2), MET
1   WRITE (6,6)
GO TO 3
2   WRITE (6,7)
CONTINUE
3   WRITE (6,8) N
WRITE (6,9) MAX
WRITE (6,10) IPRINT
WRITE (6,11) M
WRITE (6,12) XSTRT(1)
IF (MET.EQ.1) WRITE (6,13) (I,XSTRT(I),I=2,N)
IF (MET.EC.1) WRITE (6,14) (I,EPS(I),I=1,N)
IF (MET.EQ.2) WRITE (6,15) (I,EPS(I),I=2,N)
WRITE (6,16) EPS1
WRITE (6,17) EST
RETURN
4   WRITE (6,18)
CALL EXIT
5   FORMAT (1H1,*INPUT DATA*,/,1X,10(*-*),//,1X,*FOLLOWING METHODS HAVE
1E BEEN CALLED*,/)
6   FORMAT (1H0,*NEW FLETCHER METHOD*)
7   FORMAT (1H0,*FLETCHER-POWELL METHOD*)
8   FORMAT (1H0,1X,*NUMBER OF INDEPENDENT VARIABLES*,36(*.*),*N=*,15,
1/)
9   FORMAT (1H0,*MAXIMUM NUMBER OF ALLOWABLE ITERATIONS*,27(*.*),*MAX=*
1*,15,/)
10  FORMAT (1H0,*INTERMEDIATE OUTPUT TO BE PRINTED EVERY IPRINT ITERAT
1IONS*,5(*.*),*IPRINT=*,15,/)
11  FORMAT (1H0,*STARTING POINT TO BE SAME FOR ALL THE METHODS IF M=1*
1*,15(*.*),*M=*,15,/)
12  FORMAT (1H0,*STARTING VALUE FOR VECTOR X(I)*,29(*.*),*XSTRT(1)=*,1E16.8)
13  FORMAT (1H0,59X,*XSTRT(*,I2,*),*,E16.8)
14  FORMAT (1H0,,1H0,*TEST QUANTITIES TO BE USED IN NEW FLETCHER METH
100*12(*.*),*EPS(1)=*,E16.8)
15  FORMAT (1H0,51X,*EPS(*,I2,*),*,E16.8)
16  FORMAT (1H0,,1H0,*TEST QUANTITY TO BE USED IN FLETCHER-POWELL MET
1H0*14(*.*),*EPS1=*,E16.8)
17  FORMAT (1H0,,1H0,*ESTIMATE OF LOWER BOUND ON FUNCTION TO BE MINIM
117*14(*.*),*EST=*,E16.8)
18  FORMAT (1H0,*NONE OF THE OPTIMIZATION METHODS HAVE BEEN CALLED*,/,
11X,*PLEASE CHECK THE VALUE OF MET*,/,1X,*REMAINDER*,/,1X,*MET=1*
2 NEW FLETCHER METHOD WOULD BE CALLED*,/,1X,*MET=2 FLETCHER-POWE
3LL METHOD WOULD BE CALLED*)
END

```

```

*****  

SUBROUTINE FINAL (X,F,N)
    PRINTS THE RESULTS
    FOR THE OPTIMIZATION PROCESS
DIMENSION X(1)
COMMON /BLK/ KO
WRITE (6,3)
IF (KO.EQ.0) GO TO 1
WRITE (6,4)
GO TO 2
1   WRITE (6,5)
CONTINUE
2   WRITE (6,6) F
WRITE (6,7) (I,X(I),I=1,N)
RETURN
3   FORMAT (1H1)
FORMAT (1H0,4X,*FOLLOWING IS THE OPTIMUM SOLUTION*,/,41X,*-----*
1-----*)
4   FORMAT (1H0,44X,*RESULTS AT LAST ITERATION*,/45X,*-----*
1-----*)
5   FORMAT (1H0,/,48X,*F =*,E16.8,/)
6   FORMAT (1H0,44X,*X(*,I2,*),*,E16.8)
END

```

```

*****  

SUBROUTINE WRITE1 (N)
    PRINTS THE INTERMEDIATE RESULTS

```



## APPENDIX IV

### The Chebyshev polynomials

The Chebyshev polynomial of nth order,  $T_n$ , is defined

$$T_n(t) = \cos(n \arccos t), -1 \leq t \leq 1$$

$T_n$  has n zeros

$$t_k = \cos(2k-1) \frac{\pi}{2n},$$

$$= 1 - 2 \sin^2(k-0.5) \frac{\pi}{2n}, \quad k=1,2,\dots,n,$$

and  $n+1$  extrema

$$\begin{aligned} t_k &= \cos \frac{k\pi}{n} \\ &= 1 - 2 \sin^2 \frac{k\pi}{2n}, \quad k=0,1,\dots,n \end{aligned}$$

over the interval  $[-1,1]$ .

$T_n$  may be transformed on any closed interval  $[a,b]$  by introducing the change in variable

$$x = \frac{b-a}{2} t + \frac{b+a}{2}$$

where  $x \in [a,b]$ .

The zeros of  $T_n$  in the transformed interval are in

$$x_k = a + (b-a) \sin^2(k-0.5) \frac{\pi}{2n}, \quad k=1,2,\dots,n,$$

and extrema in

$$x_k = a + (b-a) \sin^2 \frac{k\pi}{2n}, \quad k=0,1,\dots,n.$$

The initial approximation for the Test 1 and Test 2 given in the Chapter V is obtained by rational interpolation where the zeros of the  $(k+1)$ st Chebyshev polynomial are used as supporting points. As a first trial for the points of the independent parameter the extrema of  $T_{n+1}$  were taken. The listing of a FORTRAN program which calculates the initial approximation and the starting point set is presented below.

```

C      DIMENSION A(10),G(10),Y(10),ASTRT(10),DET(100),X(50)
C
100 FORMAT(8I10)
200 FORMAT(//15X,*MATRIX IS SINGULAR IN SUBROUTINE SIMQ*)
300 FORMAT(8F10.5)
READ(5,100)N,L,M
READ(5,300)XA,XB
PI=3.14159
XAB=XB-XA
X(2)=XA
II=1
J=N+1
C
C      COMPUTATION OF AN INITIAL GUESS X(I) OF THE INDEPENDENT
C      VARIABLE IN THE EXTREME POINTS OF THE CHEBYSHEV POLYNOMIAL
C      T(N+1)
C
C      CONST=PI/(2.*J)
DO 10 I=1,J
D=SIN(I*0.5)*CONST
X(I+2)=D*XAB+XA
10 CONTINUE
X(N+3)=XB
C
C      CALCULATION OF THE STARTING VALUES OF THE COEFFICIENTS OF
C      THE RATIONAL FUNCTION BY RATIONAL INTERPOLATION IN THE
C      ZEROS OF THE CHEBYSHEV POLINOMIAL T(K+1)
C
K=L+1+1
15 CONST=PI/(2.*K)
DO 20 I=1,K
Q=SIN((I-0.5)*CONST)
Y(I)=Q*XAB+XA
G(I)=FCT(Y(I))
20 CONTINUE
DO 130 I=1,K
DET(I)=1.
IT=I
IF(L) 115,115,105
105 DO 110 J=1,L
IA=IT
IT=IT+K
110 DET(IT)=DET(IA)*Y(I)
115 IT=IT+K
IF(M-1)130,117,117
117 DET(IT)=-G(I)*Y(I)
IF(M-2)130,118,118
118 DO 120 J=2,M
IB=IT
IT=IT+K
120 DET(IT)=DET(IB)*Y(I)
130 CONTINUE
CALL SIMQ(DET,G,K,KS)
IF(KS)134,134,150
134 L1=L+1
DO 135 I=1,L1
A(I)=G(I)
135 IF(M)137,137,138
138 DO 136 J=1,M
I=L1+J
A(I)=G(I)
136 CONTINUE
137 CONTINUE
DO 30 I=1,K
ASTRT(I)=G(I)
30 CONTINUE
GO TO 151
150 WRITE(6,200)
151 CONTINUE
CALL EXIT

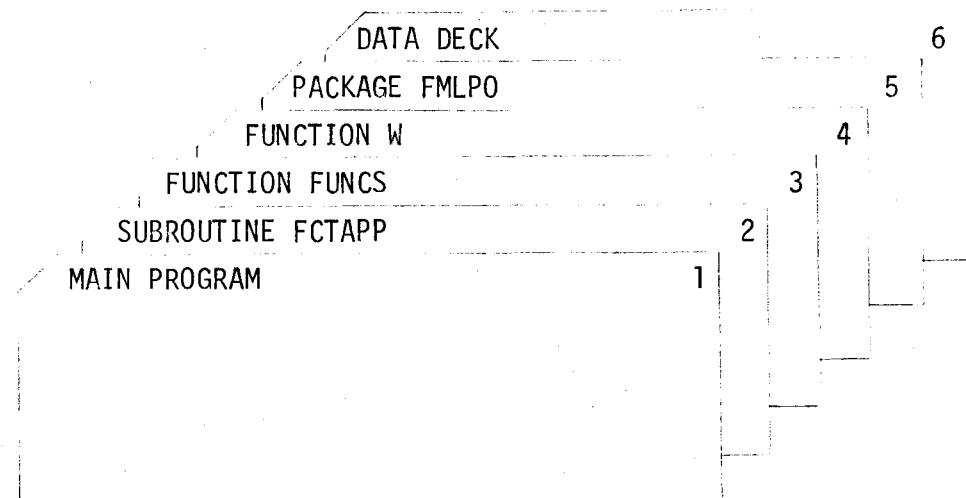
```

## APPENDIX V

### Usage of FMLPO ( Function Minimization for Least pth Objectives )

Purpose To minimize the objective function of  $k$  variables defined as the generalized discrete least  $p$ th objective (6.13), using gradient methods.

How to use Set the input deck as follows:



#### 1. Main program

Write the main program as indicated below:

- (i) Dimension the following arrays

A( $K$ ), ASTRT( $K$ ), G( $K$ ), Y( $K$ ), PY( $K$ ), DUM1( $K$ ), DUM2( $K$ ), GRAD( $K$ ),  
EPS( $K$ ), H( $M$ ), XX(3, NINT), NUMB( $N$ ), INUMB( $N$ ), X( $N$ ), X1( $N$ ),  
ERROR( $N$ ), EHELP( $N$ ), AP( $N$ ), IPA(ITER)

where

K is the number of variable parameters,

M = K·(K+7)/2,

NINT is the number of intervals,

N is the total discrete point set of independent parameter ,  
from all intervals, and

ITER is the maximum number of allowable usage of the optimiz-  
ation method.

- (ii) Call the subroutine FMLPO as follows:

```
CALL FMLPO (A, ASTRT, G, Y, PY, DUM1, DUM2, EPS, H, GRAD,  
NUMB, XX, X, X1, ERROR, EHELP, AP, INUMB, IPA)
```

- (iii) Add STOP and END cards.

## 2. Subroutine FCTAPP

Subroutine which defines the approximating function in each interval and calculates its gradients with respect to variable vector  $\underline{a}$ .

Write subroutine FCTAPP as follows:

- (i) SUBROUTINE FCTAPP (X, K, A, APP, GRAD, IINT, INDIC)  
DIMENSION A(1), GRAD(1).

where X, K, A, IINT and INDIC are input, and APP and GRAD are output variables.

INDIC may have a value 1 or 2 and indicates whether the approximating function in the IINT interval  $APP=F(\underline{A}, X)$ , or all its gradients  $GRAD(I)=\frac{\partial F(\underline{A}, X)}{\partial A(I)}$ ,  $I=1, K$  should be calculated, respectively.

(ii) Write the approximating function APP and all its gradients GRAD(I), I = 1, K for each interval IINT = 1, NINT.

(iii) Add RETURN and END cards.

### 3. Function FUNCS

Function subprogram FUNCS defines upper and lower specified function  $S_u(x)$  and  $S_\ell(x)$ , respectively, in various intervals.

FUNCTION FUNCS (X, IINT)

GO TO (1,2,...,NINT), IINT

1      FUNCS =  $S_{u/\ell}(x)$

RETURN

2

.

.

NINT    FUNCS =  $S_{u/\ell}(x)$

RETURN

END

where IINT = 1, NINT is a current interval.

Note If the upper and lower specified functions are defined for the same set of the independent variable x, consider the common interval (or subinterval) twice.

### 4. Function W

Function subprogram W defines an upper and lower positive weighting function  $W_u(x)$  and  $W_\ell(x)$ , respectively, in various intervals.

Write subprogram W as follows:

```
FUNCTION W(X,IINT)
GO TO (1,2,...,NINT), IINT
1    W = wu/l(x)
RETURN
2
.
.
.
NINT    W = wu/l(x)
RETURN
END
```

where IINT = 1,2,...,NINT is an interval under consideration.

An example which shows how to write these subroutines is the same as in Test 2 in Chapter VI, but with constraints on parameter  $a_2$  such that

$$0 \leq a_2 \leq 1.$$

The listings which the user has to write follow.

```

SUBROUTINE FCTAPP(X,K,A,APP,GRAD,IINT,INDIC)
C
C          SUBROUTINE WHICH CALCULATES APPROXIMATING
C          FUNCTION AND ITS GRADIENTS WITH RESPECT TO
C          VARIABLE PARAMETERS
C
DIMENSION A(1),GRAD(1)
GO TO(100,200),INDIC
100 GO TO (1,1,2,2),IINT
1 APP=A(3)/A(2)*EXP(-A(1)*X)*SIN(A(2)*X)
RETURN
2 APP=A(2)
RETURN
200 GO TO (3,3,4,4),IINT
3 HP1=1./A(2)*EXP(-A(1)*X)
HP2=HP1*SIN(A(2)*X)
HP3=HP1*COS(A(2)*X)
GRAD(1)=-APP*X
GRAD(2)=-1./A(2)*APP+A(3)*X*HP3
GRAD(3)=HP2
RETURN
4 GRAD(1)=0.
GRAD(2)=1.
GRAD(3)=0.
RETURN
END

```

\*\*\*\*\*

FUNCTION FUNCS(X,IINT)

```

C
C          FUNCTION SUBPROGRAM WHICH DEFINES
C          UPPFR AND LOWER
C          SPECIFIED FUNCTION
C
1 GO TO(1,1,2,3),IINT
1 FUNCS=3./20.*EXP(-X)+1./52.*EXP(-5.*X)-EXP(-2.*X)/65.*(
3.*SIN
*(2.*X)+11*COS(2.*X))
RETURN
2 FUNCS=1.0
RETURN
3 FUNCS=0.0
RETURN
END

```

.....

FUNCTION W(X,IINT)

C  
C       FUNCTION SUBPROGRAM WHICH DEFINES  
C       UPPER AND LOWER  
C       WEIGHTING FUNCTION

C  
GO TO (1,1,1,1),IINT  
1 W=1.  
RRETURN  
END

## 6. Data deck

Parameters to be supplied as data are defined below:

K	The number of independent variable parameters $\alpha$ .
NINT	The total number of upper and lower intervals.
NUMB(I), I=1, NINT	Number of subintervals in each interval of independent parameters.
XX(1,I), I=1, NINT	The left end point of ith interval.
XX(2,I), I=1, NINT	The right end point of ith interval.
XX(3,I), I=1, NINT	Numbers in floating point which supply informations on the specified function in the ith interval:  set XX(3,I)=1. for the upper specification and  set XX(3,I)=-1. for the lower specification.
IREAD	Integer which denotes whether or not the discrete set of points in each interval will be read.  If IREAD=0 the discrete set of points will be set equidistantly in each interval with NUMB(I) sub-intervals in the ith interval.  If IREAD=1 the discrete point set will be read.

$X(I), I=1, \sum_{I=1}^{NINT} (\text{NUMB}(I)+1)$	Discrete point set of the independent parameter.
KSI	The artificial margin $\xi$ .
ASTRT(I), I=1, K	Starting values for a vector containing the K variable parameters.
IGRDCH	Gradients to be checked if IGRDCH=1; it should be set to 0 if gradients are not to be checked.
MET	Optimization method to be called:  if MET=1 Fletcher method will be called;  if MET=2 Fletcher-Powell method will be called.
MAX	Maximum number of permissible iterations.
ITER	Has already been defined in the main program as a length of the working array.
IPA(I), I=1, ITER	Vector containing the values of p for different least pth objective.
IOPT	Denotes how many times the optimization is repeated with different starting points and/or different optimization techniques.
IPRINT	Intermediate output is printed out every IPRINT iterations; it should be set to 0 if no intermediate output is desired.

IDATA	Input data is printed out if IDATA=1; it should be set to 0 if input data is not to be printed out.
EST	Minimum estimated value of the objec- tive function.
EPS(I), I=1, K	Small test quantities used by the Fletcher method.
EPS1	Small test quantity used by the Fletcher-Powell method.
DIF	Small test quantity used by the sub- routine FMLPO.

Setting up the data deck

Conditions	Number of cards	Parameters	Type	Format
-	1	K,NINT,IOPT,ITER,IREAD, IGRDCH	INTEGER	6I10
-	I=1, NINT	NUMB(I), (XX(J,I), J=1,3)	1 INTEGER, 3 REALS	I10, 3E16.8
IREAD=1	As many as required by NUMB(I), I=1, NINT	X(J), J=1 (NUMB(I)+1)	REAL	5E16.8
-	1	EST, DIF, KSI	REAL	3E16.8
+ -	As many as required by K	ASTRT(I), I=1,K	REAL	5E16.8
MET=1	1 As many as required by K	MET, MAX, IPrint, IDATA EPS(I), I=1,K	INTEGER REAL	4I10 5E16.8
MET=2	1	EPSI	REAL	5E16.8
-	As many as required by ITER	IPA(I), I=1, ITER	INTEGER	8I10

Recommended values for some of the parameters

MAX = 100;

EPS(I), I=1,K, each  $10^{-6}$ ,

DIF =  $10^{-4}$

EST      A lower bound of the minimum value of the objective function may be obtained from physical reasons. If the true minimum is not known, for the case when the specification is violated EST=0 is convenient, and when the specification is satisfied choose EST sufficiently negative.

For the example for which listings were shown the data deck is presented below. The Fletcher-Powell optimization method is called.

3	4	1	5	0
50	0.F	00	10.F	00
50	0.F	00	10.F	00
0	2.0F	1	2.0E	1
0	3.0F	1	3.0E	1
-1.0F	00	1.0F	-4	20.0E
1.	F00	1.	E00	1.
2	100		1	
1.0F	-6		1	
2	10	100	1000	10000

## APPENDIX VI

### **Listing of the Package FMLPO**

```

CCCCC
      SUBROUTINE FMLPO (A,XSTART,G,Y,PY,DUM1,DUM2,EPS,H,GRAD,NUMB,XX,X,X1
1,ERPOS,EHELP,AP,INUMB,IPA)
      SUBROUTINE WHICH COORDINATES THE OTHER
      SUBROUTINES IN THE PACKAGE FMLPO

      EXTERNAL FUNCS,W,FCT,FUNGT
      LOGICAL CONVUNITH
      DIMENSION A(1),G(1),Y(1),PY(1),XSTART(1),DUM1(1),DUM2(1),EPS
1(1),H(1),GRAD(1),NUMB(1),XX(3,1),X(1),X1(1),ERPOR(1),EHELP
2(1),AP(1),INUME(1),IPA(1)
      COMMON /BLK1/KO
      EPR(2)=EPSNP(Z,IINT,FCT,W,A,N1,GRAD,APP,PSI,XX,1)
      UNITH=.TRUE.
      READ (5,36) N1,NINT,IOPT,ITER,IREAD,IGROCH
      DO 1 I=1,NINT
      READ (5,38) NUMB(I),(XX(J,I),J=1,3)
      CONTINUE
      K=0
      IF (IREAD.EQ.0) IREAD=2
      GO TO 2,4, IREAD
      DO 3 J=1,NINT
      K=K+1
      KL=K+NUMB(J)
      READ (5,40) (X(I),I=K,KL)
      K=KL
      READ (5,43) EST,EIF,PSI
      WRITE (6,37) PSI
      WRITE (6,44) EST
      READ (5,45) (XSTART(I),I=1,N1)
      DO 5 I=1,N1
      A(I)=XSTART(I)
      K=0
      GO TO 6,9, IREAD
      DO 8 J=1,NINT
      IINT=J
      WRITE (6,49) IINT
      K=K+1
      KL=K+NUMB(J)
      DO 7 I=K,KL
      EROR(I)=ERR(X(I))
      L=I-K+1
      WRITE (6,48) L,X(I),ERR(I)
      EHELP(I)=ERR(I)+XX(3,J)
      AP(I)=APP
      CONTINUE
      K=KL
      CONTINUE
      GO TO 11
      DO 10 J=1,NINT
      IINT=J
      WRITE (6,49) IINT
      L=NUMB(J)+1
      IF (NUMB(J).EQ.0) Z=XX(1,J)
      DO 12 I=1,L
      IF (NUMB(J).GT.0) Z=XX(1,J)+(XX(2,J)-XX(1,J))*(I-1)/NUMB(J)
      ER=ERR(Z)
      WRITE (6,48) I,Z,ER
      K=K+1
      EROR(K)=ERR(Z)
      EHELP(K)=ERR(K)*XX(3,J)
      X(K)=Z
      AP(K)=APP
      EMAX=EHELP(1)
      DO 13 I=2,M
      EMAX=AMAX1(EMAX,EHELP(M))
      CONTINUE
      WRITE (6,46)
      WRITE (6,47)
      WRITE (6,45) EMAY
      CALL FPRO (FCT,W,A,N1,K,GRAD,APP,PSI,2,NUMB,XX,X,X1,ERPOR,EHELP,AP
1,EMAX,N,INUMB,NINT,IP)
      WRITE (6,42)
      WRITE (6,43)
      WRITE (6,48) (J,X1(J),ERROR(J),J=1,N)

      DATA FOR THE OPTIMALITY
      FOR THE OPTIMIZATION METHOD USED

      DO 35 K=1,IOPT
      KRP=1
      IF (K-1).LT.14,14,13
      READ (5,40) (XSTART(I),I=1,N1)
      READ (5,36) MET,MAX,IP,INT,INDATA
      IF (MET.EQ.1) READ (5,40) EPS(I),I=1,N1
      IF (MET.EQ.2) READ (5,40) EPS1
      READ (5,36) (IPA(I),I=1,ITER)
      DO 34 KK=1,ITER

      OPTIMIZATION

      IP=IPA(KK)
      IF (KK.GT.1) FF=F
      IF (KR.EQ.0) GO TO 15
      DO 16 I=1,N1
      A(I)=XSTART(I)

```

```

16    CONTINUE
17    IF (IGRICH,NE.1) GO TO 17
18    CALL GFJGJK (N1,A,G,PY,Y,GRAD,APP,PSI,NUMB,XX,X,X1,EROR,EHELP,AP,
1EMAX,N,INUMB,NINT,IP,DUM1)
19    IF (KRDAT,EQ.0) GO TO 18
20    CALL INPUT (MET,MAX,N1,IPRINT,IODATA,EPSS1,EST,EPS,XSTR)
21    IF (MET,EQ.0) MET=4
22    INDEX=0
23    GO TO (19,24,30,29), MET
24    IF (IPRINT,EQ.1) GO TO 20
25    CALL WRITER (1)
26    CALL SECOND (T1)
27    IF (KR,NE.0) GO TO 22
28    DO 21 I=1,N1
29    A(I)=DUM1(I)
30    CONTINUE
31    CALL FMFPG (N1,A,F,G,H,UNITH,EST,EPSS1,MAX,IPRINT,IFXIT,GRAD,APP,PSI
1NU14B,XX,X,X1,ERRCR,EHELP,AP,EMAX,N,INUMB,NINT,IP)
32    DO 23 I=1,N1
33    DUM2(I)=A(I)
34    CONTINUE
35    CALL SECOND (T2)
36    CALL FINAL (A,F,N1)
37    T=T2-T1
38    WRITE (6,39) T
39    GO TO 23
40    IF (IPRINT,EQ.0) GO TO 25
41    CALL WRITER (2)
42    CALL SECOND (T1)
43    IF (KR,NE.0) GO TO 27
44    DO 26 I=1,N1
45    A(I)=DUM2(I)
46    CONTINUE
47    CALL FMFPG (FUNGT,N1,A,F,G,EST,EPSS1,MAX,IER,H,IPRINT,GRAD,APP,PSI
1NU14B,XX,X,X1,ERRCR,EHELP,AP,EMAX,N,INUMB,NINT,IP)
48    DO 28 I=1,N1
49    DUM2(I)=A(I)
50    CONTINUE
51    CALL SECOND (T2)
52    CALL FINAL (A,F,N1)
53    T=T2-T1
54    WRITE (6,39) T
55    INDEX=INDEX+1
56
57    KP=0
58    WRITE (6,41) IP
59    WRITE (6,44)
60    KN=0
61    KO=0
62    DO 32 J=1,NINT
63    INT=J
64    WRITE (6,49) INT
65    KO=KO+1
66    KL=KO+NUMB(J)
67    DO 31 I=KO,KL
68    L=I-KO+1
69    ER=ER(X(I))
70    WRITE (6,48) L,X(I),ER
71    KN=KN+1
72    CONTINUE
73    KO=KL
74    CONTINUE
75    WRITE (6,42)
76    WRITE (6,43)
77    WRITE (6,48) (J,X1(J),ERROR(J),J=1,N)
78    WRITE (6,46)
79    WRITE (6,47)
80    WRITE (6,45) EMAX
81    IGROCH=IGROCH+2
82    IF (KK-1) 30,33,33
83    FTST=ABS ((FF-F)/FF)
84    IF (FTST.LT.DIF) GO TO 35
85    CONTINUE
86    CONTINUE
87    RETURN
C
88    FORMAT (9I10)
89    FORMAT (1H1,2X,*KSI=*,E23.12//++)
90    FORMAT (1H1,3X,16.8)
91    FORMAT (1H0,7/25X,*EXECUTION TIME IN SECONDS=*,F10.5)
92    FOR MAT (5E15.8)
93    FORMAT (1H1,12X,*P = *,I7//++)
94    FORMAT (//2IX,*ERRP USFD IN*)
95    FORMAT (40X,*OBJECTIVE FUNCTION*/)
96    FORMAT (8X,N*6X*INDEPENDENT VARIABLE*,8X,*ERRORS*)
97    FORMAT (13X,E2.1,12)
98    FORMAT (13X,E2.1,12)
99    FORMAT (//20X,*VALUE OF*)
100   FORMAT (18X,*MAXIMUM ERROR*)
101   FORMAT (19X,3X,2E23.12)
102   FORMAT (//20X,* INTERVAL *,I2/)

END

```

\*\*\*\*\*

```

FUNCTION FCT (Z,FUNCS,W,IINT,PSI,XX)
  FUNCTION SUBPROGRAM WHICH DEFINES
  MODIFIED UPPER AND LOWER
  SPECIFIED FUNCTION
  EXTERNAL FUNCS,W
  DIMENSION XX(3,1)
  FCT=FUNCS(Z,IINT)+PSI*XX(3,IINT)/W(Z,IINT)
  RETURN
END

*****



FUNCTION EPSNP (Z,IINT,FCT,W,A,N1,GRAD,APP,PSI,XX,IPOINT)
  FUNCTION SUBPROGRAM WHICH CALCULATES
  UPPER AND LOWER WEIGHTED ERROR FUNCTION
  EXTERNAL FUNCS,W,FCT
  DIMENSION A(1), GRAD(1), XX(3, 1)
  IF (IPPOINT) 1,2,1
1  CONTINUE
  CALL FCTAPP (Z,N1,A,APP,GRAD,IINT,1)
2  CONTINUE
  IF (PSI) 3,4,3
3  EPSNP=(APP-FCT(Z,FUNCS,W,IINT,PSI,XX))*W(Z,IINT)
  RETURN
4  EPSNP=(APP-FUNCS(Z,IINT))*W(Z,IINT)
  RETURN
END

*****



SUBROUTINE ERRO (FCT,W,A,N1,K,GRAD,APP,PSI,INDIC,NUMB,XX,X,X1,ERRO
1R,EHELP,AP,EMAX,N,INUMB,NINT,IP)
  SUBROUTINE WHICH SELECTS THE WEIGHTED
  ERROR FUNCTION OF INTEREST FOR
  THE OBJECTIVE FUNCTION
  EXTERNAL FUNCS,W,FCT
  DIMENSION A(1), GRAD(1), NUMB(1), XX(3,1), X(1), X1(1), EPROR(1),
1EHELP(1), AP(1), INUMB(1)
  EPROR(1)=EPSNP(Z,IINT,FCT,W,A,N1,GRAD,APP,PSI,XX,IPOINT)
  GO TO (1,9), INDIC
1  CONTINUE
  IPOINT=1
  K=0
  KL=3
  DO 7 J=1,NINT
    IINT=J
    IF (J.EQ.1) GO TO 2
    KL=KL+1
    L=NUMB(J)+1
    DO 6 I=1,L
      K=K+1
      IF (J.EQ.1) GO TO 5
      DO 4 KK=1,KL
        IF (X(K)-X(KK)) 4,3,4
        AP(K)=AP(KK)
        APP=AP(K)
        IPOINT=3
        GO TO 3
5     CONTINUE
      EPROR(K)=ERR(X(K))
      EHELP(K)=ERROR(K)*XX(3,J)
      IF (IPPOINT.NE.0) AP(K)=APP
      IPOINT=1
6     CONTINUE
7     CONTINUE
    EMAX=EHELP(1)
    DO 8 M=2,K
      EMAX=AMAX1(EMAX,EHELP(M))
8     CONTINUE
    IF (EMAX) 10,11,11
10    IP=-IABS(IP)
    GO TO 12
11    IP=IABS(IP)
    K=0
    N=0
    INUMB(1)=0
    DO 16 J=1,NINT
      IINT=J
      L=NUMB(J)+1
      DO 15 I=1,L
        K=K+1

```

```

13 IF (IP) 14,13,13
14 IF (EHELP(K)) 15,14,14
N=N+1
X1(N)=X(K)
ERROR(N)=ERROP(K)
EHELP(N)=AP(K)
CONTINUE
INUMB(J+1)=N
CONTINUE
RETURN
END

.....  

CCCCC
15 SUBROUTINE FUNGT (N1,A,IP,G,GRAD,APP,PSI,NUMB,XX,X,X1,ERROR,EHELP
1,AP,EMAX,N,INUMB,NINT,IP)
16
17 SUBROUTINE WHICH COMPUTES THE OBJECTIVE FUNCTION
18 AND ITS GRADIENTS W.R.T. THE VARIABLE PARAMETERS
19 IN THE LEAST P-TH SENSE  

20
21 EXTERNAL FUNCS,W,FCT
22 DIMENSION A(1),GRAD(1),NUMB(1),XX(3,1),X(1),X1(1),ERROR(1),
23 EHELP(1),AP(1),INUMB(1),G(1)
24 OBJP=0.
25 GRADP=0.
26 DO 1 K=1,N1
27 G(K)=0.
28 CONTINUE
29 CALL ERRO (FCT,W,A,N1,K,GRAD,APP,PSI,1,NUMB,XX,X,X1,ERROR,EHELP,AP
1,EMAX,N,INUMB,NINT,IP)
30 DO 7 I=1,N
31 Z=X1(I)
32 DEL=ERROR(I)/EMAX
33 OBJI=DEL**IP
34 GRADI=DEL**(IP-1)
35 OBJP=OBJP+OBJI
36 DO 4 J=1,NINT
37 IF (I-INUMB(J+1)) 2,2,4
38 IF (I-INUMB(J)) 4,4,3
39 IINT=J
40 GO TO 5
41 CONTINUE
42 CONTINUE
43 APP=EHELP(I)
44 CALL FCT APP (Z,N1,A,APP,GRAD,IINT,2)
45 DO 6 K=1,N1
46 GRAD(K)=GRADI*W(Z,IINT)*GRAD(K)
47 G(K)=G(K)+GRAD(K)
48 CONTINUE
49 CONTINUE
50 PR=1./IP
51 OBJ=EMAX*(OBJP**PR)
52 GRP=OBJP**(PR-1.)
53 DO 8 K=1,N1
54 G(K)=GRP*G(K)
55 CONTINUE
56 RETURN
57 END

.....  

CCCCC
58 SUBROUTINE GROSHK (N,X,G,PY,Y,GRAD,APP,PSI,NUMB,XX,XP,X1,ERROR,EHE
59 1,LP,AP,EMAX,NP,INUMB,NINT,IP,DUM1)
60
61 SUBROUTINE WHICH CHECKS THE GRADIENTS
62 W.R.T. ALL VARIABLE PARAMETERS  

63
64 DIMENSION X(1),G(1),PY(1),Y(1),GRAD(1),NUMB(1),XX(3,1),XP(1
65 1),X1(1),ERROR(1),EHELP(1),AP(1),INUMB(1),DUM1(1)
66 CALL FUNGT (N,X,F,G,GRAD,APP,PSI,NUMB,XX,XP,X1,ERROR,EHELP,AP,EMAX
1,LP,INUMB,NINT,IP)
67 DO 1 I=1,N
68 DELX=X(1)-4*X(I)
69 X(I)=X(I)+DELX
70 CALL FUNGT (Y,X,FNEW,PY,GRAD,APP,PSI,NUMB,XX,XP,X1,ERROR,EHELP,AP,
71 1,EMAX,NP,INUMB,NINT,IP)
72 Y(I)=(FNEW-F)/DELX
73 DUM1(I)=Y(I)
74 X(I)=X(I)-DELX
75 CONTINUE
76 DO 2 I=1,N
77 IF (ABS(Y(I)).LT.1.E-20) DUM1(I)=1.E-20
78 PY(I)=A3S((Y(I)-G(I))/DUM1(I))*100.
79 CONTINUE
80 WRITE (6,6)
81 WRITE (6,7)
82 WRITE (6,8) (I,X(I),I=1,N)
83 WRITE (6,9)
84 DO 3 I=1,N

```

```

      WRITE (6,10) G(I),Y(I),PY(I)
3    CONTINUE
      DO 4 I=1,N
      IF (PY(I).GT.1.) GO TO 5
4    CONTINUE
      WRITE (6,11)
      RETURN
5    WRITE (6,12)
      CALL EXIT
6
      FORMAT (1H1)
7    FORMAT (1H0,5X,*GRADIENTS CHECKING*,/,6X,18(*-*),//,6X,*GRADIENTS
1 HAVE BEEN CHECKED AT THE FOLLOWING POINT*)
8    FORMAT (1H0,5X(*,I2,*)=*,E16.8)
9    FORMAT (//,1H0,5X,*ANALYTICAL GRADIENTS*,5X,*NUMERICAL GRADIENTS*
1 ,5X,*PERCENTAGE ERROR*,/)
10   FORMAT (1H0,5X,3(E16.8,9X))
11   FORMAT (1H0,/,6X,*GRADIENTS ARE O. K.*)
12   FORMAT (1H0,/,6X,*YOUR PROGRAM HAS BEEN TERMINATED BECAUSE GRADIE
1 NTS ARE INCORRECT*,/6X,*PLEASE CHECK IT AGAIN*)
      END

```

```

*****
SUBROUTINE FMNFG (N,X,F,G,H,UNITH,FEST,EPNS,MAXFN,IPRINT,IEXIT,GRAD
1,APP,PSI,NUMB,XX,XP,X1,ERROR,EHELP,AP,EMAX,NP,INUMB,NINT,IP)
C
      PURPOSE
      TO FIND A LOCAL MINIMUM OF A FUNCTION OF SEVERAL VARIABLES
      ASSUMING THAT ITS GRADIENTS CAN BE CALCULATED EXPLICITLY
      BY THE METHOD OF FLETCHER
      THE METHOD IS DESCRIBED IN THE FOLLOWING ARTICLE
      R. FLETCHER, A NEW APPROACH TO VARIABLE METRIC ALGORITHMS,
      COMP. JOURNAL, VOL.13, 1970, PP.317-322.
C
      DIMENSION X(1), C(1), H(1), EPNS(1), GRAD(1), NUMB(1), XX(3,1), XP(
11), X1(1), ERROR(1), EHELP(1), AP(1), INUMB(1)
      LOGICAL CONV,UNITH
      COMMON /BLK/KO
      CALL SECOND (T3)
      KO=3
      CALL FUNGT (N,X,F,G,GRAD,APP,PSI,NUMB,XX,XP,X1,ERROR,EHELP,AP,EMAX
1,IP,INUMB,NINT,IP)
      IF (F.LT.FEST) GC TO 23
      NFNS=1
      ITN=0
      STEP=1.
      I0X=N
      IDG=N+N
      IH=IDG+N
      IF (.NOT.UNITH) GO TO 2
      IJ=IH+1
      DO 1 I=1,N
      DO 1 J=I,N
      F(IJ)=0.
      IF (I.EQ.J) H(IJ)=1.0
1     IJ=IJ+1
      CONV=.TRUE.
      GDX=0.
      DO 6 I=1,N
      Z=J
      IJ=IH+I
      IF (I.EQ.1) GO TO 4
      II=I-1
      DO 3 J=1,II
      Z=Z-H(IJ)*G(J)
      IJ=IJ+N-J
      CONTINUE
      DO 5 J=I,N
      Z=Z-H(IJ)*G(J)
      IJ=IJ+1
      CONTINUE
      IF ((ABS(Z).GT.EPNS(I))) CONV=.FALSE.
      H(I0X+I)=Z
      GDX=GDX+G(I)*Z
      CONTINUE
C
      IF (IPRINT.EQ.0) GO TO 7
      IF (MOD(ITN,IPRINT).NE.0) GO TO 7
      CALL SECOND (T4)
      TIME=T4-T3
      CALL WRITE2 (X,N,G,F,NFNS,ITN,TIME)
      IEXIT=1
      IF (CONV) GO TO 24
      IEXIT=2
      IF (GDX.GE.0.) GC TO 24
      Z=1.
      IF (ITN.LT.N.AND.UNITH) Z=STEP
      W=2.*(FEST-F)/GDX
      IF (W.LT.Z) Z=W
      STEP=Z
      GDX=GDX*Z

```

```

DO 9 I=1,N
H(IDX+I)=H(IDX+I)*Z
X(I)=X(I)+H(IDX+I)
CONTINUE
1 CALL FUMGT(N,X,FP,H,GRAD,APP,PSI,NLMB,XX,XP,X1,LERROR,EHELP,AP,LMA
1 IF (FP.LT.FEST) GO TO 23
1 NFNS=NFNS+1
1 IEXIT=3
1 IF (ITN.EQ.MAXFN) GO TO 24
1 GPOX=.*
1 DO 10 I=1,N
1 H(IDG+I)=H(I)-G(I)
1 GPOX=GPOX+H(I)*H(IDX+I)
1 CONTINUE
1 GDX=GPOX-GDX
1 IF (F.GT.FP-.0001*GDX) GO TO 12
1 IEXIT=4
1 IF (GPOX.LT.1.*AND.ITN.GT.N) GO TO 24
1 Z=3.*(FP-FP)+GPOX+GDX
1 W=SQRT(1.-GDX/Z*GPOX/Z)*ABS(Z)
1 Z=1.-(GPOX*W-Z)/(DGDX+2.*W)
1 IF (Z.LT.0.1) Z=.1
1 DO 11 I=1,N
1 X(I)=X(I)-H(IDX+I)
1 CONTINUE
1 GO TO 14
12 F=FP
12 DO 13 I=1,N
12 G(I)=H(I)
13 CONTINUE
13 IF (DGDX.GT.0.) GO TO 15
12 GDX=GPOX
12 Z=4.
14 STEP=Z*STEP
14 GO TO 8
15 IF (GPOX.LT.0.5*GDX) STEP=2.*STEP
15 DGHDG=0.
15 DO 19 I=1,N
15 Z=0.
15 IJ=IH+I
15 IF (I.EQ.1) GO TO 17
15 II=I-1
15 DO 16 J=1,II
15 Z=Z+H(IJ)*H(IDG+J)
15 IJ=IJ+1-N-J
16 CONTINUE
16 DO 18 J=I,N
16 Z=Z+H(IJ)*H(IDG+J)
16 IJ=IJ+1
16 CONTINUE
16 DGHDG=DGHDG+Z*H(IDG+I)
16 H(I)=Z
16 CONTINUE
16 IF (DGHDG.LT.1.0) DGHDG=DGDX*.01
16 IF (DGDX.LT.DGHDG) GO TO 21
16 W=1.0+DGHDG/DGDX
16 DO 20 I=1,N
16 H(IDX+I)=W*H(IDX+I)-H(I)
16 CONTINUE
16 DGDX=DGDX+DGHDG
16 DGHDG=DGDX
21 IJ=IH
21 DO 22 I=1,N
21 W=H(IDX+I)/DGDX
21 Z=H(I)/DGHDG
21 DO 22 J=I,N
21 IJ=IJ+1
21 H(IJ)=H(IJ)+W*H(IDX+J)-Z*H(J)
21 ITN=ITN+1
21 GO TO 2
23 IEXIT=5
24 IF (IEXIT.EQ.1) K0=1
24 IF (IPPIN.EQ.0) RETURN
24 GO TO (25,26,27,26,28), IEXIT
25 WPITE(5,30) IEXIT
25 GO TO 29
26 WPITE(6,31) IEXIT
27 WPITE(6,32) IEXIT
28 WPITE(6,33) IEXIT
29 CONTINUE
29 RETURN
C
30 FORMAT (/1HO,*IEXIT=*,I2,*CPIITERION FOR OPTIMUM (CHANGE IN VECTOR
30 1 X .LT.EPS) HAS BEEN SATISFIED*)
31 FORMAT (/1HO,*IEXIT=*,I2,*EITHER OF THE FOLLOWING THINGS HAS HAPP
31 1ENED*,/9X,*1. EFS CHOSEN IS TOO SMALL*/9X,*2. GRADIENTS ARE NOT
32 2 CORRECT*/9X,*3. MATRIX H GOES SINGULAR*)
32 FORMAT (/1HO,*IEXIT=*,I2,*MAXIMUM NUMBER OF ALLOWABLE ITERATION H
32 HAS BEEN EXCEEDED*)
33 FORMAT (/1HO,*IEXIT=*,I2,*FUNCTION VALUE LESS THAN MINIMUM ESTIMA
33 END

```

```

.....  

SUBROUTINE FMFFG (FUNGT,N,X,F,G,EST,EPS,LIMIT,IER,H,IPRINT,GRAD,AP  

1P,PSI,NUMB,XX,XP,X1,ERROR,EHELP,AP,EMAX,NP,INUMB,NINT,IP) H 1  

C  

C PURPOSE H 1  

C TO FIND A LOCAL MINIMUM OF A FUNCTION OF SEVERAL VARIABLES H 1  

C ASSUMING THAT ITS GRADIENTS CAN BE CALCULATED EXPLICITLY H 1  

C BY THE METHOD OF FLETCHER AND POWELL H 1  

C  

C THE METHOD IS DESCRIBED IN THE FOLLOWING ARTICLE H 1  

C R. FLETCHER AND M.J.D. POWELL, A RAPIDLY CONVERGENT H 1  

C DESCENT METHOD FOR MINIMIZATION, COMP. JOURNAL, H 1  

C VOL.6, 1963, PP.163-168. H 1  

C  

C COMMON /BLK/ KO H 1  

C DIMENSION H(1), X(1), G(1), GRAD(1), NUMB(1), XX(3,1), XP(1), X1(1 H 1  

C 1), ERROR(1), EHELP(1), AP(1), INUMB(1) H 1  

C  

C COMPUTE FUNCTION VALUE AND GRADIENT VECTOR FOR INITIAL ARGUMENT H 1  

C KO=0 H 1  

C CALL SECOND (T3) H 1  

C CALL FUNGT (N,X,F,G,GRAD,APP,PSI,NUMB,XX,XP,X1,ERROR,EHELP,AP,EMAX H 1  

C 1, NP, INUMB, NINT, IP) H 1  

C KOUNT=0 H 1  

C NUMF=1 H 1  

C CALL SECOND (T4) H 1  

C TIME=T4-T3 H 1  

C IF (IPRINT.EQ.0) GO TO 1 H 1  

C CALL WRITE2 (X,N,G,F,NUMF,KOUNT,TIME) H 1  

C CONTINUE H 1  

C  

C RESET ITERATION COUNTER AND GENERATE IDENTITY MATRIX H 1  

C IER=0 H 1  

C KK=0 H 1  

C N2=N+N H 1  

C N3=N2+N H 1  

C N31=N3+1 H 1  

C K=N31 H 1  

C DO 5 J=1,N H 1  

C H(K)=1. H 1  

C NJ=N-J H 1  

C IF (NJ) 6,6,3 H 1  

C 3 DO 4 L=1,NJ H 1  

C KL=K+L H 1  

C H(KL)=0. H 1  

C 4 CONTINUE H 1  

C K=KL+1 H 1  

C 5 CONTINUE H 1  

C  

C START ITERATION LOOP H 1  

C IF (KOUNT.EQ.1) GO TO 7 H 1  

C IF (KK.NE.IPOINT) GO TO 7 H 1  

C KK=0 H 1  

C CALL SECOND (T4) H 1  

C TIME=T4-T3 H 1  

C CALL WRITE2 (X,N,G,F,NUMF,KOUNT,TIME) H 1  

C 7 CONTINUE H 1  

C KOUNT=KOUNT+1 H 1  

C KK=KK+1 H 1  

C  

C SAVE FUNCTION VALUE, ARGUMENT VECTOR AND GRADIENT VECTOR H 1  

C OLOF=F H 1  

C DO 11 J=1,N H 1  

C K=N+J H 1  

C H(K)=G(J) H 1  

C K=K+N H 1  

C H(K)=X(J) H 1  

C  

C DETERMINE DIRECTION VECTOR H H 1  

C K=N3 H 1  

C T=0 H 1  

C DO 10 L=1,N H 1  

C T=T-G(L)+H(K) H 1  

C IF (L-J) 8,9,9 H 1  

C 8 K=K+N-L H 1  

C GO TO 10 H 1  

C 9 K=K+1 H 1  

C 10 CONTINUE H 1  

C H(J)=T H 1  

C 11 CONTINUE H 1  

C  

C CHECK WHETHER FUNCTION WILL DECREASE STEPPING ALONG H. H 1  

C DY=0. H 1  

C HNRM=2. H 1  

C GNRM=0. H 1  

C  

C CALCULATE DIRECTIONAL DERIVATIVE AND TESTVALUES FOR DIRECTION H 1  

C VECTOR H AND GRADIENT VECTOR G. H 1  

C DO 12 J=1,N H 1  

C HNRM=HNRM+ABS(H(J)) H 1  

C GNRM=GNRM+ABS(G(J)) H 1  

C DY=DY+H(J)*G(J) H 1

```

```

12    CONTINUE          H 92
C      REPEAT SEARCH IN DIRECTION OF STEEPEST DESCENT IF DIRECTIONAL H 93
C      DERIVATIVE APPEARS TO BE POSITIVE OR ZERO. H 94
C      IF (DY) 13,57,57 H 95
C      REPEAT SEARCH IN DIRECTION OF STEEPEST DESCENT IF DIRECTION H 96
C      VECTOR H IS SMALL COMPARED TO GRADIENT VECTOR G. H 97
13      IF (HNRM/GNRM-EPS) 57,57,14 H 98
C      SEARCH MINIMUM ALONG DIRECTION H H 99
C      SEARCH ALONG H FOR POSITIVE DIRECTIONAL DERIVATIVE H 100
14      FY=F H 101
      ALFA=2.* (EST-F)/CY H 102
      AMBDA=1. H 103
C      USE ESTIMATE FOR STEPSIZE ONLY IF IT IS POSITIVE AND LESS THAN H 104
C      1. OTHERWISE TAKE 1. AS STEPSIZE H 105
      IF (ALFA) 15,17,15 H 106
15      IF (ALFA-AMBDA) 16,17,17 H 107
      AMBDA=ALFA H 108
      ALFA=0. H 109
C      SAVE FUNCTION AND DERIVATIVE VALUES FOR OLD ARGUMENT H 110
18      FX=FY H 111
      DX=DY H 112
C      STEP ARGUMENT ALONG H H 113
      DO 19 I=1,N H 114
      X(I)=X(I)+AMBDA*H(I) H 115
19      CONTINUE H 116
C      COMPUTE FUNCTION VALUE AND GRADIENT FOR NEW ARGUMENT H 117
      CALL FUNGT (N,X,F,G,GRAD,APP,PSI,NUMB,XX,XP,X1,ERROR,EHELP,AP,EMAX H 118
      1,NP,INUMB,NINT,IP) H 119
      NUMF=NUMF+1 H 120
      FY=F H 121
C      COMPUTE DIRECTIONAL DERIVATIVE DY FOR NEW ARGUMENT. TERMINATE H 122
      SEARCH, IF DY IS POSITIVE. IF DY IS ZERO THE MINIMUM IS FOUND H 123
      DY=0. H 124
      DO 20 I=1,N H 125
      DY=DY+G(I)*H(I) H 126
20      CONTINUE H 127
      IF (DY) 21,41,24 H 128
C      TERMINATE SEARCH ALSO IF THE FUNCTION VALUE INDICATES THAT H 129
      A MINIMUM HAS BEEN PASSED H 130
      IF (FY-FX) 22,24,24 H 131
C      REPEAT SEARCH AND DOUBLE STEPSIZE FOR FURTHER SEARCHES H 132
22      AMBDA=AMBDA+ALFA H 133
      ALFA=AMBDA H 134
      END OF SEARCH LOOP H 135
C      TERMINATE IF THE CHANGE IN ARGUMENT GETS VERY LARGE H 136
      IF (HNPM*AMBDA-1.E10) 18,18,23 H 137
C      LINEAR SEARCH TECHNIQUE INDICATES THAT NO MINIMUM EXISTS H 138
23      IER=2 H 139
      GO TO 62 H 140
C      INTERPOLATE CUBICALLY IN THE INTERVAL DEFINED BY THE SEARCH H 141
      ABOVE AND COMPUTE THE ARGUMENT X FOR WHICH THE INTERPOLATION H 142
      POLYNOMIAL IS MINIMIZED H 143
24      T=0. H 144
      IF (AMBDA) 25,41,25 H 145
25      Z=3.* (FX-FY)/AMBDA+DX+DY H 146
      ALFA=AMAX1 (ARS (Z),ABS (DX),ABS (DY)) H 147
      DALFA=Z*ALFA H 148
      DALFA=DALFA*DALFA-DX/ALFA*DY/ALFA H 149
      IF (DALFA) 57,27,27 H 150
27      W=ALFA*SORT (DALFA) H 151
      ALFA=DY-DX*W H 152
      IF (ALFA) 28,29,28 H 153
28      ALFA=(DY-Z*W)/ALFA H 154
      GO TO 31 H 155
29      ALFA=(Z+DY-W)/(Z+DX+Z+DY) H 156
30      ALFA=ALFA*AMBDA H 157
      DO 31 I=1,N H 158
      X(I)=X(I)+(T-ALFA)*H(I) H 159
31      CONTINUE H 160
C      TERMINATE IF THE VALUE OF THE ACTUAL FUNCTION AT X IS LESS H 161
      THAN THE FUNCTION VALUES AT THE INTERVAL ENDS. OTHERWISE PREDUCE H 162
      THE INTERVAL BY CHOOSING ONE END-POINT EQUAL TO X AND REPEAT H 163
      THE INTERPOLATION, WHICH END-POINT IS CHOSEN DEPENDS ON THE H 164
      VALUE OF THE FUNCTION AND ITS GRADIENT AT X H 165
      NUMF=NUMF+1 H 166
      CALL FUNGT (N,X,F,G,GRAD,APP,PSI,NUMB,XX,XP,X1,ERROR,EHELP,AP,EMAX H 167
      1,NP,INUMB,NINT,IP) H 168
      IF (F-FX) 32,32,33 H 169
32      IF (F-FY) 41,41,33 H 170
      DALFA=0. H 171
33

```

```

      DO 34 I=1,N
      DALFA=DALFA+G(I)*H(I)
      CONTINUE
      IF (DALFA) 35, 38, 38
      35   IF (F-FX) 37, 36, 38
      36   IF (DX-DALFA) 37, 41, 37
      37   FX=F
      DX=DALFA
      T=ALFA
      AMBDA=ALFA
      GO TO 25
      38   IF (FY-F) 40, 39, 40
      39   IF (DY-DALFA) 40, 41, 40
      40   FY=F
      DY=DALFA
      AMBDA=AMBDA-ALFA
      GO TO 24
C
C      TERMINATE, IF FUNCTION HAS NOT DECREASED DURING LAST ITERATION
C      IF (OLDF-F+EPS) 57,42,42
C
C      COMPUTE DIFFERENCE VECTORS OF ARGUMENT AND GRADIENT FROM
C      TWO CONSECUTIVE ITERATIONS
C      DO 43 J=1,N
      41   K=N+J
      K=N+J
      H(K)=G(J)-H(K)
      K=N+K
      H(K)=X(J)-H(K)
      CONTINUE
C
C      TEST LENGTH OF ARGUMENT DIFFERENCE VECTOR AND DIPICTION VECTOR
C      IF AT LEAST N ITERATIONS HAVE BEEN EXECUTED. TERMINATE, IF
C      BOTH ARE LESS THAN EPS
      CIER=1
      IF (KOUNT-N) 47,44,44
      44   T=0.
      Z=0.
      DO 45 J=1,N
      K=N+J
      W=H(K)
      K=K+1
      T=T+ARS(H(K))
      Z=Z+W*H(K)
      45   CONTINUE
      IF (HNRM-EPS) 46,46,47
      46   IF (T-EPS) 62,62,47
C
C      TERMINATE, IF NUMBER OF ITERATIONS WOULD EXCEED LIMIT
      47   IF (KOUNT-LIMIT) 48,55,55
C
C      PREPARE UPDATING OF MATRIX H
      48   ALFA=0.
      DO 52 J=1,N
      K=J+N3
      W=0.
      DO 51 L=1,N
      KL=N+L
      W=W+H(KL)*H(L)
      IF (L-J) 49,50,50
      49   K=K+N-L
      GO TO 51
      50   K=K+1
      CONTINUE
      K=N+J
      ALFA=ALFA+W*H(K)
      H(J)=W
      51   CONTINUE
C
C      REPEAT SEARCH IN DIRECTION OF STEEPEST DESCENT IF RESULTS
C      ARE NOT SATISFACTORY
      IF (Z*ALFA) 53,2,53
C
C      UPDATE MATRIX H
      53   K=N31
      DO 54 L=1,N
      KL=N2+L
      DO 54 NJ=L,N
      NJ=N2+J
      H(K)=H(K)+H(KL)*F(NJ)/Z-H(L)*H(J)/ALFA
      K=K+1
      GO TO 6
      54   END OF ITERATION LOOP
C
C      NO CONVERGENCE AFTER LIMIT ITERATIONS
      55   IER=1
      IF (KK,NE,IPRINT) GO TO 56
      CALL WRITE2 (X,N,G,F,NUMF,KOUNT)
      56   CONTINUE
      GO TO 62
C
C      RESTORE OLD VALUES OF FUNCTION AND ARGUMENTS
      57   DO 58 J=1,N
      K=N2+J
      X(J)=H(K)
      58   CONTINUE
      CALL FUNGT (N,X,F,G,GRAD,APP,PSI,NUMB,XX,XP,X1,ERRP,EHELP,AP,EMAX)

```

H 188  
 H 189  
 H 190  
 H 191  
 H 192  
 H 193  
 H 194  
 H 195  
 H 196  
 H 197  
 H 198  
 H 199  
 H 200  
 H 201  
 H 202  
 H 203  
 H 204  
 H 205  
 H 206  
 H 207  
 H 208  
 H 209  
 H 210  
 H 211  
 H 212  
 H 213  
 H 214  
 H 215  
 H 216  
 H 217  
 H 218  
 H 219  
 H 220  
 H 221  
 H 222  
 H 223  
 H 224  
 H 225  
 H 226  
 H 227  
 H 228  
 H 229  
 H 230  
 H 231  
 H 232  
 H 233  
 H 234  
 H 235  
 H 236  
 H 237  
 H 238  
 H 239  
 H 240  
 H 241  
 H 242  
 H 243  
 H 244  
 H 245  
 H 246  
 H 247  
 H 248  
 H 249  
 H 250  
 H 251  
 H 252  
 H 253  
 H 254  
 H 255  
 H 256  
 H 257  
 H 258  
 H 259  
 H 260  
 H 261  
 H 262  
 H 263

```

1,MP,I NUMB,NINT,IF)
NUMF=NUMF+1
C
      REPEAT SEARCH IN DIRECTION OF STEEPEST DESCENT IF DERIVATIVE
C      FAILS TO BE SUFFICIENTLY SMALL
C      IF (GNRM-EPS) 61,61,59
C
      TEST FOR REPEATED FAILURE OF ITERATION
59      IF (IER) 62,60,60
60      IER=-1
61      GO TO 2
62      IER=J
63      II=IER+2
64      IF (II.EQ.2) K0=1
65      IF (IPRINT.EQ.1) RETURN
66      GO TO (63,64,65,66), II
67      WPITE (6,68) ICR
68      GO TO 67
69      WPITE (6,69) IER
70      GO TO 67
71      WPITE (6,70) IER
72      GO TO 67
73      WPITE (6,71) IER
74      RETURN
C
      FORMAT (1HO,*IER=*,I2,* ERROR IN GRADIENTS CALCULATIONS*)
59      FORMAT (1HO,*IER=*,I2,* CRITERION FOR OPTIMUM HAS BEEN SATISFIED*)
70      FORMAT (1HO,*IER=*,I2,* MAXIMUM NUMBER OF ALLOWABLE ITERATIONS HAS
1 REACHED*)
71      FORMAT (1HO,*IER=*,I2,* CHANGE IN ARGUMENTS GETS TOO LARGE, LINEAR
1 SEARCH INDICATES THAT NO MINIMUM EXISTS*)
END

```

HSBZ00Z //// END OF LIST ////

H 284
H 285
H 286
H 287
H 288
H 289
H 290
H 291
H 292
H 293
H 294
H 295
H 296
H 297
H 298
H 299
H 300
H 301
H 302
H 303
H 304
H 305
H 306
H 307
H 308
H 309
H 310
H 311
H 312
H 313
H 314
H 315

## APPENDIX VII

### Five section transmission line low pass filter

The  $\hat{A}$  matrix of a lossless transmission line is given by

$$(7.1) \quad \hat{A} = \begin{bmatrix} \cos \beta l & jZ_0 \sin \beta l \\ j \frac{\sin \beta l}{Z_0} & \cos \beta l \end{bmatrix}$$

where  $Z_0$  is the characteristic impedance of the transmission line,  $l$  is its length and

$$\beta = \frac{2\pi f}{c}$$

where  $c$  is the velocity of propagation in the medium and  $f$  is frequency [37].

The reflection coefficient is given [37], from Fig. 7-1, as

$$(7.2) \quad \rho(j\omega) = 1 + 2R_g \frac{I_g(j\omega)}{V_g(j\omega)} .$$

The problem defined in Test 3, Chapter VI is to optimize the absolute value of the reflection coefficient (7.2) for the low pass filter design. The impedances  $R_g = R_L = 1\Omega$ .

To compute the reflection coefficient and its gradients it is easier to assume that the load current is 1 instead of fixing the value of the input generator  $V_g$  [37]. In this case

$$\begin{bmatrix} V_L \\ I_L \end{bmatrix} \triangleq \begin{bmatrix} V_6 \\ I_6 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} .$$

Using the relation (7.1) we can write

$$\begin{bmatrix} V_i \\ I_i \end{bmatrix} = \begin{bmatrix} \cos\beta_i & jZ_{oi}\sin\beta_i \\ j\sin\beta_i & \cos\beta_i \\ Z_{oi} & \end{bmatrix} \begin{bmatrix} V_{i+1} \\ I_{i+1} \end{bmatrix}, \quad i=5,4,3,2,1.$$

Also from Fig. 7-1

$$I_g = -I_1$$

$$V_g = V_1 - I_g R_g.$$

The adjoint network approach was used to calculate the gradients

$$\nabla\left(\frac{I_g}{V_g}\right) = \frac{1}{V_g^2} \begin{bmatrix} g_1 \\ g_2 \\ g_3 \\ g_4 \\ g_5 \end{bmatrix}$$

where  $V_g$  is not a function of design variables and

$$g_i = \frac{1}{Z_{oi}} (V_i I_i - V_{i+1} I_{i+1}), \quad i=1,2,3,4,5$$

are the sensitivities which may be found in tables in [37].

The gradients of the  $|\rho|$  with respect to the variable characteristic impedances are given by

$$\chi|\rho| = 2R_g \cdot \text{Re}\left\{\frac{1}{\rho}\right\} \cdot \nabla\left(\frac{I_g}{V_g}\right).$$

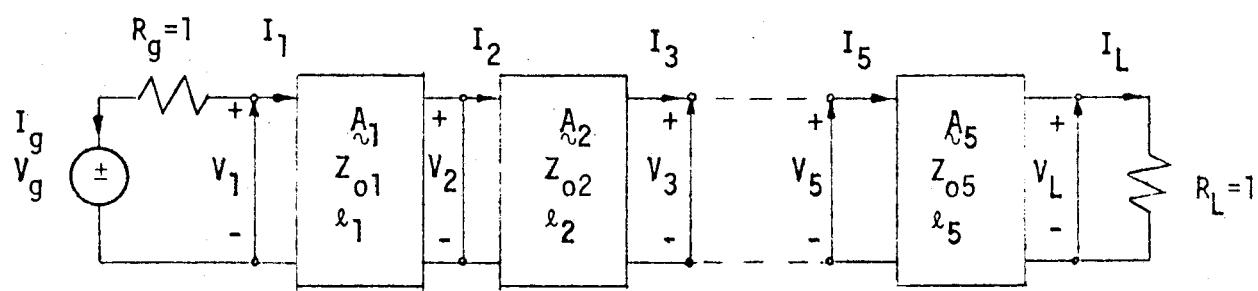


Fig. 7-1 Cascaded five section transmission-line network.

## REFERENCES

- [1] J.W. Bandler and C. Charalambous, "Practical least pth approximation with extremely large values of p", in Proc. 5th Asilomar Conf. on Circuits and Systems, (Pacific Grove, Calif. 1971), pp. 66-70.
- [2] J.W. Bandler and C. Charalambous, "Theory of generalized least pth approximation", IEEE Trans. Circuit Theory, Vol. CT-19, pp. 287-289, 1972.
- [3] G.C. Temes and D.A. Calahan, "Computer-aided network optimization, the state-of-the-art", Proc. IEEE, Vol. 55, pp. 1832-1863, 1967.
- [4] M.J. Box, "A comparison of several current optimization methods, and the use of transformations in constrained problems", Comp. J., Vol. 9, pp. 67-77, 1966.
- [5] A.V. Fiacco and G.P. McCormick, "The sequential unconstrained minimization technique for nonlinear programming, a primal-dual method", Management Sci., Vol. 10, pp. 360-366, 1964.
- [6] A.V. Fiacco and G.P. McCormick, "Computational algorithm for the sequential unconstrained minimization technique for nonlinear programming", Management Sci., Vol. 10, pp. 601-617, 1964.
- [7] A.D. Waren, L.S. Lasdon, and D.F. Suchman, "Optimization in engineering design", Proc. IEEE, Vol. 55, pp. 1885-1897, 1967.
- [8] W.I. Zangwill, Nonlinear Programming: A Unified Approach, Englewood Cliffs, N.J.: Prentice Hall, 1969.
- [9] L.S. Lasdon and A.D. Waren, "Mathematical programming for

- optimal design", Electro-Technol., Vol. 80, pp. 55-70, 1967.
- [10] D.J. Wilde, Optimum Seeking Methods, Englewood Cliffs, N.J.: Prentice Hall, 1964.
- [11] D.J. Wilde and C.S. Beightler, Foundations of Optimization, Englewood Cliffs, N.J.: Prentice Hall, 1967.
- [12] J. Kowalik and M.R. Osborne, Methods for Unconstrained Optimization Problems, Elsevier Publishing Co. Inc., New York, 1968.
- [13] M.J. Box, D. Davies and W.H. Swann, Nonlinear Optimization Techniques, Oliver and Boyd Ltd., Edinburgh, 1969.
- [14] W.C. Davidon, "Variable metric method for minimization", A.E.C. Research and Development Report, ANL-5990 (Rev.), 1959.
- [15] A.S. Householder, Principles of Numerical Analysis, McGraw-Hill, New York, 1953.
- [16] C.G. Broyden, "Quasi-Newton method and their application to function minimization", Math. of Comp., Vol. 21, pp. 368-381, 1967.
- [17] H.Y. Huang, "Unified approach to quadratically convergent algorithms for function minimization", J. Optimization Theory and Applications, Vol. 5, pp. 405-423, 1970.
- [18] R. Fletcher and M.J.D. Powell, "A rapidly convergent descent method for minimization", Comp. J., Vol. 6, pp. 163-168, 1963.
- [19] R. Fletcher, "A new approach to variable metric algorithms," Comp. J., Vol. 13, pp. 317-322, 1970.
- [20] A. Ralston, "Rational Chebyshev approximation by Remes'

- algorithms", Numer. Math., No. 7, pp. 322-331, 1965.
- [21] W.J. Cody and J. Stoer, "Rational Chebyshev approximation using interpolation", Numer. Math., No. 9, pp. 177-188, 1966.
- [22] H. Werner, J. Stoer and W. Bommers, "Rational Chebyshev approximation", Numer. Math., No. 10, pp. 289-306, 1967.
- [23] W.J. Cody, "A survey of practical rational and polynomial approximation of functions", SIAM Rev., Vol. 12, pp. 401-423, 1970.
- [24] M.R. Aaron, "The use of least squares in system design", IRE Trans. Circuit Theory, Vol. CT-3, pp. 224-231, 1956.
- [25] D.W. Marquardt, "An algorithm for least-squares estimation of nonlinear parameters", J. SIAM, Vol. 11, pp. 431-441, 1963.
- [26] G.C. Temes, "Optimization methods in circuit design", in Computer Oriented Circuit Design, F.F. Kuo and W.G. Magnusson, Jr., Eds. Englewood Cliffs, N.J.: Prentice Hall, 1969.
- [27] G.C. Temes and D.Y.F. Zai, "Least pth approximation", IEEE Trans. Circuit Theory, Vol. CT-16, pp. 235-237, 1969.
- [28] J.W. Bandler and C. Charalambous, "On conditions for optimality in least pth approximation with  $p \rightarrow \infty$ ", in Proc. 9th Allerton Conf. on Circuit and System Theory, (Urbana, Ill., 1971), pp. 404-413. J. Optimization Theory and Applications, accepted for publication.
- [29] J.W. Bandler, "Conditions for a minimax optimum", IEEE Trans. Circuit Theory, Vol. CT-18, pp. 476-479, 1971.
- [30] N.I. Achieser, Theory of Approximation (Translated by C.J. Hyman), Frederick Ungar Publ. Co., New York, 1956.
- [31] L. Fox and I.B. Parker, Chebyshev Polynomials in Numerical Analysis, Oxford University Press, 1968.

- [32] J.R. Rice, "The characterization of best nonlinear Tchebycheff approximations", Trans. Amer. Math. Soc., Vol. 96, pp. 322-340, 1960.
- [33] A.R. Curtis and M.J.D. Powell, "Necessary conditions for a minimax approximation", Comp. J., Vol. 8, pp. 358-361, 1965.
- [34] J.W. Bandler and C. Charalambous, "Practical least pth optimization of networks", IEEE Trans. Microwave Theory Tech., Vol. MTT-20, 1972.
- [35] H.J. Carlin, "Distributed Circuit Design With Transmission Line Elements", Proc. IEEE, Vol. 59, 1059-1081, 1971.
- [36] J.W. Bandler and V.K. Jha, "GRADMIN- a package for function minimization using efficient gradient methods", Department of Electrical Engineering, McMaster Univ., Hamilton, Ont., Canada, 1972, internal report.
- [37] J.W. Bandler and R.E. Seviora, "Current trends in network optimization", IEEE Trans. Microwave Theory Tech., Vol. MTT-18, pp. 1159-1170, 1970.