

**TOPICS IN OPTIMIZATION**

BY

**MICHAEL OSBORNE**

**STAN-CS-72-279**

**APRIL 1972**

**COMPUTER SCIENCE DEPARTMENT**  
**School of Humanities and Sciences**  
**STANFORD UNIVERSITY**





TOPICS IN OPTIMIZATION

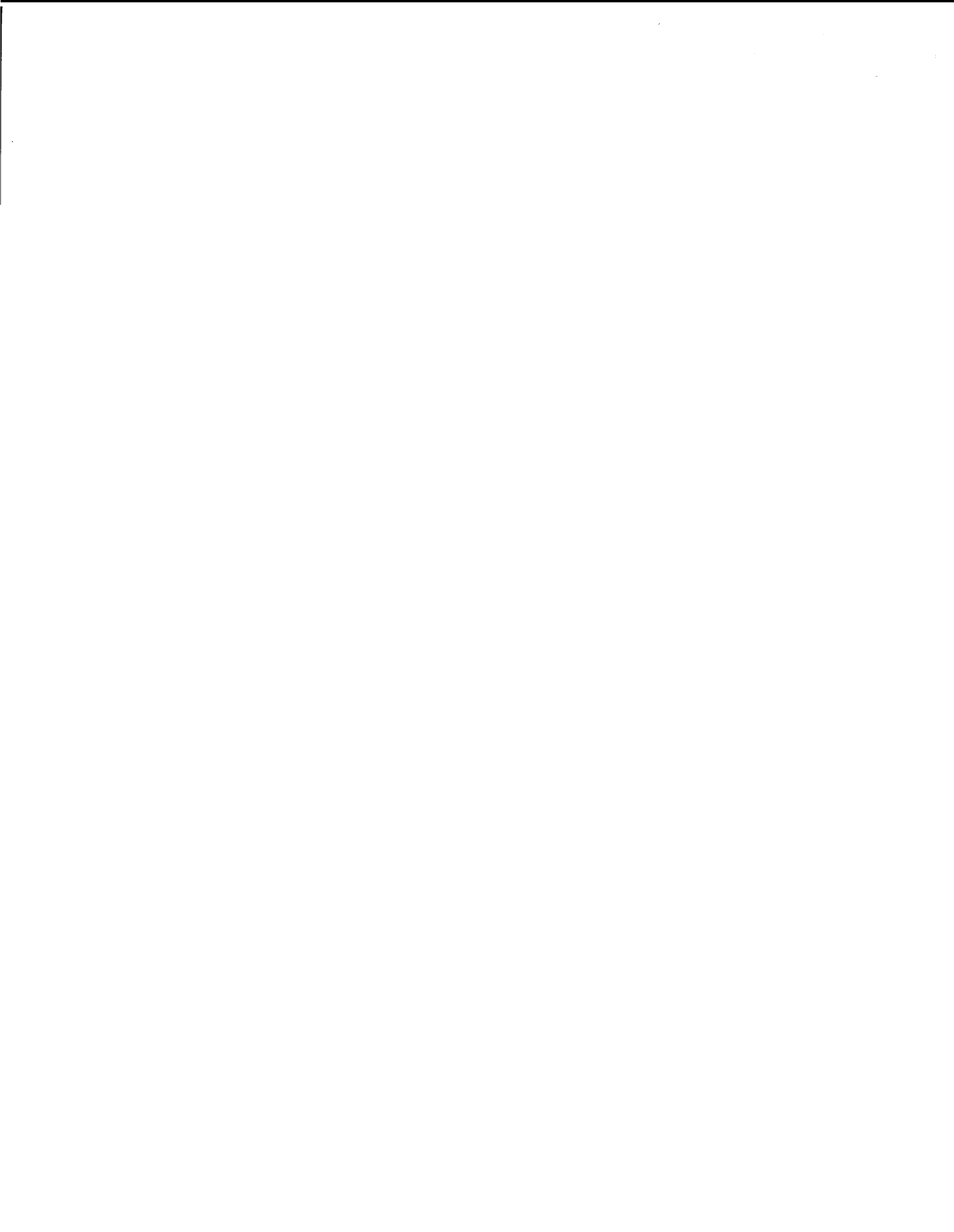
BY

Michael Osborne

STAN-CS-72-279

April 1972

Computer Science Department  
School of Humanities and Sciences  
Stanford University



TOPICS IN OPTIMIZATION

Michael Osborne  
Computer Science Department  
Stanford University

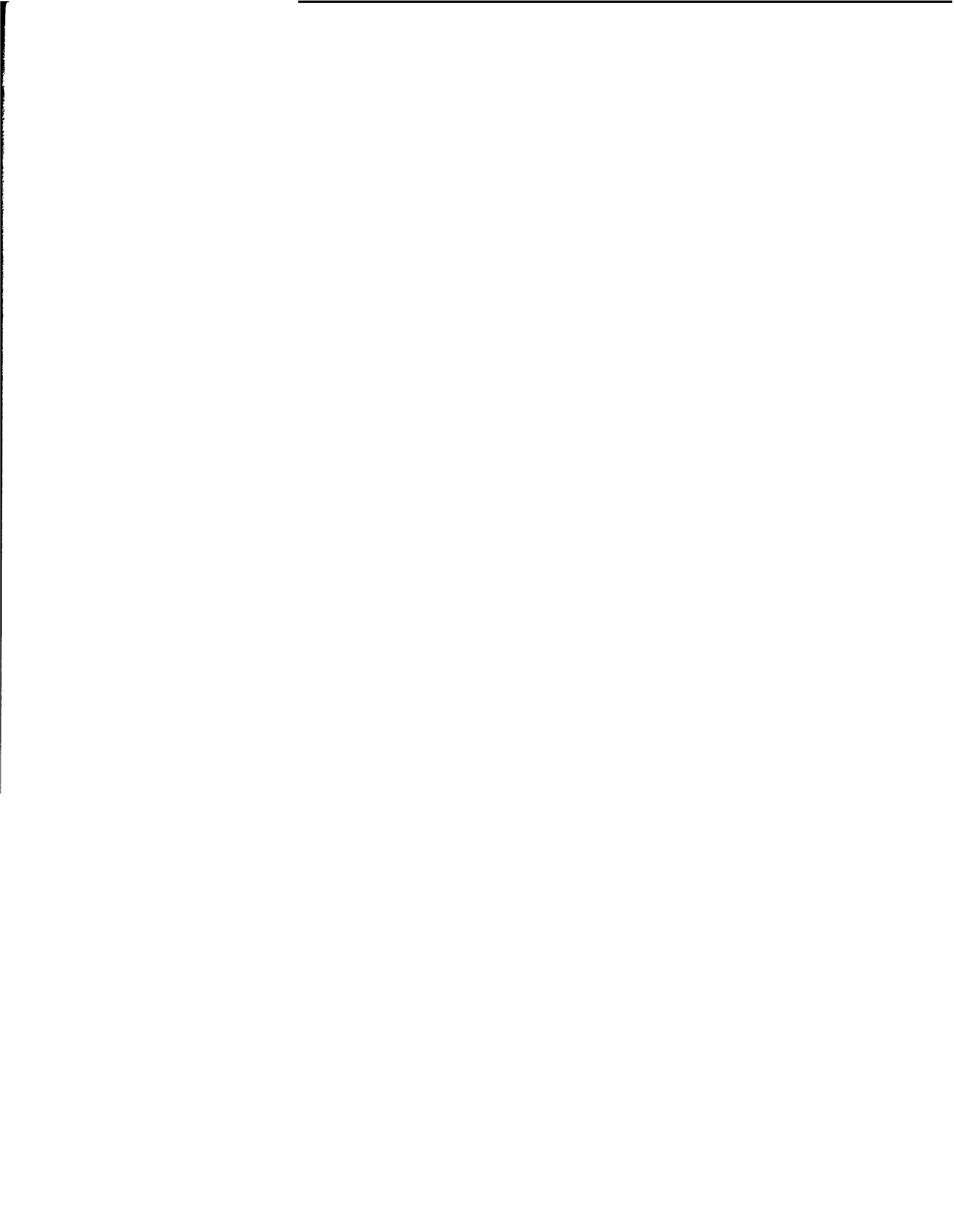
and

Australian National University  
Canberra, Australia

Abstract

These notes are based on a course of lectures given at Stanford, and cover three major topics relevant to optimization theory. First an introduction is given to those results in mathematical programming which appear to be most important for the development and analysis of practical algorithms. Next unconstrained optimization problems are considered. The main emphasis is on that subclass of descent methods which (a) requires the evaluation of first derivatives of the objective function, and (b) has a family connection with the conjugate direction methods. Numerical results obtained using a program based on this material are discussed in an Appendix. In the third section, penalty and barrier function methods for mathematical programming problems are studied in some detail, and possible methods for accelerating their convergence indicated.

This research was supported in part by the National Science Foundation under grant number 29988X, and the Office of Naval Research under contract number N-00014-67-A-0112-00029 NR 044-211 . Reproduction in whole or in part is permitted for any purpose of the United States Government.



## Introduction

These notes were prepared for a course on optimization given in the Computer Science Department at Stanford University during the fall quarter of 1971. In part they are based on lectures given during the year of study in numerical analysis funded by the United Kingdom Science Research Council at the University of Dundee, and on courses given at the Australian National University.

The choice of material has been regulated by limitations of time as well as by personal preference. Also, much material appropriate to the development of algorithms for linearly constrained optimization problems was covered in the parallel course on numerical linear algebra given by Professor Golub. Thus, despite same ambition to cover a larger range, the course eventually consisted of three main sections. These notes cover these sections and have been supplemented by brief additional comments and a list of references. A more extensive bibliography is also included. This is an amended version of a bibliography prepared by my former student Dr. D. M. Ryan.

The first section is intended to provide a solid introduction to the main results in mathematical programming (or at least to those results which appear to be the most important for the development and analysis of practical algorithms). The main aim has been to characterize local extrema, so that convexity and duality theory are not treated in any great detail. However, the material given is more than adequate for the purposes of the remaining sections. Opportunity has been taken to prevent the recent results of Gould and Tolle which provide an accessible and rather complete description of the first order conditions for an extremum. The second order conditions are also considered in detail.

The second section on unconstrained optimization is largely restricted to that subclass of descent methods which (a) requires the evaluation of first derivatives of the objective function, and (b) has some family connection with the so-called conjugate direction methods. This is an area in which there has been considerable recent activity, and here an attempt is made both to summarize significant recent developments and to indicate their algorithmic possibilities. An appendix (prepared with the help of M. A. Saunders) summarizes numerical results obtained with a program based on this material. One significant omission from this section is any detailed discussion of convergence. However, the convergence of certain algorithms (those that reset the Hessian estimate periodically or according to appropriate criteria) is an easy consequence of the material given.

In the third section, penalty and barrier function methods for non-linear programming are considered. This turns out to be a very nice application, in particular, of the results of the first section. These methods have advantages of robustness and simplicity but carry a definite cost penalty. However, attempts to remedy this situation show some promise. The material presented in this section has important connections with other areas: for example, with the method of regularization for the approximate solution of improperly posed problems.

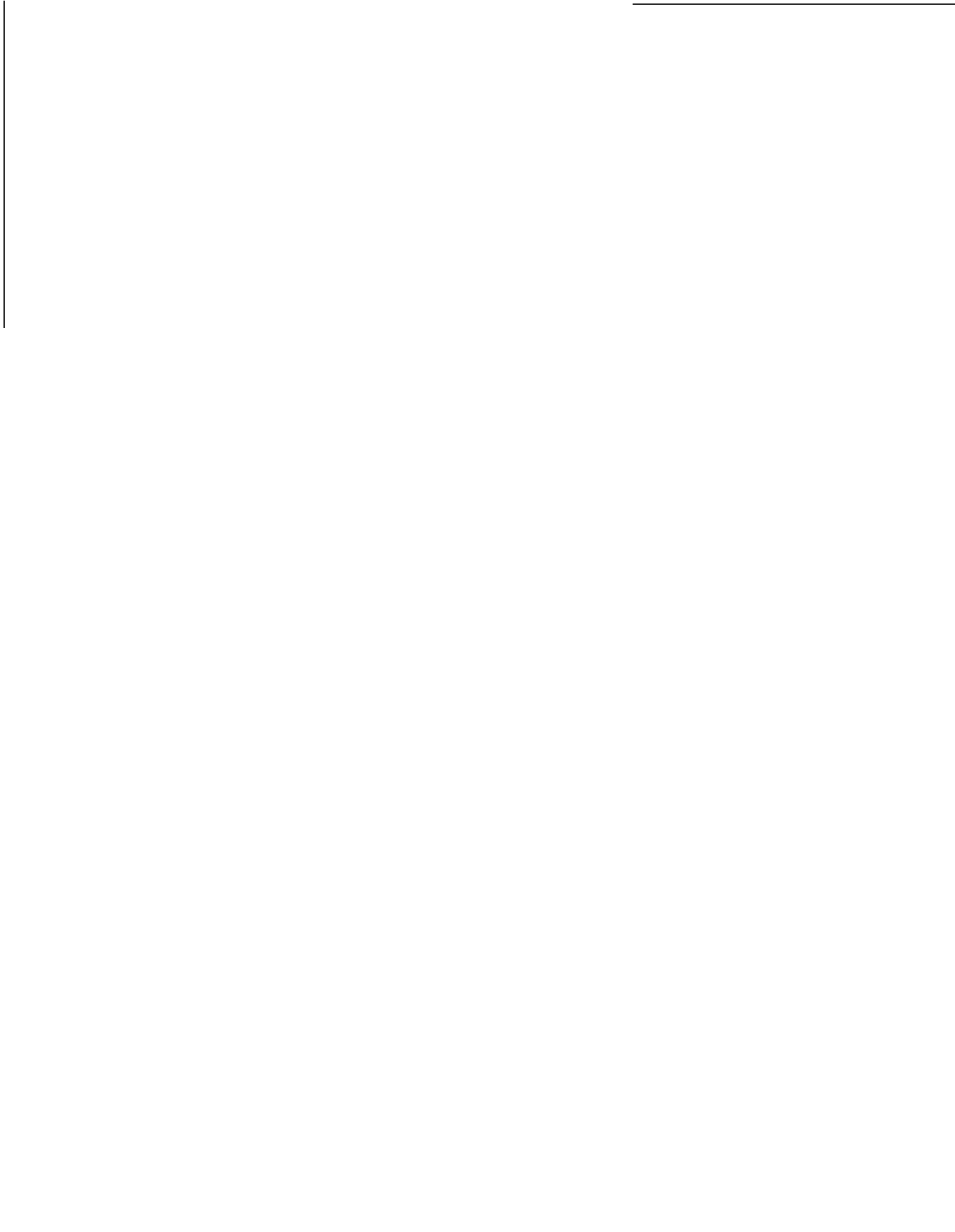
### Acknowledgments

The material presented here has benefited greatly from discussions with Roger Fletcher, Gene Golub, John Greenstadt, Shmuel Oren, Mike Powell, and Mike Saunders. The presentation of the course was shaped more by the



determination of the lecturer to cover as large a field as possible than by any consideration for the audience. In spite of this the level of continued interest was most gratifying, and it is a pleasure to single out the credit students Linda Kaufman, John Lewis and Margaret Wright in this respect.

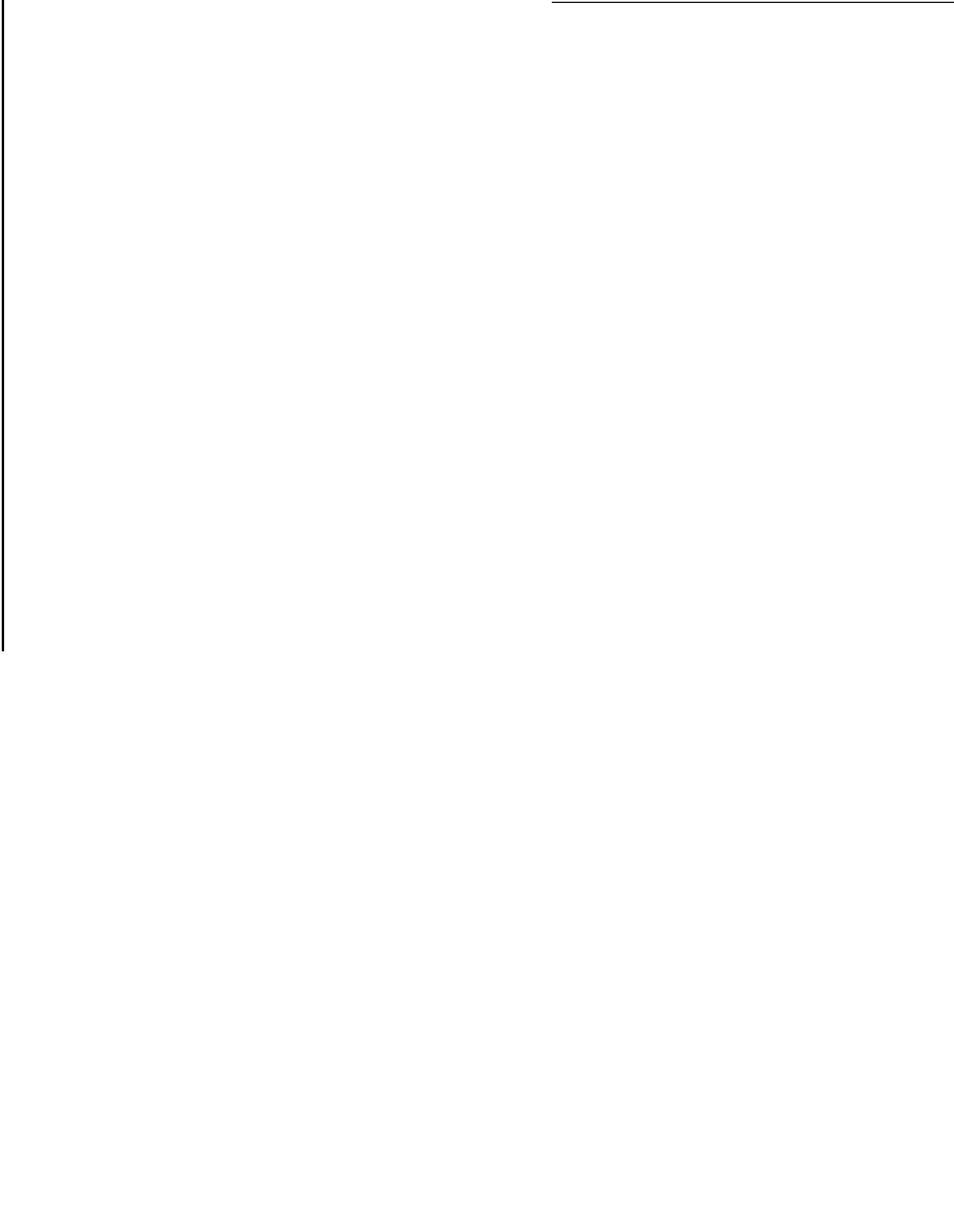
A special vote of thanks is required for George Forsythe who was responsible for the invitation to Stanford, and for completing the subsequent arrangements despite the efforts of the Australian and U.K. Post Offices. In spite of this he was still prepared to combine with John Herriot in keeping the lecturer in reasonable check, and to provide the incentive to produce these notes as a CS report.



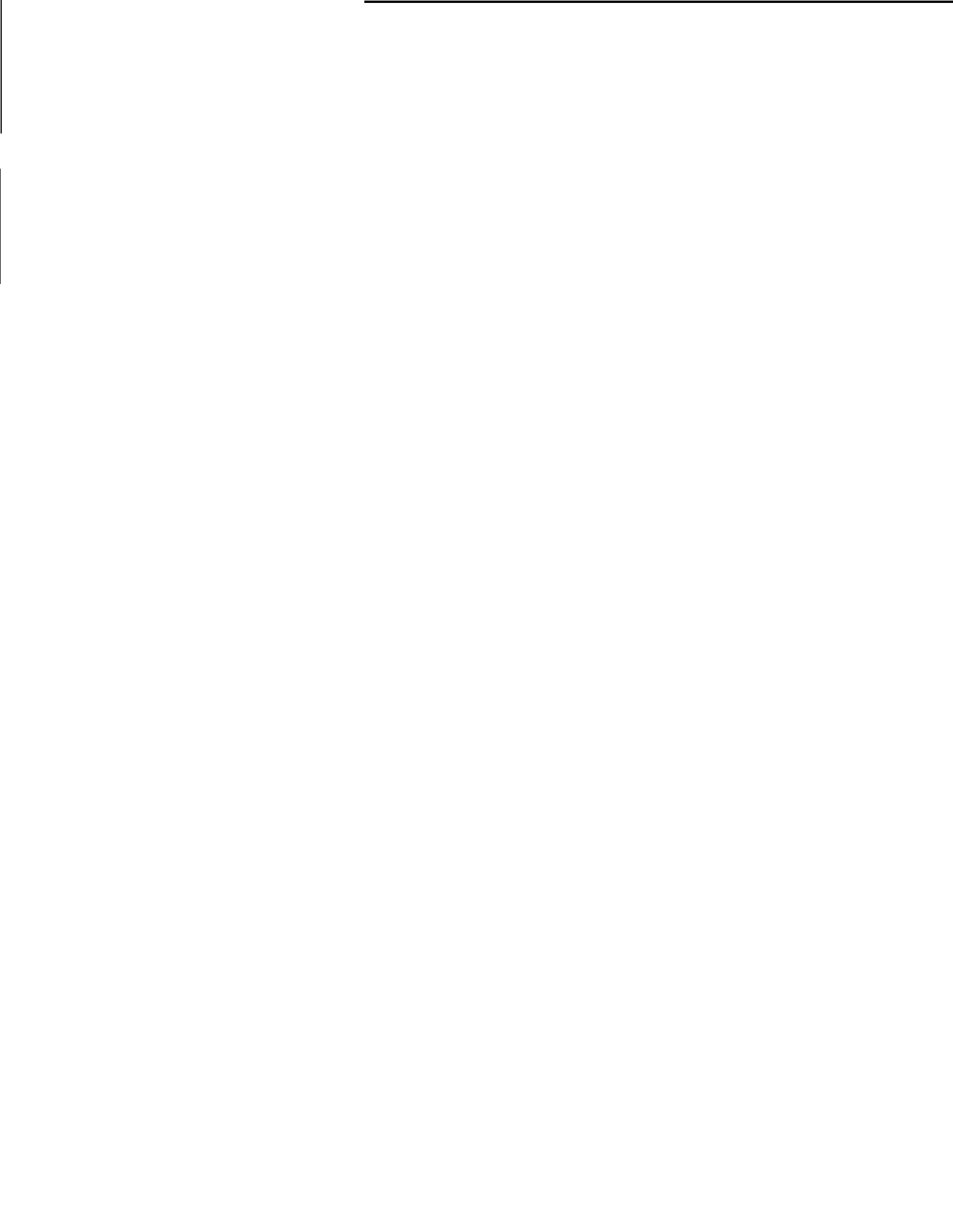
## Contents

I.	Introduction to Mathematical Programming . . . . .	5
1.	Minimum of a constrained function . . . . .	6
2.	Some properties of linear inequalities . . . . .	12
3.	Multiplier relations . . . . .	14
4.	Second order conditions . . . . .	20
5.	Convex programming problems . . . . .	25
II.	Descent Methods for Unconstrained Optimization . . . . .	36
1.	General properties of descent methods . . . . .	37
2.	Methods based on conjugate directions . . . . .	44
	Appendix. Numerical questions relating to Fletcher's algorithm . .	68
III.	Barrier and Penalty Function Methods . . . . .	78
1.	Basic properties of barrier functions . . . . .	79
2.	Multiplier relations (first order analysis) . . . . .	83
3.	Second order conditions . . . . .	86
4.	Rate of convergence results . . . . .	91
5.	Analysis of penalty function methods . . . . .	111
6.	Accelerated penalty and barrier methods . . . . .	117

Note. Equations, Theorems, etc. are numbered in sequence according to subsection. Cross references between sections are indicated explicitly. For example, equation 3 of subsection 4 in the Mathematical Programming section is referenced as MP equation 4.3 in other sections.



I. Introduction to Mathematical Programming



1. Minimum of a constrained function.

Consider a function  $f(\tilde{x})$  on  $S \subseteq E_n \rightarrow E_1$  where  $S$  is a given point set.

Definition:  $\tilde{x}^*$  is the global minimum of  $f$  on  $S$  if

$$f(\tilde{x}^*) \leq f(\tilde{x}) \quad \forall \tilde{x} \in S. \quad (1.1)$$

Remark:  $\tilde{x}^*$  exists, for example, if  $S$  is finite, or if  $S$  is compact and  $f(\tilde{x})$  continuous on  $S$ .

Definition:  $\tilde{x}^*$  is a local minimum of  $f$  on  $S$  if  $\exists \delta > 0$  such that

$$f(\tilde{x}^*) \leq f(\tilde{x}) \quad \forall \tilde{x} \in N(\tilde{x}^*, \delta) \quad (1.2)$$

where

$$N(\tilde{x}, \delta) = \{\tilde{t}; S \cap \{\tilde{t}; \|\tilde{t} - \tilde{x}\|^{+} \leq \delta\}\} . \quad (1.3)$$

If strict inequality holds in either (1.1) or (1.2) whenever  $\tilde{x} \neq \tilde{x}^*$  then the minimum is said to be isolated.

Definition:  $S$  is convex if  $\tilde{x}_1, \tilde{x}_2 \in S \Rightarrow \theta \tilde{x}_1 + (1-\theta)\tilde{x}_2 \in S$  for  $0 \leq \theta \leq 1$ .

Example: If  $S$  is convex all finite combinations of points in  $S$  is

again in  $S$ . That is  $\sum_{i=1}^m \lambda_i \tilde{x}_i \in S$  where  $\tilde{x}_i \in S$ ,  $\sum_{i=1}^m \lambda_i = 1$ ,  $\lambda_i \geq 0$ ,  $1 \leq m < \infty$ .

Definition:  $f(\tilde{x})$  is a convex function on the convex set  $S$  if

$$f(\theta \tilde{x}_1 + (1-\theta)\tilde{x}_2) \leq \theta f(\tilde{x}_1) + (1-\theta)f(\tilde{x}_2), \quad 0 < \theta < 1. \quad (1.4)$$

---

<sup>+/</sup>  $\|\tilde{t}\| = \left\{ \sum_{i=1}^n t_i^2 \right\}^{1/2}$ , the euclidean vector norm of  $\tilde{t}$ .

If strict inequality holds when  $0 < \theta < 1$  then  $f$  is strictly convex.  
 Say  $g(x)$  is concave (strictly concave) if  $-g$  is convex (strictly convex).

Lemma 1.1: If  $f(x)$  is a convex function on the convex set  $S$  then a local minimum of  $f$  is the global minimum. If  $f$  is strictly convex then the minimum is unique.

Proof: It is necessary to consider only the case  $f$  bounded below.

If  $x^*$  is a local minimum but not the global minimum  $\exists x^{**}$  such that  $f(x^{**}) < f(x^*)$ . Now, by assumption,  $\exists \delta > 0$  such that  $f(x) \geq f(x^*)$  for  $x \in N(x^*, \delta)$ . Choose  $\theta > 0$  sufficiently small for  $\theta x^{**} + (1-\theta)x^* \in N(x^*, \delta)$  then

(i)  $f(x^*) \leq f(\theta x^{**} + (1-\theta)x^*)$  as  $x^*$  is a local minimum, and

(ii)  $f(\theta x^{**} + (1-\theta)x^*) \leq \theta f(x^{**}) + (1-\theta)f(x^*)$  by convexity  
 $< f(x^*)$  unless  $f(x^*) = f(x^{**})$ .

Now assume  $x^*, x^{**}$  both are global minima and that  $f$  is strictly convex. Then

$$f(\theta x^* + (1-\theta)x^{**}) < \theta f(x^*) + (1-\theta)f(x^{**}), \quad 0 < \theta < 1$$

which gives a contradiction.  $\square$

Definition: A set  $C$  is a cone with vertex at the origin if  $x \in C \Rightarrow \lambda x \in C, \lambda \geq 0$ .  $C$  is a cone with vertex at  $p$  if  $\{x-p; x \in C\}$  is a cone with vertex at the origin.



Definition:  $\tilde{x}$  is in the tangent cone  $\mathcal{T}(S, \tilde{x}_0)$  to  $S$  at  $\tilde{x}_0$  if  
 3 sequences  $\{\lambda_n\} \geq 0$ ,  $\{\tilde{x}_n\} \rightarrow \tilde{x}_0$ ,  $\{\tilde{x}_n\} \subset S$  such that

$$\lim_{n \rightarrow \infty} \|\lambda_n (\tilde{x}_n - \tilde{x}_0) - \tilde{x}\| = 0. \quad (1.5)$$

Example: (i)  $S = \{\tilde{x}; \|\tilde{x} - \tilde{w}\| = r\}$ ,  $\mathcal{T}(S, \tilde{x}_0) = \{\tilde{x}; \tilde{x}^\top (\tilde{x}_0 - \tilde{w}) = 0\}$ .

(ii)  $S = \{\tilde{x}; \|\tilde{x} - \tilde{w}\| \leq r\}$ .  $\mathcal{T}(S, \tilde{x}_0) = E_n$  if  $\tilde{x}_0$  in interior of  $S$ ,  
 otherwise  $\mathcal{T}(S, \tilde{x}_0) = \{\tilde{x}; \tilde{x}^\top (\tilde{x}_0 - \tilde{w}) \leq 0\}$ .

Lemma 1.2:  $\mathcal{T}(S, \tilde{x}_0)$  is closed.

Proof: Consider a sequence  $\{\tilde{t}_i\} \in \mathcal{T}(S, \tilde{x}_0)$  such that  $\|\tilde{t}_i - \tilde{t}\| \rightarrow 0$ ,  $i \rightarrow \infty$ .

It is required to show that  $\tilde{t} \in \mathcal{T}(S, \tilde{x}_0)$ . Now  $\tilde{t}_i \in \mathcal{T}(S, \tilde{x}_0) \Rightarrow$

$\{\lambda_j^i\} \geq 0$ ,  $\{\tilde{x}_j^i\} \subset S$  such that  $\lim_{j \rightarrow \infty} \|\lambda_j^i (\tilde{x}_j^i - \tilde{x}_0) - \tilde{t}_i\| = 0$ . Prescribe

$\{\epsilon_i\} \downarrow 0$ . Select  $\tilde{t}_i$  such that  $\|\tilde{t}_i - \tilde{t}\| < \epsilon_i/2$ , and  $j = i(j)$  such

that  $\|\lambda_j^i (\tilde{x}_j^i - \tilde{x}_0) - \tilde{t}_i\| < \epsilon_i/2$ . Then  $\|\lambda_j^i (\tilde{x}_j^i - \tilde{x}_0) - \tilde{t}\| < \epsilon_i \Rightarrow \tilde{t} \in \mathcal{T}(S, \tilde{x}_0)$ .  $\square$

Lemma 1.3: (Necessary condition for a local minimum.) If  $f(x) \in C^1$  and

and if  $\tilde{x}_0$  is a local minimum of  $f$  on  $S$  then  $\nabla f(\tilde{x}_0)^\top \tilde{x} \geq 0$ ,

$\forall \tilde{x} \in \mathcal{T}(S, \tilde{x}_0)$ .

Proof: Let  $\tilde{x}$  be defined by sequences  $\{\lambda_n\}$ ,  $\{\tilde{x}_n\}$ . As  $\tilde{x}_0$  is a local

minimum  $\exists \delta > 0$  such that  $f(\tilde{x}^*) \geq f(\tilde{x}_0) \forall \tilde{x}^* \in N(\tilde{x}_0, \delta)$ . Consider now

the restriction of the sequences  $\{\lambda_n\}$ ,  $\{\tilde{x}_n\}$  such that  $\tilde{x}_n \in N(\tilde{x}_0, \delta)$ .

$\frac{+}{-} f \in C^1$  at  $\tilde{x}_0$  if  $f(\tilde{x}) = f(\tilde{x}_0) + \nabla f(\tilde{x}_0)^\top (\tilde{x} - \tilde{x}_0) + o(\|\tilde{x} - \tilde{x}_0\|)$ . Higher

order continuity classes are defined similarly. For example,  $f \in C^2$

if the  $o(\cdot)$  term can be estimated in the form

$$\frac{1}{2} (\tilde{x} - \tilde{x}_0)^\top \nabla^2 f(\tilde{x}_0) (\tilde{x} - \tilde{x}_0) + o(\|\tilde{x} - \tilde{x}_0\|^2).$$

We have

$$\begin{aligned} 0 &\leq f(\underline{x}_n) - f(\underline{x}_0) \\ &\leq \nabla f(\underline{x}_0)(\underline{x}_n - \underline{x}_0) + o(\|\underline{x}_n - \underline{x}_0\|) \end{aligned}$$

whence (note it is sufficient to consider  $\underline{x}$  such that  $\|\underline{x}\| = 1$ )

$$\begin{aligned} 0 &\leq \nabla f(\underline{x}_0)\lambda_n(\underline{x}_n - \underline{x}_0) + o(\lambda_n\|\underline{x}_n - \underline{x}_0\|) \\ &\leq \nabla f(\underline{x}_0)\underline{x} + o(1) \quad \text{as } n \rightarrow \infty. \end{aligned}$$

□

Example: (5.) If  $\underline{x}_0 \in S_0$  (the interior of  $S$ ) then  $\mathcal{T}(S, \underline{x}_0) = E_n$ .

Thus  $\underline{x}$  can be chosen arbitrarily- so that  $\nabla f(\underline{x}_0) = 0$ .

(ii) If  $S = \{\underline{x}; \|\underline{x} - \underline{w}\| = r\}$  then  $\mathcal{T}(S, \underline{x}_0) = \{\underline{x}; \underline{x}^T(\underline{x}_0 - \underline{w}) = 0\}$ .

In particular if  $\underline{x} \in \mathcal{T}(S, \underline{x}_0)$  then  $-\underline{x} \in \mathcal{T}(S, \underline{x}_0)$ . Thus we must have  $\nabla f(\underline{x}_0)\underline{x} = 0 \forall \underline{x}$  such that  $\underline{x}^T(\underline{x}_0 - \underline{w}) = 0$ . Thus  $\nabla f(\underline{x}_0) = \alpha(\underline{x}_0 - \underline{w})$  for some  $\alpha$ .

(iii) If  $S = \{\underline{x}; \|\underline{x} - \underline{w}\| \leq r\}$  and  $\|\underline{x}_0 - \underline{w}\| = r$  then

$\mathcal{T}(S, \underline{x}_0) = \{\underline{x}; \underline{x}^T(\underline{x}_0 - \underline{w}) \leq 0\}$ . In this case we have  $\nabla f(\underline{x}_0)\underline{x} \geq 0 \forall \underline{x}$  such that  $\underline{x}^T(\underline{x}_0 - \underline{w}) \leq 0$ . Thus  $\nabla f(\underline{x}_0) = \alpha(\underline{w} - \underline{x}_0)$  for some nonnegative  $\alpha$ .

Let  $A$  be a set in  $E_n$ .

Definition: The polar cone to  $A$  is the set  $A^* = \{\underline{x}; \underline{x}^T \underline{y} \leq 0 \forall \underline{y} \in A\}$ .

$A^*$  has the following properties.

- (i)  $A^*$  is a closed convex cone.
- (ii) If  $A_1 \subseteq A_2$  then  $A_2^* \subseteq A_1^*$ .
- (iii)  $A^{**} = A$  if and only if  $A$  is a closed convex cone.

(iv)  $A^* = (A^c)^*$  -- the polar cone of the closure of the convex hull of  $A$ . The convex hull of a set is the smallest convex set containing it. Thus  $A^c = \cap X, A \subset X, X$  convex.

(v) If  $A$  is a subspace then  $A^\perp = A^*$ .

Remark: Lemma 1.3 can be restated: 'if  $\tilde{x}^*$  is a local minimum of  $f$  on  $S$  then  $-\nabla f(\tilde{x}^*) \in \mathcal{T}(S, \tilde{x}^*)^*$ '.

Lemma 1.4: If  $y \in \mathcal{T}(S, \tilde{x}_0)^*$  then  $-y$  is the gradient of a function having a local minimum on  $S$  at  $\tilde{x}_0$ .

Remark: It is sufficient to consider the case  $\|y\| = 1, \tilde{x}_0 = 0$ .

Proof: Let  $C_e = \{\tilde{x}; \tilde{x}^T y \leq \frac{\|\tilde{x}\|}{e}\}$   $e = 1, 2, \dots$ . We first show that for each  $e, \exists \epsilon(e) > 0$  such that  $N(0, \epsilon(e)) \subset C_e$ . For assume this is not the case. Then  $\exists \{\tilde{x}_p\} \subset E_n - C_e$  with  $\tilde{x}_p \in N(0, 1/p), p = 1, 2, \dots$  such that

$$\frac{\tilde{x}_p^T y}{\|\tilde{x}_p\|} > \frac{1}{e}, p = 1, 2, \dots \quad (1.6)$$

The sequence  $\left\{ \frac{\tilde{x}_p}{\|\tilde{x}_p\|} \right\}$  is bounded and therefore contains a convergent subsequence  $\left\{ \frac{\hat{\tilde{x}}_i}{\|\hat{\tilde{x}}_i\|} \right\} \rightarrow \tilde{z}$ . By-definition  $\tilde{z} \in \mathcal{T}(S, 0)$ , but, by (1.6),

$$\tilde{z}^T y > \frac{1}{e} > 0$$

which contradicts  $y \in \mathcal{T}(S, 0)^*$ .

Now let  $\hat{\epsilon}_k = \sup\{\epsilon, N(0, \epsilon) \subset C_k\}$ . We define

$$\epsilon_1 = \min(1, \hat{\epsilon}_1) ,$$

$$\epsilon_k = \min\left(\frac{1}{2} \epsilon_{k-1}, \hat{\epsilon}_k\right) , \quad k > 1 ,$$

and

$$P(z) = 2\|z\| , \quad \|z\| \geq \epsilon_2$$

$$\frac{2\|z\|}{k-1} - \frac{\|z\| - \epsilon_{k+1}}{\epsilon_k - \epsilon_{k+1}} + \frac{2\|z\|}{k} - \frac{\epsilon_k - \|z\|}{\epsilon_k - \epsilon_{k+1}} , \quad \|z\| \in [\epsilon_{k+1}, \epsilon_k] ,$$

$$0 , \quad \|z\| \in [0, \epsilon_{k+1}] .$$

It is clear that  $\epsilon_k > 0$  and  $\epsilon_k$  monotonically decreasing. Further  $P(z) \geq 0$ ,  $P(z)$  is an increasing function of  $\|z\|$ , and

$$\|z\| \leq \epsilon_k \Rightarrow P(z) \leq \frac{2\|z\|}{k-1} .$$

Thus  $P(z) = \alpha(\|z\|)$  so that  $\nabla P(0) = 0$ .

Now let  $z = x - (\tilde{x}^T \tilde{y}) \tilde{y}$ . We show that, under appropriate conditions,  $\tilde{x}^T \tilde{y} < P(z)$ . It is sufficient to consider  $\tilde{x}^T \tilde{y} > 0$ , and in this case

$$\|\tilde{x}\| - \tilde{x}^T \tilde{y} \leq \|z\| \leq \|\tilde{x}\| + \tilde{x}^T \tilde{y} . \quad (1.7)$$

If  $\tilde{x} \in C_e$  then  $\tilde{x}^T \tilde{y} \leq \frac{\|\tilde{x}\|}{e}$ . Using (1.7) we have

$$\frac{e-1}{e} \|\tilde{x}\| \leq \|z\| \leq \frac{e+1}{e} \|\tilde{x}\| . \quad (1.8)$$

Now assume  $x \in N(0, \epsilon)$ ,  $\epsilon < \epsilon_3$ . Then  $\|x\| \in [\epsilon_{k+1}, \epsilon_k]$  for some  $k \geq 3$  whence  $\tilde{x} \in C_k$ . This gives

$$\tilde{x}^T \tilde{y} \leq \frac{\|\tilde{x}\|}{k} \leq \frac{\|\tilde{z}\|}{k-1} \quad (1.9)$$

However,  $\|\tilde{z}\| \geq \frac{k-1}{k} \epsilon_{k+1} > \epsilon_{k+2}$  whence

$$P(\tilde{z}) \geq \frac{2\|\tilde{z}\|}{k+1} \quad (1.10)$$

so that, combining (1.9) and (1.10)

$$\tilde{x}^T \tilde{y} \leq \frac{k+1}{2(k-1)} P(\tilde{z}) \leq P(\tilde{z}) \quad (1.11)$$

Thus the function

$$f(\tilde{x}) = -\tilde{x}^T \tilde{y} + P(\tilde{x} - (\tilde{x}^T \tilde{y}) \tilde{y}) \quad (1.12)$$

has a local minimum on  $S$  at  $\tilde{x}_0 = 0$ . Further  $f \in C^1$  at  $0$ , and  $\nabla f(0) = \tilde{y}$ .  $\square$

## 2. Some properties of linear inequalities.

Definition: The set  $H(\tilde{u}, \tilde{v}) = \{\tilde{x}; \tilde{u}^T \tilde{x} = \tilde{v}\}$  is a hyperplane. Note that the hyperplane separates  $E_n$  into two disjoint half spaces  $R_+ = \{\tilde{x}; \tilde{u}^T \tilde{x} \geq \tilde{v}\}$ ,  $R_- = \{\tilde{x}; \tilde{u}^T \tilde{x} < \tilde{v}\}$ .

Lemma 2.1: (lemma of separating hyperplane). Let  $S$  be a closed convex set in  $E_n$ , and let  $\tilde{x}_0 \notin S$ . Then  $\exists$  a hyperplane separating  $\tilde{x}_0$  and  $S$ .

Proof: Let  $\tilde{x}_1$  be any point in  $S$ . Then  $\min_{\tilde{x} \in S} \|\tilde{x} - \tilde{x}_0\| \leq \|\tilde{x}_1 - \tilde{x}_0\| = r$ .

The function  $\|\tilde{x} - \tilde{x}_0\|$  is continuous on the closed set  $S \cap \{\tilde{x}; \|\tilde{x} - \tilde{x}_0\| \leq r\}$  and hence the minimum is attained. Let this point be  $\tilde{x}^*$ . From

Figure 2.1 it is suggested that

$$(\tilde{x} - \tilde{x}^*)^T (\tilde{x}^* - \tilde{x}_0) = 0 \quad (2.1)$$

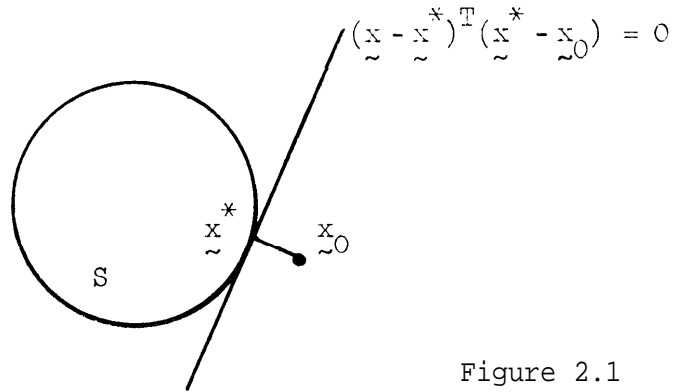


Figure 2.1

is an appropriate hyperplane. To verify this, note that  $\tilde{x}_0 \in \mathbb{R}$  so that it remains to show that  $S \subset \mathbb{R}_+$ . Let  $\tilde{x} \in S$  then for  $0 \leq \theta \leq 1$ ,

$$\|\theta \tilde{x} + (1-\theta)\tilde{x}^* - \tilde{x}_0\|^2 \geq \|\tilde{x}^* - \tilde{x}_0\|^2$$

so that

$$\theta^2 \|\tilde{x} - \tilde{x}^*\|^2 + 2\theta(\tilde{x} - \tilde{x}^*)^T (\tilde{x}^* - \tilde{x}_0) \geq 0$$

and, letting  $\theta \rightarrow 0$ ,

$$(\tilde{x} - \tilde{x}^*)^T (\tilde{x}^* - \tilde{x}_0) \geq 0$$

whence  $\tilde{x} \in \mathbb{R}_+$ .  $\square$

Definition:  $C$  is finitely generated if

$$C = \{\tilde{x}; \tilde{x} = \sum_{i=1}^p \lambda_i \tilde{c}_i, \lambda_i \geq 0, i = 1, 2, \dots, p\}. \text{ It is clear that } C$$

is a cone. It can be shown that  $C$  is closed.

Lemma 2.2: (Farkas Lemma). Let  $A$  be a  $p \times n$  matrix. If for every solution  $\tilde{y}$  of the system of linear inequalities

$$A\tilde{y} \geq 0 \tag{2.2}$$

it is true that

$$\tilde{a}^T \tilde{y} \geq 0 \quad (2.3)$$

then  $\exists \tilde{x} \geq 0$  such that  $A^T \tilde{x} = \tilde{a}$ .

Proof: Let  $C$  be the cone generated by  $\rho_i(A)$ ,  $i = 1, 2, \dots, p$ .

Then the result of Farkas lemma is that if (2.2)  $\Rightarrow$  (2.3) then  $\tilde{a} \in C$ .

We assume  $\tilde{a} \notin C$  and seek a contradiction. By Lemma 2.1 there exists

a separating hyperplane. To construct it let  $\tilde{x}^*$  be the closest point in  $C$  to  $\tilde{a}$ . Then  $\|\lambda \tilde{x}^* - \tilde{a}\|^2$  has a minimum at  $\lambda = 1$ . Differentiating and setting  $\lambda = 1$  gives

$$(\tilde{x}^* - \tilde{a})^T \tilde{x}^* = 0 \quad (2.4)$$

By (2.1) the equation of the separating hyperplane is

$$(\tilde{x} - \tilde{x}^*)^T (\tilde{x}^* - \tilde{a}) = \tilde{x}^T (\tilde{x}^* - \tilde{a}) = 0 \quad (2.5)$$

which shows that it passes through the origin.

By Lemma 2.1  $C \subset R_+$  whence

$$\tilde{v}^T A(\tilde{x}^* - \tilde{a}) \geq 0$$

for arbitrary  $\tilde{v} \geq 0$  so that

$$A(\tilde{x}^* - \tilde{a}) \geq 0, \quad (2.6)$$

but  $\tilde{a} \in R$  whence

$$\tilde{a}^T (\tilde{x}^* - \tilde{a}) < 0 \quad (2.7)$$

which gives the desired contradiction.  $\square$

Remark: Another way of looking at this result is that at most one of the following pair of systems can have a solution.

$$(i) \quad Ax = b \quad , \quad x \geq 0 \quad .$$

$$(ii) \quad A^T \tilde{y} \geq 0 \quad , \quad b^T \tilde{y} < 0 \quad .$$

This is an example of a 'theorem of the alternative'.

### 3. Multiplier relations.

We consider now the mathematical programming problem (MPP)

min  $f(x)$  subject to

$$g_i(x) \geq 0 \quad , \quad i \in I_1 \quad ,$$

$$h_i(x) = 0 \quad , \quad i \in I_2 \quad .$$

We assume that  $f$  ,  $g_i$  ,  $i \in I_1$  , and  $h_i$  ,  $i \in I_2$  , are in  $C^2$  and that the constraints on the problem are not contradictory. This corresponds to the problem discussed in Section 1 with  $S$  given by

$$S = \{x ; g_i(x) \geq 0 \quad , \quad i \in I_1 \quad , \quad h_i(x) = 0 \quad , \quad i \in I_2\} \quad . \quad (3.1)$$

At any point  $x_0 \in S$  let  $B_0$  be the index set for the constraints satisfying  $g_i(x_0) = 0$  . If  $i \in B_0$  we say that  $g_i$  is active at  $x_0$  .

Definition:  $S$  is Lagrange regular at  $x_0$  iff for every  $f$  such that

(i)  $f$  has a minimum on  $S$  at  $x_0$  , and (ii)  $f \in C^1$  at  $x_0$  (i.e.,  $f \in F_0$ )  $\exists u , v$  such that

$$(i) \quad \nabla f(x_0) = \sum_{i \in B_0} u_i \nabla g_i(x_0) + \sum_{i \in I_2} v_i \nabla h_i(x_0) \quad (3.2)$$

$$(ii) \quad u_i \geq 0 \quad , \quad i \in B_0 \quad .$$



This can also be written

$$(i) \quad \nabla f(\underline{x}_0) = \sum_{i \in I_1} u_i \nabla g_i(\underline{x}_0) + \sum_{i \in I_2} v_i \nabla h_i(\underline{x}_0) ,$$

$$(ii) \quad \underline{u}^T \underline{g}(\underline{x}_0) = 0 , \quad \text{and}$$

$$(iii) \quad \underline{u} \geq 0$$

where zero multipliers are introduced corresponding to the inactive constraints.

Remark: If (3.2) holds for  $f \in F_0$ , then  $f$  satisfies the Kuhn-Tucker conditions.

Example: It is important to realize that (3.2) need not hold. Consider the MPP

$$\min f = -x_1 ,$$

$$\text{subject to } g_1 = x_1 \geq 0 , g_2 = x_2 \geq 0 , g_3 = (1 - x_1)^3 - x_2 \geq 0 .$$

From Figure 3.1 it is clear that the minimum is attained at  $x_1 = 1$ ,  $x_2 = 0$ , and here  $g_1$  and  $g_3$  are active. We have

$$\nabla g_1 = -\nabla g_3 = \underline{e}_2$$

while

$$\nabla f = -\underline{e}_1$$

so that a relation of the form (3.2) is impossible.

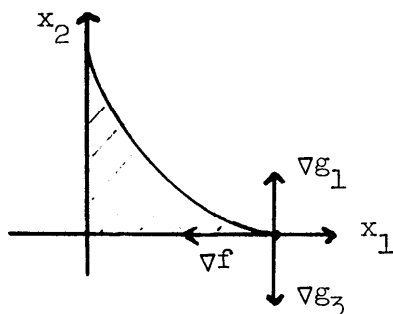


Figure 3.1

Let

$$H_0 = \{x; \nabla h_i(x_0)x = 0, i \in I_2\},$$

$$G_0 = \{x; \nabla g_i(x_0)x \geq 0, i \in B_0\}.$$

Lemma 3.1:  $S$  is Lagrange regular at  $x_0$  iff  $-\nabla f(x_0) \in (G_0 \cap H_0)^*$  for all  $f \in F_0$ .

Proof: If  $-\nabla f(x_0) \in (G_0 \cap H_0)^*$  then

$$-\nabla f(x_0)y \geq 0$$

$\forall y$  such that

$$\nabla h_i(x_0)y \geq 0, \quad i \in I_2,$$

$$-\nabla h_i(x_0)y \geq 0, \quad i \in I_2,$$

$$\nabla g_i(x_0)y \geq 0, \quad i \in B_0.$$

Thus, by Farkas Lemma,  $\nabla f(x_0)$  is a linear combination with nonnegative weights of  $\nabla g_i(x_0)$ ,  $i \in B_0$ , and  $\nabla h_i(x_0)$ ,  $-\nabla h_i(x_0)$ ,  $i \in I_2$ . Thus (3.2) holds. On the other hand, if (3.2) holds then  $\nabla f(x_0)y \geq 0$  for all  $y \in G_0 \cap H_0$ . 3

Remark: Lemma 3.1 shows the difficulty with the above example. Here

$A = \{x = -\alpha e_1, \alpha \geq 0\}$ ;  $G_0 \cap H_0 = \{x = \alpha e_1, \alpha \text{ unconstrained}\}$ . We have

$J^* = \text{right half plane}$ ,  $(G_0 \cap H_0)^*$  the  $x_2$  axis. By Lemma 1.4 for

every  $x \in J^*$  there is a function with a minimum at  $(1,0)$  and such that

$-\nabla f = x$ . Thus the conditions of Lemma 3.1 are not met in this case.

Lemma 3.2:  $(G_0 \cap H_0)^* \subseteq \mathcal{T}(S, x_0)^*$ .

Proof: This result follows if we show that  $\mathcal{T}(S, x_0) \subseteq H_0 \cap G_0$ . If

$\tilde{x} \in \mathcal{T}(S, x_0) \exists \{\tilde{x}_n\} \rightarrow \tilde{x}, \{\tilde{x}_n\} \subset S, \{\lambda_n\} \geq 0$  such that

$\{\lambda_n(x_{\tilde{x}_n} - x_0)\} \rightarrow \tilde{x}$ . We have

$$0 = h_i(x_{\tilde{x}_n}) = h_i(x_0) + \nabla h_i(x_0)(x_{\tilde{x}_n} - x_0) + o(\|x_{\tilde{x}_n} - x_0\|), \quad i \in I_2,$$

and

$$0 \leq g_i(x_{\tilde{x}_n}) = g_i(x_0) + \nabla g_i(x_0)(x_{\tilde{x}_n} - x_0) + o(\|x_{\tilde{x}_n} - x_0\|), \quad i \in B_0.$$

Multiplying by  $\lambda_n$  and repeating the argument used in Lemma 1.3 we have

$$\nabla h_i(x_0)x = 0, \quad i \in I_2, \quad \nabla g_i(x_0)x \geq 0, \quad i \in B_0,$$

so that  $\tilde{x} \in G_0 \cap H_0$ .  $\square$

Theorem 3.1: The set  $S$  is Lagrange regular at  $x_0$  iff

$$\mathcal{T}(S, x_0)^* = (G_0 \cap H_0)^*.$$

Proof: If  $\mathcal{T}(S, x_0)^* = (G_0 \cap H_0)^*$  then  $-\nabla f(x_0) \in (G_0 \cap H_0)^* \forall f \in F_0$  by Lemma 1.3. Thus (3.2) holds by Lemma 3.1. If  $S$  is Lagrange regular at  $x_0$  then by Lemma 3.1  $-\nabla f(x_0) \in (G_0 \cap H_0)^* \forall f \in F_0$ .  $\therefore$ , by Lemma 1.4,  $\mathcal{T}(S, x_0)^* \subseteq (G_0 \cap H_0)^*$ . Thus  $\mathcal{T}(S, x_0)^* = (G_0 \cap H_0)^*$  by Lemma 3.2.  $\square$

Remark: Conditions which ensure that  $S$  is Lagrange regular at  $x_0$  are called restraint conditions. Theorem 3.1 gives a necessary and sufficient restraint condition.

Corollary 3.1: (Kuhn Tucker restraint condition). If  $\nabla g_i(x_0)t \geq 0, i \in B_0$ , and  $\nabla h_i(x_0)t = 0, i \in I_2 \Rightarrow t$  is tangent at  $x_0$  to a once

differentiable arc  $x = x(\theta)$ ,  $x(0) = x_0$  contained in  $N(x_0, \delta)$  for some  $\delta > 0$  then  $S$  is Lagrange regular at  $x_0$ .

Proof: It is clear that  $t \in \mathcal{T}(S, x_0)$  for consider a sequence  $\{\theta_i\} \downarrow 0$  and define  $\{x_n\} = \{x(\theta_n)\}$ ,  $\{\lambda_n\} = \{\frac{1}{\theta_n}\}$  then

$$\{\lambda_n(x_n - x_0)\} \rightarrow \frac{dx(0)}{d\theta} = t \in \mathcal{T}(S, x_0).$$

Thus the Kuhn Tucker restraint condition implies  $(G_0 \cap H_0) \subseteq \mathcal{T}(S, x_0)$ . The result now follows from Lemma 3.2 and Theorem 3.1.  $\square$

Lemma 3.3: Let  $k_i(x) \in C^2$ ,  $k_i(x^*) = 0$ , and  $\nabla k_i(x^*)t = 0$ ,  $i = 1, 2, \dots, s < n$ . We assume  $\exists \epsilon > 0$  such that the  $\nabla k_i(x)$ ,  $i = 1, 2, \dots, s$  are linearly independent for  $\|x - x^*\| < \epsilon$ . Then  $\exists$  a smooth arc  $x = x(\theta)$ ,  $x(0) = x^*$ , such that  $k_i(x(\theta)) = 0$ ,  $i = 1, 2, \dots, s$ , for  $\|x(\theta) - x^*\| < \epsilon$  and  $\frac{dx(0)}{d\theta} = t$ .

Proof: Let  $P(x) = K^T(K K^T)^{-1}K$  where  $\rho_i(K) = \nabla k_i(x)$ ,  $i = 1, 2, \dots, s$ . Then  $x(\theta)$  can be found by integrating the differential equation

$$\frac{dx}{d\theta} = (I - P(x))t \quad (3.3)$$

subject to the initial condition  $x(0) = x^*$ .  $\square$

Remark: Let the  $k_i$  be as given in the statement of Lemma 3.3. Then the linear independence of the  $\nabla k_i$  in a region containing  $x^*$  is a consequence of the linear independence at  $x^*$ . For consider the matrix  $KK^T$ . At  $x = x^*$  this matrix is positive definite as  $K$  has rank  $s$ .

Thus the smallest eigenvalue is positive. Clearly it is a continuous function of  $x$  so that it remains positive in a small enough neighborhood of  $x^*$ , and in this neighborhood the  $\nabla k_i(x)$  are linearly independent.

Lemma 3.4: (Restraint condition A).  $S$  is Lagrange regular at  $x_0$  if the set of vectors  $\nabla g_i(x_0)$ ,  $i \in B_0$ ,  $\nabla h_i(x_0)$ ,  $i \in I_2$  are linearly independent.

Proof: This is a consequence of Corollary 3.1 and Lemma 3.3. For let  $t \in G_0 \cap H_0$ , and let  $B(t)$  be the index set such that  $\nabla g_i(x_0)t = 0$ ,  $i \in B(t)$ . Then by Lemma 3.3 a smooth arc can be constructed such that  $x = x(\theta)$ ,  $g_i(x(\theta)) = 0$ ,  $i \in B(t)$ ,  $h_i(x(\theta)) = 0$ ,  $i \in I_2$ ,  $g_i(x(\theta)) \geq 0$ ,  $i \in I_1 - B(t)$ ,  $x(\theta) \in N(x, \delta)$  for some  $\delta > 0$ , and  $\frac{dx(\theta)}{d\theta} = t$ .  $\square$

Lemma 3.5: (Restraint condition B). If  $\nabla h_i(x_0)$ ,  $i \in I_2$  are linearly independent, and if  $\exists t$  such that  $\nabla g_i(x_0)t > 0$ ,  $i \in B_0$ ,  $\nabla h_i(x_0)t = 0$ ,  $i \in I_2$ , then  $S$  is Lagrange regular at  $x_0$ .

Proof: Assume  $w \in G_0 \cap H_0$  but  $w \notin \mathcal{T}(S, x_0)$ . Prescribe  $\{\varepsilon_k\} \downarrow 0$  and set  $w_k = w + \varepsilon_k t$ . Then  $\nabla g_i(x_0)w_k > 0$ ,  $i \in B_0$ ,  $\nabla h_i(x_0)w_k = 0$ ,  $i \in I_2$ .

Now construct  $x_k = x_k(\theta)$  such that  $x_k(0) = x_0$ ,  $\frac{dx_k(\theta)}{d\theta} = w_k$ ,  $h_i(x_k(\theta)) = 0$ ,  $i \in I_2$ , for  $x_k(\theta)$  is some neighborhood of  $x_0$ . By continuity there will be a subneighborhood (say  $N(x_0, \delta_k)$  for some

$\delta_k > 0$ ) such that (i)  $\nabla g_i(x_k(\theta)) \frac{dx_k(\theta)}{d\theta} \geq 0$ ,  $i \in B_0$ , and

(ii)  $g_i(x_k(\theta)) \geq 0$ ,  $i \in I_1 - B_0$  for  $x_k(\theta) \in N(x_0, \delta_k)$ . The argument of Corollary 3.1 now gives  $w_k \in \mathcal{T}(S, x_0)$ . But, by construction,

$\{w_k\} \rightarrow w$ . Thus  $w \in \mathcal{T}(S, x_0)$  as  $\mathcal{T}$  is closed.  $\square$

4. Second order conditions.

In certain cases it is possible to further characterize local minima of  $f$  on  $S$  by looking at second derivative information.

Lemma 4.1: Let  $w(x) \in C^2$ ,  $w$  have a local minimum on  $S$  at  $\underline{x}_0$ , and  $\nabla w(\underline{x}_0) = 0$ . Then  $\underline{t}^T \nabla^2 w(\underline{x}_0) \underline{t} \geq 0 \forall \underline{t} \in \mathcal{T}(S, \underline{x}_0)$ . If  $\underline{t}^T \nabla^2 w(\underline{x}_0) \underline{t} > 0 \forall \underline{t} \in \mathcal{T}(S, \underline{x}_0)$  then  $\exists \delta > 0, m > 0$  such that  $w(x) \geq w(\underline{x}_0) + m \|x - \underline{x}_0\|^2, x \in N(\underline{x}_0, \delta)$ .

Proof: Let  $\{\underline{x}_n\}, \{\lambda_n\}$  be defining sequences for  $\underline{t} \in \mathcal{T}(S, \underline{x}_0)$ . Then for  $n$  large enough we have, as  $\nabla w(\underline{x}_0) = 0$ ,

$$0 \leq w(\underline{x}_n) - w(\underline{x}_0) = \frac{1}{2} (\underline{x}_n - \underline{x}_0)^T \nabla^2 w(\underline{x}_0) (\underline{x}_n - \underline{x}_0) + o(\|\underline{x}_n - \underline{x}_0\|^2)$$

$$\therefore 0 \leq \lambda_n^2 (w(\underline{x}_n) - w(\underline{x}_0)) = \frac{1}{2} \underline{t}^T \nabla^2 w(\underline{x}_0) \underline{t} + o(1) \quad \text{as } \underline{x}_n \rightarrow \underline{x}_0.$$

Now assume  $\underline{t}^T \nabla^2 w(\underline{x}_0) \underline{t} > 0 \forall \underline{t} \in \mathcal{T}(S, \underline{x}_0)$  and  $\exists$  no  $m > 0$  such that  $w(x) \geq w(\underline{x}_0) + m \|x - \underline{x}_0\|^2$  for  $x$  in any neighborhood of  $\underline{x}_0$ . This implies that for any integer  $q$ ,  $\exists \underline{x}_q \in S$  such that (i)  $\underline{x}_q \in N(\underline{x}_0, 1/q)$ , (ii)  $w(\underline{x}_q) - w(\underline{x}_0) < \frac{1}{q} \|\underline{x}_q - \underline{x}_0\|^2$ . Select a subsequence of the  $\underline{x}_q$  such

that  $\left\{ \frac{\underline{x}_q - \underline{x}_0}{\|\underline{x}_q - \underline{x}_0\|} \right\} \rightarrow \underline{t} \in \mathcal{T}(S, \underline{x}_0)$ . Then (ii)  $\Rightarrow \underline{t}^T \nabla^2 w(\underline{x}_0) \underline{t} < 0$  which gives a contradiction.  $\square$

Definition: The Lagrangian function associated with the MPP is given by

$$L(\underline{x}, \underline{u}, \underline{v}) = f(\underline{x}) - \sum_{i \in I_1} u_i g_i(\underline{x}) - \sum_{i \in I_2} v_i h_i(\underline{x}). \quad (4.1)$$

It will frequently be convenient to suppress the dependence of  $f$  on  $\underline{u}$  and  $\underline{v}$  in the case where these are implied by the Kuhn Tucker conditions. In this case (3.2) becomes

$$\nabla f(\underline{x}_0) = 0 \quad . \quad (4.2)$$

Lemma 4.2: Let  $S$  be Lagrange regular at  $\underline{x}_0$ ,  $f(\underline{x})$  have a local minimum on  $S$  at  $\underline{x}_0$ , and  $S_1 = \{\underline{x}; \underline{x} \in S, g_i(\underline{x}) = 0, i \in B_0\}$  then

$$\underline{t}^T \nabla^2 f(\underline{x}_0) \underline{t} \geq 0 \quad \forall \underline{t} \in \mathcal{T}(S_1, \underline{x}_0) \quad . \quad (4.3)$$

Proof: Note that  $f = f$  on  $S_1$  so that  $f$  has a minimum on  $S_1$  at  $\underline{x}_0$ .  $S_1$  is Lagrange regular at  $\underline{x}_0$ ,  $\nabla f(\underline{x}_0) = 0$ . Thus the result follows from Lemma 4.1.  $\square$

Remark: If  $S_1$  is Lagrange regular at  $\underline{x}_0$  then  $\underline{t}^T \nabla^2 f(\underline{x}_0) \underline{t} \geq 0 \quad \forall \underline{t}$  such that  $\nabla g_i(\underline{x}_0) \underline{t} = 0, i \in B_0$  and  $\nabla h_i(\underline{x}_0) \underline{t} = 0, i \in I_2$ .

Example: Consider

$$g_1 = x_1^2 + (x_2 + 1)^2 - 1 \geq 0, \quad g_2 = 1 - x_1^2 - (x_2 - 1)^2 \geq 0 \quad .$$

$S$  is illustrated diagrammatically in Figure 4.1. At  $x_1 = x_2 = 0$ ,  $\nabla g_1 = \nabla g_2 = (0, 2)$ . However  $S$  is Lagrange regular at the origin -- for example,  $e_2$  satisfies

$$\nabla g_1 e_2 > 0, \quad \nabla g_2 e_2 > 0 \quad \text{so}$$

that restraint condition B applies. In this case  $S_1$  is the single point  $\underline{x} = \theta$  so that  $\mathcal{T}(S_1, \theta)$  is null.

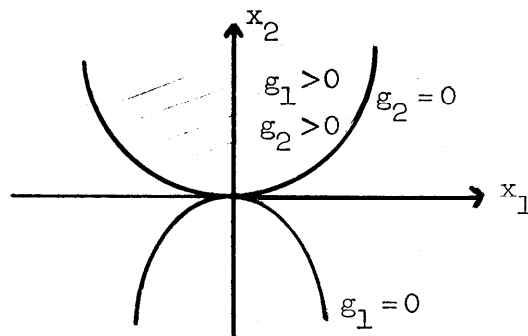


Figure 4.1

Lemma 4.1: If  $t^T \nabla^2 f(x_0) t > 0 \quad \forall t \in \mathcal{T}(S, x_0)$  such that  $\nabla g_i(x_0) t = 0$   
 $\forall i \in B_0$  such that  $u_i > 0$  then  $\exists m, \delta > 0$  such that

$$f(x) \geq f(x_0) + m \|x - x_0\|^2, \quad x \in N(x_0, \delta). \quad (4.4)$$

Proof: Assume  $\nexists m, \delta > 0$  such that (4.4) holds. Then for each integer  $q \exists x_{\sim q}$  such that (i)  $x_{\sim q} \in N(x_0, 1/q)$ , (ii)  $f(x_{\sim q}) - f(x_0) < \frac{1}{q} \|x_{\sim q} - x_0\|^2$ . Select a subsequence of the  $x_{\sim q}$  such that

$$\frac{x_{\sim q} - x_0}{\|x_{\sim q} - x_0\|} \rightarrow t \in \mathcal{T}(S, x_0). \text{ Set } G = \sum_{i \in B_0} u_i g_i(x). \text{ Then } G \geq 0 \text{ on } S,$$

$G(x_0) = 0$ , and  $f = f + G$ . For the subsequence defining  $t$  we have

$$\frac{f(x_{\sim q}) - f(x_0)}{\|x_{\sim q} - x_0\|^2} + \frac{G(x_{\sim q})}{\|x_{\sim q} - x_0\|^2} < \frac{1}{q}, \quad (4.5)$$

Thus

$$t^T \nabla^2 f(x_0) t + \limsup_{q \rightarrow \infty} \frac{G(x_{\sim q})}{\|x_{\sim q} - x_0\|^2} \leq 0 \quad (4.6)$$

A:  $G(x_{\sim q}) \geq 0$ , the second term is bounded and nonnegative. Therefore

$$0 = \lim_{q \rightarrow \infty} \frac{G(x_{\sim q})}{\|x_{\sim q} - x_0\|} = \sum_{i \in B_0} u_i \nabla g_i(x_0) t. \quad (4.7)$$

Thus

$$\nabla g_i(x_0) t = 0, \quad \forall i \in B_0 \text{ such that } u_i > 0 \quad (4.8)$$

so that (4.6) states that  $\exists t \in \mathcal{T}(S, x_0)$  such that  $t$  satisfies (4.8) and that  $t^T \nabla^2 f(x_0) t \leq 0$ . This gives a contradiction.  $\square$



Consider now the system

$$\begin{aligned}
 \nabla f(\underline{x}_0)^T &= 0 \\
 u_i g_i(\underline{x}) &= 0, \quad i = 1, 2, \dots, m, \\
 h_i(\underline{x}) &= 0, \quad i = 1, 2, \dots, p
 \end{aligned} \tag{4.9}$$

where explicit enumerations of  $I_1$  and  $I_2$  are assumed.

Definition:  $J(\underline{x}_0)$  is the Jacobian of the system (4.9) with respect to  $(\underline{x}, \underline{u}, \underline{v})$ .

$$J(\underline{x}_0) = \begin{bmatrix}
 \nabla^2 f(\underline{x}_0) & -\nabla g_1(\underline{x}_0)^T \dots -\nabla g_m(\underline{x}_0)^T & -\nabla h_1(\underline{x}_0)^T \dots -\nabla h_p(\underline{x}_0)^T \\
 u_1 \nabla g_1(\underline{x}_0) & g_1(\underline{x}_0) & \\
 \vdots & \ddots & 0 \\
 u_m \nabla g_m(\underline{x}_0) & g_m(\underline{x}_0) & \\
 \nabla h_1(\underline{x}_0) & & \\
 \vdots & 0 & 0 \\
 \nabla h_p(\underline{x}_0) & &
 \end{bmatrix} \tag{4.10}$$

Lemma 4.4: If  $J(\underline{x}_0)$  is nonsingular, then  $\underline{x}_0$  is an isolated local minimum of  $f$  on  $S$ .

Remark: Note that the condition  $J(\underline{x}_0)$  nonsingular imposes strong conditions on the problem. For example,

- (i) the active constraint gradients must be linearly independent, and

(ii) if  $g_i(x_0) = 0$  then  $u_i > 0$  (this condition is called strict complementarity).

In particular  $S_1$  is Lagrange regular at  $x_0$ .

Proof: If  $J$  is singular there is a vector  $\begin{bmatrix} y \\ a \\ \sim \\ b \\ \sim \end{bmatrix}$  satisfying

$$J \begin{bmatrix} y \\ \sim \\ a \\ \sim \\ b \\ \sim \end{bmatrix} = 0 . \quad (4.11)$$

This relation gives

- (i)  $\nabla h_i(x_0)y = 0$ ,  $i = 1, 2, \dots, p$ ,
- (ii)  $u_i \nabla g_i(x_0)y + a_i g_i(x_0) = 0$ ,  $i = 1, 2, \dots, m$ , and
- (iii)  $\nabla^2 f(x_0)y - \sum_{i=1}^m a_i \nabla g_i(x_0)^T - \sum_{i=1}^p b_i \nabla h_i(x_0)^T = 0$ .

From (ii) we see that  $u_i > 0 \Rightarrow \nabla g_i(x_0)y = 0$  while  $u_i = 0 \Rightarrow a_i = 0$ .

Now consider the problem

$$\min_y \frac{1}{2} \nabla^2 f(x_0)y$$

subject to  $\nabla g_i(x_0)y = 0$ ,  $i \in B_0$ ,  $\nabla h_i(x_0)y = 0$ ,  $i \in I_2$ , and  $\|y\|^2 = 1$ .

Clearly the constraint 'gradients are linearly independent as

$2y = v(\|y\|^2)$  is in the orthogonal complement of the set spanned by the other constraint gradients. Thus the set of feasible  $y$  is Lagrange regular at every point by restraint condition A. Let  $y_0$  minimize the objective function (the minimum exists as the constraint set is compact), then the Lagrange regularity ensures that  $\exists$  multipliers  $\lambda$ ,  $a_i$ ,  $i \in B_0$ ,  $b_i$ ,  $i \in I_2$  such that

$$2\nabla^2 \mathcal{L}(\underline{x}_0) \underline{y}_0 - 2\lambda \underline{y}_0 - \sum_{i \in B_0} a_i \nabla g_i(\underline{x}_0)^T - \sum_{i \in I_2} b_i \nabla h_i(\underline{x}_0)^T = 0 \quad (4.12)$$

whence

$$\lambda = \underline{y}_0^T \nabla^2 \mathcal{L}(\underline{x}_0) \underline{y}_0 = \min_{\underline{y}} \underline{y}^T \nabla^2 \mathcal{L}(\underline{x}_0) \underline{y} \geq 0$$

Now if  $\lambda = 0$ , (4.12) shows that conditions (i) - (iii) above are satisfied and hence  $J(\underline{x}_0)$  singular. Thus if  $J(\underline{x}_0)$  nonsingular, then  $\lambda > 0$ . In this case Lemma 4.3 shows that the minimum of the MPP is isolated.  $\square$

### 5. Convex programming problems.

If  $g_i(\underline{x})$  concave,  $i \in I_1$ , then the set  $S = \{\underline{x}; g_i(\underline{x}) \geq 0, i \in I_2\}$  is convex. The problem of minimizing a convex function on  $S$  is called a convex programming problem. In this section certain properties of this problem are studied. We require the following characterization of convex functions.

Lemma 5.1: If  $f(\underline{x}) \in C^1$  then  $f(\underline{x})$  is convex on  $S$  iff

$$f(\underline{x}) + \nabla f(\underline{x})(\underline{y} - \underline{x}) \leq f(\underline{y}), \quad \underline{x}, \underline{y} \in S \quad (5.1)$$

Proof: If  $f$  convex then, for  $0 \leq \lambda \leq 1$ ,

$$f(\underline{x} + (1-\lambda)(\underline{y} - \underline{x})) \leq \lambda f(\underline{y}) + (1-\lambda)f(\underline{x})$$

whence, if  $\lambda < 1$ ,

$$\frac{f(\underline{x} + (1-\lambda)(\underline{y} - \underline{x})) - f(\underline{x})}{1-\lambda} < \lambda f(\underline{y}) - f(\underline{x})$$

The necessity follows on letting  $\lambda \rightarrow 1$ . Now if (5.1) holds then

$$f(\lambda x + (1-\lambda)y) + \lambda \nabla f(\lambda x + (1-\lambda)y) \cdot (y-x) \leq f(y) \quad (5.2)$$

$$f(\lambda x + (1-\lambda)y) - (1-\lambda) \nabla f(\lambda x + (1-\lambda)y) \cdot (y-x) \leq f(x) \quad (5.3)$$

Multiplying (5.2) by  $(1-\lambda)$ , (5.3) by  $\lambda$  and adding gives (1.4) which demonstrates sufficiency.  $\square$

Lemma 5.2: If  $S = \{x; g_i(x) \geq 0, g_i \text{ concave}, i \in I_1\}$  has an interior point  $x^*$ , then every point of  $S$  is Lagrange regular.

Proof: Consider  $x_0 \in S$ . Let  $i \in B_0$  then Lemma 5.1 gives

$$\nabla g_i(x_0) \cdot (x^* - x_0) \geq g_i(x^*) > 0 \quad (5.4)$$

as  $g_i(x_0) = 0, i \in B_0$ . Thus restraint condition B is satisfied.  $\square$

Lemma 5.3: If  $f$  convex satisfies the Kuhn Tucker conditions at  $x_0$  then  $f$  has a minimum on  $S$  at  $x_0$ .

Proof: In this case (3.2) gives

$$\nabla f(x_0) = \sum_{i \in I_1} u_i \nabla g_i(x_0), \quad u_i^T g_i(x_0) = 0, \quad u_i \geq 0.$$

Let  $x$  be any other point of  $S$ , then

$$f(x) \geq f(x_0) - \sum_{i \in I_1} u_i g_i(x) = \mathcal{L}(x) \quad (5.5)$$

where  $\mathcal{L}(x)$  is convex on  $S$  as the  $g_i(x), i \in I_1$ , are concave. Thus

$$f(x) \geq \mathcal{L}(x_0) + \nabla \mathcal{L}(x_0) \cdot (x - x_0) = f(x_0)$$

Remark: If  $f$  has an interior the Kuhn Tucker conditions are both necessary and sufficient for a minimum of the convex programming problem.

Definition: The primal function for the convex programming problem is

$$\omega(\underline{z}) = \inf_{x \in S} f(x) \quad , \quad S_{\underline{z}} = \{x; g(x) \geq \underline{z}\} \quad (5.6)$$

Note that if  $\underline{z}_1 \geq \underline{z}_2$  then  $S_{\underline{z}_1} \subseteq S_{\underline{z}_2}$  so that  $\omega(\underline{z}_1) \geq \omega(\underline{z}_2)$  and that if  $S$  has an interior then  $S_{\underline{z}}$  nonempty for  $\underline{z} \geq 0$  and small enough.

Lemma 5.3:  $\omega(\underline{z})$  is convex.

Proof: If  $\underline{x}_1 \in S_{\underline{z}_1}, \underline{x}_2 \in S_{\underline{z}_2}$  then, by concavity of  $g_i, i \in I_1,$

$$g(\lambda \underline{x}_1 + (1-\lambda)\underline{x}_2) \geq \lambda \underline{z}_1 + (1-\lambda)\underline{z}_2 \quad , \quad 0 \leq \lambda \leq 1 .$$

Thus  $\lambda \underline{x}_1 + (1-\lambda)\underline{x}_2 \in S_{\lambda \underline{z}_1 + (1-\lambda)\underline{z}_2}$ . We have

$$\begin{aligned} \omega(\lambda \underline{z}_1 + (1-\lambda)\underline{z}_2) &\leq \inf_{\substack{\underline{x}_1 \in S_{\underline{z}_1}, \\ \underline{x}_2 \in S_{\underline{z}_2}}} f(\lambda \underline{x}_1 + (1-\lambda)\underline{x}_2) \\ &\leq \inf_{\substack{\underline{x}_1 \in S_{\underline{z}_1}, \\ \underline{x}_2 \in S_{\underline{z}_2}}} (\lambda f(\underline{x}_1) + (1-\lambda)f(\underline{x}_2)) \text{ by convexity} \\ &\leq \lambda \inf_{\underline{x}_1 \in S_{\underline{z}_1}} f(\underline{x}_1) + (1-\lambda) \inf_{\underline{x}_2 \in S_{\underline{z}_2}} f(\underline{x}_2) \\ &\leq \lambda \omega(\underline{z}_1) + (1-\lambda)\omega(\underline{z}_2) \quad . \quad \square \end{aligned}$$

Definition: The dual function is

$$\phi(\underline{z}^*) = \inf_{x \in \Omega} f(x) - g^T(x) \underline{z}^* \quad , \quad \underline{z}^* \geq 0 \quad (5.7)$$

where  $\Omega$  is the region on which  $f, -g_i, i \in I_1,$  are convex.

Lemma 5.4:  $\phi(z^*)$  is concave.

Proof: Let  $0 \leq \lambda \leq 1$ , and  $z_1^*, z_2^* \geq 0$ , then

$$\begin{aligned} \phi(\lambda z_1^* + (1-\lambda)z_2^*) &= \inf_x \{f(x) - g^T(x)(\lambda z_1^* + (1-\lambda)z_2^*)\} \\ &= \inf_x \{\lambda(f - g^T z_1^*) + (1-\lambda)(f - g^T z_2^*)\} \\ &\geq \lambda \inf_x (f - g^T z_1^*) + (1-\lambda) \inf_x (f - g^T z_2^*) \\ &\geq \lambda \phi(z_1^*) + (1-\lambda)\phi(z_2^*) \quad \square \end{aligned}$$

Lemma 5.5: Let  $\Gamma = \{z; \exists x \in \Omega \text{ such that } g(x) \geq z\}$ . Then

$$\phi(z^*) = \inf_{z \in \Gamma} (w(z) - z^T z^*) \quad (5.8)$$

Proof:

$$\begin{aligned} \phi(z^*) &= \inf_x (f(x) - g^T(x)z^*) , \\ &\leq \inf_{x \in S_z} (f(x) - z^T z^*) , \\ &= w(z) - z^T z^* \end{aligned}$$

$$\therefore \phi(z^*) \leq \inf_{z \in \Gamma} (w(z) - z^T z^*) \quad (5.9)$$

Now let  $g(x_1) = z_1$ . Then

$$\begin{aligned}
f(\underline{x}_1) - \underline{g}^T(\underline{x}_1) \underline{z}^* &\geq \inf_{\substack{\underline{x} \in S \\ \underline{z}_1}} (f(\underline{x}) - \underline{z}_1^T \underline{z}^*) \\
&\geq \omega(\underline{z}_1) - \underline{z}_1^T \underline{z}^* \\
&\geq \inf_{\substack{\underline{z} \in \Gamma}} (\omega(\underline{z}) - \underline{z}^T \underline{z}^*)
\end{aligned}$$

$$\therefore \inf_{\substack{\underline{x}_1}} (f(\underline{x}_1) - \underline{g}^T(\underline{x}_1) \underline{z}^*) \geq \inf_{\substack{\underline{z} \in \Gamma}} (\omega(\underline{z}) - \underline{z}^T \underline{z}^*) \quad . \quad (5.10)$$

The result follows from the inequalities (5.9) and (5.10).  $\square$

Theorem 5.1: (Duality theorem). (i)  $\sup_{\substack{\underline{z}^* \geq 0}} \phi(\underline{z}^*) < \inf_{\substack{\underline{x} \in S}} f(\underline{x})$  .

(ii) If  $S$  has an interior, and  $\exists \underline{x}_0$  such that the Kuhn Tucker conditions are satisfied, then  $\exists \underline{z}^*$  maximizing  $\phi(\underline{z}^*)$  and equality holds in (i).

Proof: From Lemma 5.5 we have that

$$\phi(\underline{z}^*) \leq \omega(\underline{\theta}) = \inf_{\substack{\underline{x} \in S}} f(\underline{x})$$

holds for each  $\underline{z}^* \geq 0$  . Thus

$$\sup_{\substack{\underline{z}^* \geq 0}} \phi(\underline{z}^*) \leq \inf_{\substack{\underline{x} \in S}} f(\underline{x}) \quad . \quad (5.11)$$

If  $\exists \underline{x}_0$  such that the Kuhn Tucker conditions are satisfied then  $\underline{x}_0$  minimizes  $f$  on  $S$  . Defining  $\underline{z}^* = \{u_1, \dots, u_m\}^T$  where the  $u_i > 0$  are the multipliers in the Kuhn Tucker conditions we see that

$$\phi(\underline{z}^*) = f(\underline{x}_0) \quad . \quad \square$$

Corollary 5.1: (Wolfe's form of the duality theorem). Consider the primal problem minimize the convex function  $f(\underline{x})$  subject to the concave constraints  $g_i(\underline{x}) \geq 0$ ,  $i = 1, 2, \dots, m$ , and the dual problem maximize  $\mathcal{L}(\underline{x}, \underline{u})$  subject to  $\nabla_{\underline{x}} \mathcal{L} = 0$ ,  $\underline{u} \geq 0$ . If a solution to the primal exists then the dual problem has a solution and the objective function values are equal.

Remark: (i) The linear programming problem

$$\min_{\underline{x}} \underline{a}^T \underline{x} \quad \text{subject to} \quad \underline{A} \underline{x} - \underline{b} \geq 0 \quad (5.12)$$

is a special case of a convex programming problem as linear functions have the special property of being both convex and concave -- this is an immediate consequence of Lemma 5.1. This property of linear constraints permits the previous discussion to be extended to permit linear equality constraints. Note that if the linear equality constraints are not to be contradictory, then their gradients must be linearly independent.

(ii) If the restraint condition B is satisfied at  $\underline{x}_0$ , and  $f(\underline{x})$  has a minimum on S at  $\underline{x}_0$  then  $\underline{x}_0$  also solves the linear programming problem

$$\min f(\underline{x}_0) + \nabla f(\underline{x}_0) (\underline{x} - \underline{x}_0)$$

subject to

$$(i) \quad g_i(\underline{x}_0) + \nabla g_i(\underline{x}_0) (\underline{x} - \underline{x}_0) \geq 0, \quad i \in I_1, \quad \text{and}$$

$$(ii) \quad h_i(\underline{x}_0) + \nabla h_i(\underline{x}_0) (\underline{x} - \underline{x}_0) \geq 0, \quad ,$$

$$-h_i(\underline{x}_0) - \nabla h_i(\underline{x}_0) (\underline{x} - \underline{x}_0) \rightarrow 0, \quad i \in I_2, \quad ,$$

as the Kuhn Tucker conditions are both necessary and sufficient for a



solution to the linear programming problem. That the converse need not be true is readily seen from the example  $\min -x$  subject to  $1 - x^2 - y^2 \geq 0$  which has a minimum at  $x = 1, y = 0$ . The associated linear programming problem is  $\min -x$  subject to  $1 - x \geq 0$  which has the solution  $(1, y)$  for any  $y$ . Thus additional conditions are required if the converse is to hold (for example, Lemmas 4.3 or 4.4 could be used).

Example: (i) (Duality in linear programming). Consider the primal problem

$$\text{minimize } \underline{\underline{a}}^T \underline{\underline{x}} \quad \text{subject to } \underline{\underline{A}} \underline{\underline{x}} - \underline{\underline{b}} \geq 0 .$$

The corresponding dual is

$$\text{maximize } \underline{\underline{b}}^T \underline{\underline{u}} \quad \text{subject to } \underline{\underline{A}}^T \underline{\underline{u}} - \underline{\underline{a}} = \underline{\underline{0}} \quad \underline{\underline{u}} \geq 0 .$$

If the primal has a solution then so does the dual and the objective function values are equal.

(ii) (The cutting plane algorithm).

(a) Consider the set  $S = \{\underline{\underline{x}}; g_i(\underline{\underline{x}}) \geq 0 \text{ and } g_i \text{ concave, } i \in I_1\}$ .

If  $\underline{\underline{x}}^* \notin S$  then  $g_i(\underline{\underline{x}}^*) < 0$  for at least one  $i$ . Let  $\alpha$  satisfy  $g_\alpha(\underline{\underline{x}}^*) \leq g_i(\underline{\underline{x}}^*)$ ,  $i \in I_1$ . Consider the half space  $U = \{\underline{\underline{x}}; g_\alpha(\underline{\underline{x}}^*) + \nabla g_\alpha(\underline{\underline{x}}^*) \cdot (\underline{\underline{x}} - \underline{\underline{x}}^*) \geq 0\}$ . Then  $\underline{\underline{x}}^* \notin U$ . Now if  $g_i(\underline{\underline{x}}) \geq 0$  then, as  $g_\alpha$  concave,

$$g_\alpha(\underline{\underline{x}}^*) + \nabla g_\alpha(\underline{\underline{x}}^*) \cdot (\underline{\underline{x}} - \underline{\underline{x}}^*) \geq g_\alpha(\underline{\underline{x}}) \geq 0 .$$

Thus  $g_i(\underline{\underline{x}}) \geq 0 \Rightarrow \underline{\underline{x}} \in U$  so that

$$S_\alpha = \{\underline{\underline{x}}; g_\alpha(\underline{\underline{x}}) \geq 0\} \subseteq U .$$

We have  $S \subseteq S_\alpha \subseteq U$ . Thus the hyperplane  $g_\alpha(x^*) + \nabla g_\alpha(x^*)(x - x^*) = 0$  separates  $x^*$  and  $S$ .

(b) The convex programming problem minimize  $f(x)$  subject to  $x \in S$  is equivalent to the problem minimize  $x_{n+1}$  subject to  $x \in S$ ,  $x_{n+1} - f(x) \geq 0$  where  $x_{n+1}$  is a new independent variable (note that the new constraint is concave). This equivalence follows from the Kuhn Tucker conditions by noting that the new constraint must be active. Thus a convex programming problem can be replaced by the problem of minimizing a linear objective function subject to an enlarged constraint set.

(c) Consider the problem of minimizing  $c^T x$  subject to  $x \in S$  and  $S$  bounded. In particular we assume that  $S \subseteq R_0 = \{x; Ax - b \geq 0\}$ . We can now state the cutting plane algorithm

- (0)  $i = 0$ .
- (i) Let  $x_i$  minimize  $c^T x$  subject to  $x \in R_i$ .
- (ii) Determine  $\alpha$  such that  $g_\alpha(x_i) \leq g_j(x_i)$ ,  $j \in I_1$ .
- (iii) If  $g_\alpha(x_i) > 0$  go to (v).
- (iv) Set  $R_{i+1} = R_i \cap \{x; g_\alpha(x_i) + \nabla g_\alpha(x_i)(x - x_i) \geq 0\}$ ,  
 $i := i+1$ , go to (i)
- (v) stop.

Note that step (i) requires the solution of a linear programming problem.

(d) The cutting plane algorithm generates a sequence of points  $x_i$  with the property that

$$c^T x_0 \leq c^T x_1 \leq \dots \leq c^T x_i \leq \dots \leq \min_{x \in S} c^T x$$

as  $R_0 \supseteq R_1 \supseteq \dots \supseteq S$ . Thus the sequence  $\{c^T x_i\}$  is increasing and

bounded above and therefore convergent. Let  $x^*$  be a limit point of the  $\{x_i\}$ . Then  $x^* \in S$  and therefore solves the convex programming problem. To prove this, assume  $x^* \notin S$ . Then

$$\min_i g_i(x^*) = g_\alpha(x^*) = -\lambda < 0 .$$

Let a subsequence  $\{x_j\} \rightarrow x^*$ , then,  $\exists k$  such that

$$(i) \quad \|x_k - x^*\| < \frac{\lambda}{2C} , \text{ and}$$

$$(ii) \quad g_\alpha(x_k) < -\frac{\lambda}{2}$$

where  $C \geq \|\nabla g_i(x)\|$ ,  $x \in R_0$ ,  $i \in I_1$ .

Let

$$\min_i g_i(x_k) = g_\beta(x_k) .$$

Then  $g_\beta(x_k) < -\frac{\lambda}{2}$ . Now  $x^*$  a limit point of  $\{x_i\} \Rightarrow x^* \in \cap R_i$ . In particular,  $x^* \in R_{k+1}$  whence

$$g_\beta(x_k) + \nabla g_\beta(x_k)(x^* - x_k) \geq 0 .$$

But

$$\|\nabla g_\beta(x_k)(x^* - x_k)\| < \frac{\lambda}{2C} C < \frac{\lambda}{2}$$

so that

$$\begin{aligned} g_\beta(x_k) + \nabla g_\beta(x_k)(x^* - x_k) &< -\frac{\lambda}{2} + \|\nabla g_\beta(x_k)(x^* - x_k)\| \\ &< 0 \end{aligned}$$

which gives a contradiction.

## Notes

1. For properties of tangent cones, see Hestenes. Luenberger discusses polar cones (which he calls negative conjugate cones) on pp. 157-159. Lemma 1.4 is due to Gould and Tolle. The proof is due to Nashed et al.
2. Hestenes is a good general reference for this section and includes a proof that a finitely generated cone is closed. The proof given here of Farkas Lemma is standard (see for example Vajda's paper). An extensive list of alternative theorems is given in Mangasarian.
3. The main result is due to Gould and Tolle. The treatment of the other restraint conditions follows Fiacco and McCormick.
4. The treatment of second order conditions is based on Hestenes. Similar material is given in Fiacco and McCormick.
5. The treatment of duality is based on Luenberger. A related treatment is given by Whittle who is good value on applications. Vajda is a good reference for the mathematical programming application. Wolfe's papers in both the Abadie books discuss various aspects of the cutting plane method.

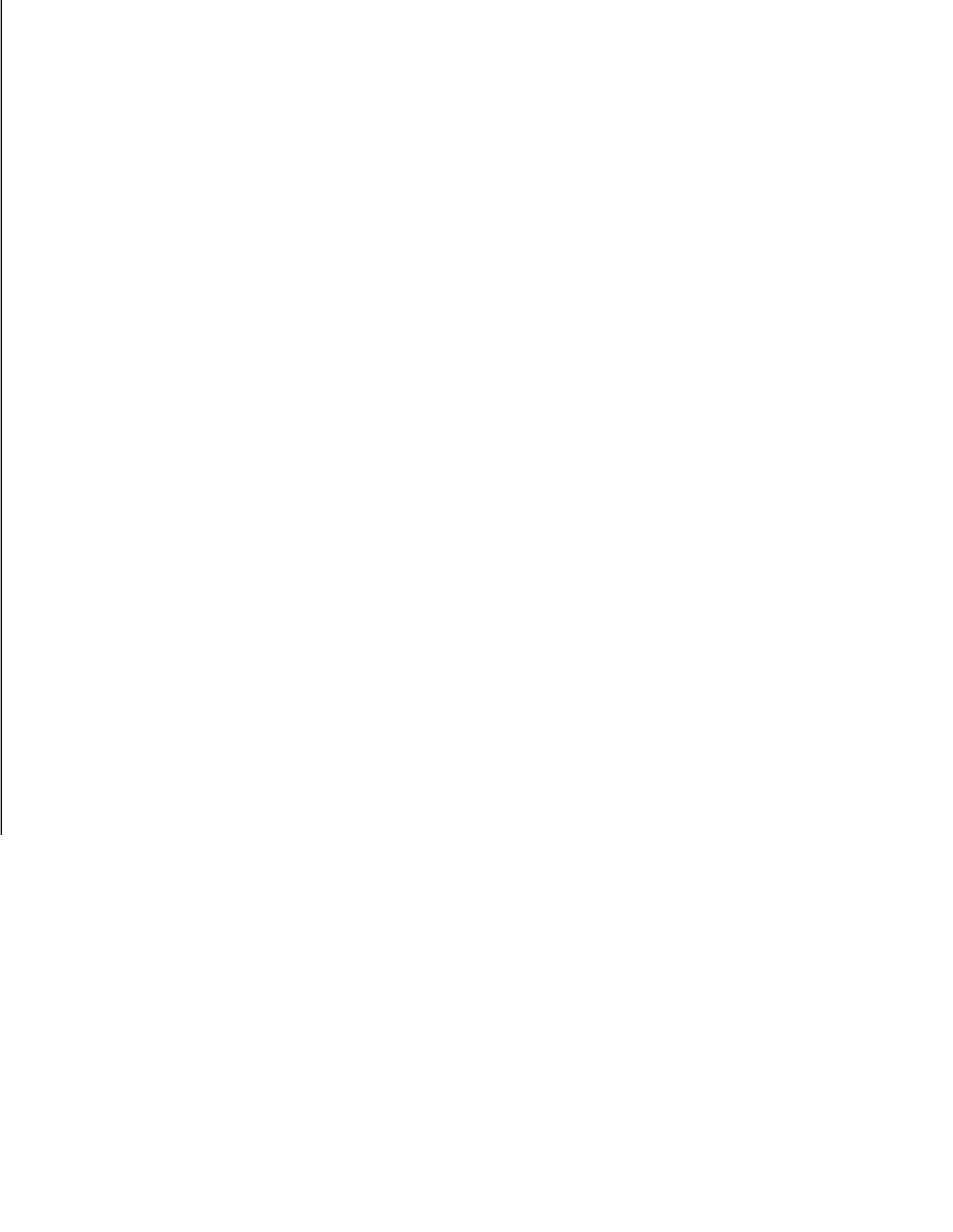
## References

- J. Abadie (Editor) (1967): Nonlinear Programming.  
(1970): Integer and Nonlinear Programming. North Holland.
- M. S. Bazaraa, J. J. Goode, M. Z. Nashed, and C. M. Shetty (1971):  
Nonlinear programming without differentiability in Banach Spaces:  
Necessary and Sufficient Constraint Qualification. To appear in  
Applicable Analysis.

- A. V. Fiacco and G. P. McCormick (1968): Nonlinear Programming and Sequential Unconstrained Minimization Techniques. Wiley.
- F. J. Gould and J. W. Tolle (1971): A necessary and sufficient constraint qualification for constrained optimization, SIAM Applied Maths., 20, pp. 164-172.
- M. R. Hestenes (1966): Calculus of Variations and Optimal Control Theory. Wiley.
- O. L. Mangasarian (1969): Nonlinear Programming. McGraw-Hill.
- Peter Whittle (1971): Optimization under Constraints. Wiley.



## II. Descent Methods for Unconstrained Minimization





1. General properties of descent methods.

The class of descent methods for minimizing an unconstrained function  $F(x)$  solve the problem iteratively by means of a sequence of one dimensional minimizations. The main idea is illustrated in Figure 1.1. At the current point  $\tilde{x}_i$  a direction  $\tilde{t}_i$  is provided, and the closest minimum to  $\tilde{x}_i$  of the function

$$G_i(\lambda) = F(\tilde{x}_i + \lambda \tilde{t}_i)$$

sought. At  $\tilde{x}_{i+1}$  we have

$$G'_i(\lambda_i) = \nabla F(\tilde{x}_{i+1}) \tilde{t}_i = 0 \quad (1.1)$$

where  $\tilde{x}_{i+1} = \tilde{x}_i + \lambda_i \tilde{t}_i$ .

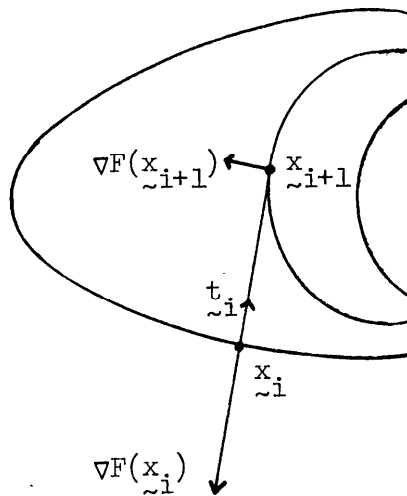


Figure 1.1

Definition: A step in which  $\tilde{x}_{i+1}$  is determined by satisfying the above conditions is said to satisfy the descent condition. We consider  $\tilde{t}_i$  a profitable search direction if  $F(\tilde{x}_i + \lambda \tilde{t}_i)$  decreases initially as  $\lambda$  increases from zero. This condition is formalized as follows.

Definition: (i) The vector  $\underline{t}_i$  is downhill for minimizing  $F$  at  $\underline{x}_i$  if  $\nabla F(\underline{x}_i)\underline{t}_i < 0$ . (ii) The sequence of unit vectors  $\{\underline{t}_i\}$  is downhill for minimizing  $F$  at the sequence of points  $\{\underline{x}_i\}$  if  $\exists \delta > 0$ , independent of  $i$ , such that  $\nabla F(\underline{x}_i)\underline{t}_i \leq -\delta \|\nabla F(\underline{x}_i)\|$ .

Example: The sequence of vectors  $\{-\nabla F(\underline{x}_i) / \|\nabla F(\underline{x}_i)\|\}$  satisfies the downhill condition with  $\delta = 1$ . In this case we say that  $\underline{t}_i$  is in the direction of steepest descent.

An estimate of the value of  $\lambda$  minimizing  $G_i$  is readily given.

We have

$$0 = \nabla F(\underline{x}_{i+1})\underline{t}_i = \nabla F(\underline{x}_i)\underline{t}_i + \lambda_i \underline{t}_i^T \nabla^2 F(\bar{\underline{x}}_i)\underline{t}_i$$

where  $\bar{\underline{x}}_i = \underline{x}_i + \lambda_i \underline{t}_i$  is an appropriate mean value. Thus

$$\lambda_i = \frac{-\nabla F(\underline{x}_i)\underline{t}_i}{\underline{t}_i^T \nabla^2 F(\underline{x}_i)\underline{t}_i} \geq \frac{\delta \|\nabla F(\underline{x}_i)\|}{\|\nabla^2 F(\bar{\underline{x}}_i)\|} \quad (1.2)$$

Theorem 1.1: (Ostrowski's descent theorem). Let  $R = \{\underline{x}; F(\underline{x}) \leq F^*\}$ , and assume that  $F$  bounded below and  $\underline{t}_i^T \nabla^2 F(\underline{x}_i)\underline{t}_i \leq K \|\underline{t}_i\|^2$ ,  $\underline{x}_i \in R$ .

Define

$$\underline{x}_{i+1} = \underline{x}_i + \frac{\delta \|\nabla F(\underline{x}_i)\|}{K} \underline{t}_i, \quad \text{and}$$

$$\nabla F(\underline{x}_i)\underline{t}_i \leq -\delta \|\nabla F(\underline{x}_i)\|, \quad \|\underline{t}_i\| = 1,$$

for  $i = 1, 2, \dots$  where  $\delta > 0$ . Then  $\{F(\underline{x}_i)\}$  converges, and the limit points of  $\{\underline{x}_i\}$  are stationary values of  $F$ .

Proof: As  $\{t_i\}$  downhill then  $\{x_i\} \subset R$ . Expanding by the mean value theorem we obtain

$$F(x_{i+1}) = F(x_i) + \frac{\delta \|\nabla F(x_i)\|}{K} \nabla F(x_i) t_i + \frac{1}{2} \left( \frac{\delta \|\nabla F(x_i)\|}{K} \right)^2 t_i^T \nabla^2 F(\bar{x}_i) t_i$$

where  $\bar{x}_i$  is a mean value. We have

$$\begin{aligned} F(x_{i+1}) &\leq F(x_i) - \frac{(\delta \|\nabla F(x_i)\|)^2}{K} + \frac{1}{2} \left( \frac{\delta \|\nabla F(x_i)\|}{K} \right)^2 K \\ &\leq F(x_i) - \frac{\delta^2 \|\nabla F(x_i)\|^2}{K} \end{aligned} \quad (1.3)$$

Thus the sequence  $\{F(x_i)\}$  is decreasing and bounded below and therefore convergent. Further, from (1.3),

$$\|\nabla F(x_i)\| \leq \frac{1}{\delta} \sqrt{2K(F(x_i) - F(x_{i+1}))} \quad (1.4)$$

$$\rightarrow 0, \quad i \rightarrow \infty.$$

Thus  $\nabla F(x^*) = 0$  if  $x^*$  is a limit point of  $\{x_i\}$ .  $\square$

Remark: By (1.2) the step taken in the direction  $t_i$  underestimates the step to the minimum of  $G_i$ . Thus (1.3) holds if the descent condition is satisfied so that the conclusions of the theorem are valid also in this case.

Theorem 1.2: (Goldstein's descent theorem). Let  $R = \{x; F(x) \leq F^*\}$  be bounded, and assume  $F \in C^1$  and bounded below on  $R$ . Define

$$\Delta(\underline{x}_i, \lambda) = F(\underline{x}_i) - F(\underline{x}_i + \lambda \underline{t}_i) \quad ,$$

$$\psi(\underline{x}_i, \lambda) = - \frac{\Delta(\underline{x}_i, \lambda)}{\lambda \nabla F(\underline{x}_i) \underline{t}_i}$$

where  $\{\underline{t}_i\}$  downhill, and the  $\{\underline{x}_i\}$  are generated by the algorithm

- (i)  $\underline{x}_{i+1} = \underline{x}_i$  if  $\Delta(\underline{x}_i, \lambda) = 0$  .
- (ii) If  $\psi(\underline{x}_i, 1) < \sigma$  where  $0 < \sigma < 1/2$   
then choose  $\lambda_i$  such that  $\sigma \leq \psi(\underline{x}_i, \lambda_i) \leq 1 - \sigma$  ,  
else choose  $\lambda_i = 1$  .
- (iii)  $\underline{x}_{i+1} = \underline{x}_i + \lambda_i \underline{t}_i$  .

Then the limit points of  $\{\underline{x}_i\}$  are stationary points of  $F$  .

Proof:  $\Delta(\underline{x}_i, \lambda) = -\lambda \nabla F(\underline{x}_i) \underline{t}_i + o(h)$  .

Thus  $\Delta(\underline{x}_i, \lambda) = 0 \Rightarrow \|\nabla F(\underline{x}_i)\| = 0$  as  $\{\underline{t}_j\}$  downhill so that  $\underline{x}_i$  is a stationary point. Otherwise  $\nabla F(\underline{x}_i) \underline{t}_i < 0$  so that  $\psi(\underline{x}_i, \lambda) = 1 + o(1)$

whence  $\psi(\underline{x}_i, 0) = 1$  . Also the boundedness of  $R$  implies that

$\Delta(\underline{x}_i, \lambda) < 0$  for some  $\lambda$  large enough so that, as  $\psi(\underline{x}_i, \lambda)$  is continuous,

$\lambda_i$  can be found to satisfy condition (ii) of the algorithm. Note that

$\{\underline{x}_i\} \subset R$  . We have

$$\begin{aligned} F(\underline{x}_i) - F(\underline{x}_{i+1}) &= -\lambda_i \psi(\underline{x}_i, \lambda_i) \nabla F(\underline{x}_i) \underline{t}_i \\ &\geq \lambda_i \sigma \delta \|\nabla F(\underline{x}_i)\| \quad . \end{aligned} \tag{1.5}$$

Thus  $\{F(\underline{x}_i)\}$  decreasing and bounded below and therefore convergent. To show that the limit points of  $\{\underline{x}_i\}$  are stationary values of  $F$  consider

the subsequence  $\{x_{\tilde{\mu}_i}\} \rightarrow x^*$  and assume  $\|\nabla F(x^*)\| > \epsilon > 0$ . Then  $\|\nabla F(x_{\tilde{\mu}_i})\| > \epsilon$  for  $i > i_0$ . This implies that  $\inf_i \lambda_{\mu_i} = \lambda_0 > 0$  as otherwise  $\sup_i \psi(x_{\tilde{\mu}_i}, \lambda_{\mu_i}) = 1$  contradicting  $\psi(x_i, \lambda_i) \leq 1 - a$ . Thus

$$\|\nabla F(x_{\tilde{\mu}_i})\| \leq \frac{F(x_{\tilde{\mu}_i}) - F(x_{\tilde{\mu}_i+1})}{\lambda_0 \sigma \delta} . \quad (1.6)$$

The right hand side  $\rightarrow 0$  as  $i \rightarrow \infty$  which establishes a contradiction.  $\square$

Remark: There are two aspects of this theorem which are of particular interest. (i) It is necessary to assume only that  $f \in C^1$  in  $R$ . However, the boundedness of  $R$  is used explicitly. (ii) The algorithm for determining the step length  $\lambda_i$  is readily implemented. A value of  $\lambda$  satisfying condition (ii) of the algorithm will be said to satisfy the Goldstein condition.

Theorem 1.3: (i) Let the vector sequence in the Goldstein algorithm be defined by

$$s_i = -A_i \nabla F(x_{\tilde{\mu}_i})^T, \quad t_i = s_i / \|s_i\| \quad (1.7)$$

where  $A_i$  is positive definite, bounded, and  $\kappa(A_i) = \frac{\lambda_{\min}(A_i)}{\lambda_{\max}(A_i)} \geq \omega > 0$ ,

$i = 1, 2, \dots$ . Then  $\{t_i\}$  is downhill with constant  $\delta = \omega$ .

(ii) Assume that  $\{x_{\tilde{\mu}_i}\} \rightarrow x^*$ , and that  $\|A_i^{-1} - \nabla^2 F(x_{\tilde{\mu}_i})\| = o(1)$ , then  $\lambda_i = \|s_i\|$  satisfies the Goldstein condition for  $i$  large enough.

(iii) The ultimate rate of convergence of the algorithm is superlinear for this choice of  $\lambda_i$ .

Proof: (i)

$$\begin{aligned}
\nabla F(\tilde{x}_i) t_i &= - \frac{\nabla F(\tilde{x}_i) A_i \nabla F(\tilde{x}_i)^T}{\|\nabla F(\tilde{x}_i) A_i\|} \\
&\leq \frac{\lambda_{\min}(A_i) \|\nabla F(\tilde{x}_i)\|}{\lambda_{\max}(A)} \\
&\leq -\omega \|\nabla F(\tilde{x}_i)\| .
\end{aligned} \tag{1.8}$$

$$\begin{aligned}
(ii) \quad \psi(\tilde{x}_i, \lambda) &= - \frac{F(\tilde{x}_i) - F(\tilde{x}_i + \lambda t_i)}{\lambda \nabla F(\tilde{x}_i) t_i} \\
&= \frac{\lambda \nabla F(\tilde{x}_i) t_i + \frac{\lambda^2}{2} t_i^T \nabla^2 F(\bar{x}_i) t_i}{\lambda \nabla F(\tilde{x}_i) t_i}
\end{aligned}$$

where  $\bar{x}_i$  is a mean value dependent on  $\lambda$ . Now, writing

$$\nabla^2 F(\bar{x}_i) = A_i^{-1} + E_i$$

and noting that  $\|E_i\| \rightarrow 0$  as  $i \rightarrow \infty$ , we have

$$\psi(\tilde{x}_i, \lambda) = 1 - \frac{\lambda}{2} \left\{ \frac{1}{\|s_i\|} + \frac{t_i^T E_i t_i}{\nabla F(\tilde{x}_i) t_i} \right\}$$

so that

$$\begin{aligned}
|\psi(\tilde{x}_i, \lambda) - (1 - \frac{\lambda}{2\|s_i\|})| &\leq \frac{\lambda}{2\|s_i\|} \frac{\|E_i\| \|s_i\|}{\omega \|\nabla F(\tilde{x}_i)\|} \\
&\leq \frac{\lambda}{2\|s_i\|} \frac{\|E_i\| \|A_i\|}{\omega} .
\end{aligned} \tag{1.9}$$

In particular

$$\left| \psi(\tilde{x}_i, \|\tilde{s}_i\|) - \frac{1}{2} \right| \leq \frac{1}{2} \frac{\|\tilde{E}_i\| \|\tilde{A}_i\|}{\omega} \rightarrow 0, \quad i \rightarrow \infty.$$

(iii) Another application of the mean value theorems gives

$$\begin{aligned} \tilde{s}_i &= -A_i \nabla F(\tilde{x}_i)^\top = -A_i (\nabla F(\tilde{x}_i) - \nabla F(\tilde{x}^*))^\top \\ &= -A_i (\nabla^2 F(\tilde{x}_i) (\tilde{x}_i - \tilde{x}^*) + o(\|\tilde{x}_i - \tilde{x}^*\|)) \\ &= -(\tilde{x}_i - \tilde{x}^*) \cdot \cdot \cdot \|\tilde{x}_i - \tilde{x}^*\| \end{aligned} \quad (1.10)$$

Thus

$$\begin{aligned} \tilde{x}_{i+1} - \tilde{x}^* &= \tilde{x}_i + \gamma_i \tilde{s}_i - \tilde{x}^* \\ &= (1 - \gamma_i) (\tilde{x}_i - \tilde{x}^*) + o(\|\tilde{x}_i - \tilde{x}^*\|) \end{aligned} \quad (1.11)$$

From (1.11) the choice  $\gamma_i = 1$  ( $\lambda_i = \|\tilde{s}_i\|$ ) gives superlinear convergence. C1

Remark: Theorem 1.3 shows that if  $\nabla^2 F$  is positive definite in the neighbourhood of an unconstrained minimum, then it is possible to have algorithms with superlinear convergence without the necessity of satisfying the descent condition.. It is not generally considered economic to compute the second partial derivatives of  $F$ , and considerable emphasis has been placed on developing approximations to the inverse Hessian using only first derivative information. Although the steepest descent direction is initially in the direction of most rapid decrease of the function it gives in general only linear convergence.

2. Methods based on conjugate directions.

The problem of minimizing a positive definite quadratic form is an important special case of the general unconstrained optimization problem. In particular it is frequently used as a model problem for the development of new algorithms. It is argued that in a neighborhood of the minimum, a general function having a positive definite Hessian at the minimum will be well represented by a quadratic form so that methods which work well in this particular case should work well in general.

Let  $F$  be given by

$$F(\underline{x}) = a + \underline{b}^T \underline{x} + \frac{1}{2} \underline{x}^T C \underline{x} \quad (2.1)$$

where  $C$  is a positive definite, necessarily symmetric matrix. We have

$$\nabla F(\underline{x}) = \underline{b}^T + \underline{x}^T C \quad (2.2)$$

Consider now a descent step from  $\underline{x}_i$  in the direction  $\underline{t}_i$ . The descent condition gives

$$0 = \nabla F(\underline{x}_{i+1}) \underline{t}_i = \underline{t}_i^T (C(\underline{x}_i + \lambda \underline{t}_i) + \underline{b})$$

whence

$$\lambda = - \frac{\underline{t}_i^T \underline{g}_i}{\underline{t}_i^T C \underline{t}_i} \quad (2.3)$$

where  $\underline{g}_i = \nabla F(\underline{x}_i)$ . To calculate the change in the value of  $F$  in a descent step we have



$$\begin{aligned}
F(\tilde{x}_i + \lambda_i \tilde{t}_i) - F(\tilde{x}_i) &= \lambda_i \tilde{b}_i^T \tilde{t}_i + \lambda_i \tilde{x}_i^T C \tilde{t}_i + \frac{1}{2} \lambda_i^2 \tilde{t}_i^T C \tilde{t}_i \\
&= \lambda_i \tilde{g}_i^T \tilde{t}_i + \frac{1}{2} \lambda_i^2 \tilde{t}_i^T C \tilde{t}_i \\
&= -\frac{1}{2} \frac{(\tilde{g}_i^T \tilde{t}_i)^2}{\tilde{t}_i^T C \tilde{t}_i} .
\end{aligned} \tag{2.4}$$

Example: (Linear convergence of the method of steepest descent).

Let  $F = \frac{1}{2} \tilde{x}^T C \tilde{x}$ . Then (2.4) gives

$$\begin{aligned}
F(\tilde{x}_{i+1}) - F(\tilde{x}_i) &= -\frac{1}{2} \frac{(\tilde{x}_i^T C \tilde{x}_i)^2}{\tilde{x}_i^T C^3 \tilde{x}_i} \\
&= -\frac{1}{2} \frac{(\tilde{w}_i^T \tilde{w}_i)^2}{\tilde{w}_i^T C \tilde{w}_i}
\end{aligned}$$

where  $\tilde{w}_i = C \tilde{x}_i$ .

We have  $F(\tilde{x}_i) = \frac{1}{2} \tilde{w}_i^T C^{-1} \tilde{w}_i$  so that

$$F(\tilde{x}_{i+1}) = \left( 1 - \frac{1}{2} \frac{(\tilde{w}_i^T \tilde{w}_i)^2}{\tilde{w}_i^T C \tilde{w}_i \tilde{w}_i^T C^{-1} \tilde{w}_i} \right) F(\tilde{x}_i)$$

The Kantorovich inequality gives

$$\frac{(\tilde{w}_i^T \tilde{w}_i)^2}{\tilde{w}_i^T C^{-1} \tilde{w}_i \tilde{w}_i^T C \tilde{w}_i} \geq \frac{4 \sigma_1 \sigma_n}{(\sigma_1 + \sigma_n)^2}$$

where  $\sigma_1$  and  $\sigma_n$  are the smallest and largest eigenvalues of  $C$  respectively, whence

$$F(\underline{x}_{i+1}) \leq \left( \frac{\sigma_n - \sigma_1}{\sigma_n + \sigma_1} \right)^2 F(\underline{x}_i)$$

which shows that the rate of convergence of steepest descent is at least linear.

To show that it is exactly linear consider the particular case in which

$$\underline{x}_i = \alpha_1^i \underline{v}_1 + \alpha_n^i \underline{v}_n$$

where  $\underline{v}_1$  and  $\underline{v}_n$  are the normalized eigenvectors associated with  $\sigma_1$  and  $\sigma_n$  respectively. We have

$$\underline{x}_{i+1} = \alpha_1^{i+1} \underline{v}_1 + \alpha_n^{i+1} \underline{v}_n = (1 - \lambda_i \sigma_1) \alpha_1^i \underline{v}_1 + (1 - \lambda_i \sigma_n) \alpha_n^i \underline{v}_n$$

with (from (2.3))

$$\lambda_i = \frac{(\alpha_1^i)^2 \sigma_1^2 + (\alpha_n^i)^2 \sigma_n^2}{(\alpha_1^i)^2 \sigma_1^3 + (\alpha_n^i)^2 \sigma_n^3},$$

so that

$$\alpha_1^{i+1} = \frac{(\alpha_n^i)^2 (\sigma_n - \sigma_1) \sigma_n^2}{(\alpha_1^i)^2 \sigma_1^3 + (\alpha_n^i)^2 \sigma_n^3} \alpha_1^i$$

and

$$\alpha_n^{i+1} = - \frac{(\alpha_1^i)^2 (\sigma_n - \sigma_1) \sigma_1^2}{(\alpha_1^i)^2 \sigma_1^3 + (\alpha_n^i)^2 \sigma_n^3} \alpha_n^i.$$

In particular

$$\frac{\alpha_1^{i+1}}{\alpha_n^{i+1}} = - \frac{\sigma_n^2}{\sigma_1^2} \frac{\alpha_n^i}{\alpha_1^i} = \frac{\alpha_1^{i-1}}{\alpha_n^{i-1}}$$

so that the ratios  $\alpha_1^i / \alpha_n^i$  assume just two values for all  $i$  (depending on  $i$  even or odd). Now

$$|\alpha_{1+1}^i| \geq \left(1 - \frac{\sigma_1}{\sigma_n}\right) \frac{\gamma^i}{(\alpha_1^i)^2 + (\alpha_n^i)^2} |\alpha_1^i|,$$

and

$$|\alpha_n^{i+1}| \geq \left(1 - \frac{\sigma_1}{\sigma_n}\right) \left(\frac{\sigma_1}{\sigma_n}\right)^2 \frac{(\alpha_1^i)^2}{(\alpha_1^i)^2 + (\alpha_n^i)^2} |\alpha_n^i|$$

so that

$$\begin{aligned} |\alpha_1^{i+1}| + |\alpha_n^{i+1}| &\geq \left(1 - \frac{\sigma_1}{\sigma_n}\right) \left(\frac{\sigma_1}{\sigma_n}\right)^2 \frac{\alpha_1^i \alpha_n^i}{(\alpha_1^i)^2 + (\alpha_n^i)^2} (|\alpha_1^i| + |\alpha_n^i|) \\ &\geq \left(1 - \frac{\sigma_1}{\sigma_n}\right) \left(\frac{\sigma_1}{\sigma_n}\right)^2 \gamma (|\alpha_1^i| + |\alpha_n^i|) \end{aligned}$$

where  $\gamma = \min \left\{ \frac{\frac{|\alpha_1^i|}{\alpha_n^i}}{1 + \frac{|\alpha_1^i|}{\alpha_n^i}}, \frac{\frac{|\alpha_1^{i+1}|}{\alpha_n^{i+1}}}{1 + \frac{|\alpha_1^{i+1}|}{\alpha_n^{i+1}}} \right\} < 1$ , and  $\gamma$  is independent

of  $i$ . This inequality shows that the rate of convergence of steepest descent is linear.

Definition: Directions  $t_{\sim 1}, t_{\sim 2}$  are conjugate with respect to  $C$  if

$$t_{\sim 1}^T C t_{\sim 2} = 0 \quad (2.5)$$

In what follows it will frequently be convenient to speak about a

'direction of search' without intending to imply that its norm is unity. However, the null vector is excluded from any set of mutually conjugate directions. It is clear that any set of mutually conjugate directions are linearly independent.

Example: The eigenvectors of  $C$  are conjugate. The property of being both conjugate and orthogonal specializes the eigenvectors.

Lemma 2.1: Let  $\tilde{t}_1, \dots, \tilde{t}_n$  be a set of mutually conjugate directions (with respect to  $C$ ). Starting from  $x_1$  let  $x_2, x_3, \dots, x_{n+1}$  be points produced by descent steps applied to (2.1). Then

$$g_i^T \tilde{t}_j = 0, \quad j = 1, 2, \dots, i-1. \quad (2.6)$$

Proof: The descent condition gives  $g_i^T \tilde{t}_{i-1} = 0$  so it is necessary only to verify the result for  $j < i-1$ . We have

$$\begin{aligned} g_i^T \tilde{t}_s &= (Cx_i + b)^T \tilde{t}_s \\ &= (Cx_{s+1} + b + \sum_{k=s+1}^{i-1} \lambda_k C \tilde{t}_k)^T \tilde{t}_s \\ &= g_{s+1}^T \tilde{t}_s + \sum_{k=s+1}^{i-1} \lambda_k \tilde{t}_k^T C \tilde{t}_s, \quad s = 1, 2, \dots, i-2 \\ &= 0 \quad . \quad \square \end{aligned}$$

Corollary 2.1: The minimum of a positive definite quadratic form can be found by making at most one descent step along each of  $n$  mutually conjugate directions.

Proof: From Lemma 2.1 we have  $\underline{g}_{n+1}^T \underline{t}_{\sim i} = 0$ ,  $i = 1, 2, \dots, n$ .

Thus  $\underline{g}_{n+1}$  is orthogonal to  $n$  linearly independent directions and therefore vanishes identically.  $\square$

Remark: A method which minimizes a quadratic form in a finite number of steps is said to have a quadratic termination property.

Example: The sequence of vectors

$$\begin{aligned} \underline{t}_{\sim 1} &= -\underline{g}_{\sim 1}, \\ \underline{t}_{\sim i} &= -\underline{g}_{\sim i} + \frac{\|\underline{g}_{\sim 1}\|^2}{\|\underline{g}_{\sim i-1}\|^2} \underline{t}_{\sim i-1}, \quad i = 2, \dots, n \end{aligned} \quad (2.7)$$

are conjugate. The algorithm based on this choice is called the method of conjugate gradients.

We now consider the generation of sequences of conjugate directions to provide a basis for a descent calculation. To do this we note that the minimum of (2.1) is at  $\underline{x} = -\underline{C}^{-1}\underline{b}$  so that if we minimize in the direction  $\underline{t} = -\underline{C}^{-1}(\underline{C}\underline{x}_1 + \underline{b}) = -\underline{C}^{-1}\nabla F(\underline{x}_1)$  then the minimum is found in a single step. In general  $\underline{C}^{-1}$  is not known in advance, so that we are lead to consider processes in which each step consists of two parts  
(i) a descent calculation in the direction

$$\underline{t}_{\sim i} = -\underline{H}_i \underline{g}_{\sim i} \quad (2.8)$$

where  $\underline{H}_i$  is the current estimate of  $\underline{C}^{-1}$ , and (ii) the calculation of a correction to  $\underline{H}_i$  which serves both the purposes of making the  $\underline{t}_{\sim i}$  conjugate and making  $\underline{H}_i$  approach  $\underline{C}^{-1}$ . It is convenient in what follows to assume that the  $\underline{H}_i$  are symmetric. This seems a natural condition given the symmetry of  $\underline{C}$  but is in fact not necessary.

If we assume that  $t_{\tilde{s}}$ ,  $s < i$ , are mutually conjugate then the condition that each be conjugate to  $t_{\tilde{i}}$  is

$$t_{\tilde{i}}^T C t_{\tilde{s}} = g_{\tilde{i}}^T H_i C t_{\tilde{s}} = 0, \quad s < i,$$

and, by Lemma 2.1, this is certainly satisfied if

$$H_i C t_{\tilde{s}} = \rho_{\tilde{s}} t_{\tilde{s}}, \quad s < i.$$

We write this equation in the equivalent form (multiplying both sides by  $\lambda_s$ )

$$H_i y_{\tilde{s}} = \rho_{\tilde{s}} d_{\tilde{s}}, \quad s < i, \quad (2.9)$$

where

$$d_{\tilde{i}} = x_{\tilde{i+1}} - x_{\tilde{i}}, \quad y_{\tilde{i}} = g_{\tilde{i+1}} - g_{\tilde{i}}. \quad (2.10)$$

Consider the symmetric updating formula

$$H_{i+1} = H_i + \xi_i d_i d_i^T + \eta_i H_i y_i y_i^T H_i - \zeta_i (d_i y_i^T H_i + H_i y_i d_i^T) \quad (2.11)$$

where  $\xi_i$ ,  $\eta_i$ ,  $\zeta_i$  are to be determined (or prescribed). We have

$$H_{i+1} y_{\tilde{s}} = H_i y_{\tilde{s}} = \rho_{\tilde{s}} d_{\tilde{s}}, \quad s < i, \quad \text{provided (2.9) holds as}$$

$$d_{\tilde{i}}^T y_{\tilde{s}} = \lambda_i \lambda_s t_{\tilde{i}}^T C t_{\tilde{s}} = 0, \quad \text{and} \quad y_{\tilde{i}}^T H_i y_{\tilde{s}} = \lambda_s y_{\tilde{i}}^T H_i C t_{\tilde{s}} = \rho_{\tilde{s}} \lambda_s y_{\tilde{i}}^T d_{\tilde{s}} = \rho_{\tilde{s}} \lambda_s \lambda_i d_{\tilde{i}}^T C d_{\tilde{s}} = 0. \quad \text{Thus (2.9) is satisfied for } i := i+1 \text{ if}$$

$$0 = 1 + \eta_i (y_{\tilde{i}}^T H_i y_{\tilde{i}}) - \zeta_i (d_{\tilde{i}}^T y_{\tilde{i}}), \quad (2.12)$$

and

$$\rho_i = \xi_i (d_{\tilde{i}}^T y_{\tilde{i}}) - \zeta_i (y_{\tilde{i}}^T H_i y_{\tilde{i}}). \quad (2.13)$$

If  $\xi_i$  and  $\eta_i$  are expressed in terms of  $\rho_i$  and  $\zeta_i$  from (2.12) and (2.13) we have

$$\eta_i = -\frac{1}{\tilde{y}_i^T H_i \tilde{y}_i} + \zeta_i \frac{\tilde{d}_i^T \tilde{y}_i}{\tilde{y}_i^T H_i \tilde{y}_i},$$

and

$$\xi_i = \frac{\rho_i}{\tilde{d}_i^T \tilde{y}_i} + \zeta_i \frac{\tilde{y}_i^T H_i \tilde{y}_i}{\tilde{d}_i^T \tilde{y}_i},$$

so that equation (2.11) becomes

$$\begin{aligned} H_{i+1} &= H_i + \rho_i \frac{\tilde{d}_i \tilde{d}_i^T}{\tilde{d}_i^T \tilde{y}_i} - \frac{H_i \tilde{y}_i \tilde{y}_i^T H_i}{\tilde{y}_i^T H_i \tilde{y}_i} \\ &\quad + \zeta_i \left\{ \frac{\tilde{y}_i^T H_i \tilde{y}_i}{\tilde{d}_i^T \tilde{y}_i} \tilde{d}_i \tilde{d}_i^T + \frac{\tilde{d}_i^T \tilde{y}_i}{\tilde{y}_i^T H_i \tilde{y}_i} H_i \tilde{y}_i \tilde{y}_i^T H_i - \tilde{d}_i \tilde{y}_i^T H_i - H_i \tilde{y}_i \tilde{d}_i^T \right\} \\ &= D(\rho_i, H_i) + \zeta_i \tau_i \tilde{v}_i \tilde{v}_i^T \end{aligned} \quad (2.14)$$

where

$$\tilde{v}_i = \tilde{d}_i - \frac{1}{\tau_i} H_i \tilde{y}_i, \quad (2.15)$$

and

$$\tau_i = \frac{\tilde{y}_i^T H_i \tilde{y}_i}{\tilde{d}_i^T \tilde{y}_i}. \quad (2.16)$$

Example: The particular case  $\rho_i = 1$ ,  $\zeta_i = 0$ ,  $i = 1, 2, \dots$  gives the variable metric or DFP formula which is the most frequently used member of the family.

The class of formulae described by (2.14) generate recursively a set of conjugate directions so that the first of our aims is satisfied. It still remains to show the relationship between the  $H_i$  and  $C^{-1}$ . To do

this note that (2.9) can be written (introducing the symmetric square root  $C^{1/2}$  of the positive definite matrix  $C$ ).

$$C^{1/2} H_i C^{1/2} t_{\sim s} = \rho_s C^{1/2} t_{\sim s}, \quad s = 1, 2, \dots, i-1,$$

or, more briefly,

$$\hat{H}_i t_{\sim s} = \rho_s t_{\sim s}, \quad s = 1, 2, \dots, i-1. \quad (2.9a)$$

Defining the matrix  $\hat{T}$  by  $\kappa_i(\hat{T}) = \frac{t_{\sim s}}{(t_{\sim s}^T C t_{\sim s})^{1/2}}$ ,  $i = 1, 2, \dots, n$ , and the

diagonal matrix  $P$  by  $P_{ii} = \rho_i$ ,  $i = 1, 2, \dots, n$ , we can write (2.9a) in the case  $i = n+1$  in the form

$$H_{n+1} T = T P. \quad (2.9b)$$

Now  $\hat{T}$  is an orthogonal matrix so that

$$\hat{H}_{n+1} = \hat{T} P \hat{T}^T,$$

whence

$$H_{n+1} = C^{-1/2} \hat{T} P \hat{T}^T C^{-1/2}. \quad (2.17)$$

In particular, if  $P = \rho I$ ,

$$H_{n+1} = \rho C^{-1}. \quad (2.18)$$

Remark: Remember the motivation for developing the recursion (2.14) is the search for efficient descent directions. Specifically we are looking not only for conjugate directions but also for good estimates of the inverse Hessian. This indicates that  $\rho = 1$  is the natural choice (or at least  $\rho = \text{constant}$ ), and almost all published methods use  $\rho = 1$ . However, from (2.17), the choice of  $\rho$  variable may well have scaling advantages in the initial phases of a computation with a general objective function. Presumably the strategy for choosing  $\rho$  should make  $\rho \rightarrow \text{constant}$  to ensure a fast rate of ultimate convergence.



Lemma 2.2: Provided the descent condition is satisfied,

$$H_{i+1} \underline{g}_{i+1} \parallel \underline{v}_i \text{ or null.}$$

Remark: In what follows it is convenient to drop the  $i$  subscripts. Quantities subscripted  $i+1$  will be starred. In what follows we assume  $\rho$  is constant.

Proof: We have (using the descent condition, the definition of  $t$ , and  $d = \lambda t$ )

$$\begin{aligned} D(\rho, H) \underline{g}^* &= (H + \rho \frac{d d^T}{\underline{y}^T \underline{y}} - \frac{H \underline{y} \underline{y}^T H}{\underline{y}^T H \underline{y}}) \underline{g}^* \\ &= H \underline{y} + H \underline{g} - \frac{H \underline{y} \underline{y}^T H (\underline{y} + \underline{g})}{\underline{y}^T H \underline{y}} \\ &= \text{cm}^{-1} (d - \frac{\underline{y}^T d}{\underline{y}^T H \underline{y}} H \underline{y}) \end{aligned}$$

Whence

$$H^* \underline{g}^* = - \left( \frac{1}{\lambda} + \zeta \tau \underline{v}^T \underline{g}^* \right) \underline{v} \quad \text{cl} \quad (2.19)$$

Remark: (i) The condition that  $H^* \underline{g}^* = 0$  when  $\underline{v} \neq 0$  gives a condition which determines  $\zeta$ . We have

$$\underline{v}^T \underline{g}^* = - \frac{1}{\tau} \underline{g}^{*T} H \underline{y} = - \frac{1}{\tau} \underline{g}^{*T} H \underline{g}^*$$

so that (from (2.19))

$$\zeta = \frac{1}{\lambda \underline{g}^{*T} H \underline{g}^*} \quad (2.20)$$

Provided this value of  $\zeta$  is excluded from consideration, then  $\tilde{x}^*$  is independent of  $\zeta$ . Note that this result is true for a general function as no properties specific to a quadratic form have been used in its derivation.

(ii) We can only have  $\tilde{y} = 0$  with  $\tilde{d}$  and  $\tilde{g}^*$  nonnull if  $H$  is singular, and in this case  $H^*$  is also singular and the null space of  $H^*$  is at least as large as that of  $H$ . This follows from (2.15) which can vanish only if (a)  $H\tilde{g}^* = 0$  and  $1 + \frac{1}{\lambda\tau} = 0$ , or (b)  $H\tilde{g}$  and  $H\tilde{g}^*$  are parallel. Now if  $H$  is singular  $\exists \tilde{w}, \tilde{w}^T H = 0$ . Thus  $\tilde{w}^T \tilde{d} = 0$ , and hence  $\tilde{w}^T H^* = 0$ .

Clearly it is important that  $H_i$  positive definite  $\Rightarrow H_{i+1}$  positive definite,  $i = 1, 2, \dots$  in order that premature termination should be avoided ( $H^* \tilde{g}^* = 0$  and  $H^*$  positive definite  $\Rightarrow \tilde{g}^* = 0$  whence  $\tilde{x}^*$  is a stationary point). Conditions which ensure this are given in the following lemma (due to Powell).

Lemma 2.3: If  $0 < \rho, \tau < \infty$ ,  $H$  positive semidefinite, and  $H H^+ \tilde{v} = \tilde{v}$  (where  $H^+$  is the generalized inverse of  $H$ ), then  $H^*$  is positive semidefinite, and the null space of  $H^*$  is equal to that of  $H$  provided

$$\zeta > \frac{-\tilde{y}^T \tilde{d}}{(\tilde{d}^T H^+ \tilde{d})(\tilde{y}^T H \tilde{y}) - (\tilde{d}^T \tilde{y})^2} \quad (2.21)$$

Proof: We first note the identity

$$D(\rho, H) = (I + \tilde{u}\tilde{y}^T)H(I + \tilde{y}\tilde{u}^T) \quad (2.22)$$

where

$$\underline{\underline{u}} = \frac{1}{\sqrt{(\underline{\underline{d}}^T \underline{\underline{y}})(\underline{\underline{y}}^T H \underline{\underline{y}})}} \left\{ \sqrt{\rho} \underline{\underline{d}} - \frac{1}{\sqrt{\tau}} H \underline{\underline{y}} \right\} . \quad (2.23)$$

and

$$\det(I + \underline{\underline{u}} \underline{\underline{u}}^T) = 1 + \underline{\underline{y}}^T \underline{\underline{u}} = \sqrt{\frac{\rho}{\tau}} \quad (2.24)$$

so that, by the assumptions,  $I + \underline{\underline{u}} \underline{\underline{u}}^T$  is nonsingular. Now  $\underline{\underline{y}}^T \underline{\underline{v}} = 0$  so that  $H^*$  can be written

$$H^* = (I + \underline{\underline{u}} \underline{\underline{u}}^T) (H + \zeta \tau \underline{\underline{v}} \underline{\underline{v}}^T) (I + \underline{\underline{y}} \underline{\underline{y}}^T) . \quad (2.25)$$

Thus the problem reduces to considering  $H + \zeta \tau \underline{\underline{v}} \underline{\underline{v}}^T$ . We have

$$H + \zeta \tau \underline{\underline{v}} \underline{\underline{v}}^T = H (I + \zeta \tau H^+ \underline{\underline{v}} \underline{\underline{v}}^T) .$$

The null spaces of  $H$  and  $H^*$  will agree provided  $I + \zeta \tau H^+ \underline{\underline{v}} \underline{\underline{v}}^T$  is nonsingular. The condition for singularity is

$$\begin{aligned} 0 &= \det(I + \zeta \tau H^+ \underline{\underline{v}} \underline{\underline{v}}^T) \\ &= 1 + \zeta \tau \underline{\underline{v}}^T H^+ \underline{\underline{v}} . \end{aligned}$$

Noting that  $H H^+ H = H$ , and  $H H^+ \underline{\underline{v}} = \underline{\underline{v}} \Rightarrow H H^+ \underline{\underline{d}} = \underline{\underline{d}}$  we have

$$\begin{aligned} 1 + \zeta \tau \underline{\underline{v}}^T H^+ \underline{\underline{v}} &= 1 + \zeta \tau \left( \underline{\underline{d}}^T H^+ \underline{\underline{d}} - 2 \frac{\underline{\underline{d}}^T \underline{\underline{y}}}{\tau} + \frac{\underline{\underline{y}}^T H \underline{\underline{y}}}{\tau^2} \right) \\ &= 1 + \zeta \tau \left( \underline{\underline{d}}^T H^+ \underline{\underline{d}} - \frac{(\underline{\underline{d}}^T \underline{\underline{y}})^2}{\underline{\underline{y}}^T H \underline{\underline{y}}} \right) \end{aligned}$$

and this vanishes provided

$$\zeta = \frac{-\underline{\underline{y}}^T \underline{\underline{d}}}{(\underline{\underline{d}}^T H^+ \underline{\underline{d}})(\underline{\underline{y}}^T H \underline{\underline{y}}) - (\underline{\underline{y}}^T \underline{\underline{d}})^2} .$$

The stated result is a consequence of this and the observation that decreasing  $\zeta$  below this value will make  $H^*$  indefinite.  $\square$

Remark: (i) The condition on  $\tau$  is automatically satisfied if  $H$  is positive definite and the descent condition is satisfied for then  $\tilde{d}^T \tilde{y} = -\tilde{g}^T \tilde{d} = \lambda \tilde{g}^T H \tilde{g}$ . However the lemma does not require that the descent condition be satisfied and remains valid even though the exact minimum in the direction  $\tilde{t}$  is not found. In this case the condition on  $\tau$  is necessary.

Corollary 2.2: If  $H_1$  positive definite, and  $H_{i+1} = D(\rho, H_i)$ ,  $i = 1, 2, \dots$  then provided the descent condition is satisfied for  $i = 1, 2, \dots$  then  $H_{i+1}$  is positive definite.

Proof: This is a consequence of (2.22) and the above remark which shows that if  $H$  is positive definite, and if the descent condition is satisfied, then  $I + \tilde{u}\tilde{y}^T$  is nonsingular.  $\square$

Theorem 2.1: (Dixon's equivalence theorem). If (i) the formula (2.14) is used to generate descent directions, (ii)  $\zeta_i$  satisfies (2.21) for  $i = 1, 2, \dots$  and  $H_1$  is positive definite, and (iii) the descent condition is satisfied in each descent step, then the sequence of points generated by the algorithm depends only on  $F$ ,  $H_1$ ,  $\rho$ , and  $\tilde{x}_1$  and is independent of  $\zeta_i$ ,  $i = 1, 2, \dots$ .

Remark: It is important to note that  $F$  is not restricted to be a quadratic form in this result.

Proof: Let  $D_1 = H_1$ ,  $D_i = D(\rho, D_i)$ ,  $i = 2, 3, \dots$ . We show that if  $H_i = D_i + \alpha \tilde{d}_i \tilde{d}_i^T$ , then  $H_{i+1} = D_{i+1} + \beta \tilde{d}_{i+1} \tilde{d}_{i+1}^T$ . By Lemma 2.2 we have  $H^* = D(\rho, H) + \gamma \tilde{d}^* \tilde{d}^{*T}$ . Now

$$\begin{aligned}
 D(\rho, H) &= D + \rho \frac{\tilde{d} \tilde{d}^T}{\tilde{d}^T \tilde{y}} - \frac{(D + \alpha \tilde{d} \tilde{d}^T) \tilde{y} \tilde{y}^T (D + \alpha \tilde{d} \tilde{d}^T)}{\tilde{y}^T (D + \alpha \tilde{d} \tilde{d}^T) \tilde{y}} + \alpha \tilde{d} \tilde{d}^T \\
 &= D + \rho \frac{\tilde{d} \tilde{d}^T}{\tilde{d}^T \tilde{y}} - \frac{D \tilde{y} \tilde{y}^T D + \alpha (\tilde{y}^T \tilde{d}) (D \tilde{y} \tilde{d}^T + \tilde{d} \tilde{y}^T D) + \alpha^2 (\tilde{d}^T \tilde{y})^2 \tilde{d} \tilde{d}^T}{\tilde{y}^T D \tilde{y} + \alpha (\tilde{y}^T \tilde{d})^2} + \alpha \tilde{d} \tilde{d}^T \\
 &= D^* - \frac{-D \tilde{y} \tilde{y}^T D \alpha \frac{(\tilde{y}^T \tilde{d})^2}{\tilde{y}^T D \tilde{y}} + \alpha (\tilde{y}^T \tilde{d}) (D \tilde{y} \tilde{d}^T + \tilde{d} \tilde{y}^T D) - \alpha (\tilde{y}^T D \tilde{y}) \tilde{d} \tilde{d}^T}{\tilde{y}^T D \tilde{y} + \alpha (\tilde{y}^T \tilde{d})^2} \\
 &= D^* + \frac{\alpha (\tilde{y}^T D \tilde{y})}{\tilde{y}^T D \tilde{y} + \alpha (\tilde{y}^T \tilde{d})^2} \left( \tilde{d} - \frac{\tilde{y}^T \tilde{d}}{\tilde{y}^T D \tilde{y}} D \tilde{y} \right) \left( \tilde{d} - \frac{\tilde{y}^T \tilde{d}}{\tilde{y}^T D \tilde{y}} D \tilde{y} \right)^T. \quad (2.26)
 \end{aligned}$$

By Lemma 2.2,  $\tilde{d}^* \parallel D(\rho, H) \tilde{g}^*$ . By (2.26)  $D(\rho, H) \tilde{g}^* \parallel D^* \tilde{g}^*$ . Thus

$\tilde{d}_j \parallel D_j \tilde{g}_j$ ,  $j = 1, 2, \dots, i \Rightarrow \tilde{d}_{i+1} \parallel D_{i+1} \tilde{g}_{i+1}$ . But the case  $j = 1$  is a consequence of Lemma 2.2 so the result follows by induction.  $\square$

Example: Equivalence results for a wide class of conjugate direction algorithms applied to a given positive definite quadratic form can be demonstrated by noting that at the  $i$ -th stage we find the minimum in the translation to  $\tilde{x}_1$  of the subspace spanned by  $\tilde{t}_1, \dots, \tilde{t}_i$ , and that this subspace is also spanned by  $H_1 \tilde{g}_1, \dots, H_i \tilde{g}_i$ . Thus  $\tilde{x}_{i+1}$  depends only on

$\tilde{x}_1, \dots, \tilde{x}_i$  and not on the particular updating formula for the inverse Hessian estimate. If  $H_{i-1} = I$  this equivalence extends to the conjugate gradient algorithm (2.7).

Lemma 2.4: If the descent condition is satisfied at each stage then the sequence  $\tilde{g}_i^T D_i \tilde{g}_i$ ,  $i = 1, 2, \dots$  is strictly decreasing provided  $D_1$  is positive definite.

Proof: We have (as  $\tilde{g}^{*T} d = 0$ )

$$\begin{aligned} \tilde{g}^{*T} D \tilde{g}^* &= \tilde{g}^{*T} D \tilde{g}^* - \frac{(\tilde{g}^{*T} D y)^2}{y^T D y}, \\ &= \tilde{g}^{*T} D \tilde{g}^* - \frac{(\tilde{g}^{*T} D \tilde{g}^*)^2}{\tilde{g}^{*T} D \tilde{g}^* + \tilde{g}^T D \tilde{g}}, \\ &= \frac{(\tilde{g}^{*T} D \tilde{g}^*)(\tilde{g}^T D \tilde{g})}{\tilde{g}^{*T} D \tilde{g}^* + \tilde{g}^T D \tilde{g}}. \end{aligned}$$

Thus

$$\frac{1}{\tilde{g}^{*T} D \tilde{g}^*} = \frac{1}{\tilde{g}^{*T} D \tilde{g}^*} + \frac{1}{\tilde{g}^T D \tilde{g}}. \quad (2.27)$$

By Corollary 2.2, the  $D_i$  are positive definite so that the desired result follows from (2.27).  $\square$

Remark: This result indicates a potential defect of the DFP algorithm. For if the choice of  $D_1$  is poor in the sense that it leads to too small a value of  $\tilde{g}_1^T D_1 \tilde{g}_1$  then the algorithm has no mechanism to correct this, and must initially generate a sequence of directions which are nearly orthogonal to the gradient. This must also happen if, for any

reason, an abnormally small value of  $\underline{g}^T D \underline{g}$  is generated at some stage. A possible cause of such behaviour is poor scaling of the problem.

Lemma 2.5: Corresponding to the formula (2.14) for updating H there is a similar formula for updating  $H^{-1}$ . Specifically we have

$$H^{*-1} = D(\rho, H)^{-1} + \gamma \mu \underline{w} \underline{w}^T \quad (2.28)$$

where

$$D(\rho, H)^{-1} = H^{-1} + \left( \frac{1}{\rho} + \mu \right) \frac{\underline{y} \underline{y}^T}{\underline{d}^T \underline{y}} - \frac{1}{\underline{d}^T \underline{y}} \left( \underline{y} \underline{d}^T H^{-1} + H^{-1} \underline{d} \underline{y}^T \right), \quad (2.29)$$

$$\mu = \frac{\underline{d}^T H^{-1} \underline{d}}{\underline{d}^T \underline{y}}, \quad (2.30)$$

$$\underline{w} = \underline{y} - \frac{1}{\mu} H^{-1} \underline{d}, \quad (2.31)$$

and  $\gamma$  is related to  $\zeta$  by

$$\gamma = - \frac{\zeta \tau \mu}{1 + \zeta \tau \underline{v}^T H^{-1} \underline{v}}. \quad (2.32)$$

Proof: From (2.22) we have

$$D(\rho, H)^{-1} = \left( I - \sqrt{\frac{\tau}{\rho}} \underline{y} \underline{u}^T \right) H^{-1} \left( I - \sqrt{\frac{\tau}{\rho}} \underline{u} \underline{y}^T \right)$$

and (2.29) follows from this by an elementary calculation. From (2.25)

$$\begin{aligned} H^{*-1} &= \left( I - \sqrt{\frac{\tau}{\rho}} \underline{y} \underline{u}^T \right) \left( H + \zeta \tau \underline{v} \underline{v}^T \right)^{-1} \left( I - \sqrt{\frac{\tau}{\rho}} \underline{u} \underline{y}^T \right), \\ &= \left( I - \sqrt{\frac{\tau}{\rho}} \underline{y} \underline{u}^T \right) \left( H^{-1} - \frac{\zeta \tau}{1 + \zeta \tau \underline{v}^T H^{-1} \underline{v}} H^{-1} \underline{v} \underline{v}^T H^{-1} \right) \left( I - \sqrt{\frac{\tau}{\rho}} \underline{u} \underline{y}^T \right) \end{aligned} \quad (2.33)$$

Now

$$\begin{aligned}
 (\mathbf{I} - \sqrt{\frac{\tau}{\rho}} \underline{\underline{y}} \underline{\underline{u}}^T) \mathbf{H}^{-1} \underline{\underline{v}} &= \mathbf{H}^{-1} \underline{\underline{d}} - \left( \frac{1}{\tau} + \sqrt{\frac{\tau}{\rho}} \underline{\underline{u}}^T \mathbf{H}^{-1} \underline{\underline{d}} - \frac{1}{\sqrt{\rho\tau}} \underline{\underline{u}}^T \underline{\underline{y}} \right) \underline{\underline{y}} \\
 &= \mathbf{H}^{-1} \underline{\underline{d}} - \mu \underline{\underline{y}} \\
 &= -\mu \underline{\underline{w}} \quad , \tag{2.34}
 \end{aligned}$$

so that (2.28) is a direct consequence of (2.33) and (2.34).  $\square$

Remark: If we take  $\underline{\underline{y}} = -\frac{1}{\underline{\underline{d}}^T \underline{\underline{d}}}$  in (2.28) then we obtain

$$\begin{aligned}
 G(\rho, \mathbf{H}^{-1}) &= \mathbf{H}^{-1} + \frac{1}{\rho} \frac{\underline{\underline{y}} \underline{\underline{y}}^T}{\underline{\underline{d}}^T \underline{\underline{d}}} - \frac{\mathbf{H}^{-1} \underline{\underline{d}} \underline{\underline{d}}^T \mathbf{H}^{-1}}{\underline{\underline{d}}^T \mathbf{H}^{-1} \underline{\underline{d}}} \\
 &= (\mathbf{I} + \underline{\underline{z}} \underline{\underline{d}}^T) \mathbf{H}^{-1} (\mathbf{I} + \underline{\underline{d}} \underline{\underline{z}}^T) \tag{2.35}
 \end{aligned}$$

where

$$\underline{\underline{z}} = \frac{1}{\sqrt{(\underline{\underline{d}}^T \underline{\underline{y}})(\underline{\underline{d}}^T \mathbf{H}^{-1} \underline{\underline{d}})}} \left\{ \frac{1}{\sqrt{\rho}} \underline{\underline{y}} - \frac{1}{\sqrt{\mu}} \mathbf{H}^{-1} \underline{\underline{d}} \right\} . \tag{2.36}$$

We have

$$\begin{aligned}
 D(\rho, \mathbf{H}^{-1})^{-1} &= G(\rho, \mathbf{H}^{-1}) + \frac{\mu}{\underline{\underline{y}}^T \underline{\underline{d}}} \underline{\underline{w}} \underline{\underline{w}}^T \\
 &= (\mathbf{I} + \underline{\underline{z}} \underline{\underline{d}}^T) (\mathbf{H}^{-1} + \frac{\mu}{\underline{\underline{y}}^T \underline{\underline{d}}} \underline{\underline{w}} \underline{\underline{w}}^T) (\mathbf{I} + \underline{\underline{d}}^T \underline{\underline{z}}) \tag{2.37}
 \end{aligned}$$

as  $\underline{\underline{d}}^T \underline{\underline{w}} = 0$  .



To summarize these results we have the following:

$$(i) \quad D(\rho, H) = (I + \frac{yy^T}{\rho})H(I + \frac{yu^T}{\rho}) \quad ,$$

$$G(\rho, H^{-1}) = (I + \frac{zd^T}{\rho})H^{-1}(I + \frac{dz^T}{\rho})$$

$$(ii) \quad D(\rho, H)^{-1} = G(\rho, H^{-1}) + \frac{\mu}{y^T d} \frac{ww^T}{\rho}$$

$$G(\rho, H^{-1})^{-1} = D(\rho, H) + \frac{\tau}{y^T d} \frac{vv^T}{\rho}$$

	update formula	update formula for inverse
$D(\rho, H)$	$H + \rho \frac{dd^T}{d^T y} - \frac{Hyy^T H}{y^T H y}$	$H^{-1} + (\frac{1}{\rho} + \mu) \frac{yy^T}{d^T y} - \frac{1}{d^T y} (y d^T H^{-1} + H^{-1} d y^T)$
$G(\rho, H^{-1})$	$H^{-1} + \frac{1}{\rho} \frac{yy^T}{y^T d} - \frac{H^{-1} d d^T H^{-1}}{d^T H^{-1} d}$	$H + (\rho + \tau) \frac{dd^T}{d^T y} - \frac{1}{d^T y} (d y^T H + H y d^T)$

$D(\rho, H)$  ,  $G(\rho, H^{-1})$  have been called dual formulae by Fletcher.

**Lemma 2.6:** Let  $A$  be a symmetric matrix,  $A = T\Lambda T^T$  where  $\Lambda$  diagonal ( $\Lambda_{ii} = \lambda_i$  ,  $i = 1, 2, \dots, n$ ) , and  $T$  orthogonal. Let  $\lambda_i^*$  ,  $i = 1, 2, \dots, n$  be the eigenvalues of  $A + \sigma \frac{aa^T}{a^T a}$  , then either  $\sigma > 0$  and  $\lambda_i \leq \lambda_i^* \leq \lambda_{i+1}$  ,  $i = 1, 2, \dots, n$  , or  $\sigma < 0$  and  $\lambda_{i-1} \leq \lambda_i^* \leq \lambda_i$  .

Proof: We have

$$\begin{aligned} \det\{A + \sigma \underline{\underline{a}} \underline{\underline{a}}^T - \lambda I\} &= \prod_{i=1}^n (\lambda_i - \lambda) \det\{I + \sigma(\Lambda - \lambda I)^{-1} (\underline{\underline{T}}^T \underline{\underline{a}}) (\underline{\underline{T}}^T \underline{\underline{a}})^T\} \\ &= \prod_{i=1}^n (\lambda_i - \lambda) (1 + \sigma (\underline{\underline{T}}^T \underline{\underline{a}})^T (\Lambda - \lambda I)^{-1} (\underline{\underline{T}}^T \underline{\underline{a}})) \\ &= \prod_{i=1}^n (\lambda_i - \lambda) \left\{ 1 + \sigma \sum_{i=1}^n \frac{v_i^2}{\lambda_i - \lambda} \right\}, \quad v_i = \rho_i (\underline{\underline{T}}^T \underline{\underline{a}}), \end{aligned}$$

and the desired result is an easy consequence of this expression.  $\square$

In the following theorem we consider specifically the minimization of a positive definite quadratic form. We assume that the initial estimate of the Hessian  $H_1$  is positive definite, and we make use of the following sequences of updates for the current Hessian estimates

- (a)  $H_{i+1} = D(\rho, H_i)$ ,  $i = 1, 2, \dots$ , and
- (b)  $\hat{H}_{i+1} = G(\rho, H_i^{-1})^{-1}$ ,  $i = 1, 2, \dots$ .

Further we do not assume that the descent condition is satisfied.

Theorem 2.2: (i) Let  $K_i = C^{1/2} H_i C^{1/2}$ , and let the eigenvalues of  $K_i$  ordered in increasing magnitude be  $\lambda_j^{(i)}$ ,  $j = 1, 2, \dots, n$ . Then if  $\lambda_j^{(1)} \geq \rho$  then  $\lambda_j^{(1)} \geq \lambda_j^{(2)} \geq \dots \geq \rho$ , while if  $\lambda_j^{(1)} \leq \rho$  then  $\lambda_j^{(1)} \leq \lambda_j^{(2)} \leq \dots \leq \rho$  for  $j = 1, 2, \dots, n$ . (ii) Let  $\hat{K}_i = C^{1/2} \hat{H}_i C^{1/2}$ , and let the eigenvalues of  $\hat{K}_i$  be  $\hat{\lambda}_j^{(i)}$ ,  $j = 1, 2, \dots, n$ . If  $\hat{\lambda}_j^{(1)} \geq \rho$  then  $\hat{\lambda}_j^{(1)} \geq \hat{\lambda}_j^{(2)} \geq \dots \geq \rho$ , while if  $\hat{\lambda}_j^{(1)} \leq \rho$  then  $\hat{\lambda}_j^{(1)} \leq \hat{\lambda}_j^{(2)} \leq \dots \leq \rho$  for  $j = 1, 2, \dots, n$ .

Remark: This result is important because it shows that we have a 'weak' convergence result for these Hessian estimates when minimizing a positive definite quadratic form even when the descent conditions is not satisfied at each step.

Proof: Noting that  $C^{1/2} \underline{d} = C^{-1/2} \underline{y} = \underline{a}$ , we can write the formula for updating  $K$  as

$$K^* = K + \rho \frac{\underline{a} \underline{a}^T}{\underline{a}^T \underline{a}} - \frac{K \underline{a} \underline{a}^T K}{\underline{a}^T K \underline{a}}$$

We can break this into the two operations

$$J = K - \frac{K \underline{a} \underline{a}^T K}{\underline{a}^T K \underline{a}}, \quad \text{and}$$

$$K^* = J + \rho \frac{\underline{a} \underline{a}^T}{\underline{a}^T \underline{a}}.$$

Note that  $J$  has a zero eigenvalue, and that  $\underline{a}$  is the corresponding eigenvector. By Lemma 2.6 we have  $\lambda_1(J) = 0$ , and  $\lambda_{j-1} \leq \lambda_j(J) \leq \lambda_j$  for  $j = 2, 3, \dots, n$ . The rank one modification which takes  $J$  into  $K^*$  changes the zero eigenvalue to  $\rho$  and leaves the other eigenvalues of  $J$  unchanged. Assume that  $\lambda_j(J) \leq \rho \leq \lambda_{j+1}(J)$  then reordering the eigenvalues in increasing order of magnitude we have  $\lambda_k^* = \lambda_{k+1}(J)$ ,  $k = 1, 2, \dots, j-1$ ,  $\lambda_j^* = \rho$ ,  $\lambda_k^* = \lambda_k(J)$ ,  $k = j+1, \dots, n$ . This establishes the first part of the theorem. The second is demonstrated in similar fashion by noting that  $K^{-1}$  satisfies a formally similar update relation. This establishes the result for the eigenvalues of  $\hat{K}^{-1}$  and hence for their reciprocals.  $\square$

Remark: Note that both  $H_i$  and  $\hat{H}_i$  are positive definite  $i = 2, 3, \dots$  if  $H_1$  is positive definite. In this case the result does not depend on the descent conditions being satisfied.

Theorem 2.3: Let  $H$  be positive definite, and consider a step  $d$  in the direction  $-Hg$ . Let  $H^* = D(\rho, H)$ ,  $\hat{H} = G(\rho, H)^{-1} =$

$$D(\rho, H) + \frac{\tau}{\tilde{y}^T d} \tilde{v} \tilde{v}^T, \text{ and } H_\theta = \theta \hat{H} + (1-\theta)H^* = D(\rho, H) + \frac{\theta\tau}{\tilde{y}^T d} \tilde{v} \tilde{v}^T. \text{ Let}$$

$K = C^{1/2} H C^{1/2}$ , and define  $K^*$ ,  $\hat{K}$ ,  $K_\theta$  similarly. Let the eigenvalues of  $K$ ,  $K^*$ ,  $\hat{K}$ ,  $K_\theta$  be  $\lambda_j$ ,  $\lambda_j^*$ ,  $\hat{\lambda}_j$ , and  $\lambda_j^\theta$  respectively,  $j = 1, 2, \dots, n$ . Let  $0 \leq \theta \leq 1$ . If  $\lambda_j \geq \rho$  then  $\lambda_j \geq \hat{\lambda}_j \geq \lambda_j^\theta \geq \lambda_j^* \geq \rho$  while if  $\lambda_j \leq \rho$  then  $\lambda_j \leq \lambda_j^* \leq \lambda_j^\theta \leq \hat{\lambda}_j \leq \rho$ . If  $\theta \notin [0, 1]$  then  $\lambda_j^\theta$  need not lie in the interval defined by  $\lambda_j$  and  $\rho$ .

Proof: It follows from the definition of  $H^*$ ,  $\hat{H}$ , and  $H_\theta$  and Lemma 2.6 that  $\lambda_j^* \leq \lambda_j^\theta \leq \hat{\lambda}_j$ ,  $j = 1, 2, \dots, n$ , provided  $0 \leq \theta \leq 1$ . The first part of the result is now a consequence of Theorem 2.2. To show that  $\lambda_j^\theta$  need not lie in the interval defined by  $\lambda_j$  and  $\rho$ , consider the example

$$C = \begin{bmatrix} 1+\epsilon & \sqrt{\epsilon} \\ \sqrt{\epsilon} & \epsilon \end{bmatrix}, \quad H = I, \quad \rho = 1, \quad a = \begin{bmatrix} 0 \\ 1 \end{bmatrix}.$$

We have  $\lambda_1 = \eta$ ,  $\lambda_2 = 1+2\epsilon - \eta$  where  $\eta = \frac{1}{2}(1+2\epsilon - \sqrt{1+4\epsilon})$ . Thus  $\eta$  is positive and  $O(\epsilon^2)$ . In this case we have

$$K=C, \quad \tau = \frac{a^T K a}{a^T a} = \epsilon, \quad C^{1/2} \tilde{v} = a - \frac{1}{\epsilon} K a = - \begin{bmatrix} 1/\sqrt{\epsilon} \\ 0 \end{bmatrix}.$$

It is readily verified that  $K^* = \begin{bmatrix} \epsilon & 0 \\ 0 & 1 \end{bmatrix}$ , so that  $K_{-\epsilon} = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}$ ,

$K_{1+\epsilon} = \begin{bmatrix} 1+2\epsilon & 0 \\ 0 & 1 \end{bmatrix}$ . In both cases eigenvalues lie outside the prescribed interval. In the first case we have  $0 < \eta$ , and in the second  $1+2\epsilon > 1+2\epsilon - \eta$ .

Remark: This result shows that  $\hat{H}$  gives the best improvement in the eigenvalues  $< \rho$ , while  $H^*$  has a similar property for those  $> \rho$ . This suggests an algorithm in which a choice is made between updating  $H$  to  $\hat{H}$  or  $H^*$  depending on some appropriate criterion. Fletcher suggests that if  $\tau > 1$  (that is,  $\underline{y}^T H \underline{y} > \underline{y}^T C^{-1} \underline{y}$ ) then  $H^*$  should be used, while if  $\tau < 1$  then  $\hat{H}$  is chosen. He has used this criterion in an implementation of Goldstein's algorithm, and has reported satisfactory results.

## Notes

1. For Ostrowski's theorem see his book 'Solution of Equations' (2nd edition) or Kowalik and Osborne. Goldstein's theorem is from his paper 'On steepest descent' in SIAM Control, 1965. Theorem 1.3 is abstracted from Goldstein and Price, 'An Effective Algorithm for Minimization', Num. Math. 1967.
2. For background material see Kowalik and Osborne. The form of the update for the inverse Hessian is due to Powell 'Recent Advances in Unconstrained Optimization' to appear in Math. Prog. It is a specialization of a form derived in Huang, 'Unified approach to quadratically terminating algorithms for function minimization', JOTA, 1970. The form (2.14) and the result of Lemma 2.2 are probably due (in the case  $\rho = 1$ ) to Fletcher 'A new approach to variable metric algorithms', Comp. J., 1970, and Broyden, 'Convergence of a class of double rank minimization algorithms', JIMA, 1970. Lemma 2.3 is due to Powell (to be published). The product update form (2.22) is due to Greenstadt (to be published). Dixon's paper containing Theorem 2.1 is to appear in Math. Prog. The significance of (2.27) for the successful performance of the DFP algorithm was noted in Powell's survey paper already cited. Attention was drawn to the dual updating formulae by Fletcher. This material together with Theorems 2.2 and 2.3 are included in his paper already cited.

## References

- C. G. Broyden (1970): The convergence of a class of double rank minimization algorithms, J. Inst. Math. Applic., 6, pp. 76-90.

- R. Fletcher (1970): A new approach to variable metric algorithms, Computer J., 13, pp. 317-322.
- A. A. Goldstein (1965): On steepest descent, SIAM J. Control, 3, pp. 147-151.
- A. A. Goldstein and J. Price (1967): An effective algorithm for minimization, Num. Math., 10, pp. 184-189.
- N. Y. Huang (1970): Unified approach to quadratically convergent algorithms for function minimization, J.O.T.A., 5, pp. 405-423.
- J. Kowalik and M. R. Osborne (1968): Methods for Unconstrained Optimization Problems, Elsevier.
- A. M. Ostrowski (1966): Solutions of Equations and Systems of Equations (second edition), Academic Press.
- M. J. D. Powell (1970): Recent advances in unconstrained optimization, Report TP.430, AERE Harwell.

## APPENDIX Numerical Questions Relating to Fletcher's Algorithm

### 1. Implementation

In this section we consider two questions relating to the implementation of Fletcher's algorithm. These are

- (i) an appropriate strategy for determining  $\lambda$  to satisfy the Goldstein condition, and
- (ii) the use of the product updating formulae for the inverse Hessian estimate.

In his program Fletcher uses a cubic line search to determine  $\lambda$ . Here we use a somewhat simpler procedure which has the advantage of requiring only additional function values. Also we work with the Choleski decomposition of the inverse Hessian estimate. This has certain numerical advantages which have been outlined by Gill and Murray<sup>\*/</sup>. In particular, it is possible to ensure the positive definiteness of  $H_i$ , and this can be lost through the effect of accumulated rounding error when direct evaluation of the updating formulae is used. Another possible advantage of the Choleski decomposition is that we can work with an estimate of the Hessian (that is  $H^{-1}$ ) rather than with  $H$  as division by a triangular matrix does not differ greatly in cost to multiplication. We felt this could well be an advantage in problems with singular or near singular Hessians, in which case  $H$  would be likely to contain large numbers.

To implement the line search we note that by Theorem 1.2 we should test first if  $\psi(\tilde{x}_i, \tilde{t}_i, \|\tilde{s}_i\|) = \psi(\tilde{x}_i, \tilde{s}_i, 1)$  satisfies the Goldstein condition. This requires the evaluation of  $F(\tilde{x}_i + \tilde{s}_i)$ , and this, together with the

---

<sup>\*/</sup> NPL Mathematics Division, Report 97, 1970.



known values  $F(\underline{x}_i)$  and  $F'(\underline{x}_i) = \nabla F(\underline{x}_i) \underline{s}_i$ , gives sufficient information to determine a quadratic interpolating polynomial to  $I$ ? . We write this as

$$P(\lambda) = F(\underline{x}_i) + F'(\underline{x}_i)\lambda + A\lambda^2 \quad (\text{A.1})$$

where  $A$  is to be determined by setting  $P(1) = F(\underline{x}_i + \underline{s}_i)$ . This gives

$$\begin{aligned} A &= F(\underline{x}_i + \underline{s}_i) - F(\underline{x}_i) - F'(\underline{x}_i) \\ &= F'(\underline{x}_i)(\psi(\underline{x}_i, \underline{s}_i, 1) - 1) \quad . \end{aligned} \quad (\text{A.2})$$

The minimum of  $P(h)$  is given by

$$\lambda = - \frac{F'(\underline{x}_i)}{2A} = \frac{1}{2(1 - \psi(\underline{x}_i, \underline{s}_i, 1))} . \quad (\text{A.3})$$

To test if this is an appropriate value we compute  $\psi(\underline{x}_i, \underline{s}_i, \lambda)$ . This gives

$$\psi(\underline{x}_i, \underline{s}_i, \lambda) = \frac{1}{2} + \frac{1}{2} \left\{ 1 - \frac{\frac{1}{2} F''(\underline{x}_i + \bar{\lambda} \underline{s}_i)}{A} \right\} \quad (\text{A.4})$$

where  $\bar{\lambda}$  is a mean value. Thus, if  $F$  is quadratic and  $\psi(\underline{x}_i, \underline{s}_i, 1) < \sigma$  then  $\lambda$  given by equation (A.3) satisfies the Goldstein condition for any allowable  $\sigma$  (normally  $\sigma$  is chosen small -- say  $10^{-4}$ ): For nonquadratic  $F$  the test is satisfied if the relative error in estimating  $\frac{1}{2} F''(\underline{x}_i + \bar{\lambda} \underline{s}_i)$  by  $A$  is not too large.

This analysis provides the basis for our method which is given below.

#### Algorithm

- (i) Calculate  $\|\underline{s}_i\|$ , set  $w = \min(1, \|\underline{s}_i\|)$ ,  $\lambda = 1$ .
- (ii) Evaluate  $\psi = \psi(\underline{x}_i, \underline{s}_i, \lambda)$ .

(iii) If  $\lambda < \sigma$  then begin  $p\lambda = \lambda$  ,

$$\lambda = \frac{p\lambda}{2(1-\psi)} ,$$

go to (ii), end.

(iv) If  $\lambda \leq 1-\sigma$  got0 EXIT .

If  $\lambda \geq 1$  then begin if  $\lambda \geq 1/w$  go to EXIT,

$$\lambda = 2\lambda , \text{ end.}$$

$$\text{else } \lambda = .5(\lambda + p\lambda) .$$

go to (ii).

#### Remark

(i) Numerical experience has shown that the value of  $\lambda$  predicted in (iii) can be too small, and that an additional instruction

$$\text{If } \lambda < s*p\lambda \text{ then } \lambda = s*p\lambda$$

should be included. A value for  $s$  of about .1 has proved satisfactory (1/8 was used in the numerical experiments reported in the next section).

(ii) It is readily verified that  $\lim_{\lambda \rightarrow 0} \psi(\tilde{x}_i, \tilde{s}_i, \lambda) = 1$  . Thus the algorithm can be expected to return a value of  $\lambda$  satisfying the Goldstein condition unless  $\psi$  exhibits rather pathological behavior.

We write the Choleski decomposition of  $H$  as

$$H = R^T R \tag{A.5}$$

where  $R$  is an upper triangular matrix. Thus we require to find  $R^*$  such that

$$R^{*T} R^* = H^* \tag{A.6}$$

where  $H^*$  is given by either

$$(i) H^* = (I + \underline{\underline{u}}\underline{\underline{y}}^T)R^T R(I + \underline{\underline{y}}\underline{\underline{u}}^T), \text{ or}$$

$$(ii) H^* = (I + \underline{\underline{u}}\underline{\underline{y}}^T)\{R^T R + \zeta \tau \underline{\underline{v}}\underline{\underline{v}}^T\}(I + \underline{\underline{y}}\underline{\underline{u}}^T).$$

The second case can be reduced to the first if we write

$$\hat{R}^T \hat{R} = R^T R + \zeta \tau \underline{\underline{v}}\underline{\underline{v}}^T \quad (A.7)$$

To calculate  $\hat{R}$  note that

$$\begin{aligned} R^T R + \zeta \tau \underline{\underline{v}}\underline{\underline{v}}^T &= [R^T \mid \sqrt{\zeta \tau} \underline{\underline{v}}] \begin{bmatrix} R \\ \sqrt{\zeta \tau} \underline{\underline{v}}^T \end{bmatrix} \\ &= [R^T \mid \sqrt{\zeta \tau} \underline{\underline{v}}] Q^T Q \begin{bmatrix} R \\ \sqrt{\zeta \tau} \underline{\underline{v}}^T \end{bmatrix} \end{aligned} \quad (A.8)$$

where  $Q$  is orthogonal. Thus we seek an orthogonal matrix  $Q$  such that

$$Q \begin{bmatrix} R \\ \sqrt{\zeta \tau} \underline{\underline{v}}^T \end{bmatrix} = \begin{bmatrix} \hat{R} \\ 0 \end{bmatrix} \quad (A.9)$$

Let  $W(i, j, \{p, q\})$  be the plane rotation such that  $W(i, j, \{p, q\})A$  combines the  $i$ -th and  $j$ -th rows of  $A$ , and reduces  $A_{pq}$  to zero. It is necessary that  $p$  be either  $i$  or  $j$ . Then  $Q$  is given explicitly by

$$Q = \prod_{i=n}^1 W(i, n+1, \{n+1, i\}) \quad (A.10)$$

It is readily verified that the zero introduced by each transformation is preserved by the subsequent transformations provided they are carried out in the order indicated.

Consider now the problem of constructing the Choleski decomposition of  $S^T S$  where  $S = T + ab^T$ , and  $T$  is upper triangular. This corresponds to our problem with the identifications  $T = R$  or  $\hat{R}$ ,  $a = Ry$  or  $\hat{R}y$ , and  $b = u$ . In this case the decomposition is done in two stages. Our method uses ideas due independently to Stoer, Golub, and Gill and Murray.

(i) We determine an orthogonal matrix  $Q_1$  such that

$$Q_1 a = \|a\| e_n \quad (A.11)$$

If we set

$$Q_1 = \prod_{i=1}^{n-1} W(i, n, \{i, *\}) \quad (A.12)$$

where the  $*$  indicates that the rotation is defined by being applied to zero an element of a vector, then  $Q_1 S = Q_1 T + \|a\| e_n b^T$  differs from an upper triangular matrix only in having possible nonzero elements in the last row.

(ii) To complete the determination of  $R^*$  we sweep out the elements in the first  $(n-1)$  places in the last row of  $Q_1 S$  by plane rotations. Thus  $R^*$  is given by

$$R^* = Q_2 (Q_1 T + \|a\| e_n b^T) \quad (A.13)$$

where

$$Q_2 = \prod_{i=n-1}^1 W(i, n, \{n, i\}) \quad (A.14)$$

It will be seen that the updating of the Choleski factorization can be carried out very cheaply. Depending on the update formula used, the major cost is either  $2n$  or  $3n$  plane rotations. It should be noted that

$$y^T H y = \|Ry\|^2 = \|\tilde{a}\|^2 \quad (\text{A.15})$$

is required in the update formula. Thus  $\tilde{a}$  can be already available for S .

## 2. Numerical Results

In this section we report the results of numerical experiments carried out to test some of our techniques. We consider four line search strategies:

- (i) a standard cubic interpolation procedure with  $\lambda = 1$  as initial search interval,
- (ii) a standard cubic interpolation procedure with  $\lambda$  given by the step to the minimum in the previous line search,
- (iii) a strategy for satisfying the Goldstein condition in which  $\lambda$  is reduced by the factor  $1/8$  if  $\psi < \sigma$  , and
- (iv) the method for satisfying the Goldstein condition given in the previous section.

Product form updating for the Choleski factorization of both H and  $H^{-1} = G$  has been implemented, and the results obtained for each are given.

The problems considered include:

- (i) **Hilbert:** Minimization of a quadratic form with matrix given by the Hilbert matrix of order 5 . Here

$$F = \frac{1}{2} \sum_{i=1}^5 \sum_{j=1}^5 \frac{(x_i - 1)(x_j - 1)}{i+j-1} ,$$

and the starting point is given by

$$x_i = -4/i \quad , \quad i = 1, 2, \dots, 5 .$$

(ii) Banana(n) : The Banana function in the cases  $n = 2$  (the Rosenbrock function) and  $n = 8$ . Here

$$F = \sum_{i=1}^{n-1} \{100(x_{i+1} - x_i^2)^2 + (1 - x_i)^2\} ,$$

and the starting point is given by

$$x_1 = -1.2 \quad \text{if } i \text{ odd, otherwise } x_1 = 1 .$$

(iii) Woods: Here

$$\begin{aligned} F = & 100(x_1 - x_2)^2 + (1 - x_1)^2 + 90(x_4 - x_3)^2 \\ & + (1 - x_3)^2 + 10.1((1 - x_2)^2 + (1 - x_4)^2) \\ & + 19.8(1 - x_2)(1 - x_4) , \end{aligned}$$

and the starting point is

$$\mathbf{x}^T = \{-3, -1, -3, -1\} .$$

(iv) Singular: Powell's singular function is designed to test the performance of algorithms on a function with a singular Hessian at the solution. Here

$$F = (x_1 + 10x_2)^2 + 5(x_3 - x_4)^2 + (x_2 - 2x_3)^4 + 10(x_1 - x_4)^4 ,$$

and the starting point is

$$\mathbf{x}^T = \{3, -1, 0, 1\} .$$

(v) Helix: Here we define

$$R = \{x_1^2 + x_2^2\}^{1/2} ,$$

$$\begin{aligned} T = & \text{if } x_1 \geq 0 \text{ then } \frac{1}{2\pi} \arctan \frac{x_2}{x_1} \\ & \text{if } x_1 < 0 \text{ then } -\frac{1}{2\pi} \arctan \left( \frac{x_2}{x_1} \right) + .5 , \end{aligned}$$

and set

$$F = 100((x_3 - 10T)^2 + (R-1)^2) + x_3^2.$$

The starting vector is

$$x^T = \{-1, 0, 0\}.$$

Numerical results are given in Table A.1. For historical reasons, the test for terminating the calculations was based on the size of  $\|s_i\|$  ( $\|s_i\| \leq \text{EPS}/n$  with  $\text{EPS} = 10^{-8}$ ). This proved reasonably satisfactory for all cases except the singular function -- in fact in all other cases the ultimate convergence was clearly superlinear, and the results were accordingly only marginally affected by the size of **EPS**. In the case of singular the convergence test proved difficult to satisfy in most cases (indicated by \* in Table A.1), and these computations were terminated by the number of iterations exceeding the specified limit. However, in all cases the answers were correct to at least six decimal places. There is some variation in the H and G columns. This shows the effect of rounding error, as these would be identical in exact arithmetic. The most interesting case is the H column in both cases of the Banana (8) when satisfying the Goldstein condition. In these cases both H and G formulae produce very similar results until the 10-th iteration at which point the H formulae produce much larger reductions in F than do the G. However, this progress is not maintained and at the 20-th iteration (in the case of the line search algorithm of Section 3) the H matrix becomes singular and the iteration is terminated. A restart procedure could have been used at this stage.

The numerical results indicate that the new algorithm is promising. In general, although more iterations are required, we make significantly

fewer function evaluations in comparison with the routine using a standard line search. As only one derivative evaluation is required in each iteration, the real saving can be considerable. We note that on the basis of the evidence presented it is not possible to draw conclusions as to the relative values of the H and G algorithms. However, that both manage to produce very comparable results provides some evidence of their stability.

The program which gave the results presented here is coded in ALGOL W for the IBM 360/67 at Stanford University. The calculations were carried out using long precision (14 hexadecimal digits). A FORTRAN version of the program has been developed at the Australian National University.



	cubic line search						Goldstein algorithm					
	$\lambda^*$ from prev. tt			$\lambda^I = 1$			$(1/8)^P$			parabolic interpolation		
	H	G	H	H	G	H	H	G	H	H	G	
Hilbert	6	6	6	6	6	40	40	40	39	41	39	41
Banana (2)	23	27	24	24	24	61	61	78	46	56	44	58
Banana (8)	74	75	74	235	72	NC	NC	143	NC	NC	80	114
Woods	27	27	30	102	29	94	114	112	62	84	61	86
Singular	50	50	32	129	60	100	*	111	69	83	76	89
Helix	NC	45	24	70	23	65	NC	NC	35	47	35	48

Table A.1 Number of iterations and number of function values in numerical experiments



### III. Barrier and Penalty Function Methods



1. Basic properties of barrier functions.

Consider the inequality constrained problem (ICP)

$$\min f(x)$$

$$\text{subject to } g_i(x) \geq 0, \quad i = 1, 2, \dots, m \quad (i \in I_1),$$

where we assume (as before) that  $f, g_i, i \in I_1$ , are in  $C^2$ . We also assume that  $S = \{x; g_i(x) \geq 0, i \in I_1\}$  is compact, has a nonvoid interior  $S_0$ , and satisfies the regularity condition that every neighborhood of points of  $S$  contains points of  $S_0$  (this precludes  $S$  having 'whiskers'). If  $x \in S$  and  $g_i(x) = 0$  for some  $i$  then it is assumed that  $x \notin S_0$ .

Definition:  $\phi(g(x))$  is a barrier function for  $S$  if the following conditions are satisfied.

(i)  $\phi > 0, x \in S$ . If  $X$  closed set,  $X \subset S_0$ , then  $\phi \in C^2$  on  $X$ .

(ii)  $\phi \rightarrow \infty, g_i \rightarrow 0, i \in I_1$ .

(iii)  $\frac{\partial \phi}{\partial g_i} < 0$  if  $g_i < \rho_i$  where the  $\rho_i, i \in I_1$ , are fixed positive constants.

(iv)  $|\frac{\partial \phi}{\partial g_i}|$  bounded on  $N(x, \delta)$  if  $g_i > 0$  on  $N(x, \delta)$ .

Example: (i)  $\phi = \sum_{i=1}^m 1/g_i(x)$  (inverse barrier function),

(ii)  $\phi = \sum_{i=1}^m (\log(1+g_i(x)) - \log g_i(x))$ .

Remark: In the second example the term with argument  $1+g_i(x)$  merely ensures that the positivity condition is satisfied. It could be replaced by a bound  $k_i$  for  $\log(1+g_i(x))$  on  $S$  if this is known. In

practice it is of no consequence. The barrier function

$$\phi = \sum_{i=1}^m (k_i - \log g_i(\tilde{x})) \quad \text{is called the log barrier function.}$$

Definition:  $T(\tilde{x}, r)$  is a barrier objective function if

$$T(\tilde{x}, r) = f(\tilde{x}) + r\phi(\tilde{g}(\tilde{x})) \quad (1.1)$$

where  $r > 0$  .

Lemma 1.1:  $\exists \tilde{x} = \tilde{x}(r) \in S_0$  such that  $T(\tilde{x}(r), r) = \min_{\tilde{x} \in S} T(\tilde{x}, r)$  .

Proof:  $T(\tilde{x}, r)$  is bounded below on  $S$  , and  $T(\tilde{x}, r) \rightarrow +\infty$  as  $\tilde{x} \rightarrow \partial S$  .  
cl

Lemma 1.2: Let  $\{r_j\} \downarrow 0$  , and let  $\tilde{x}(r_j) = \tilde{x}_j$  . Then

- (i)  $\{T(\tilde{x}_j, r_j)\}$  is strictly decreasing,
- (ii)  $\{f(\tilde{x}_j)\}$  is nonincreasing, and
- (iii)  $\{\phi(\tilde{x}_j)\}$  is nondecreasing.

Proof: Let  $r_i < r_j$ , then

$$\begin{aligned} f(\tilde{x}_i) + r_i\phi(\tilde{g}(\tilde{x}_i)) &\leq f(\tilde{x}_j) + r_i\phi(\tilde{g}(\tilde{x}_j)) , \\ &< f(\tilde{x}_j) + r_j\phi(\tilde{g}(\tilde{x}_j)) , \\ &\leq f(\tilde{x}_i) + r_j\phi(\tilde{g}(\tilde{x}_i)) . \end{aligned}$$

This demonstrates (i). Subtracting the inside and outside inequalities gives

$$(r_j - r_i)\phi(\tilde{g}(\tilde{x}_i)) \geq (r_j - r_i)\phi(\tilde{g}(\tilde{x}_j))$$

which gives (ii). From the first inequality we have

$$\begin{aligned} f(\underline{x}_i) &\leq f(\underline{x}_j) + r_i(\phi(g(\underline{x}_j)) - \phi(g(\underline{x}_i))) \\ &\leq f(\underline{x}_j) \quad . \quad \square \end{aligned}$$

Remark: If  $T(\underline{x}, r)$  is strictly convex, then all inequalities are strict.

Theorem 1.1: The sequence  $\{T_1(\underline{x}_i, r_i)\}$  converges, and

$$\lim_{i \rightarrow \infty} T(\underline{x}_i, r_i) = \min_{\underline{x} \in S} f(\underline{x}) .$$

Proof: By Lemma 1.2,  $\{T(\underline{x}_i, r_i)\}$  is decreasing and bounded below and hence convergent. Let  $f^* = \min_{\underline{x} \in S} f(\underline{x})$ , then

$$T(\underline{x}_i, r_i) \geq f(\underline{x}) \geq f^*$$

whence

$$\lim_{i \rightarrow \infty} T(\underline{x}_i, r_i) \geq f^* . \quad (1.2)$$

Now let  $\varepsilon > 0$  be given. Choose  $\bar{\underline{x}} \in S_0$  such that  $f(\bar{\underline{x}}) - f^* < \varepsilon/2$  (this is possible because of the regularity condition on  $S$ ), and choose  $r_i$  such that  $r_i \phi(g(\bar{\underline{x}})) < \varepsilon/2$ . Then

$$\min_{\underline{x}} T(\underline{x}, r_i) \leq T(\bar{\underline{x}}, r_i) < f^* + \varepsilon$$

whence

$$\lim_{i \rightarrow \infty} T(\underline{x}_i, r_i) \leq f^* . \quad \square \quad (1.3)$$

Corollary 1.1: The limit points of  $\{\underline{x}_i\}$  are local minima of the ICP.

Remark: The generality of these results should be noted. For example, we have not required  $S$  to be Lagrange regular at the limits points of  $\{\underline{x}_i\}$ .

Definition:  $Q(\underline{x}, \underline{r})$  is a separable barrier objective function if

$$Q(\underline{x}, \underline{r}) = f(\underline{x}) + \sum_{i=1}^m r_i \phi_i(g_i(\underline{x})) = f(\underline{x}) + \underline{r}^T \underline{\phi}(\underline{x}) \quad (1.4)$$

where  $r_i > 0$ , and  $\phi_i$  is a barrier function for  $S_i = \{\underline{x}; g_i(\underline{x}) \geq 0\}$ ,  $i = 1, 2, \dots, m$ .

The previous results are readily extended to this case and are summarized in the following theorem.

Theorem 1.2: Let  $r_i > r_{i+1}$ ,  $i = 1, 2, \dots$ , and  $\lim_{i \rightarrow \infty} r_i = \underline{\theta}$ . Then

- (i)  $\min_{\underline{x} \in S} Q(\underline{x}, \underline{r}_k)$  is attained for some  $\underline{x}_k \in S_0$ ,
- (ii)  $\{Q(\underline{x}_k, \underline{r}_k)\}$  is strictly decreasing,  $\{f(\underline{x}_k)\}$  is nonincreasing, and
- (iii)  $\lim_{k \rightarrow \infty} Q(\underline{x}_k, \underline{r}_k) = f^*$ , and the limit points of  $\{\underline{x}_k\}$  are local minima of the ICP.

Remark: Given a sequence of positive vectors tending to zero then it is possible to select a subsequence which is strictly decreasing. Conclusions (i) and (iii) remain valid in this more general case.



## 2. Multiplier relations (first order analysis)

In this section we assume sequences  $\{r_k\} \downarrow 0$ ,  $\{\tilde{x}_k\} \rightarrow \tilde{x}^*$ . The condition that  $T(\tilde{x}, r_k)$  is stationary at  $r_k$  gives

$$\begin{aligned} \nabla T(\tilde{x}_k, r_k) &= \nabla f(\tilde{x}_k) + \sum_{i=1}^m r_k \frac{\partial \phi}{\partial g_i} \nabla g_i(\tilde{x}_k) \\ &= \nabla f(\tilde{x}_k) - \sum_{i=1}^m u_i^k \nabla g_i(\tilde{x}_k) = 0 \end{aligned} \quad (2.1)$$

where  $u_i^k = -r_k \frac{\partial \phi}{\partial g_i} (g(\tilde{x}_k))$ . Note that  $u_i^k = o(r_k) \rightarrow 0$ ,  $k \rightarrow \infty$ ,

if  $i \notin B_0$ , and  $u_i^k \geq 0$  for  $i \in B_0$  and  $k \geq k_0$  by the conditions defining barrier functions. Equation (2.1) is formally similar to the multiplier relations given earlier (MP(3.2)), and it is comparatively straightforward to deduce these relations from (2.1) in certain special cases. We assume that  $B_0 = \{1, 2, \dots, t\}$ , that the rank of the system of vectors  $\nabla g_i(\tilde{x}^*)$ ,  $i \in B_0$  is  $s \leq t$ , and  $\nabla g_1(\tilde{x}^*), \dots, \nabla g_s(\tilde{x}^*)$  are linearly independent. We define matrices  $C_1(x)$ ,  $C_2(x)$  by  $\kappa_i(C_1) = \nabla g_i(x)^T$ ,  $i = 1, \dots, s$ , and  $\kappa_i(C_2) = \nabla g_{s+i}(x)^T$ ,  $i = 1, 2, \dots, t-s$ , and vectors  $u_1^{(k)T} = \{u_1^k, \dots, u_s^k\}$ ,  $u_2^{(k)T} = \{u_{s+1}^k, \dots, u_t^k\}$ .

Lemma 2.1: If  $\{u_2^k\}$  is bounded then the Kuhn-Tucker conditions hold at  $\tilde{x}^*$ .

Proof: From (2.1) we have

$$\nabla f(\tilde{x}_k)^T = C_1(\tilde{x}_k) u_1^{(k)} + C_2(\tilde{x}_k) u_2^{(k)} + o(r_k) \quad (2.2)$$

The linear dependence of the set of vectors  $\nabla g_i(\tilde{x}^*)$ ,  $i \in B_0$ , gives

$$C_2(\tilde{x}^*) = C_1(\tilde{x}^*)R \quad (2.3)$$

so that (2.2) can be written

$$\nabla f(\tilde{x}_k)^T = C_1(\tilde{x}_k) \{u_1^{(k)} + Ru_2^{(k)}\} + \{C_2(\tilde{x}_k) - C_1(\tilde{x}_k)R\}u_2^{(k)} + o(r_k) \quad (2.4)$$

Provided  $k$  is large enough the rank of  $C_1(\tilde{x}_k)$  will be  $s$  (see MP Remark following Lemma 3.3). Thus

$$u_1^{(k)} + Ru_2^{(k)} = \{C_1(\tilde{x}_k)^T C_1(\tilde{x}_k)\}^{-1} C_1(\tilde{x}_k)^T \{\nabla f(\tilde{x}_k)^T + (C_2(\tilde{x}_k) - C_1(\tilde{x}_k)R)u_2^{(k)} + o(r_k)\} \quad (2.5)$$

As  $u_2^{(k)}$  bounded we conclude from (2.5)

(i)  $u_1^{(k)}$  bounded, and

(ii)  $\lim_{k \rightarrow \infty} u_1^{(k)} + Ru_2^{(k)} = \{C_1(\tilde{x}^*)^T C_1(\tilde{x}^*)\}^{-1} C_1(\tilde{x}^*)^T \nabla f(\tilde{x}^*)$ .

As  $\{u_1^{(k)}\}$ ,  $\{u_2^{(k)}\}$  are bounded and nonnegative (at least for  $k$  large enough) this property is shared by the limit points of the sequences.

Consider subsequences tending to  $u_1^*$ ,  $u_2^*$  respectively. From (2.2)

we have

$$\nabla f(\tilde{x}^*)^T = C_1(\tilde{x}^*)u_1^* + C_2(\tilde{x}^*)u_2^*$$

or

$$\nabla f(\tilde{x}^*) = \sum_{i \in B_0} u_i^* \nabla g_i(\tilde{x}^*) \quad (2.6)$$

Thus the Kuhn-Tucker conditions are satisfied.  $\square$

Corollary 2.1: If the Kuhn-Tucker conditions do not hold at  $\tilde{x}^*$ ,

then  $\{u_1^k\}$ ,  $\{u_2^k\}$  are unbounded.

Corollary 2.2: If restraint condition A holds then the  $\nabla g_i(x^*)$  are linearly independent for  $i \in B$ . In this case  $\{u_2^k\}$  is null, and  $\{u_1^k\}$  converges. If restraint condition A holds then the multipliers in the Kuhn-Tucker conditions are uniquely determined.

Lemma 2.2: If restraint condition B holds then  $\{u_2^{(k)}\}$  is bounded.

Remark: By MP Lemma 5.2, this implies that  $\{u_2^{(k)}\}$  is bounded for the convex programming problem provided S has an interior.

Proof: If restraint condition B holds then  $\exists \delta$  such that  $\nabla g_i(x^*) \delta > 0$ ,  $i = 1, 2, \dots, t$ . From (2.2) we have

$$\sum_{i=1}^t u_i^k \nabla g_i(x_k) \delta = \nabla f(x_k) \delta + O(r_k). \quad (2.7)$$

As  $\nabla g_i(x) \delta$  is a continuous function, we must have  $u_i^k > 0$  and  $\nabla g_i(x_k) \delta > 0$ ,  $i = 1, 2, \dots, t$ , provided  $k$  large enough. Thus (2.7) gives

$$\sum_{i=1}^t u_i^k < \frac{\nabla f(x_k) \delta + O(r_k)}{\min_i \nabla g_i(x_k) \delta}. \quad (2.8)$$

This relation shows that the  $u_i^k$  are bounded as  $k \rightarrow \infty$ .  $\square$

Remark: The results of the first section showed that convergence of barrier function algorithms can be proved under very few assumptions. The results of this section show that valuable structural information on the problem is available as a by-product of the computation. Note that the condition that the  $u_i^k$  be bounded is a weaker restraint condition than either A or B.

3. Second order conditions.

Consider now a barrier function  $\phi$  and a sequence  $\{r_k\} \downarrow 0$  such that  $\{\underline{x}_k\} \rightarrow \underline{x}^*$ ,  $\{\underline{u}_k\} \rightarrow \underline{u}^*$ . It is convenient to assume the following properties which are satisfied by all barrier functions of practical interest.

$$(i) \quad u_i^k = -r_k \frac{\partial \phi}{\partial g_i}(\underline{x}_k) > 0, \quad r_k \frac{\partial^2 \phi}{\partial g_i^2}(\underline{x}_k) > 0, \quad i = 1, 2, \dots, m, \forall k.$$

$$(ii) \quad r_k \frac{\partial^2 \phi}{\partial g_i^2}(\underline{x}_k) \rightarrow +\infty, \quad k \rightarrow \infty, \quad i \in B. \quad (\text{But see Example 4(ii) p. 100 for qualification.})$$

$$(iii) \quad a_{ij} = 0 \quad \text{if } i \neq j.$$

Lemma 3.1: If the matrices  $\nabla^2 T(\underline{x}_k, r_k)$  are positive definite for  $k \geq k_0$  and the  $\nabla g_i(\underline{x}^*)$ ,  $i \in B_0$ , are linearly independent then  $\underline{v}^T \nabla^2 T(\underline{x}^*, \underline{u}^*) \underline{v} > 0$ ,  $\forall \underline{v} \neq 0$  such that  $\nabla g_i(\underline{x}^*) \underline{v} = 0$ ,  $\forall i \in B_0$ .

Proof: Differentiating  $T(\underline{x}, r_k)$  gives

$$\nabla^2 T(\underline{x}_k, r_k) = \nabla^2 f(\underline{x}_k, \underline{u}_k) + r_k \sum_{i=1}^m \frac{\partial^2 \phi}{\partial g_i^2} \nabla g_i(\underline{x}_k)^T \nabla g_i(\underline{x}_k) \quad (3.1)$$

which can be written

$$\nabla^2 T(\underline{x}_k, r_k) = \nabla^2 f(\underline{x}_k, \underline{u}_k) + C_1(\underline{x}_k) D_k C_1(\underline{x}_k)^T + \sum_{i \in I_1 - B_0} r_k \frac{\partial^2 \phi}{\partial g_i^2} \nabla g_i(\underline{x}_k)^T \nabla g_i(\underline{x}_k)$$

Where

$$(D_k)_{ii} = r_k \frac{\partial^2 \phi}{\partial g_i^2}(\underline{x}_k), \quad i = 1, 2, \dots, t.$$

Let

$$P_k = C_1(x_k)C_1(x_k)^+ \quad (3.2)$$

then, for arbitrary nonzero  $v$  such that  $(I - P_k)v \neq 0$

$$\begin{aligned} 0 &< v^T(I - P_k)\nabla^2 T(x_k, r_k)(I - P_k)v \\ &= v^T(I - P_k)\nabla^2 f(x_k, u_k)(I - P_k)v + o(1) \end{aligned}$$

as

$$(I - P_k)C_1(x_k) = 0.$$

The desired result follows from this on letting  $k \rightarrow \infty$ .  $\square$

Corollary 3.1: If in addition to the conditions of Lemma 3.1 we have also strict complementarity then the second order sufficiency conditions (the conditions of MP Lemma 4.3) hold at  $x^*$ .

Remark: The problem of generalizing this result to the case where the active constraint gradients are not linearly independent is the following.

In general, when  $k < \infty$ ,  $\text{rank}[C_1(x_k)|C_2(x_k)] > s$ . Thus

$$V_k = \{v; \nabla g_i(x_k)v = 0, \forall i \in B_0\} \subset U_k = \{v; v = (I - P_k)u, u \in E_n\}.$$

We have

$$\lim_{k \rightarrow \infty} U_k = V^* = \{v; \nabla g_i(x^*)v = 0, \forall i \in B_0\}.$$

It is not difficult to construct examples in which  $\lim_{k \rightarrow \infty} V_k \subset V^*$ .

Consider  $C_1(x_k) = e_1$ ,  $C_2(x_k) = e_1 + e_2^T(x - x^*)e_2$ . Then

$V_k = \{v; e_1^T v = e_2^T v = 0\} = \lim_{k \rightarrow \infty} V_k \subset V^* = \{v; e_1^T v = 0\}$ . The

argument of Lemma 3.1 shows only that  $\frac{1}{V} \frac{2}{V} f(x^*, u^*)_V > 0$  for

$v \in \lim_{k \rightarrow \infty} V_k$ .

Lemma 3.2: Let

$$w = U + \gamma V V^T,$$

$$N = \{t; \|t\| = 1, V^T t = 0\},$$

$$M = \{u; \|u\| = 1, u \in N^\perp\},$$

$$v = \min_{t \in N} t^T U t > 0, \quad \sigma = \min_{u \in M} u^T U u, \quad \mu = \min_{u \in M} \|V u\| > 0,$$

$$\eta = \min_{t \in N, u \in M} t^T U u, \quad \rho = \min(0, \eta)$$

then  $W$  is positive definite provided

$$\gamma > \frac{\rho^2 + v^2 - \sigma v}{v \mu} \quad (3.3)$$

Proof: Any unit vector  $w$  can be written

$$w = \alpha u + \beta t \quad \text{where } u \in M, t \in N, \text{ and } \alpha^2 + \beta^2 = 1.$$

Thus

$$\begin{aligned} w^T W w &= \alpha^2 u^T U u + 2\alpha\beta u^T U t + \beta^2 t^T U t + \gamma \alpha^2 \|V u\|^2 \\ &\geq \alpha^2 (\sigma + \gamma \mu^2 - v) + 2|\alpha| (1 - \alpha^2)^{1/2} \rho + v \\ &\geq \alpha^2 (\sigma + \gamma \mu^2 - v) + 2|\alpha| \rho + v \end{aligned} \quad (3.4)$$

as  $\rho \leq 0$ . Provided  $\gamma > \frac{\nu - \sigma}{\mu^2}$  (a weaker condition than (3.3)), the

right hand side has a minimum at  $\alpha = \frac{-\rho}{\sigma + \gamma\mu^2 - \nu} > 0$ . The value at

this minimum is  $\frac{-\rho^2}{\sigma + \gamma\mu^2 - \nu} + \nu$  which is positive if (3.3) holds.  $\square$

Corollary 3.2: If  $W = U + VD V^T$  where D diagonal, and if the conditions of Lemma 3.2 hold, then W is positive definite provided

$\min_i D_{ii} > \frac{\rho^2 + \nu^2 - \sigma\nu}{\nu\mu^2}$ . If D is positive definite the result holds

provided the smallest eigenvalue of D satisfies this inequality.

Proof: We have

$$\begin{aligned} \tilde{w}^T VD V^T \tilde{w} &= \alpha^2 \tilde{u}^T VD V^T \tilde{u} = \alpha^2 \left\{ \sum_{i=1}^t D_{ii} (\rho_i (V^T \tilde{u}))^2 \right\} \\ &\geq \min_i D_{ii} \alpha^2 \|V^T \tilde{u}\|^2. \end{aligned}$$

The result now follows as from equation (3.4) above.  $\square$

Lemma 3.3: If the second order sufficiency conditions hold at  $\tilde{x}^*$  then  $\nabla^2 T(\tilde{x}_k, \tilde{r}_k)$  is positive definite for k large enough.

Proof: We have

$$V^* \subseteq U^* = \{ \tilde{v} ; \nabla g_{i_0}(\tilde{x}^*) \tilde{v} = 0, \forall i \in B_0 \text{ such that } u_i^* > 0 \}.$$

Thus the second order sufficiency conditions imply that

$\tilde{v}^T \nabla^2 \mathcal{L}(\tilde{x}^*, \tilde{u}^*) \tilde{v} > 0$ ,  $\forall \tilde{v} \in V^*$  such that  $\tilde{v} \neq 0$ . From Corollary 3.2

it follows that  $\exists D_0$  such that  $\nabla^2 \mathcal{L}(\tilde{x}^*, \tilde{u}^*) + [C_1(\tilde{x}^*) | C_2(\tilde{x}^*)] D [C_1(\tilde{x}^*) | C_2(\tilde{x}^*)]^T$

is positive definite for  $D \geq D_0$ . By continuity this implies that the corresponding matrix evaluated at  $x_k$  is positive definite for  $k$  large enough. The desired result follows from this as  $(D_k)_{ii} \rightarrow \infty$ ,  $i = 1, 2, \dots, t$  as  $k \rightarrow \infty$ .  $\square$

Lemma 3.4: If  $U$  is nonsingular,  $D$  diagonal, and  $V$  of full rank, then the system of linear equations

$$[U + VD V^T] \tilde{x} = V \tilde{y} \quad (3.5)$$

has the solution

$$x = U^{-1} V (I + M)^{-1} M \tilde{y} \quad (3.6)$$

where  $M = (V^T U^{-1} V)^{-1} D^{-1}$  provided  $I + M$  is nonsingular. A sufficient condition for  $I + M$  nonsingular is  $\|M\| < 1$  which is satisfied if  $\min_i |D_{ii}|$  is large enough.

Proof: The result follows on substituting (3.6) into (3.5).  $\square$

Remark: From (3.6) it follows that

$$x \sim U^{-1} V (V^T U^{-1} V)^{-1} D^{-1} \tilde{y} \quad (3.7)$$

as  $\min_i |D_{ii}| \rightarrow \infty$ .

Corollary 3.3: If the right hand side of equation (3.5) is  $\tilde{z}$ , a general vector, then the solution is given by

$$\begin{aligned} \tilde{x} = & U^{-1} V (I + M)^{-1} M (V^T U^{-1} V)^{-1} V^T U^{-1} \tilde{z} \\ & + U^{-1} (I - V (V^T U^{-1} V)^{-1} V^T U^{-1}) \tilde{z} . \end{aligned} \quad (3.8)$$



4. Rate of convergence results.

In this section, rate of convergence estimates for barrier function algorithms are considered. Unless stated otherwise the conditions imposed in Section 3 are assumed, together with the condition that  $\|\nabla g_i(x^*)\| \neq 0$ ,  $i \in B_0$ .

Lemma 4.1: Provided  $\{u_k\}$  is bounded then

$$f(x_k) - f(x^*) = \sum_{i=1}^t u_i^k g_i(x_k) + O(\max\{\|x_k - x^*\|^2, r_k \|x_k - x^*\|\}) \quad (4.1)$$

Proof: The result follows by taking the scalar product of (2.2) with  $x_k - x^*$  and identifying with terms in the Taylor series expansion.  $\square$

Definition: We say that  $u_k$  is  $SO(v_k)$  (strict order  $v_k$ ) provided

- (i)  $u_k = O(v_k)$ , and
- (ii)  $\exists k_0 < \infty$  and  $\mu > 0$  such that  $|u_k| \geq \mu |v_k|$  for  $k \geq k_0$ .

Remark: (4.1) gives an error estimate provided the remainder term is small. A sufficient condition for this is

$$f(x_k) - f(x^*) = SO(\|x_k - x^*\|) \quad (4.2)$$

This implies that for at least one  $i$ ,  $u_i^k g_i(x_k) = SO(\|x_k - x^*\|)$ .

If (4.2) does not hold then for  $i = 1, 2, \dots, t$  either

- (i)  $u_i^k \rightarrow 0$ ,  $k \rightarrow \infty$ , or
- (ii)  $g_i(x_k) = o(\|x_k - x^*\|)$ .

If (ii) holds then the approach of  $x_k$  to  $x^*$  is tangential to the surface  $g_i(x) = 0$  at  $x = x^*$ .

Lemma 4.2: If the ICP is convex, then

(i)  $x_k, u_k$  are dual feasible, and

$$(ii) f(x_k) - f(x^*) < \sum_{i=1}^m u_i^k g_i(x_k) .$$

Proof: The dual feasibility is a consequence of (2.1) and assumption (i) of Section 3. This follows directly from Wolfe's form of the duality theorem (MP Corollary 5.1). We have

$$f(x_k, u_k) = f(x_k) - \sum_{i=1}^m u_i^k g_i(x_k) \leq f(x^*, u^*) = f(x^*) \leq f(x_k) \quad (4.3)$$

which demonstrates the second part of the desired result.  $\square$

Example: For  $i \in B_0$  let  $g_i(x_k) = SO(\|x_k - x^*\|)$ ,  $u_i^* > 0$ .

(a) inverse barrier function. We have

$$u_i^k = r_k / g_i(x_k)^2$$

whence

$$g_i(x_k) = \sqrt{r_k / u_i^k} .$$

This gives .

$$\|x_k - x^*\| = O(r_k^{1/2}) .$$

(b) log barrier functions. In this case

$$u_i^k = r_k / g_i(x_k)$$

which gives

$$g_i(\underline{x}_k) = r_k / u_i^k$$

and

$$\|\underline{x}_k - \underline{x}^*\| = O(r_k) .$$

Thus the strict order condition permits us to deduce a rate of convergence result. We now show that the SO condition is equivalent to the condition of strict complementarity for the inverse and log barrier functions. In these cases the remark following Lemma 4.1 gives us a geometric interpretation of strict complementarity.

To discuss this equivalence, consider the following system of equations which define  $\underline{x}_k$  and  $\underline{u}_k$  as functions of  $r_k$ .

$$\nabla f(\underline{x}_k) - \sum_{i=1}^m u_i^k \nabla g_i(\underline{x}_k) = 0 ,$$

and

$$u_i^k \frac{\partial \phi}{\partial g_i}(\underline{x}_k) = - r_k , \quad i = 1, 2, \dots, m . \quad (4.4)$$

If the Jacobian  $H(\underline{x}, \underline{u})$  of this system with respect to  $\underline{x}, \underline{u}$  or an appropriate transform of it is nonsingular then we can study the behavior of  $\underline{x}(r), \underline{u}(r)$  as a function of  $r$  by integrating the system of differential equations

$$H(\underline{x}, \underline{u}) \begin{bmatrix} \frac{d\underline{x}}{dr} \\ \frac{d\underline{u}}{dr} \end{bmatrix} = - \begin{bmatrix} 0 \\ \underline{e} \end{bmatrix} \quad (4.5)$$

where  $\underline{e}^T = \{1, 1, \dots, 1\}$ . We have

$$H(\tilde{x}, \tilde{u}) = \begin{bmatrix} -\frac{2}{V} \mathcal{L}(\tilde{x}, \tilde{u}) & -\nabla_{g_1}(\tilde{x})^T \dots & \nabla_{g_m}(\tilde{x})^T \\ \frac{-u_1}{\left(\frac{\partial \phi}{\partial g_1}\right)^2} \frac{\partial^2 \phi}{\partial g_1^2} \nabla_{g_1}(\tilde{x}) & 1/\frac{\partial \phi}{\partial g_1} & \\ \vdots & \ddots & \\ \frac{-u_m}{\left(\frac{\partial \phi}{\partial g_m}\right)^2} \frac{\partial^2 \phi}{\partial g_m^2} \nabla_{g_m}(\tilde{x}) & & 1/\frac{\partial \phi}{\partial g_m} \end{bmatrix} \quad (4.6)$$

Let D be the diagonal matrix

$$D = \begin{bmatrix} I & & & \\ & w_1 & & \\ & & \ddots & \\ & & & w_m \end{bmatrix} \quad (4.7)$$

where

$$w_i = -\left(\frac{\partial \phi}{\partial g_i}\right)^2 / \frac{\partial^2 \phi}{\partial g_i^2} \quad (4.8)$$

then

$$DH = \begin{bmatrix} \nabla^2 \mathcal{L}(\tilde{x}, \tilde{u}) & -\nabla_{g_1}(\tilde{x})^T & \dots & -\nabla_{g_m}(\tilde{x})^T \\ u_1 \nabla_{g_1}(\tilde{x}) & -\frac{\partial \phi}{\partial g_1} / \frac{\partial^2 \phi}{\partial g_1^2} & & \\ \vdots & & \ddots & \\ u_m \nabla_{g_m}(\tilde{x}) & & & -\frac{\partial \phi}{\partial g_m} / \frac{\partial^2 \phi}{\partial g_m^2} \end{bmatrix} \quad (4.9)$$

Provided  $\frac{\partial \phi}{\partial g_i} / \frac{\partial^2 \phi}{\partial g_i^2} \rightarrow 0$ ,  $k \rightarrow \infty$ ,  $i \in B_0$  then DH will be non-singular for  $k$  large enough provided  $J(x^*)$  is nonsingular (see MP Lemma 4.4). This implies

- (i) the second order sufficiency conditions hold,
- (ii) the active constraint gradients are linearly independent, and
- (iii) strict complementarity holds.

In this case

$$DH \begin{bmatrix} \frac{dx}{dr} \\ \frac{du}{dr} \end{bmatrix} = - \begin{bmatrix} 0 \\ w \end{bmatrix} \quad (4.10)$$

whence

$$\begin{bmatrix} x_k - x^* \\ u_k - u^* \end{bmatrix} = - \int_0^{r_k} (DH)^{-1} \begin{bmatrix} 0 \\ w \end{bmatrix} dr \quad (4.11)$$

provided the components of  $w$  are integrable.

Example: (i) inverse barrier function.

We have

$$\frac{\partial \phi}{\partial g_i} / \frac{\partial^2 \phi}{\partial g_i^2} = - \frac{g_i}{2}, \quad \left( \frac{\partial \phi}{\partial g_i} \right)^2 / \frac{\partial^2 \phi}{\partial g_i^2} = \frac{1}{2g_i}.$$

In particular, it follows that  $DH(x_k)$  is nonsingular for  $k$  large enough, while  $-w_1 = \frac{1}{2} \sqrt{u_i / r}$ . We have

$$\begin{bmatrix} \frac{dx}{d\tilde{r}} \\ \frac{du}{d\tilde{r}} \end{bmatrix} = \frac{1}{2r^{1/2}} (\text{DH})^{-1} \begin{bmatrix} 0 \\ \sqrt{u_1} \\ \vdots \\ \sqrt{u_m} \end{bmatrix} \quad (4.12)$$

whence, changing to  $r^{1/2}$  as independent variable,

$$\begin{bmatrix} \frac{dx}{d(r^{1/2})} \\ \frac{du}{d(r^{1/2})} \end{bmatrix} = (\text{DH})^{-1} \begin{bmatrix} 0 \\ \sqrt{u_1} \\ \vdots \\ \sqrt{u_m} \end{bmatrix} . \quad (4.13)$$

Thus we have

$$\begin{bmatrix} x(r) - x^* \\ u(r) - u^* \end{bmatrix} = r^{1/2} [\text{DH}(x^*)]^{-1} \begin{bmatrix} 0 \\ \sqrt{u_1^*} \\ \vdots \\ \sqrt{u_m^*} \end{bmatrix} + o(r^{1/2}) . \quad (4.14)$$

(ii) log barrier function. In this case

$$\frac{\partial \phi}{\partial g_i} / \frac{\partial^2 \phi}{\partial g_i^2} = -g_i , \quad \left( \frac{\partial \phi}{\partial g_i} \right)^2 / \frac{\partial^2 \phi}{\partial g_i^2} = 1$$

so that (4.10) becomes

$$\begin{bmatrix} \frac{dx}{d\tilde{r}} \\ \frac{du}{d\tilde{r}} \end{bmatrix} = J^{-1} \begin{bmatrix} 0 \\ \varepsilon \end{bmatrix} . \quad (4.15)$$

In particular,  $\tilde{x}(r)$ ,  $\tilde{u}(r)$  inherit the differentiability properties of  $f$  and  $g_i$ ,  $i = 1, \dots, m$  for  $r$  small enough (if  $f, g \in C^p$  then  $\tilde{x}, \tilde{u} \in C^{p-1}$ ). Thus

$$\begin{bmatrix} \tilde{x}(r) - \tilde{x}^* \\ \tilde{u}(r) - \tilde{u}^* \end{bmatrix} = r J(\tilde{x}^*)^{-1} \begin{bmatrix} 0 \\ \tilde{\varepsilon} \end{bmatrix} + o(r^2) . \quad (4.16)$$

Remark: These results can also be derived by differentiating (2.1) with respect to  $r$ . We obtain

$$\nabla^2 T(\tilde{x}(r), r) \frac{d\tilde{x}}{dr} = - \sum_{i=1}^m \frac{\partial \phi}{\partial g_i} \nabla g_i(\tilde{x}(r))^T . \quad (4.17)$$

In this case Lemma 3.3 guarantees that  $\nabla^2 T(\tilde{x}(r), r)$  is positive definite for  $r$  small enough, and Corollary 3.3 can be used to give the solution to (4.17). Note that (4.17) results if  $\frac{d\tilde{u}}{dr}$  is eliminated from (4.10).

We can now proceed to the main result.

Theorem 4.1: Provided  $J(\tilde{x}^*)$  is nonsingular, then the strict complementarity and strict order conditions are equivalent for the inverse and log barrier functions.

Proof: The argument is essentially the same for both barrier functions, but is simplest for the log function. Thus only this case is considered here.

For the log penalty function  $u_i^k = r_k / g_i(\tilde{x}_k)$  so that

$$\begin{aligned}
f(\tilde{x}_k) - f(\tilde{x}^*) &= \sum_{i \in B_0} \frac{r_k}{g_i(\tilde{x}_k)} g_i(\tilde{x}_k) + o(\|\tilde{x}_k - \tilde{x}^*\|) \\
&= \tau r_k + o(\|\tilde{x}_k - \tilde{x}^*\|) .
\end{aligned} \tag{4.18}$$

If strict complementarity does not hold then for at least one  $i \in B_0$

$$\lim_{k \rightarrow \infty} \frac{r_k}{g_i(\tilde{x}_k)} = 0 .$$

As  $g_i(\tilde{x}_k) = \nabla g_i(\tilde{x}^*) (\tilde{x}_k - \tilde{x}^*) + o(\|\tilde{x}_k - \tilde{x}^*\|)$  is  $o(\|\tilde{x}_k - \tilde{x}^*\|)$  at most, it follows that  $r_k = o(\|\tilde{x}_k - \tilde{x}^*\|)$  and hence, from (4.18), that  $f(\tilde{x}_k) - f(\tilde{x}^*) = o(\|\tilde{x}_k - \tilde{x}^*\|)$ . Thus the SO condition does not hold.

If strict complementarity does hold then asymptotically for small  $r$

$$\frac{r}{\|\nabla g(\tilde{x}^*)\| \|\tilde{x}(r) - \tilde{x}^*\|} \sim \leq \frac{r}{g_i(\tilde{x}(r))} \rightarrow u_i^* > 0 , \quad i \in B_0 .$$

Thus  $r = o(\|\tilde{x}(r) - \tilde{x}^*\|)$ . As  $J(\tilde{x}^*)$  is nonsingular (4.16) holds and this implies (as  $r = o(\|\tilde{x}(r) - \tilde{x}^*\|)$ ) that

$$\|\tilde{x}(r) - \tilde{x}^*\| \leq Kr + o(r^2)$$

for some  $K > 0$ . This shows that  $r = SO(\|\tilde{x}(r) - \tilde{x}^*\|)$  so that, by (4.18), the strict order condition is satisfied.  $\square$

Remark: The above argument shows that if strict complementarity does not hold then the strict order condition cannot. The condition that  $J(\tilde{x}^*)$  be nonsingular is required only for the second part of the theorem.



Example: (i) minimize  $x_1 + x_2$   
 subject to  $-x_1^2 + x_2 \geq 0$ ,  $x_1 \geq 0$

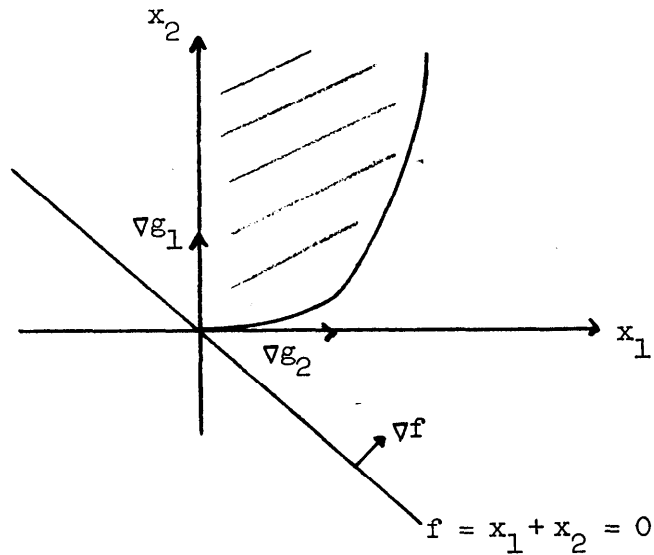


Figure 4.1

From Figure 4.1 it is clear that the minimum is  $f = 0$  at  $x_1 = x_2 = 0$ , and that strict complementarity holds.

(a) inverse barrier function

$$T = x_1 + x_2 + r \left\{ \frac{1}{x_2 - x_1^2} + \frac{1}{x_1} \right\}$$

$$\frac{\partial T}{\partial x_1} = 0 = 1 + r \left\{ \frac{2x_1}{(x_2 - x_1^2)^2} - \frac{1}{x_1^2} \right\}$$

$$\frac{\partial T}{\partial x_2} = 0 = 1 + r \left\{ \frac{-1}{(x_2 - x_1^2)^2} \right\}$$

This gives a pair of equations for  $x_1$  and  $x_2$  as function of  $r$ . We have

$$x_1 = r^{1/2} - r + o(r^{3/2}),$$

$$x_2 = r^{1/2} + r + o(r^{3/2}).$$

(b) log barrier function

$$T = x_1 + x_2 - r \{ \log(x_2 - x_1^2) + \log x_1 \}$$

$$\frac{\partial T}{\partial x_1} = 0 = 1 - r \left\{ \frac{-2x_1}{x_2 - x_1^2} + \frac{1}{x_1} \right\}$$

$$\frac{\partial T}{\partial x_2} = 0 = 1 - r \left\{ \frac{1}{x_2 - x_1^2} \right\}.$$

Solving for  $x_1$  and  $x_2$  as functions of  $r$  gives

$$x_1 = r - 2r^2 + o(r^3),$$

$$x_2 = r + r^2 + o(r^3).$$

(ii) minimize  $x_2$

subject to  $x_2 - x_1^2 > 0$ ,  $x_1 \geq 0$ .

In this case the minimum is again  $f = 0$  at  $x_1 = x_2 = 0$ . However,  $\nabla f(0) = e_2^T$  is orthogonal to  $\nabla g_2(0) = e_1^T$ . Thus, as both constraints are active at zero, strict complementarity does not hold. Note that the constraint  $g_2 = x_1 \geq 0$  is redundant, and that the barrier function trajectory is tangential to the constraint surface  $g_1 = 0$ . Note also that the rate of convergence is reduced, and that  $r_k \frac{\partial \phi}{\partial g^2}$  does not tend to  $\infty$  for the constraint with the zero multiplier.

(a) inverse barrier function.

$$T = x_2 + r \left\{ \frac{1}{x_2 - x_1^2} + \frac{1}{x_1} \right\}$$

$$\frac{\partial T}{\partial x_1} = 0 = r \left\{ \frac{2x_1}{(x_2 - x_1^2)^2} - \frac{1}{x_1^2} \right\}$$

$$\frac{\partial T}{\partial x_2} = 0 = 1 + r \left\{ \frac{-1}{(x_2 - x_1^2)^2} \right\}$$

whence

$$x_1 = (r/2)^{1/3}, \quad u_1 = 1, \quad r \frac{\partial^2 \phi}{\partial g_1^2} = \frac{2}{r^{1/2}},$$

$$x_2 = 3r/2, \quad u_2 = (4r)^{1/3}, \quad r \frac{\partial^2 \phi}{\partial g_2^2} = 4.$$

(b) log barrier function.

$$T = x_2 - r \{ \log(x_2 - x_1^2) + \log x_1 \}$$

$$\frac{\partial T}{\partial x_1} = 0 = -r \left\{ \frac{-2x_1}{x_2 - x_1^2} + \frac{1}{x_1} \right\}$$

$$\frac{\partial T}{\partial x_2} = 0 = 1 - r \left\{ \frac{1}{x_2 - x_1^2} \right\}$$

whence

$$x_1 = (r/2)^{1/2}, \quad u_1 = 1, \quad r \frac{\partial^2 \phi}{\partial g_1^2} = \frac{1}{r}$$

$$x_2 = 3r/2, \quad u_2 = (2r)^{1/2}, \quad r \frac{\partial^2 \phi}{\partial g_2^2} = 2.$$

(iii) minimize  $-x_1$   
 subject to  $(1-x_1)^3 - x_2 \geq 0$  ,  $x_2 \geq 0$  .

This is the example used in MP Section 3 (see Figure 3.1). The optimum is  $f = -1$  at  $x_1 = 1$  ,  $x_2 = 0$  . The Kuhn Tucker conditions do not hold at this point.

(a) inverse barrier function.

$$T = -x_1 + r \left\{ \frac{1}{(1-x_1)^3 - x_2} + \frac{1}{x_2} \right\}$$

$$\frac{\partial T}{\partial x_1} = 0 = -1 + r \left\{ \frac{3(1-x_1)^2}{((1-x_1)^3 - x_2)^2} \right\}$$

$$\frac{\partial T}{\partial x_2} = 0 = r \left\{ \frac{1}{((1-x_1)^3 - x_2)^2} - \frac{1}{x_2^2} \right\}$$

whence

$$x_1 = 1 - 2^{1/2} 3^{1/4} r^{1/4} ,$$

$$x_2 = 2^{1/2} 3^{3/4} r^{3/4} .$$

In this case  $u_1(r) = u_2(r) = \frac{1}{2 \cdot 3^{3/2} \cdot r^{1/2}}$

(b) log barrier function.

$$T = -x_1 - r \{ \log((1-x_1)^3 - x_2) + \log x_2 \}$$

$$\frac{\partial T}{\partial x_1} = 0 = -1 - r \left\{ \frac{-3(1-x_1)^2}{(1-x_1)^3 - x_2} \right\}$$

$$\frac{\partial \Gamma}{\partial x_2} = 0 = -r \left\{ \frac{-1}{(1-x_1)^3 - x_2} + \frac{1}{x_2} \right\}$$

whence

$$x_1 = 1 - 6r \quad ,$$

$$x_2 = 108r^3 \quad .$$

In this case  $u_1(r) = u_2(r) = \frac{1}{108r^2}$  .

The above examples confirm the predictions of our analysis, and for a given fixed sequence of  $r_k$  values effective convergence is attained more rapidly (i.e., for earlier members of the sequence) with the log barrier function.

Now let  $\phi$  be a barrier function. Then

$$\phi_1 = \log(\sigma + \phi) \quad , \quad \sigma > 1 \tag{4.19}$$

is a barrier function. Let  $\underline{x}_k$  minimize  $f + r_k \phi$  ,  $\hat{\underline{x}}_k$  minimize  $f + r_k \phi_1$  . Then comparing corresponding Lagrange multiplier estimates gives

$$\frac{\frac{\partial \phi}{\partial g_i}(\hat{\underline{x}}_k)}{(\sigma + \phi(\hat{\underline{x}}_k)) \frac{\partial \phi}{\partial g_i}(\underline{x}_k)} \rightarrow 1 \quad , \quad r_k \rightarrow 0$$

whence

$$\frac{\partial \phi}{\partial g_i}(\underline{x}_k) = o\left(\frac{\partial \phi}{\partial g_i}(\hat{\underline{x}}_k)\right) \quad \text{as } r_k \rightarrow 0 \quad .$$

Essentially this says that  $g_i(\hat{\underline{x}}_k) \rightarrow 0$  more rapidly than  $g_i(\underline{x}_k)$  , so that a faster rate of convergence is anticipated for the  $\phi_1$  barrier function.

Consider now the sequence of barrier functions defined recursively by

$$\begin{aligned}\phi_j^{(1)} &= \log(k_j - \log(g_j(x))) \quad , \\ \phi_j^{(i)} &= \log(\sigma + \phi_j^{(i-1)}) \quad , \quad i = 2, 3, \dots, \sigma > 1 \quad , \\ \phi_j^{(i)} &= \sum_{j=1}^m \phi_j^{(i)} \quad .\end{aligned}\tag{4.20}$$

In this case the error estimate is

$$f(x_k^{(i)}) - f(x^*) = -r_k \sum_{j=1}^t \frac{\partial \phi_j^{(i)}}{\partial g_j} g_j(x_k) = \sum_{j=1}^t \frac{r_k}{k_j - \log(g_j)} \prod_{s=2}^{i-1} \frac{1}{\sigma + \phi_j^{(s)}}\tag{4.21}$$

The right hand side of (4.21) tends to zero as  $i \rightarrow \infty$  , and this suggests that increasingly rapid rates of convergence can be obtained by using barrier functions associated with large values of  $i$  .

However, an even more interesting result is possible. This shows that in certain circumstances it is possible to choose a barrier function having the property that the solution to the ICP is approximated arbitrarily closely by the result of a single unconstrained minimization, without requiring  $r$  to be taken arbitrarily small. Let

$$T^{(i)}(\underline{x}, r) = f(\underline{x}) + r \sum_{j=1}^{\bar{m}} \phi_j^{(i)}(g_j(\underline{x})) \quad , \quad \text{and}$$

$$Q(\underline{x}, \lambda) = f(\underline{x}) + \sum_{j=1}^m \lambda_j (k_j - \log(g_j(\underline{x}))) \quad .$$

Theorem 4.2: Let  $Q(\underline{x}, \underline{\lambda})$  have a unique stationary value (necessarily a minimum) in  $S_0$  for each  $\underline{\lambda} > 0$ , and let  $\underline{x}^{(i)}$  minimize  $T^{(i)}(\underline{x}, \underline{r})$  for  $i = 1, 2, \dots$  and fixed  $\underline{r}$ . Then the limit points of  $\{\underline{x}^{(i)}\}$  are local minima of the ICP.

Remark: Note that  $\underline{r}$  does not have to be small in this result.

Proof: If  $\underline{x}^{(i)}$  minimizes  $T^{(i)}(\underline{x}, \underline{r})$  then

$$\begin{aligned} \nabla f(\underline{x}^{(i)}) - r \sum_{j=1}^m \frac{1}{k_j - \log(g_j(\underline{x}^{(i)}))} \prod_{s=2}^{i-1} \frac{1}{\sigma + \phi_j^{(s)}(g_j(\underline{x}^{(i)}))} \frac{1}{g_j(\underline{x}^{(i)})} \nabla g_j(\underline{x}^{(i)}) \\ = 0 \quad , \end{aligned} \quad (4.22)$$

and this expression has the form

$$\nabla Q(\underline{x}^{(i)}, \underline{\lambda}^{(i)}) = 0 \quad (4.23)$$

where  $\lambda_j^{(i)}$  has the numerical value

$$\lambda_j^{(i)} = \frac{r}{k_j - \log(g_j(\underline{x}^{(i)}))} \prod_{s=2}^{i-1} \frac{1}{\sigma + \phi_j^{(s)}(g_j(\underline{x}^{(i)}))} \quad , \quad j = 1, 2, \dots, m \quad . \quad (4.24)$$

Thus the  $\{\underline{x}^{(i)}\}$  also correspond to a sequence minimizing  $Q(\underline{x}, \underline{\lambda}^{(i)})$  by the assumed uniqueness of these stationary values. Now, as  $\sigma > 1$ ,  $\phi_j^{(s)} > 0$ ,  $s = 1, 2, \dots, i-1$ ,  $\lambda_j^{(i)}$  can be made arbitrarily small for each  $j$  by choosing  $i$  large enough. The desired result is thus a consequence of the remark following Theorem 1.2.  $\square$

Remark: The conditions of the theorem are satisfied if  $f(x)$  convex,  $g_i(x)$ ,  $i = 1, 2, \dots, m$  concave, and strict convexity / concavity holds for at least one of these functions.

In what follows it is convenient to use the superscript  $i$  to indicate the appropriate member of the log barrier function sequence (4.20).

Lemma 4.3:

$$\frac{\partial^2 \phi^{(i)}}{\partial g_j^2} \bigg/ \frac{\partial \phi^{(i)}}{\partial g_j} = - \frac{\rho_j^{(i)}}{g_j}, \quad i = 1, 2, \dots, \quad (4.25)$$

where

$$\rho_j^{(i)} = \rho_j^{(i-1)} + g_j \frac{\partial \phi^{(i)}}{\partial g_j} < \rho_j^{(i-1)}, \quad (4.26)$$

and

$$\rho_j^{(i)} \rightarrow 1, \quad g_j \rightarrow 0. \quad (4.27)$$

Proof: Let  $\phi^{(0)} = -\log g_j$ . then

$$\frac{\partial^2 \phi^{(0)}}{\partial g_j^2} \bigg/ \frac{\partial \phi^{(0)}}{\partial g_j} = - \frac{1}{g_j} \quad (4.28)$$

so that  $\rho_j^{(0)} = 1$ . Now, differentiating the relation

$$\frac{\partial \phi^{(i+1)}}{\partial g_j} = \frac{1}{\sigma + \phi^{(i)}} \frac{\partial \phi^{(i)}}{\partial g_j} \quad (4.29)$$

gives

$$\frac{\partial^2 \phi^{(i+1)}}{\partial g_j^2} = - \frac{1}{(\sigma + \phi^{(i)})^2} \left( \frac{\partial \phi^{(i)}}{\partial g_j} \right)^2 + \frac{1}{\sigma + \phi^{(i)}} \frac{\partial^2 \phi^{(i)}}{\partial g_j^2}$$



so that

$$\begin{aligned} \frac{\partial^2 \phi^{(i+1)}}{\partial g_j^2} / \frac{\partial \phi^{(i+1)}}{\partial g_j} &= - \frac{\partial \phi^{(i+1)}}{\partial g_j} + \frac{\partial^2 \phi^{(i)}}{\partial g_j^2} / \frac{\partial \phi^{(i)}}{\partial g_j} \\ &= - \frac{1}{g_j} \left( \rho^{(i)} + g_j \frac{\partial \phi^{(i+1)}}{\partial g_j} \right) \end{aligned}$$

This demonstrates (4.25) and (4.26), (4.27) follows on noting that  $\rho_j^{(0)} = 1$ , and that, from (4.29),

$$g_j \frac{\partial \phi^{(i)}}{\partial g_j} \rightarrow 0, \quad g_j \rightarrow 0, \quad i = 1, 2, \dots \quad \square$$

A consequence of this lemma is that, provided  $J(\underline{x}^*)$  nonsingular, then  $D^{(i)}_{H^{(i)}}$  is nonsingular for  $\underline{x}(r)$  sufficiently close to  $\underline{x}^*$ .

Lemma 4.4: Let  $J(\underline{x}^*)$  be nonsingular, and  $I_1 = B_0$ , then the SO condition is satisfied.

Proof: We have from (4.29) that

$$w_j^{(i)} = - \left( \frac{\partial \phi^{(i)}}{\partial g_j} \right)^2 / \frac{\partial^2 \phi^{(i)}}{\partial g_j^2} = \frac{g_j}{\rho_j^{(i)}} \frac{\partial \phi^{(i)}}{\partial g_j}, \quad (4.31)$$

so that  $w_j^{(i)} \rightarrow 0$  as  $r \rightarrow 0$  for  $j \in B_0$ . Now

$$\left\| \begin{bmatrix} \underline{x}_k - \underline{x}_k^* \\ \underline{u}_k - \underline{u}_k^* \end{bmatrix} \right\| \leq \left\| [D^{(i)}_{H^{(i)}}(\underline{x}_k)]^{-1} \begin{bmatrix} I \\ P^{(i)} \end{bmatrix} \begin{bmatrix} 0 \\ u_1^k g_1 \\ \vdots \\ u_m^k g_m \end{bmatrix} \right\|$$

+ smaller terms

where  $P^{(i)}$  is diagonal, and  $P_{jj}^{(i)} = 1/\rho_j^{(i)} \rightarrow 1$ ,  $r_k \rightarrow 0$ ,  
 $j = 1, 2, \dots, m$ . This result implies that for  $k$  large enough

$$\|\tilde{x}_k - x^*\| \leq K \sum_{j \in B_0} u_j^k g_j(\tilde{x}_k).$$

The SO condition is an immediate consequence of this inequality.  $\square$

Remark: If  $B_0 \subset I_1$  then  $w_j^{(i)}$  need not tend to zero for  $j \notin B_0$ . Thus eventually the largest components of  $w^{(i)}$  will be those associated with the inactive constraints. This implies that  $\|\tilde{x}_k - x^*\| = o(r_k)$ .

But  $Q_j - Y_{j \in B_0}$ , is  $o(r_k)$  which suggests that, in general, the SO condition does not apply. This case should be contrasted with the log and inverse cases where the contributions of the inactive constraints do not dominate in  $w$  (in the inverse case the active constraints dominate). We note that the SO condition is only sufficient for (4.1) to provide an error estimate and numerical experience indicates that it is applicable in the calculations with the log sequence. However, the above discussion suggests that to **attain** the maximum rate of convergence with the members of the log sequence, the inactive constraints should be identified and discarded. A possible way to do this automatically is by the use of a separable barrier objective function

$$Q(\tilde{x}, \tilde{p}_k) = f(\tilde{x}) + r_k \sum_{i=1}^m u_i^{k-1} \phi(g_i(\tilde{x})) \quad (4.32)$$

where  $\tilde{p}_k = r_k^{k-1} u$ ,  $r_k$  is the usual barrier parameter, and  $u_i^{k-1}$  is the multiplier estimate obtained from the previous minimization. This objective function has the property of forcing the multiplier estimates

for the inactive constraints to zero at a very fast rate. We have

$$\begin{aligned}
 u_i^{k+1} &= -r_{k+1} u_i^k \frac{\partial \phi}{\partial g_i}(\tilde{x}_{k+1}) \\
 &\leq r_{k+1} \rho_i u_i^k \leq \dots \leq \rho_i^{k+1} \prod_{j=1}^{k+1} r_j u_i^0
 \end{aligned} \tag{4.33}$$

where  $\rho_i$  is a bound for  $\frac{\partial \phi}{\partial g_i}(\tilde{x}_k)$ ,  $k = 1, 2, \dots$ .

This choice can also be favorable in the case of nonstrict complementarity. Consider the previous example

$$\min x_2 \quad \text{subject to } x_2 - x_1^2 \geq 0, \quad x_1 \geq 0.$$

Set  $Q = x_2 - r_k \{ \log(x_2 - x_1^2) + u_{k-1} \log x_1 \}$ . Then  $\nabla Q = 0$  gives

$$x_2 - x_1^2 = r_k, \quad 2x_1^2 = r_k u_{k-1},$$

so that

$$u_k = \frac{r_k u_{k-1}}{x_1} = 2^{1/2} r_k^{1/2} u_{k-1}^{1/2}.$$

Setting  $r_k = \alpha^{-k}$ ,  $u_k/2 = \beta_k$  reduces this to

$$\beta_k = \frac{1}{2} \beta_{k-1} - \frac{k}{2}.$$

The solution to this difference equation satisfying the initial condition  $\beta_0 = 0$  is

$$\beta_k = -k + (1 - (\frac{1}{2})^k).$$

From this it follows that  $u_k = O(r_k)$ , and hence that  $x_1 = O(r_k)$ .

Thus, for this example, we are able to obtain results as favorable as those in which strict complementarity holds.

Example. Show that the error estimate (4.1) is valid in this case.

There is a penalty to pay for the generality of the barrier function algorithms, and this is a significant burden of calculation associated with each of the successive unconstrained minimizations. This can be explained (at least in part) by looking at the Hessian of the barrier objective function. Experience (in part supported by theoretical results) indicates that the condition number of the Hessian is a good indicator of the degree of difficulty of an unconstrained optimization problem when it is solved by descent methods.

On the assumptions that the second order sufficiency conditions hold at  $\mathbf{x}^*$ , and that the active constraint gradients are linearly independent, then it is possible to deduce fairly complete information on the eigenvalues and eigenvectors of  $\nabla^2 T(\mathbf{x}_k, \mathbf{r}_k)$  from (3.8).

(i) There are  $n-t$  eigenvectors associated with eigenvalues of  $\nabla^2 T(\mathbf{x}_k, \mathbf{r}_k)$  that are  $O(1)$  as  $\mathbf{r}_k \rightarrow 0$ . The smallest eigenvalue tends to

$$m = \min_{\mathbf{v}} \frac{\mathbf{v}^T \nabla^2 f(\mathbf{x}^*, \mathbf{u}^*) \mathbf{v}}{\mathbf{v}^T \mathbf{v}}, \quad \forall \mathbf{v} \text{ such that } \nabla g_i(\mathbf{x}^*) \mathbf{v} = 0, \quad \forall i \in B_0.$$

(ii) There are  $t$  eigenvectors associated with eigenvalues of  $\nabla^2 T(\mathbf{x}_k, \mathbf{r}_k)$  which tend to  $\infty$  as  $\mathbf{r}_k \rightarrow 0$ . These eigenvectors are asymptotic to vectors of the form  $C_1(\mathbf{x}^*) \mathbf{y}_i$  where  $\mathbf{y}_i$  are eigenvectors of the problem

$$[C_1(\mathbf{x}^*)^T C_1(\mathbf{x}^*) - \mu_i A] \mathbf{y}_i = 0$$

where  $A$  is a diagonal matrix,  $A_{ii} = \frac{\partial^2 \phi}{\partial g_i^2} / \max_{1 \leq j \leq t} \frac{\partial^2 \phi}{\partial g_j^2}$ ,  $i = 1, \dots, t$

The corresponding eigenvalues tend to  $\infty$  like  $\mu_i r_k \max_{1 \leq j \leq t} \frac{\partial^2 \phi}{\partial g_j^2}$

-- that is, like  $\mu_i \frac{u_\alpha^*}{g_\alpha(x_k)}$  where  $\alpha$  is the maximizing index.

This shows that the condition number of  $\nabla^2 T(x_k, r_k)$  tends to  $\infty$  like  $1/g_\alpha(x_k)$  or like  $1/\|x_k - x_k^*\|$  if the SO condition is satisfied.

In this latter case we have shown that our measure of the cost of a barrier function calculation depends in the main on the accuracy desired rather than on the choice of barrier functions. However, our estimates for the log family indicate that these will be **somewhat** more expensive than the above estimate except when **all** constraints are active.

Note that the device introduced to force more effective elimination of the inactive constraints does not force the Hessian to be worse conditioned in the case that strict **complementarity** does not obtain, at least in the examples that have been worked out. The use of this device would appear to be an important improvement in barrier function algorithms.

##### 5. Analysis of penalty function methods.

Consider now the equality constrained problem (ECP)

$$\min_{x \in S} f(x), \quad S = \{x; h_i(x) = 0, i = 1, 2, \dots, q (i \in I_2)\}. \quad (3.1)$$

It is assumed that  $S$  is nonempty, and that (3.1) has a bounded minimum (say  $\bar{f}$ ).

Remark: The inequality constraint  $g(\underline{x}) \geq 0$  can be written as the equality constraint

$$h(\underline{x}) = \min(0, g(\underline{x})) = 0 \quad (5.2)$$

so that formally the ICP is a special case of an ECP. However  $h(\underline{x})$  given by (5.2) can have discontinuous first derivatives.

Definition:  $F(\underline{x}, \lambda)$  is a penalty objective function if

$$F(\underline{x}, \lambda) = f(\underline{x}) + \lambda \sum_{i=1}^q \psi(h_i(\underline{x})) \quad (5.3)$$

where  $\psi(h)$  is a monotonic increasing function of  $|h|$ , and  $\psi(0) = 0$ .

Example: Let  $\psi(h) = |h|^{1+\alpha}$  then  $\psi$  is a penalty function if  $\alpha > -1$ . If  $g(\underline{x})$  is concave then, from (5.2), so is  $h(\underline{x})$ , and  $\psi(h)$  is convex provided  $\alpha \geq 0$ . If  $\alpha < 0$  then  $\frac{d\psi}{dh}$  is unbounded as  $h \rightarrow 0$ .

Theorem 5.1: Let  $\{\lambda_i\} \uparrow \infty$ , and  $\underline{x}_r$  minimize  $F(\underline{x}, \lambda_r)$ . Then

$\{f(\underline{x}_r)\}$  nondecreasing,  $\{F(\underline{x}_r, \lambda_r)\}$  strictly increasing unless  $\underline{x}_r \in S$ ,

and  $\left\{ \sum_{i=1}^q \psi(h_i(\underline{x}_r)) \right\}$  nonincreasing. If  $\{\underline{x}_r\} \rightarrow \underline{x}^*$  then  $\underline{x}^*$  solves the ECP.

Proof: Let  $\lambda_r < \lambda_s$ . Then, provided  $\underline{x}_r, \underline{x}_s \notin S$ ,

$$F(\underline{x}_r, \lambda_r) \leq F(\underline{x}_s, \lambda_r) < F(\underline{x}_s, \lambda_s) \leq F(\underline{x}_r, \lambda_s).$$

Thus (compare Lemma 1.2) the results for the sequences follow as before.

We have

$$\min_{\underline{x}} F(\underline{x}, \lambda) \leq \min_{\underline{x} \in S} F(\underline{x}, \lambda) = \min_{\underline{x} \in S} f(\underline{x}) = \bar{f}. \quad (5.4)$$

Thus the  $F(\underline{x}_r, \lambda_r)$  are bounded, and hence  $x^* \in S$ . Now

$$\underline{x}^* \in S \Rightarrow f(x^*) \geq \bar{f},$$

but, by (5.4),

$$f(\underline{x}_r) \leq F(\underline{x}_r, \lambda_r) \leq \bar{f}$$

so that

$$\lim_{r \rightarrow \infty} f(\underline{x}_r) \leq \bar{f}.$$

Thus  $f(x^*) = \bar{f}$ , and  $x^*$  solves the ECP.  $\square$

Remark: In the more general case in which  $\{\underline{x}_r\}$  is bounded it follows, by restricting attention to convergent subsequences, that all limit points of  $\{\underline{x}_r\}$  solve the ECP.

Theorem 5.2: Let  $\{\underline{x}_r\} \rightarrow x^*$  and assume  $\psi(h_i)$  continuously differentiable, and  $\nabla h_i(x^*)$ ,  $i \in I_2$ , linearly independent. Define  $\underline{u}_r$  by

$$u_i^r = -\lambda_r \frac{\partial \psi}{\partial |h_i|} \text{sgn}(h_i), \quad i = 1, 2, \dots, q \quad (5.5)$$

then  $\{\underline{u}_r\} \rightarrow u^*$ , the vector of Lagrange multipliers for the ECP.

Proof: Define the matrix  $B(\underline{x}_r)$  by  $\kappa_i(B(\underline{x}_r)) = \nabla h_i(\underline{x}_r)^T$ ,  $i = 1, 2, \dots, q$ . The condition that  $\underline{x}_r$  minimize  $F(\underline{x}_r, \lambda_r)$  gives

$$0 = \nabla F(\underline{x}_r, \lambda_r) = \nabla f(\underline{x}_r) + \lambda_r \sum_{i=1}^q \frac{\partial \psi}{\partial |h_i|} \text{sgn}(h_i) \nabla h_i(\underline{x}_r) \quad (5.6)$$

so that

$$\nabla f(\underline{x}_r)^T = B(\underline{x}_r) \underline{u}_r. \quad (5.7)$$

Now  $B(\underline{x}_r)$  has full rank for  $\|\underline{x}_r - \underline{x}^*\|$  small enough. Thus

$$\begin{aligned} \underline{u}_r &= (B(\underline{x}_r)^T B(\underline{x}_r))^{-1} B(\underline{x}_r)^T \nabla f(\underline{x}_r)^T \\ &\rightarrow (B(\underline{x}^*)^T B(\underline{x}^*))^{-1} B(\underline{x}^*)^T \nabla f(\underline{x}^*)^T = \underline{u}^* . \quad \square \end{aligned} \quad (5.8)$$

Remark: (i) If strict complementarity holds so that  $|u_i^*| > 0$ ,  $i = 1, 2, \dots, q$ , then the convergence of the Lagrange multiplier estimates implies (from (5.5)) that  $\text{sgn}(h_i)$  is constant for  $i$  large enough as  $\frac{\partial \psi}{\partial |h_i|} > 0$ ,  $\underline{x} \in S$ . Thus the minimizing sequence approaches  $S$  'from one side'. In this sense  $S$  acts like a barrier.

(ii) Note that  $h_i(\underline{x}) = \min(0, g_i(\underline{x})) = 0$  identically in a neighborhood of  $\underline{x}^*$  if  $g_i(\underline{x}^*) > 0$ . Thus  $\nabla h_i(\underline{x}^*) = \underline{0}$  so that, in this trivial sense, the constraint gradients are not linearly independent. However, if strict complementarity holds, then a multiplier result can be proved for the active constraints (do this!). In fact, the strict complementarity restriction can be relaxed somewhat.

Theorem 5.3: If the conditions of Theorem 5.2 hold, and, in addition,

$$\{\underline{u}_r\} \rightarrow \underline{u}^*, \text{ and } \lambda_r \frac{d^2 \psi}{dh_i^2} \rightarrow \infty, \quad i = 1, 2, \dots, q, \text{ as } r \rightarrow \infty, \text{ then the}$$

second order sufficiency conditions hold at  $\underline{x}^*$  if and only if

$$\nabla^2 F(\underline{x}_r, \lambda_r)$$
 is positive definite for  $r$  sufficiently large.

Proof: This is essentially the same as that of Lemmas 3.1 and 3.2.  $\square$

Example: Derive the analogues of Lemmas 3.1 and 3.3 which apply when the ECP is obtained by transforming an ICP by means of (5.2).



Remark: The condition that  $\lambda_r \frac{d^2\psi}{dh_i^2} \rightarrow \infty$  is related to strict

complementarity. Consider  $\psi = |h_i|^{1+\alpha}$ ,  $\alpha > 0$ . Then

$$\frac{d\psi}{dh_i} = (1+\alpha) |h_i|^\alpha \text{sgn}(h_i), \quad -\lambda_r \frac{d\psi}{dh_i} = u_i^r \quad (5.9)$$

so that

$$\lambda_r \frac{d^2\psi}{dh_i^2} = (1+\alpha)\alpha |h_i|^{\alpha-1} = \frac{\alpha |u_i^r|}{|h_i|} \quad (5.10)$$

Thus  $\lambda_r \frac{d^2\psi}{dh_i^2} \rightarrow \infty$ ,  $h_i \rightarrow 0$  if  $|u_i| > 0$ . Strict complementarity

is of particular importance for equality constraints derived from inequality constraints by (5.2). In this case, the one sided convergence implied by the multiplier relations is needed if we are to be able to talk about second derivatives at all.

The parallel development of the treatment of the ECP by penalty function methods and the treatment of the ICP by barrier functions can be completed by discussing convergence rates of penalty function algorithms in much the same way as we treated the barrier function case. For example, multiplying (5.6) by  $\tilde{x}_r - \tilde{x}^*$  gives

$$f(\tilde{x}^*) - f(\tilde{x}_r) = \sum_{i=1}^q u_i^r h_i(\tilde{x}_r) = o(\|\tilde{x}_r - \tilde{x}^*\|^2) \quad (5.11)$$

The assumption that the SO condition is satisfied can now be used to provide estimates. From (5.9)

$$(1+\alpha)\lambda_r |h_i|^\alpha \text{sgn}(h_i) = -u_i^r$$

so that

$$|h_i| = \left\{ \frac{|u_i^r|}{(1+\alpha)\lambda_r} \right\}^{1/\alpha} \quad (5.12)$$

This suggests a rate of convergence of  $O((\frac{1}{\lambda_r})^{1/\alpha})$ , which contrasts favorably with the estimates obtained for the barrier function algorithms. In particular as  $\alpha \rightarrow 0$ , (5.12) suggests that the convergence rate becomes arbitrarily great. However, the results of the previous section also indicate that the condition number of the Hessian will become arbitrarily large as  $\alpha \rightarrow 0$ . The next result provides information on the limiting case  $\alpha = 0$ .

Theorem 5.4: In the ICP let  $f(x)$  be convex, and  $g_i(x)$ ,  $i \in I_1$ , concave. Let  $w$  be an infeasible point,  $x_0$  an interior point of  $S$ ,  $a = \min_{i \in I_1} g_i(x_0)$ ,  $b = f(x_0) - f(x^*)$ , and  $\lambda_0 = (b+1)/a$ . Then  $x^*$

minimizes

$$F(x, \lambda) = f(x) - \lambda \sum_{i=1}^p \min(0, g_i(x)) \quad (5.13)$$

provided  $\lambda \geq \lambda_0$ .

Remark: It is necessary to demonstrate the result only for  $\lambda = \lambda_0$ . For all larger  $\lambda$  it is then a consequence of Theorem 5.1.

Proof: Let  $v$  be the boundary point of  $S$  on the join of  $w$  and  $x_0$ , and  $B_v$  be the index set of constraints active at  $v$ . Define

$$s(x) = f(x) - \lambda_0 \sum_{i \in B_v} g_i(x), \quad (5.14)$$

then

$$\begin{aligned}
s(\underline{x}_0) &= f(\underline{x}_0) - \lambda_0 \sum_{i \in B_V} g_i(\underline{x}_0) \\
&\leq f(\underline{x}_0) - (b+1) \\
&= f(\underline{x}^*) - 1 \\
&\leq f(\underline{v}) = s(\underline{v}) = F(\underline{v}, \lambda_0).
\end{aligned} \tag{5.15}$$

As  $s(\underline{x})$  is convex and  $\underline{v}$  is on the join of  $\underline{x}_0$  and  $\underline{w}$ ,  $\exists \theta$ ,  $0 < \theta < 1$ , such that

$$\begin{aligned}
s(\underline{v}) &= \theta s(\underline{x}_0) + (1-\theta) s(\underline{w}) \\
&\leq \theta s(\underline{v}) + (1-\theta) s(\underline{w})
\end{aligned}$$

whence

$$s(\underline{v}) \leq s(\underline{w}) . \tag{5.16}$$

Now  $s(\underline{w}) \leq F(\underline{w}, \lambda_0)$  so that, from (5.15) and (5.16)

$$F(\underline{v}, \lambda_0) \leq F(\underline{w}, \lambda_0) .$$

Thus  $\min_{\mathbb{R}} F(\underline{x}, \lambda_0)$  must be attained at a feasible point.  $\square$

## 6. Accelerated penalty and barrier methods.

The problems of poor conditioning of the computational problem and (comparatively) slow convergence make it worthwhile to search for methods for accelerating the convergence of the penalty/barrier function algorithms. Consider the (generalized) penalty objective function

$$P(\underline{x}, \underline{W}, \underline{\eta}) = f(\underline{x}) + \sum_{i=1}^q W_{ii} \psi(h_i(\underline{x}) + \eta_i) \tag{6.1}$$

where  $W$  is the diagonal matrix of penalty parameters, and the  $\eta_i$  are further parameters to be used in the acceleration process.

At a minimum of  $P$ ,  $\tilde{x}(W, \eta)$  satisfies

$$\nabla f(x) - \sum_{i=1}^q u_i(W, \eta) \nabla h_i(x) = 0 \quad (6.2)$$

where  $u_i(W, \eta) = -W_{ii} \frac{\partial \psi}{\partial h_i}$ . Provided  $\|\tilde{x}(W, \eta) - x^*\|$  and  $\|u(W, \eta) - u^*\|$  are sufficiently small, the second order sufficiency conditions hold at  $x^*$ , and  $\nabla h_i(x^*)$ ,  $i \in I_2$ , are linearly independent, then (by Theorem 5.3)  $x^*$  also solves the ECP

$$\min_{\tilde{x} \in S_{W, \eta}} f(x) \quad , \quad S_{W, \eta} = \{x ; h_i(x) = h_i(\tilde{x}(W, \eta)) , i \in I_2\} .$$

One sequential strategy for making  $\tilde{x}(W, \eta) \rightarrow x^*$  is to force  $\lambda = \min_i W_{ii} \uparrow \infty$ . However, the parameter vector  $\eta$  is also available, and we ask is it possible to adjust it to make

$$h_i(\tilde{x}(W, \eta)) = 0 \quad , \quad i \in I_2 \quad . \quad (6.3)$$

Let  $\frac{\partial x}{\partial \eta}$  be the matrix with components  $\frac{\partial x_i}{\partial \eta_j}$ ,  $i = 1, \dots, n$ ,  $j = 1, 2, \dots, q$ . If  $\nabla^2 P(x(W, \eta), W, \eta)$  is nonsingular, then, by the implicit function theorem, we can solve (6.2) for  $\tilde{x} = \tilde{x}(\eta)$  holding  $W$  fixed. We have

$$\nabla^2 P \frac{\partial x}{\partial \eta} = - \sum_{i=1}^q W_{ii} \frac{\partial^2 \psi}{\partial h_i \partial \eta_i} \nabla h_i^T$$

where all quantities are evaluated at  $\tilde{x}(W, \eta)$ . Defining the diagonal

matrix  $V$  by  $V_{ii} = W_{ii} \frac{\partial^2 \psi}{\partial h_i^2} = W_{ii} \frac{\partial^2 \psi}{\partial h_i \partial \eta_i}$ ,  $i = 1, 2, \dots, q$ , and the matrix  $B$  by  $\kappa_i(B) = \nabla h_i^T$ ,  $i = 1, 2, \dots, q$ , then (6.4) can be written

$$(\nabla^2 \mathcal{L} + BVB^T) \frac{\partial x}{\partial \eta} = -BV \quad (6.5)$$

Choose  $V_0$  to make  $U = \nabla^2 \mathcal{L} + BV_0B^T$  positive definite, and set  $V_1 = V - V_0$ . Then, by (3.6), if  $\min_i V_{ii}$  is sufficiently large

$$\frac{\partial x}{\partial \eta} \sim -U^{-1}B(B^T U^{-1}B)^{-1} + o(V^{-1}) \quad (6.6)$$

This relation can be justified if

- (i) the second order sufficiency conditions hold at  $\tilde{x}^*$  and  $\nabla h_i(\tilde{x}^*)$ ,  $i \in I_2$ , are linearly independent,
- (ii)  $\|\tilde{x} - \tilde{x}^*\|$  and  $\|u - u^*\|$  are sufficiently small,
- (iii)  $\min_i V_{11} \rightarrow \infty$  as  $\min_i W_{ii} \rightarrow \infty$  for  $\tilde{\eta} = \tilde{\theta}$ , and
- (iv)  $\min_i W_{11}$  sufficiently large.

Consider now the use of Newton's method for solving (6.3). This suggests that a correction  $\delta \eta$  to  $\tilde{\eta}$  be found by solving

$$\nabla_{\tilde{\eta}} h \frac{\partial x}{\partial \eta} \delta \tilde{\eta} = B^T \frac{\partial x}{\partial \eta} \delta \tilde{\eta} = -h \quad (6.7)$$

But, by (6.6),

$$B^T \frac{\partial x}{\partial \eta} \sim -I + o(V^{-1}) \quad (6.8)$$

so that

$$\delta \tilde{\eta} = \tilde{h} + o(V^{-1}) \quad (6.9)$$

Thus we expect the simple correction  $\delta\eta = h$  to approximate arbitrarily closely to a second order process provided  $V$  is sufficiently large.

Algorithm (ECP)

- (i) Initialize  $\eta^{(1)}$ ,  $W^{(1)}$ .
- (ii) Minimize  $P(x, W^{(k)}, \eta^{(k)})$  to determine  $x_{\sim k}$ ,  $u_{\sim k}$ .
- (iii) IF  $\sum_{i=1}^q u_i^k h_i(x_{\sim k}) < \text{TOL}$  THEN STOP.
- (iv) FOR  $I = 1$  STEP 1 UNTIL  $Q$  DO  
 IF  $\text{ABS}(h_i(x_{\sim k})) < \text{DECR} * \text{ABS}(h_i(x_{\sim k-1}))$   
 THEN  $\eta_i^{k+1} = \eta_i^k + h_i(x_{\sim k})$   
 ELSE  $W_{ii}^{(k+1)} = \text{EXP} * W_{ii}^{(k)}$   

$$\eta_i^{k+1} = \left( \frac{\partial \psi}{\partial h_i} \right)^{-1} \left( - \frac{u_i^k}{W_{ii}^{(k+1)}} \right).$$
- (v)  $K := K+1$ .
- (vi) GO TO (ii).

Remark: The idea behind the algorithm is that the correction (6.9)

is used whenever the convergence of  $h_i$  to zero is satisfactory.

Otherwise it is assumed that  $W_{ii}$  is too small and it is increased

accordingly.  $\eta_i$  is modified at the same time to ensure that

$u_i^k \rightarrow u_i^*$  as  $h_i \rightarrow 0$ .  $\left( \frac{\partial \psi}{\partial h} \right)^{-1}$  indicates the inverse function to  $\frac{\partial \psi}{\partial h}$ .

For the ICP we consider the modified barrier function

$$R(x, W^{(k)}, \eta^{(k)}) = f(x) + \sum_{i=1}^m W_{ii}^{(k)} u_i^{k-1} \phi(\mathbf{g}_i(x) + \eta_i^k) \quad (6.10)$$

where  $\underline{u}_{k-1}$  is the vector of multiplier estimates from the previous minimization, and  $W^{(k)}$  is now the diagonal matrix of barrier parameters. We note that, in the particular case in which all constraints are active, the previous analysis is applicable, at least formally, and suggests a correction

$$\delta \underline{\eta}^{(k)} = \underline{\eta}^{(k)} + \underline{g}(\underline{x}_k) \quad (6.11)$$

with order of magnitude departure from a second order iteration of

$$O\left(\left(\min_i W_{ii}^{(k)} u_i^{k-1} \frac{\partial^2 \phi}{\partial g_i^2}\right)^{-1}\right). \quad \text{However, we require automatic selection}$$

of the active constraints if we are to make use of this result, and it is important to note that this is provided naturally in the algorithm by the options

- (i) if  $g_i \rightarrow 0$  at a satisfactory rate then  $\eta_i^{k-1} = \eta_i^k + g_i$ , and
- (ii) if the convergence rate is too slow, then decrease the barrier parameter.

This second option can be expected to apply to the inactive constraints, and will drive the contribution to (6.4) from this source rapidly to zero by (4.33). Note that the boundedness of the barrier terms requires that  $g_i + \eta_i$  be positive. If  $\underline{\eta}^{(1)}$  is set to zero then (6.11) ensures that this condition will be met initially. Provided strict complementarity holds, the convergence of the multiplier estimates will ensure that it must hold ultimately. Of course, the calculation must be started from a feasible point.

Algorithm (ICP)

- (i) Initialize  $\eta^{(1)}$ ,  $W^{(1)}$ ,  $u_0$ .
- (ii) Minimize  $R(x, W^{(k)}, \eta^{(k)})$  to determine  $\tilde{x}_k$ ,  $\tilde{u}_k$ .
- (iii) IF  $\sum_{i=1}^m u_i^k g_i(\tilde{x}_k) < \text{TOL}$  THEN STOP.
- (iv) FOR I = 1 STEP 1 UNTIL M DO  
 IF  $\text{ABS}(g_i(\tilde{x}_k)) < \text{DECR} * \text{ABS}(g_i(\tilde{x}_{k-1}))$   
 THEN  $W_{ii}^{(k+1)} = -1 / \frac{\partial \phi}{\partial g_i}$ ,  
 $\eta_i^{(k+1)} = \eta_i^{(k)} + g_i(\tilde{x}_k)$   
 ELSE  $W_{ii}^{(k+1)} = \text{DECR} * W_{ii}^k$ ,  
 $\eta_i^{(k+1)} = \left( \frac{\partial \phi}{\partial g_i} \right)^{-1} \left( - \frac{1}{W_{ii}^{(k+1)}} \right)$ .
- (v)  $K := K+1$ .
- (vi) GO TO (ii).

Remark: As in the previous algorithm  $\left( \frac{\partial \phi}{\partial g} \right)^{-1}$  denotes the inverse function to  $\frac{\partial \phi}{\partial g}$ . For example, if  $\phi = -\log g$  then  $\eta = W$ .

Consider now another modified penalty function for the ECP

$$s(\tilde{x}) = f(\tilde{x}) - \tilde{u}(\tilde{x})^T \tilde{h}(\tilde{x}) + \tilde{h}(\tilde{x})^T W \tilde{h}(\tilde{x}) \quad (6.12)$$

where the matrix  $W$  is positive definite.

Lemma 6.1: If the second order sufficiency conditions hold for  $x = \tilde{x}^*$ , and the  $\nabla h_i(\tilde{x}^*)$ ,  $i \in I_2$ , are linearly independent then  $S(x)$  has a



local minimum at  $\underline{x} = \underline{x}^*$  provided  $\underline{u}(\underline{x}) \rightarrow \underline{u}^*$  as  $\underline{x} \rightarrow \underline{x}^*$  and the smallest eigenvalue of  $W$  is large enough.

Proof: We have

$$\begin{aligned} \nabla S(\underline{x}^*) &= \nabla f(\underline{x}^*) - \underline{u}(\underline{x}^*)^T \nabla h(\underline{x}^*) \\ &\quad - \underline{h}(\underline{x}^*)^T (\nabla \underline{u}(\underline{x}^*) - 2W \nabla h(\underline{x}^*)) \\ &= \nabla f(\underline{x}^*) - \underline{u}^{*T} \nabla h(\underline{x}^*) \\ &= 0 \end{aligned}$$

as  $\underline{u}^*$  is the vector of Lagrange multipliers for the ECP. Thus  $S$  has a stationary point for  $\underline{x} = \underline{x}^*$ . Now

$$\begin{aligned} \nabla^2 S(\underline{x}^*) &= \nabla^2 \underline{f}(\underline{x}^*, \underline{u}^*) - \nabla \underline{u}(\underline{x}^*)^T \nabla h(\underline{x}^*) \\ &\quad - \nabla h(\underline{x}^*)^T \nabla \underline{u}(\underline{x}^*) \\ &\quad + 2 \nabla h(\underline{x}^*)^T W \nabla h(\underline{x}^*) \end{aligned} \tag{6.13}$$

where terms which vanish at  $\underline{x}^*$  have been ignored. Corollary 3.2 can now be applied to show  $\nabla^2 S(\underline{x}^*)$  is positive definite. We set  $V = \nabla h(\underline{x}^*)^T$ ,  $U = \nabla^2 \underline{f}(\underline{x}^*, \underline{u}^*) - \nabla \underline{u}(\underline{x}^*)^T \nabla h(\underline{x}^*) - \nabla h(\underline{x}^*)^T \nabla \underline{u}(\underline{x}^*)$ , and note that

$$\min_{\substack{\underline{v}^T \underline{t} = 0 \\ \|\underline{t}\| = 1}} \underline{t}^T U \underline{t} = \min_{\substack{\underline{v}^T \underline{t} = 0 \\ \|\underline{t}\| = 1}} \underline{t}^T \nabla^2 \underline{f}(\underline{x}^*, \underline{u}^*) \underline{t} = m > 0$$

as the second order sufficiency conditions hold at  $\underline{x}^*$ .  $\square$

Theorem 6.1: Let the conditions of Lemma 6.1 hold at  $x^*$  and set

$$\tilde{u}(\tilde{x}) = B(\tilde{x}) + \nabla f(\tilde{x})^T \quad (6.14)$$

where  $B(\tilde{x}) = \nabla h(\tilde{x})^T$ . Then (6.12) has a local minimum at  $x^*$  provided the smallest eigenvalue of  $W$  is large enough.

Proof: This result is an immediate consequence of Lemma 6.1. As the  $\nabla h_i(x^*)$ ,  $i \in I_2$ , are linearly independent,  $B(x^*)^+$  is a bounded operator for  $\|x - x^*\|$  small enough. Thus  $\tilde{u}(\tilde{x}) \rightarrow u^*$  as  $x \rightarrow x^*$ .  $\square$

Remark: (i) By using (6.14) we can construct a penalty function which is differentiable in a neighborhood of  $x^*$  (contrast with (5.13)) and which has a local minimum at  $x = x^*$  for sufficiently large but finite values of the penalty parameter. However, (6.14) requires first derivatives of the problem functions so that minimization of (6.12) with a method that requires first derivatives of  $S$  will require second derivatives of the problem functions. Two cases have been considered (Fletcher).

$$(i) S(x) = f(x) - h(\tilde{x})^T B(\tilde{x})^+ \nabla f(\tilde{x})^T + \sigma \|h(\tilde{x})\|^2, \quad \text{and} \quad (6.15)$$

$$(ii) S(x) = f(x) - h(\tilde{x})^T B(\tilde{x})^+ \nabla f(\tilde{x})^T + \sigma \|(B(\tilde{x})^+)^T h(\tilde{x})\|^2, \quad (6.16)$$

where  $\sigma$  is a penalty parameter.

(ii) There is a close connection between the penalty function (6.15) and the algorithm based on (6.1) in the case  $\psi(h) = h^2$ . At a minimum of  $P$  we have (as  $VP = 0$ )

$$2W(\tilde{h} + \theta) = -B^+(\nabla f)^T. \quad (6.17)$$

Thus the correction formula corresponds to updating the Lagrange multiplier estimate by (6.14) at the end of each unconstrained minimization rather than continuously which the use of S requires.

(iii) Note that  $S(\underline{x})$  can be interpreted as a **Lagrangian**. For example, in the case  $S(\underline{x})$  is given by (6.16),

$$S(\underline{x}) = \mathcal{L}(\underline{x}, \underline{w}(\underline{x})) \quad (6.18)$$

where

$$\underline{w}(\underline{x}) = B(\underline{x})^+ \nabla f(\underline{x})^T - \sigma B(\underline{x})^+ (B(\underline{x})^+)^T \underline{h}(\underline{x}) \quad (6.19)$$

**Lemma 6.2:**  $\underline{w}(\underline{t})$  defined by (6.19) is the vector of Lagrange multipliers for the problem

$$\underset{\underline{x}}{\text{minimize}} \quad f(\underline{t}) + \nabla f(\underline{t})^T (\underline{x} - \underline{t}) + \frac{\sigma}{2} \|\underline{x} - \underline{t}\|^2 \quad (6.20)$$

subject to the linear constraints

$$\underline{h}(\underline{t}) + \nabla \underline{h}(\underline{t})^T (\underline{x} - \underline{t}) = 0 \quad (6.21)$$

provided this minimum exists.

Proof: Any point satisfying the constraints (6.21) has the form

$$\underline{x} = \underline{t} - (B(\underline{t})^T)^+ \underline{h}(\underline{t}) + A(\underline{t}) \underline{z} \quad (6.22)$$

where  $B(\underline{t})^T A(\underline{t}) = 0$ . The multiplier relation for (6.20), (6.21) is

$$\nabla f(\underline{t}) + \sigma (\underline{x} - \underline{t})^T = \underline{u}^T B(\underline{t})^T \quad (6.23)$$

so that  $\underline{u}$  can be taken as (substituting (6.22) into (6.23))

$$\underline{u} = B(\underline{t})^+ \{ \nabla f(\underline{t})^T - \sigma (B(\underline{t})^T)^+ \underline{h}(\underline{t}) \} \quad \square$$

If (6.22) is substituted into (6.20) the problem becomes one of minimizing w.r.t.  $\tilde{z}$

$$\nabla f A \tilde{z} + \frac{\sigma}{2} \tilde{z}^T A^T A \tilde{z}$$

whence

$$\tilde{z} = - \frac{1}{\sigma} (A^T A)^{-1} A^T \nabla f^T \quad (6.24)$$

Thus  $\sigma$  plays a role in ensuring that  $\tilde{x}(t)$ , the minimum of (6.20), cannot deviate far from  $t$  (cf. remark (ii) following MP Corollary 5.1).

Example: (i) The Lagrangian interpretation provides a method for generalizing the above discussion to inequality constraints. Consider the problem  $\min_x \mathcal{J}(x, w(x))$  where  $\tilde{w}(t)$  is the vector of multipliers for

the problem

$$\begin{aligned} \min_X & f(t) + \nabla f(t)(x-t) + \frac{\sigma}{2} \|x-t\|^2 \\ & \text{subject to } g(t) + \nabla_{NH} g(t)(x-t) \geq 0 \end{aligned}$$

Under what conditions does  $\mathcal{J}$  have an unconstrained minimum at  $\tilde{x}^*$ . What role does strict complementarity play in this problem?

(ii) (6.12) can be generalized to other penalty functions and to barrier functions (cf. Remark (ii) above). How much of the above analysis goes through? What modifications are required? Evaluate the resulting algorithms.

## Notes

- 1., 2. See Fiacco and McCormick's book. Also the paper 'Penalty function methods for mathematical programming problems', J. Math. Anal. and Applic. (1970), by Osborne and Ryan.
3. Fiacco and McCormick were the first to draw attention to the importance of these (as they were to much of the material in this section).
4. The log family is due to Osborne and Ryan. The importance of the conditioning of the Hessian to Walter Murray. Rate of convergence formulae have also been developed by F. A. Lootsma, (Thesis, also survey paper at Dundee conference).
5. Fiacco and McCormick. The exact penalty function is due to Zangwill.
6. The algorithm for the ECP is due to Powell in the case  $\psi = h^2$  (Harwell report, also Proceedings of Keele Conference). The exact penalty function  $S(x)$  is due to Fletcher who has developed it together with his student Shirley Lill and described it in several Harwell reports. The extension to inequality constraints (example (i)) is also due to Fletcher.

## References

- A. V. Fiacco and G. P. McCormick (1968): Nonlinear Programming: Sequential Unconstrained Minimization Techniques, Wiley.
- R. Fletcher and S. A. Lill (1970): A class of methods for non-linear programming, in Proceedings of Nonlinear Programming Symposium, Madison.

- F. A. Lootsma (1970): Boundary properties of penalty functions for constrained minimization, Phillips Res. Repts. Suppl., No. 3.
- W. Murray (1969): Constrained Optimization, Nat. Phys. Lab. Rep. Ma79.
- M. R. Osborne and D. M. Ryan (1970): On penalty functions for non-linear programming problems, J. Math. Anal. Appl., 31, pp. 559-578.
- M. J. Powell (1969): A Method for non-linear constraints in minimization problems, in Optimization (editor: R. Fletcher), Academic Press.

## Bibliography

- Abadie, J. (editor), (1967) Nonlinear Programming, (1970) Integer and Nonlinear Programming, North Holland.
- Ablow, C. M. and Brigham, G. (1955). An analog solution of programming problems. *Op. Res.*, 3, pp 388-394.
- Akaike, H. (1959). On a successive transformation of probability distribution and its application to the analysis of the optimal gradient method. *Ann. Inst. Stat. Math.*, Tokyo, 11, p 1.
- Allran, R. R. and Johnsen, S. E. J. (1970). An algorithm for solving nonlinear programming problems subject to nonlinear inequality constraints. *Comp. J.*, 13(2), pp 171-177.
- Arrow, K. J. and Hurwicz, L. (1956). Reduction of constrained maxima to saddle point problems; in "Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability", ed. J. Neyman. University of California Press, pp 1-20.
- Arrow, K. J., Hurwicz, L. and Uzawa, H. (1958). "Studies in Linear and Nonlinear Programming". Stanford University Press.
- Arrow, K. J., Hurwicz, L. and Uzawa, H. (1961). Constraint qualifications in maximization problems. *Nav. Res. Log. Q.*, 8, pp 175-191.
- Bard, Y. (1968). On a numerical instability of Davidon-like methods. *Maths. Comp.*, 22 (103), pp 665-666.
- Bard, Y. (1970). Comparison of gradient methods for the solution of nonlinear parameter estimation problems. *SIAM J. Numer. Anal.*, 7 (1), pp 157-186.
- Bartels, R. H., Golub, G. H., and Saunders, M. A. (1970), Numerical techniques in Mathematical Programming, in Nonlinear Programming, Academic Press.
- Bazaraa, M. S. Goode, J. J., Nashed, M. Z., and Shetty, C. M. (1971). Nonlinear programming without differentiability in Banach spaces: Necessary and sufficient constraint qualification. *Applicable Analysis*, to appear.

- Beale, E. M. L. (1955). On minimizing a convex function subject to linear inequalities. *J. Roy. Stat. Soc., Ser. B*, **17** (2), pp 173-184.
- Beale, E. M. L. (1959). Numerical methods; in "Nonlinear Programming", ed. J. Adadie. North-Holland Publishing Co., p 150.
- Bellmore, M., Greenberg, H. J. and Jarvis, J. J. (1970). Generalized penalty function concepts in mathematical optimization. *Op. Res.*, **18** (2), pp 229-252.
- Beltrami, E. J. (1969). A constructive proof of the Kuhn-Tucker multiplier rule. *J. Math. Anal. Appl.*, **26**, pp 297-306.
- Beltrami, E. J. (1970). "An Algorithmic Approach to Nonlinear Analysis and Optimization". Academic Press.
- Beltrami, E. J. and McGill, R. (1966). A class of variational problems in search theory and the maximum principle. *Op. Res.*, **14** (2), pp 267-278.
- Box, M. J. (1966). A comparison of several current optimization methods and the use of transformations in constrained problems. *Comp. J.*, **9**, pp 67-77.
- Box, M. J., Davies, D. and Swann, W. H. (1969). "Nonlinear Optimization Techniques". ICI Monograph, Oliver and Boyd.
- Bracken, J. and McCormick, G. P. (1968). "Selected Applications of Nonlinear Programming". Wiley.
- Broyden, C. G. (1965). A class of methods for solving nonlinear simultaneous equations. *Math.- Comp.*, 19 (92), pp 577-593.
- Broyden, C. G. (1967). Quasi-Newton methods and their application to function minimization. *Math. Comp.*, 21 (19), pp 368-381.
- Broyden, C. G. (1970a). The convergence of a class of double-rank minimization algorithms; 1. general considerations. *J. Inst. Math. Appl.*, 6 (1) pp 76-90.
- Broyden, C. G. (1970b). The convergence of single-rank quasi-Newton methods. *Math. Comp.*, 24 (110), pp 365-382.



- Bui Trong Lieu and Huard, P. (1966). La méthode des centres dans un espace topologique. *Numer. Math.*, 8 (1), pp 56-67.
- Butler, T. and Martin, A. V. (1962). On a method of Courant for minimizing functionals. *J. Math. Phys.*, 41, pp 291-299.
- Camp, G. D. (1955). Inequality-constrained stationary value problems. *J. Op. Res. Soc. Am.*, 3, pp 548-550.
- Carrol, c. w. (1961). The created response surface technique for optimizing nonlinear restrained systems. *Op. Res.*, 9 (2), pp 169-184.
- Charnes, A. and Lemke, C. E. (1954). Minimization of nonlinear separable convex functionals. *Nav. Res. Log. Q.*, 1, pp 301-302.
- Cheney, E. W. and Goldstein, A. A. (1959). Newton's method for convex programming and Tchebycheff approximation. *Numer. Math.*, 1, pp 254-268.
- Colville, A. R. (1968). A comparative study on nonlinear programming codes. IBM New York Scientific Centre Tech. Rep. No. 320-2949.
- Courant, R. (1943). Variational methods for the solution of problems of equilibrium and vibrations. *Bull. Am. Math. Soc.*, 49, pp 1-23.
- Crockett, J. B. and Chernoff, H. (1955). Gradient methods of maximization. *Pac. J. Math.*, 5, pp 33-50.
- Curry, H. B. (1944). The method of steepest descent for non-linear minimization problems. *Quart. Appl. Math.*, 2 (3), pp 250-261.
- Daniel, J. W. (1967a). The conjugate gradient method for linear and nonlinear operator. equations. *SIAM J. Numer. Anal.*, 4, pp 10-26.
- Daniel, J. W. (1967b). Convergence of the conjugate gradient method with computationally convenient modifications. *Numer. Math.*, 10, pp 125-131.
- Dantzig, G. B. (1951). Maximization of a linear function of variables subject to linear inequalities: in "Activity Analysis of Production and Allocation", ed. T. C. Koopmans. Cowles Commission Monograph No. 13, Wiley.

- Dantzig, G. B. (1963). Variable metric method for minimization. U.S.A.E.C. Research and Development Rep. ANL-5990, Argonne National Laboratory.
- Davidon, W. C. (1968). Variance algorithm for minimization. *Comp. J.*, 10, pp 406-410.
- Dorn, W. S. (1960). Duality in quadratic programming. *Quart. Appl. Math.*, 18, pp 155-162.
- Edelbaum, T. N. (1962). Theory of maxima and minima: in "Optimization Techniques", ed. G. Leitmann. Academic Press.
- El-Hodiri, Mohamed A. (1971). Constrained Extrema, Lecture Notes in Operations Research and Mathematical Systems 56, Springer-Verlag.
- Evans, J. P. and Gould, F. J. (1970). Stability in non-linear programming. *Op. Res.*, 18 (1), pp 107-118.
- E'laure, P. and Huard, P. (1966). Résultats nouveaux relatifs à la méthode des centres. Quatrième Conference de Recherche Opérationnelle, Cambridge, Mass.
- Fiacco, A. V. (1967). Sequential unconstrained minimization methods for nonlinear programming. Ph.D. Dissertation, Northwestern University, Ill.
- Fiacco, A. V. (1970). Penalty methods for mathematical programming in  $E^n$  with general constraint sets. *J. Opt. Th. Appl.*, 6 (3), pp 252-268.
- Fiacco, A. V. and McCormick, G. P. (1963). Programming under non-linear constraints by unconstrained minimization: a primal-dual method. Tech. Paper RAC-TP-96, Research Analysis Corporation.
- Fiacco, A. V. and McCormick, G. P. (1964a). The sequential unconstrained minimization technique for nonlinear programming: a primal-dual method. *Man. Sci.*, 10 (2), pp 360-366.
- Fiacco, A. V. and McCormick, G. P. (1964b). Computational algorithm for the sequential unconstrained minimization technique for nonlinear programming. *Man. Sci.*, 10 (2), pp 601-617.

- Fiacco, A. V. and McCormick, G. P. (1966). Extensions of SUMT for nonlinear programming: equality constraints and extrapolation. *Man. Sci.*, 12 (11), pp 816-829.
- Fiacco, A. V. and McCormick, G. P. (1967). The slacked unconstrained minimization technique for convex programming. *SIAM J. Appo. Math.*, 15 (3), pp 505-515.
- Fiacco, A. V. and McCormick, G. P. (1968). "Nonlinear Programming: Sequential Unconstrained Minimization Techniques". Wiley.
- Fletcher, R. (1965). Function minimization without evaluating derivatives - a review. *Comp. J.*, 8, pp 33-41.
- Fletcher, R. (1968). Programming under linear equality and inequality constraints. ICI Rep. No. MSDH/68/19.
- Fletcher, R. (Ed.) (1969a). "Optimization". Academic Press.
- Fletcher, R. (1969b). A review of methods for unconstrained optimization: in "Optimization", ed. R. Fletcher. Academic Press.
- Fletcher, R. (1969c). A new approach to variable metric algorithms. AERE Tech. Rep. No. 383.
- Fletcher, R. (1970). A general quadratic programming algorithm. AERE Tech. Rep. No. 401.
- Fletcher, R. (1970). A new approach to variable metric algorithms, *Computer J.*, 13, pp 317-322.
- Fletcher, R. and McCann, A. P. (1969). Acceleration techniques for nonlinear programming: in "Optimization", ed. R. Fletcher. Academic Press. .
- Fletcher, R. and Powell, M. J. D. (1963). A rapidly convergent descent method for minimization. *Comp. J.*, 6, pp 163-168.
- Fletcher, R. and Reeves, C. M. (1964). Function minimization by conjugate gradients. *Comp. J.*, 7, pp 149- 154.
- Forsythe, G. E. (1967). On the asymptotic directions of the s-dimensional optimum gradient method. Tech. Rep. No. CS61, Stanford University.

- Forsythe, G. E. and Motzkin, T. S. (1951). Asymptotic properties of the optimum gradient method (abstract). Bull. Am. Math. Soc., **57**, p 183.
- E'risch, K. R. (1955). The logarithmic potential method of convex programming. Memorandum May 13, 1955, University Institute of Economics, Oslo.
- Geoffrion, A. M. (1966). Strictly concave parametric programming, part 1: basic theory. Man. Sci., 13, pp 244-253.
- Geoffrion, A. M. (1967). Strictly concave parametric programming, part 2: additional theory and computational considerations. Man. Sci., 13 (5), pp 359-370.
- Gill, P. E. and Murray, W. (1970). Stable implementation of unconstrained optimization algorithms, NPL report Maths. 97.
- Goldfarb, D. (1969a). Sufficient conditions for the convergence of a variable metric algorithm: in "Optimization", ed. R. Fletcher. Academic Press.
- Goldfarb, D. (1969b). Extension of Davidon's variable metric method to maximization under linear inequality and equality constraints. SIAM J. Appl. Math., **17**, pp 739-764.
- Goldfarb, D. (1970). A family of variable-metric methods derived by variational means. Math. Comp., **24**, pp 23-26.
- Goldfarb, D. and Lapidus, L. (1968). Conjugate gradient method for nonlinear programming problems with linear constraints. I. and E. C. Fundamentals, **7** (1), pp 142-151.
- Goldstein, A. A. (1962). Cauchy's method of iminimization. Numer. Math., **4**, pp 146-150.
- Goldstein, A. .A. (1965). On steepest descent, SIAM J. Control, **3**, pp 147-151.
- Goldstein, A. A. (1968). Constructive Real Analysis, Harper and Row.
- Goldstein, A. A. and Price, J. (1967). An effective algorithm for minimization, Num. Math., **10**, pp. 184-189.

- Golub, G. H., and Saunders, M. A. (1969). Numerical methods for solving linear least squares problems, Computer Science Dept. Tech. Rep. CS134, Stanford University.
- Gould, F. J. and Tolle, J. W. (1971). A necessary and sufficient constraint qualification, SIAM J. Appl. Maths. 20, pp 164-172.
- Greenstadt, J. (1967). On the relative efficiencies of gradient methods. Math. Comp., 21, pp 360-367.
- Greenstadt, J. (1970). Variations on variable-metric methods. Math. Comp., 24, pp 1-25.
- Haarhoff, P. C. and Buys, J. D. (1970). A new method for the optimization of a nonlinear function subject to nonlinear constraints. Comp. J., 13 (2), pp 178-184.
- Hadley, G. (1964). "Nonlinear and Dynamic Programming". Addison-Wesley.
- Hartley, H. O. and Hocking, R. R. (1963). Convex programming by tangential approximation. Man. Sci., 9, pp 600-612.
- Hestenes, M. R. (1966). "Calculus of Variations and Optimal Control Theory". Wiley.
- Huang, H. Y. (1969). Unified approach to quadratically convergent algorithms for function minimization. J.O.T.A., 5, pp 405-425.
- Huang, H. Y. and Levy, A. V. (1969). Numerical experiments on quadratically convergent algorithms for function minimization. Aero-Astronautics Rep. No. 66, Rice University.
- Huard, P. (1964). Résolution de-programmes mathématiques à contraintes non-linéaires par la méthode des centres. Note E.D.F., HR5690/317.
- Huard, P. (1967). Resolution of mathematical programming with nonlinear constraints by the method of centres; in "Nonlinear Programming", ed. J. Adadie. North-Holland Pub. Co., Amsterdam.
- Huard, P. (1968). Programmation mathématique convexe. R.I.R.O., 7, pp 43-59.

- John, F. (1948). Extremum problems with inequalities as subsidiary conditions: in "Studies and Essays". Courant Anniversary Volume, Interscience, pp 187-204.
- Kelley, H. J. (1962). Methods of gradients: in "Optimization Techniques", ed. G. Leitmann. Academic Press.
- Kelley, J. E., Jr. (1960). The cutting plane method for solving convex programs. J. Soc. Ind. Appl. Math., 8 (4), pp 703-712.
- Kowalik, J. (1966). Nonlinear programming procedures and design optimization. Acta Polytechnica Scandinavica, 13, Trondheim.
- Kowalik, J. and Osborne, M. R. (1968). "Methods for Unconstrained Problems". American Elsevier.
- Kowalik, J., Osborne, M. R. and Ryan, D. M. (1969). A new method for constrained optimization problems. Op. Res., 17 (6), pp 973-983.
- Kuhn, H. W. and Tucker, A. W. (1951). Non-linear Programming: in "Proceedings of the Second Berkeley Symposium on Mathematical Statistics and Probability", ed. J. Neyman. University of California Press, pp 481-493.
- Kunzi, H. P. and Oettli, W. (1969). Nichtlineare Optimierung, Lecture Notes in Operations Research and Mathematical Systems 16, Springer-Verlag.
- Lemke, C.E. (1962). A method for solution of quadratic programming problems. Man. Sci., 8, pp 442-453.
- Lootsma, F. A. (1967). Logarithmic programming: a method of solving nonlinear-programming problems. Philips Res. Rep., 22 (3), pp 329-344.
- Lootsma, F. A. (1968a). Extrapolation in logarithmic programming. Philips Res. Res. Rep., 23, pp 108-116.
- Lootsma, F. A. (1968b). Constrained optimization via penalty functions. Philips Res. Rep., 23, pp 408-423.
- Lootsma, F. A. (1969). Hessian matrices of penalty functions for solving constrained-optimization problems. Philips Res. Resp., 24, pp 322-331.

- Lootsma, F. A. (1970). Boundary properties of penalty functions for constrained minimization. Philips Res. Resp. Supp. No. 3.
- Lootsma, F. A. (editor), (1972). Numerical. Methods for Nonlinear Optimization, Academic Press.
- Luenburger, D. G. (1969). "Optimization by Vector Space Methods". Wiley.
- Mangasarian, O. L. (1969). Nonlinear Programming. McGraw-Hill.
- Mangasarian, O. L. and Fromowitz, S. (1967). The Fritz John necessary optimality conditions in the presence of equality and inequality constraints. J. Math. Anal. Appl., 17 (1), pp 37-47.
- McCormick, G. P. (1969). The rate of convergence of the reset Davidon variable metric method. MRC Tech. Rep. No. 1012, University of Wisconsin.
- McCormick, G. P. and Pearson, J. D. (1969). Variable metric methods and unconstrained optimization: in "Optimization", ed. R. Fletcher. Academic Press.
- Meyers, G. E. (1968). Properties of the conjugate-gradient and Davidon methods. J. Opt. Th. Appl., 2 (4), pp 209-219.
- Morrison, D. D. (1968). Optimization by least squares. SIAM J. Numer. Anal., 5 (1), pp 83-88.
- Murray, W. (1969a). Indefinite quadratic programming. Nat. Phys. Lab. Rep. Ma 76.
- Murray, W. (1969b). Behaviour of Hessian matrices of barrier and penalty functions arising in optimization. Nat. Phys. Lab. Rep. Ma 77.
- Murray, W. (1969c). Constrained Optimization. Nat. Phys, Lab. Rep. Ma 79.
- Murtagh, B. A. and Sargent, R. W. H. (1969). A constrained minimization method with quadratic convergence; in "Optimization, ed. R. Fletcher, Academic Press.
- Murtagh, B. A. and Sargent, R. W. H. (1970). Computational experience with quadratically convergent minimization methods. Comp. J., 13 (2), pp 185-194.

- Nelder, J. A. and Mead, R. (1965). A simplex method for function minimization. *Comp. J.*, **7**, pp 308-313.
- Osborne, M. R. and Ryan, D. M. (1970a). On penalty function methods for nonlinear programming problems. *J. Math. Anal. Appl.*, **31** (3), pp 559-578.
- Osborne, M. R. and Ryan, D. M. (1970b). An algorithm for nonlinear programming. Tech. Rep. No. 35, Computer Centre, A.N.U.
- Osborne, M. R. and Ryan, D. M. (1972). A hybrid method for nonlinear programming, in *Numerical Methods for Nonlinear Optimization* (ed. F. A. Lootsma), Academic Press.
- Ostrowski, A. M. (1966). *Solution of Equations and Systems of Equations* (2nd Edition), Academic Press.
- Parisot, G. R. (1961). Resolution numérique approchée du problème de programmation linéaire par application de la programmation logarithmique. *Rev. Fr. Recherche Opérationnelle*, **20**, pp 227-259.
- Pearson, J. D. (1969). On variable metric methods of minimization. *Comp. J.*, **12**, pp 171-178.
- Pietrzykowski, T. (1962). Application of the steepest descent method to concave programming; in "Proceedings of IFIPS Congress", Munich, North-Holland Pub. Co., pp 185-189.
- Polak, E. and Ribière, G. (1969). Note sur la convergence de méthodes de directions conjuguées. *R.I.R.O.*, **3**, pp 35-43.
- Pomontale, T. (1965). A new method for solving conditioned maxima problems. *J. Math. Anal. Appl.*, **10**, pp 216-220.
- Powell, M. J. D. (1967). A method for nonlinear constraints in minimization problems. AERE Tech. Rep. No. 310.
- Powell, M. J. D. (1968). A survey of numerical methods for unconstrained optimization. AERE Tech. Rep. No. 340.
- Powell, M. J. D. (1969a). Rank one methods for unconstrained optimization. AERE Tech. Rep. No. 372.



- Powell, M. J. D. (1969b). On the convergence of the variable metric algorithm. AERE Tech. Rep. No. **382**.
- Powell, M. J. D. (1970). A new algorithm for unconstrained optimization. AERE Tech. Rep. No. 393.
- Rosen, E. M. (1966). A review of quasi-Newton methods in nonlinear equation solving and unconstrained optimization; in "Proceedings 21st National Conference of the A.C.M.". Thompson Book Co., pp 37-41.
- Rosen, J. B. (1960). The gradient projection method for nonlinear programming, part 1: linear constraints. J. Soc. Ind. Appl. Math., **8** (1), pp 181-217.
- Rosen, J. B. and Suzuki, S. (1965). Construction of nonlinear programming test problems. Comm. A.C.M., **8**, p 113.
- Rosenbrock, H. H. (1960). An automatic method for finding the greatest or least value of a function. Comp. J., 3, pp 175-184.
- Rudin, W. (1953). "Principles of Mathematical Analysis". McGraw-Hill, N. Y.
- Schmit, L. A. and Fox R. L. (1965). Integrated approach to structural synthesis, AIAA J., 3 (6), pp 1104-1112.
- Shah, B. V., Buehler, R. J. and Kempthorne, O. (1964). Some algorithms for minimizing a function of several variables. J. Soc. Ind. Appl. Math., **12**, pp 74-92.
- Spang, H. A. (1962). A review of minimization techniques for nonlinear functions. SIAM Review, 4 (4), pp 343-365.
- Stoer, J. (1971). On the numerical solution of constrained least squares problems, SIAM J. Numerical Analysis, **8**, pp 382-411.
- Strong, R. E. (1965). A note on the sequential unconstrained minimization technique for non-linear programming. Man. Sci., 12 (1), pp 142-144.
- Tremolières, R. (1968). La méthode des centre á troncature variable. Thkse, Paris.

- Whittle, P. (1971). "Optimization under Constraints", Wiley.'
- Wolfe, P. (1959). The simplex method for quadratic programming.  
Econometrica, **27** (3), pp 382-398.
- Wolfe, P. (1961). A duality theorem for nonlinear programming.  
Quart. Appl. Math., **19** (3), pp 239-244.
- Zangwill, W. I. (1967). Nonlinear programming via penalty functions.  
Man. Sci., **13** (5), pp 344-358.
- Zangwill, W. I. (1969). "Nonlinear Programming: A Unified Approach".  
Prentice-Hall.
- Zoutendijk, G. (1960). "Methods of Feasible Directions". Elsevier.