# CONSTRUCTIVE APPROACHES TO APPROXIMATE SOLUTIONS OF OPERATOR EQUATIONS AND CONVEX PROGRAMMING

by

Henry Wolkowicz

A thesis submitted to the Faculty of Graduate Studies and Research, McGill University, in partial fulfillment of the requirements for the Degree of Doctor of Philosophy.

> Department of Mathematics McGill University Montreal, Quebec August, 1978.

To my mother

and

the memory of my father.

#### ABSTRACT

i

The research in this thesis lies in two related areas of applied mathematics: approximation and optimization. In the area of approximation, new classes of iterative methods are introduced to calculate best approximate solutions of operator equations in Banach spaces. Also, Kantorovich's approximation theory is extended to include, possibly inconsistent, operator equations. As a special case, the convergence of a Galerkin type method is established. In the second part of the thesis, the geometry of optimality conditions for nondifferentiable convex optimization problems is studied. Necessary and sufficient conditions, under which the Kuhn-Tucker theory is valid, are stated. The results are used to formulate a numerical algorithm and to calculate various objects which have recently appeared in the theory of optimization.

# RÉSUMÉ

La recherche contenue dans cette thèse porte sur deux domaines voisins des mathématiques appliquées: l'approximation et l'optimisation. En théorie de l'approximation, on présente une nouvelle classe des méthodes itératives pour calculer les meilleures solutions approximatives d'équations d'opérateurs dans des espaces de Banach. On généralise aussie, la théorie d'approximation de Kantorovich au cas des équations d'opérateurs possiblement incohérentes. Dans ce contexte la preuve de convergence d'une méthode de type Galerkin apparaît comme un cas particulier. Dans la deuxième partie, on étudie la géométrie des conditions d'optimalité pour des problèmes nondifférentiables d'optimisation convexe. On énonce égalément des conditions nécessaires et suffisantes pour que la théorie de Kuhn-Tucker soit valide. Ces résultats servent à formuler un algorithme numérique pour calculer certaines expressions qui on fait leur apparition récemment en théorie de l'optimisation.

# ACKNOWLEDGEMENT

I would like to express my thanks to my supervisor, Professor S. Zlobec, whose personal interest in my work led to my studying mathematics. I would also like to thank Professor Zlobec for the help and advice received throughout my graduate career. His imaginative ideas have been the main influence in my thesis and in my research in general.

I would like to thank Professor G.P.H. Styan and my wife, Gail, for their continuous encouragement and support. In particular, I am grateful to Professor Styan for many enjoyable and helpful conversations and to my wife, Gail, for help in both discussing the contents and proofreading this manuscript.

Finally, I would like to thank Marlene Pyykko and Suzanna Clarkson for their careful and expert typing and Louis Paul Rivest for translating the abstract into French.

This research was supported by the National Research Council of Canada, le Gouvernement du Québec and McGill University.

iii

Abstract	i
Résumé	ii
Acknowledgement	iii

PART A. OPERATOR EQUATIONS	
Chapter I. INTRODUCTION	1
Chapter II. PRELIMINARIES	
1. Notation	6
2. Approximate and Best Approximate Solutions	7
Chapter III. ITERATIVE METHODS	
1. Stationary Methods	12
2. Nonstationary Methods	22
3. Examples	24
Chapter IV. KANTOROVICH'S THEORY FOR SINGULAR EQUATIONS	
1. Condition for Consistency of the Exact Equation	30
2. Solvability of the Approximate Equation	37
3. Error Estimates	45
4. Convergence Criteria	49
Chapter V. GALERKIN'S METHOD FOR BEST APPROXIMATE SOLUTIONS	
1. Description of the Method	57
2. An Example	65

68

References

C

TANI D. CONVENTINOURIE	PART	Β.	CONVEX	PROGRAMMING
------------------------	------	----	--------	-------------

Chapter	I. INTRODUCTION	75
Chapter	II. PRELIMINARIES	
1.	Notation	77
2.	The Convex Program	80
3.	Cones of Directions and Faithfully Convex Functions	81
4.	Subdifferentiability .	85
5.	The 'Badly Behaved' Constraints	87
6.	Polar Sets and Closedness Criteria	93
7.	Some Special Cones	98
8.	The Cone of Tangents	101
Chapter	r III. CHARACTERIZATIONS OF OPTIMALITY	
1.	Introduction	102
2.	A Basic Lemma	103
3.	Gould and Tolle Optimality Criteria	111
4.	Some Choices for the Cone G	114
Chapter	r IV. CONSTRAINT QUALIFICATIONS AND REGULARIZATION	TECHNIQUES
1.	Kuhn-Tucker Points	120
2.	Regular Points and Slater's Condition	120
3.	Weakest Constraint Qualifications	122
4.	Regularization	126
5.	Strongest and Weakest Optimality Conditions	128
Chapter	r V. THE METHOD OF REDUCTION	
1.	Introduction	134

C

2.	Calculating the Cone of Directions of Constancy	136
3.	Calculating the Sets $P^{=}$ and $D_{p}^{=}$	144
4.	The Method of Reduction	159
5.	Applications	169
6.	Examples	175
Appendi	x. COMPUTER PROGRAM	183
Referen	ces	198

.

.

C

PART A

 $\bigcirc$ 

C

# OPERATOR EQUATIONS

## I. INTRODUCTION

1

1.

A solution of a consistent operator equation

(1.1) 
$$Ax = b$$
,

where A is a bounded linear operator from a Banach space X into itself and b is an element of X, can be calculated in two ways. One can use a simple iterative scheme set up in X, e.g. Krasnosel'skii et al. [33, Chapter 1], or an extension to Banach spaces of various well-known matrix iterative schemes, as suggested by e.g. Petryshyn [47], [49], Kammerer and Plemmons [29], and Gudder and Neumann [22]. The other way is to approximate the original equation (1.1) by a sequence of equations

$$\bar{Ax} = \bar{b},$$

which are possibly easier to handle, and use appropriate error analysis. The latter approach is generally more successful. One of the first theories which studies the relationship between (1.1) and (1.2) was given by Kantorovich [30] and elaborated in the book by Kantorovich and Akilov [31]. Kantorovich's theory has been developed only for consistent equations. In particular, it is concerned with the following problems:

(i) Find conditions under which the consistency of (1.1)

implies the consistency of (1.2).

(ii) If both (1.1) and (1.2) are consistent, estimate the distance between their solutions. 2

- (iii) Find conditions under which the solutions of a sequence of approximate equations (1.2) converge to the solution of the equation (1.1).
- (iv) Estimate the norm of A in terms of the norm of  $\overline{A}$  and vice versa.

Kantorovich's approximation theory is rather general and, therefore, it is in principle applicable in many consistent situations, including the study and numerical treatment of infinite systems of linear equations, integral equations, ordinary differential equations and boundary value problems.

Various approximation theories have been recently developed and applied to particular problems by different authors, many of whom use the Kantorovich theory as a starting point. For instance, Thomas [59] refines some of Kantorovich's ideas and applies them to develop an approximation theory for the Nyström method of solving integral equations. Phillips [50] and Prenter [51] formulate approximation theories for the collocation method, while Ikebe [26] works with the Galerkin method. (For more details on these and other approaches for solving integral equations, see e.g. Houstis and Papatheodorou [25], Delves and Walsh [19] and Atkinson [7].)

Anselone [4] and Anselone and Moore [5] use the notion of collectively compact operators to formulate a different error analysis. Moore and Nashed [40] further developed the ideas of Anselone and Moore for possibly inconsistent operator equations in Banach spaces. They use the notions of generalized inverses of linear operators on Banach spaces and "best approximate" solutions of linear operator equations. Furthermore, they get, in special cases, some results in the perturbation theory of rectangular matrices obtained earlier by Ben-Israel [10] and Stewart [57].

3

An approximation theory for general, possibly inconsistent, linear equations in Hilbert spaces has been studied using the classical approach of Kantorovich (rather than the one of Moore and Nashed) by Zlobec [65]. One of the objectives of Part A is to continue the latter approach and formulate Kantorovich's theory for general, possibly inconsistent, linear equations in Banach spaces. The basic idea here is to establish and explore a relationship between approximate solutions of (1.1) and (1.2) and then use this relationship as a source for formulating various specific schemes for calculating approximate solutions of (1.1).

In the iterative computation of approximate solutions, as well as in Kantorovich's theory for singular equations, we will often use the concept of the generalized inverse of an operator. Some basic results on generalized inverses in Banach spaces are

I.1

summarized in Chapter II. In Chapter III, stationary and nonstationary iterative schemes are set up in Banach spaces for calculating both an approximate solution and the generalized inverse. This section extends from Hilbert to Banach spaces some results from the book by Ben-Israel and Greville [12, Chapter 8]. In Chapter IV, conditions for the consistency of Ax = y, for every y in a given subspace, are stated in terms of an approximate equation. Various error estimates are obtained as special cases. Kantorovich's theory for general linear equations is also formulated. The results are formulated in such a way that a comparison with the corresponding results for the nonsingular case from [31] is easily made. The most important results in this chapter are Theorem IV.3.1, which gives an error estimate, and Theorem IV.4.1, which gives conditions for convergence of approximate schemes. Using Kantorovich's theory, in Chapter V, a Galerkin-type method for calculating the best approximate solution is stated and its convergence is established for a class of operator equations in Banach spaces.

Situations where inconsistent linear operator equations arise are numerous and they include: integral equations in the theory of elasticity, potential theory and hydromechanics, e.g. Muskhelishvili [42], the integral formulation of the interior Neumann problem for the Laplacian, e.g. Kammerer and Nashed [27] and Atkinson [7], the eigenvalue problem in the case of a nonhomogeneous integral equation

4

when the associated homogeneous equation has a nontrivial solution, e.g. Kammerer and Nashed [27], and boundary value problems, e.g. Langford [34], Varga [62], Bramble and Shatz [16-18] and Serbin [54]. They also appear in the numerical solution of differential equations, for instance in the collocation method when the number of collocation points is bigger than the number of coefficients to be determined.

points is bigger than the number of coefficients to be determined, e.g. Krasnosel'skii et al. [33] and in the numerical solution of nonlinear equations where the Fréchet derivative is singular, e.g. Boggs [15] and Gay [20]. If the number of collocation points is smaller than the number of coefficients, then, if consistent, the approximate equation (1.2) has infinitely many solutions and one may again be interested in calculating the best approximate solution. Under- and over-determined initial value problems have been studied by Lovass-Nagy and Powers [38]. In the finite dimensional case, the under- and over-determined systems appear frequently in statistics, e.g. Rao and Mitra [52], see also Ben-Israel and Greville [12], Abdelmalek [1], [2] and Anderson [3].

5

#### II. PRELIMINARIES

1. Notation X, Y, <del>X</del>, <del>Y</del> real or complex Banach spaces  $\ell(X,Y)$ the set of all linear operators from X into Y  $\ell_{h}(X,Y)$ the set of all bounded linear operators from X into Y  $\ell(X)$  and  $\ell_{h}(X)$ the sets  $\ell(X,X)$  and  $\ell_h(X,X)$ , respectively 21 the Banach space of absolutely convergent sequences A the operator norm of A A restricted to the set S <sup>A</sup>|s σ(A) the spectrum of A ρ(A) the spectral radius of A R(A)the range space of A N(A) the null space of A A\* the adjoint of A.

For the above notions and their properties, see e.g. Taylor [58].

 $R{A,B} = \{Z \in \ell_b(X,Y): Z = AUB \text{ for some } U \in \ell(X,Y)\}, \text{ e.g.}$ 

Ben-Israel [9]

Mthe closure of the set MP<br/>Mthe projection onto the set M, II.2

 $M \oplus N$  the direct sum of M and N, II.2

M<sup>C</sup> the topological complement of M, II.2

II.1

6

x*	the best approximate solution, II.2
у*	an approximate solution, II.2
D(A)	the domain of the operator $\ensuremath{A}$ , II.2
A <sup>+</sup>	the generalized inverse of A, II.2

## 2. Approximate and Best Approximate Solutions

In order to formulate iterative methods for calculating approximate solutions and develop Kantorovich's theory for general, possibly inconsistent, operator equations in Banach spaces, we employ the following notions.

A linear operator  $P \in l(X)$ , is called a <u>projection</u> (of X) if  $P^2 = P$ . If R(P) = M, then we denote P by  $P_M$  and call it the <u>projection of X onto M</u>. Every projection  $P_M$  decomposes X into two <u>algebraic complements</u>,  $M = R(P_M)$  and  $N = R(I - P_M)$ . This implies X = M + N and we write  $N = M^C$ . If M and N are both closed, then we say that M has a <u>topological complement</u> in X and write

$$(2.1) X = M \oplus N.$$

For an example of decomposition (2.1), the reader is referred to Nashed's paper [43, p.327]. Recall that a closed subspace M of

X has a topological complement if and only if there exists a continuous projection  $P_M$  (of X), e.g. Taylor [58, p.241]. However, not every closed subspace has a topological complement, as shown by Murray [41] in 1937.

Consider  $A \in l_b(X,Y)$ . We shall assume that there exist continuous projections,  $P_{N(A)} \in l_b(X)$  and  $P_{R(A)} \in l_b(Y)$ . (In particular, such an A must have a closed range.)  $P_{N(A)}$  determines the complement  $N(A)^c = (I - P_{N(A)})X$ . Similarly,  $P_{R(A)}$ determines the complement  $R(A)^c = (I - P_{R(A)})Y$ . Hence, X = $N(A) \oplus N(A)^c$  and  $Y = R(A) \oplus R(A)^c$ . When  $A \in l_b(X,Y)$  and projections  $P_{N(A)} \in l_b(X)$  and  $P_{R(A)} \in l_b(Y)$  are given, then the system

$$(2.2) A^{\dagger} A = P_{N(A)} c$$

$$(2.3) AA^+ = P_{\mathcal{R}(A)}$$

(2.4) 
$$A^{+}P_{R(A)} = A^{+}$$

always has a unique solution  $A^+ \in \ell_b(Y,X)$ , called the <u>generalized</u> <u>inverse</u> of A (relative to the projections  $P_{N(A)}$  and  $P_{R(A)}$ ). The operator  $A^+$  then establishes a one-to-one correspondence between R(A) and  $N(A)^c$ , i.e.  $A^+|_{R(A)} = (A|_{N(A)}c)^{-1}$ , e.g. Nashed [43], Kammerer and Plemmons [29]. Note that, by the closed graph theorem,  $A^+$  is bounded when R(A) is closed.

If z is any vector in N(A), then  $y^* = A^+b + z$  is called an <u>approximate solution of the equation</u> Ax = b (relative to  $P_{R(A)}$ ), e.g. Gudder and Neumann [22], while  $x^* = A^+b$  is called the <u>best approximate solution of the equation</u> Ax = b (relative to  $P_{N(A)}$  and  $P_{R(A)}$ ), e.g. Moore and Nashed [40]. We see that  $\{A^+b + z : z \in N(A)\}$  is the set of all solutions of the projectional equation  $Ax = P_{R(A)}b$ , while  $A^+b$  is the unique one which lies in  $N(A)^{c}$ .

<u>Remark 2.1.</u> If R(A) is not closed but the complement  $M = \overline{R(A)}^{C}$  exists, then A has a unique unbounded generalized inverse  $A^{+} \in l(\mathcal{D}(A^{+}), X)$  relative to the complements M and  $N(A)^{C}$ , where  $\mathcal{D}(A^{+}) = R(A) + M$ . Many of the results herein can be extended to include this case as well as the converse case when A is a densely defined unbounded operator with closed range, see e.g. the approach in Nashed [44].

<u>Remark 2.2.</u> The term "best approximate" solution is used by Newman and Odell [46] under different circumstances. There  $\hat{x}$ is a "best approximate" solution of Ax = b, where  $A:X \rightarrow Y$ ,  $b \in Y$  if, for every  $x \in X$  with  $x \neq \hat{x}$ , either

$$\|A\hat{\mathbf{x}} - \mathbf{b}\| < \|A\mathbf{x} - \mathbf{b}\|$$

 $\mathbf{or}$ 

$$\|A\hat{x} - b\| = \|Ax - b\|$$
 and  $\|\hat{x}\| < \|x\|$ .

(This corresponds to the notion of <u>best least squares solution</u> in the case of Hilbert spaces.) In order to avoid possible ambiguity, we shall refer to the above  $\hat{x}$  as the "X, Y-best approximate"

9

solution of the equation Ax = b. If the norms on X and Y are strictly convex, then an "X, Y-best approximate" solution exists. If they are not strictly convex, then an "X, Y-best approximate" solution may not exist. In order to find  $\hat{x}$ , we need the notion of an X-projection (also called a "metric projection" by Blather, Morris and Wulbert in [14]). Suppose that S is a subspace of X. Then the mapping  $E_c$  is the X-projection onto S if, for every  $x \in X$ ,  $E_{S}x$  solves the minimization problem min ||x - y||. In general, the mapping  $E_{c}$  is not linear. An y∈S instance in which  $E_{S}$  is linear is when S and S<sup>C</sup> have a basis and the norm in X is a "TK norm" with respect to these bases, e.g. Singer [55]. In Hilbert spaces, E<sub>S</sub> corresponds to the orthogonal projection  $P_{S}$ . When the "X, Y-best approximate" solution  $\hat{x}$  exists, then  $\hat{x} = Bb$ , where  $B = (I - E_{N(A)})A^{\dagger}E_{R(A)}$ and  $A^+$  is any generalized inverse of A with respect to some  $P_{N(A)}c$  and  $P_{R(A)}$ . (Note that B need not be linear.) Thus, we see that when  $E_{N(A)}$  and  $E_{R(A)}$  are linear, one may choose  $P_{N(A)C} = I - E_{N(A)}$  and  $P_{R(A)} = E_{R(A)}$  in which case the "X, Ybest approximate" solution  $\hat{x} = Bb$  coincides with the "best approximate" solution  $x^* = A^+b$ .

Suppose that Y (but not necessarily X) is a Hilbert space, A  $\in l_b(X,Y)$ , A has closed range and X = N(A)  $\oplus N(A)^C$ . Then one may choose  $P_{R(A)} = E_{R(A)}$ , which is now the orthogonal projection on R(A), i.e.  $R(A)^C = R(A)^{\perp}$ , and write Y =  $R(A) \oplus R(A)^C$ . If  $A^+$  is the generalized inverse of A with respect to  $P_{N(A)}c$ and  $P_{R(A)}$ , the best approximate solution  $x^* = A^+b$  is the unique <u>least squares solution</u> of Ax = b in  $N(A)^{c}$ , i.e.  $x^*$ solves the problem

$$(2.5) \qquad \min_{\mathbf{x} \in \mathbf{X}} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|$$

and among all solutions of (2.5), it is the only one in  $N(A)^{c}$ , e.g. Kammerer and Plemmons [29]. The set of all least squares solutions corresponds to the set of all approximate solutions of Ax = b with respect to  $R(A)^{c} = R(A)^{\perp}$ . If both X and Y are Hilbert spaces and  $A \in \ell_{b}(X,Y)$  has closed range, we may choose  $P_{R(A)} = E_{R(A)}$  and  $P_{N(A)c} = I - E_{N(A)}$ . These are now the orthogonal projections, i.e.  $N(A)^{c} = R(A^{*})$ ,  $R(A)^{c} = N(A^{*})$ and  $\hat{x} = Bb = A^{+}b = x^{*}$  is the <u>best least squares solution</u> of the equation Ax = b. This means that  $x^{*}$  is the only solution in  $N(A)^{c}$  of the minimization problem (2.5) and among all solutions of (2.5) it is the unique one of smallest norm. For a detailed discussion of the generalized inverse and best least squares solution in Hilbert spaces, the reader is referred to the book by Ben-Israel and Greville [12].

11

#### III. ITERATIVE METHODS

#### 1. Stationary Methods

In order to calculate an approximate solution y\* of the operator equation (I.1.1)

Ax = b,

where  $A \in \ell_b(X,Y)$ ,  $P_{N(A)} c \in \ell_b(X)$  and  $P_{R(A)} \epsilon \ell_b(Y)$ , one can use the following iterative scheme:

(1.1) 
$$x_{k+1} = x_k - BAx_k + Bb$$
,  $k = 0, 1, ...,$ 

where

$$B \in R\{P_{N(A)}^{c}, P_{R(A)}\}$$
.

This scheme has been suggested for calculating the best least squares solution in Hilbert spaces in [65], see also [12, p.356].

<u>Theorem 1.1.</u> Let  $A \in l_b(X,Y)$ ,  $b \in Y$ ,  $P_{N(A)^C} \in l_b(X)$ ,  $P_{R(A)} \in l_b(Y)$  and  $B \in R\{P_{N(A)^C}, P_{R(A)}\}$  be given. Then the sequence  $\{x_k\}$ , generated by (1.1), converges, for any  $x_0 \in X$ , to the approximate solution  $y^* = x^* + P_{N(A)}x_0$  of Ax = b, if  $\rho(P_{N(A)^C} - BA) < 1$ . Moreover, if  $x_0 \in N(A)^C$ , then  $y^* = x^*$ , i.e. we obtain the best approximate solution.

Proof. We find that

$$x_{k+1} - x^* = (I - BA)x_k + Bb - x^*$$
, by (1.1)

$$= (I - BA) x_{k} + BP_{R(A)} b - x^{*}, \text{ since } B \in R\{P_{N(A)} c, P_{R(A)}\}$$

$$= (I - BA) (x_{k} - x^{*}), \text{ since } P_{R(A)} b = Ax^{*}$$

$$= (I - BA)^{k+1} (x_{0} - x^{*}), \text{ by iteration}$$

$$= (I - BA)^{k+1} (P_{N(A)} x_{0} + P_{N(A)} c x_{0} - x^{*})$$

$$(1.2) = (P_{N(A)} c - BA)^{k+1} (x_{0} - x^{*}) + P_{N(A)} x_{0}, \text{ since } x^{*} \text{ is in } N(A)^{c}$$

But  $\rho(P_{N(A)}c - BA) = \lim_{n \to \infty} \sup \|(P_{N(A)}c - BA)^n\|^{1/n} < 1$ , by the property of  $\rho$  (e.g. Taylor [58]) and the assumption. Therefore, there exists a real number s and a positive integer  $n_0$  such that

$$\|(P_{N(A)}c - BA)^n\|^{1/n} \le s \le 1$$
, for all  $n \ge n_0$ .

Hence,  $\|(P_{N(A)}c - BA)^n\| \le s^n \ne 0$  as  $n \ne \infty$ . This implies  $(P_{N(A)}c - BA)^n \ne 0$  as  $n \ne \infty$ . Thus  $x_k$  converges to  $x^* + P_{N(A)}x_0$ , by (1.2).

<u>Remark 1.1.</u> Necessary conditions for convergence of  $x_k$  to  $x^* + P_{N(A)} x_0$ , for every  $x_0 \in X$ , are

$$\rho(P_{N(A)}c - BA) \le 1 \quad [not \quad \rho(P_{N(A)}c - BA) < 1 !]$$

and

 $P_{N(A)}c$  - BA has no eignevalue  $\lambda$  such that  $|\lambda| = 1$ .

<u>Proof.</u> When  $x_k \rightarrow x^* + P_{N(A)} x_0$  for any  $x_0 \in X$ , then

$$(P_{N(A)}c - BA)^{k}(x_{0} - x^{*}) \rightarrow 0 \text{ as } k \rightarrow \infty,$$

by (1.2). Hence, for every  $x \in X$ ,

$$\sup_{k\geq 1} \| (P_{N(A)}^{c} - BA)^{k} x \| < \infty.$$

Now, by the Banach-Steinhaus theorem, there exists M > 0 such that  $\|(P_{N(A)}c - BA)^k\| \le M$ , k = 1, 2, ... But

$$\left[\rho(P_{N(A)}c - BA)\right]^{k} = \rho\left[\left(P_{N(A)}c - BA\right)^{k}\right]$$
, by the spectral mapping theorem

$$\leq \|(P_{N(A)} c - BA)^k\| \leq M.$$

Therefore,  $\rho(P_{N(A)}c - BA) \leq 1$ . In order to prove the second necessary condition, we observe that if  $0 \neq x \in X$  and  $|\lambda| = 1$ , such that  $(P_{N(A)}c - BA)x = \lambda x$ , then  $x \in N(A)^{c}$  and  $(P_{N(A)}c - BA)^{k}x = \lambda^{k}x \neq 0$  as  $k \neq \infty$ , contradicting  $x_{k} \neq x^{*}$ , by (1.2).

Example 1.1. The above remark is demonstrated by the operator  $A \in \ell_b(\ell_1)$  defined on  $x = (x_i)$  by

 $(Ax)_{1} = 0$  and  $(Ax)_{i} = (1 - 2^{-i})x_{i}$ , i = 2, 3, ...

It is clear that

$$N(A) = \{x \in \ell_1 : x_i = 0, i = 2, 3, ...\}$$

and

$$R(A) = \{x \in \ell_1 : x_1 = 0\}.$$

If we choose  $N(A)^{c} = R(A)$ ,  $R(A)^{c} = N(A)$  and B = 2A, then  $P_{N(A)c}$  is defined by

$$(P_{N(A)}c^{x})_{1} = 0$$
 and  $(P_{N(A)}c^{x})_{i} = x_{i}$ ,  $i = 2, 3, ...,$ 

while  $BA = 2A^2$  is defined by

$$(BAx)_1 = 0$$
 and  $(BAx)_i = 2(1 - 2^{-i})^2 x_i$ ,  $i = 2, 3, ...$ 

Therefore

$$\rho(P_{N(A)}c - BA) = \sup_{i \ge 2} \{ |1 - 2(1 - 2^{-1})^{2} | \}$$
$$= \sup_{i \ge 2} \{ |-1 + 2^{2-i} - 2^{1-2i} | \}$$
$$= 1.$$

But, for every  $x \in X$ ,

.

$$((I - BA)^{k}x)_{i} = \begin{cases} x_{1}, & \text{if } i = 1 \\ (-1 + 2^{2-i} - 2^{1-2i})^{k}x_{i}, & \text{if } i = 2, 3, \dots \end{cases}$$

Now, let  $\varepsilon > 0$  be given. Then, since

 $x \in \ell_1$  and  $|-1+2^{2-i}-2^{1-2i}| < 1$  for i = 2, 3, ...,

there exists integers  $N \ge 2$  and K such that

$$\sum_{i=N+1}^{\infty} |x_i| < \varepsilon \quad \text{and} \quad \sum_{i=2}^{N} |-1+2^{2-i}-2^{1-2i}|^k |x_i| < \varepsilon,$$

when  $k \ge K$ . Therefore

$$\|(I - BA)^{k}x - P_{N(A)}x\| = \sum_{i=2}^{\infty} |((I - BA)^{k}x)_{i}| < 2\varepsilon,$$

when  $k \ge K$ . This implies that

$$(I - BA)^{K} x \neq P_{N(A)} x$$
, for any  $x \in \ell_{1}$ ,

which, in turn, proves convergence of the iterative scheme, by (1.2).

<u>Remark 1.2.</u> It is a consequence of Remark 1.1 that Theorem 1.1 gives not only a sufficient but also a necessary condition for convergence, if X is finite dimensional. In the case of a Hilbert space X and a normal operator  $T = P_{N(A)}c - BA$ , one can show that  $T^{T}x \neq 0$ , for all  $x \in X$ , if and only if  $\rho(T) \leq 1$  and the spectrum of T has no mass on the unit circle |z| = 1. For, by the spectral theorem, e.g. Rudin [53],

$$T^{n}x = \int_{\sigma(T)} \lambda^{n} dE_{\lambda}x$$

And since, on  $\sigma(T)$ ,  $|\lambda^n| \le 1$  and  $\lambda^n \to 0$  almost everywhere, the Lebesgue dominated convergence theorem implies that  $T^n_x \to 0$ . This proves sufficiency. Now suppose that the spectrum of T has mass on the unit circle. By Remark 1.1, to prove necessity we need only find  $x \in X$  such that  $T^n x \neq 0$ . By the assumption, there exists a Borel subset  $\omega$  of the unit circle such that, if the subspace M is the range of the projection  $E_{\lambda}(\omega)$ , then  $M \neq \{0\}$ . But then T is invariant on M and  $\sigma(T_{|M}) = \omega$  is a subset of the unit circle. This implies that T is unitary on M and thus

$$T''x \neq 0$$
, for all  $0 \neq x \in M$ .

Specifying

$$P_{N(A)}c = P_{R(A^*)}$$
 and  $B = \alpha A^*$ 

one can show, see e.g. Petryshyn [47], that

$$\rho(\mathbf{T}) = \|\mathbf{P}_{\mathcal{R}(\mathbf{A}^*)} - \alpha \mathbf{A}^* \mathbf{A}\| \le 1$$

if and only if

$$0 \le \alpha \le \frac{2}{\|A^*A\|} .$$

One now establishes the following characterization of convergence: For any  $x_0 \in X$ , the sequence  $\{x_k\}$  converges to the least squares solution  $y^* = x^* + P_{N(A)} x_0$  if and only if

(i) 
$$0 < \alpha \le \frac{2}{\|A^*A\|}$$

and

(ii) ||A\*A|| is not an eigenvalue of A\*A if  $\alpha = \frac{2}{||A*A||}$ .

We recall that Petryshyn gives only the sufficient condition

$$0 < \alpha < \frac{2}{\|A^*A\|} .$$

Specifying B in (1.1), one obtains various iterative schemes for computing an approximate solution. In particular, if one splits A = M + N, chooses M<sup>+</sup> with respect to the continuous projections  $P_{R(M)}$  and  $P_{N(M)}c$  such that  $P_{R(M)} = P_{R(M)}P_{R(A)}$ and  $P_{N(M)}c = P_{N(A)}cP_{N(M)}c$  and specifies  $B = \omega M^{+}$ ,  $\omega \neq 0$ , then (1.1) becomes

(1.3) 
$$x_{k+1} = [(1 - \omega)I - \omega M^{\dagger}N]x_{k} + \omega M^{\dagger}b$$

Further, for  $\omega = 1$ , (1.3) becomes

(1.4) 
$$x_{k+1} = -M^{\dagger}Nx_{k} + M^{\dagger}b$$

If both A and M are invertible, then (1.3) and (1.4) become, respectively,

(1.5) 
$$x_{k+1} = [(1 - \omega)I - \omega M^{-1}N]x_k + \omega M^{-1}b$$

and

(1.6) 
$$x_{k+1} = -M^{-1}Nx_k + M^{-1}b$$
.

The scheme (1.5) has been studied by Pertryshyn [48], who calls it the "Extrapolated Jacobi Method". The scheme (1.6) is the wellknown Jacobi method. Other methods can be obtained by the splitting

A = D + S + Q with B = 
$$\left(\frac{1}{\omega}D + S\right)^+$$
,  $\omega \neq 0$  where  $\frac{1}{\omega}D + S \in \ell_b(X,Y)$   
and  $P_{R}\left(\frac{1}{\omega}D + S\right) = P_{R}\left(\frac{1}{\omega}D + S\right)^{P_{R}(A)}$ ,  $P_{N}\left(\frac{1}{\omega}D + S\right)^{C} = P_{N}(A)^{C}P_{N}\left(\frac{1}{\omega}D + S\right)^{C}$ .

Then (1.1) becomes

(1.7) 
$$x_{k+1} = (D + \omega S)^{+} [(1 - \omega)D - \omega Q]x_{k} + \omega (D + \omega S)^{+}b.$$

If both A and  $D + \omega S$  are invertible, the scheme (1.7) becomes

(1.8) 
$$x_{k+1} = (D + \omega S)^{-1} [(1 - \omega)D - \omega Q] x_k + \omega (D + \omega S)^{-1} b$$

which is known as the "Successive Over-Relaxation Method" (abbreviated SOR method). Specifying  $\omega = 1$  in (1.8) one obtains

(1.9) 
$$x_{k+1} = -(D+S)^{-1}Qx_k + (D+S)^{-1}b$$
,

which is known as the Gauss-Seidel method. In the case of an  $n \times n$ invertible matrix  $A = (a_{ij})$ , one frequently specifies

$$D = M = (a_{ii}), \quad i = 1,...,n$$

$$S = (a_{ij}), \quad i > j, \quad i = 1,...,n; \quad j = 1,...,n-1$$

$$Q = A - D - S,$$

$$N = A - M$$

in (1.5), (1.6), (1.8), and (1.9). Properties of these schemes for systems with invertible matrices (and linear operators) have been studied extensively, see e.g. Varga [60], [61] and Petryshyn [47], [48]. Scheme (1.4) has been studied by Berman and Plemmons [13] for systems with singular matrices and by Gudder and Neumann [22] for singular operator equations in Hilbert space. For some other schemes using splittings of A, see e.g. Hadjidimos [23] and Meijerink and van der Vorst [39].

<u>Remark 1.3.</u> One calls the splitting A = M + N a <u>proper</u> <u>splitting</u>, if R(A) = R(M),  $R(A)^{C} = R(M)^{C}$ , N(A) = N(M) and  $N(A)^{C} = N(M)^{C}$ . Note that to obtain a proper splitting in the case of Hilbert spaces and orthogonal complements, one need only check that N(A) = N(M) and R(A) = R(M), e.g. [13] and [22]. The proper splittings are not only useful in iterative calculation of least squares solutions, but they also play an important role in the Kantorovich approximation theory (see Chapter V).

One can slightly modify (1.1) in order to compute  $A^+$ , the generalized inverse of A relative to given projections,  $P_{N(A)}$  and  $P_{R(A)}$ .

<u>Theorem 1.2.</u> Let  $A \in \ell_b(X,Y)$ ,  $P_{N(A)} \in \ell_b(X)$  and  $P_{R(A)} \in \ell_b(Y)$ . If  $B \in R\{P_{N(A)}c, P_{R(A)}\}$ , then the sequence  $\{X_k\}$ , generated by

(1.10)  $X_{k+1} = X_k - BAX_k + B$ , k = 0, 1, 2, ...,

converges to  $A^{\dagger} + P_{N(A)}X_0$  for all  $X_0 \in \ell_b(Y,X)$ , if  $\rho(P_{N(A)}c - BA) < 1$ . Moreover, if  $R(X_0) \subset N(A)^c$ , then we obtain the generalized inverse  $A^{\dagger}$ . Proof. Here

$$X_{k+1} - A^{+} = X_{k} - BAX_{k} + B - A^{+}$$
  
= (I - BA) X<sub>k</sub> + BP<sub>R(A)</sub> - P<sub>N(A)</sub>c A<sup>+</sup>  
= (I - BA) (X<sub>k</sub> - A<sup>+</sup>) , by (II.2.3)  
= (P<sub>N(A)</sub>c - BA)<sup>k+1</sup> (X<sub>0</sub> - A<sup>+</sup>) + P<sub>N(A)</sub>X<sub>0</sub>.

The rest of the proof is analogous to the proof of Theorem 1.1  $\hfill \Box$ 

Note that, for  $b \in Y$ ,

$$(A^{+} + P_{N(A)}X_{0})b = A^{+}b + P_{N(A)}X_{0}b$$

is an approximate solution of Ax = b. Moreover, when  $X_0$  satisfies the condition

$$X_0 = X_0^P R(A)$$
,

we see that  $A^+ + P_{N(A)} X_0$  is the generalized inverse of A relative to (i) the same  $P_{R(A)}$  as  $A^+$ , and, by (II.2.2), (ii) the new projection on N(A),

$$P_{N(A)} + P_{N(A)} X_0^A$$
.

If (1.10) is modified as follows:

(1.11) 
$$X_{k+1} = X_k - BX_k + B$$
,  $k = 0, 1, 2, ...,$ 

III.2

and if one chooses

(1.12) 
$$B \in R\{P_{R(A)}, P_{R(A)}\}$$
 and  $X_0 \in \ell_b(Y)$ ,

such that  $R(X_0) \subset R(A)$ , then the sequence generated has the property

$$X_{k+1} - P_{\mathcal{R}(A)} = (P_{\mathcal{R}(A)} - B)^{k+1} (X_0 - P_{\mathcal{R}(A)}), \quad k = 0, 1, 2, \dots$$

Hence one concludes that whenever  $\rho(P_{R(A)} - B) < 1$ , the sequence  $\{X_k\}$ , generated by (1.11) and (1.12), converges to the projection  $P_{R(A)}$ . If X and Y are Hilbert spaces and  $P_{R(A)}$  is the orthogonal projection on R(A), then one can choose

$$X_0 = AZ_1$$
 and  $B = AZ_2A^*$ ,

where  $Z_1 \in {}^{\ell}_{b}(Y,X)$  and  $Z_2 \in {}^{\ell}_{b}(X,X)$ . In particular, one can specify  $Z_1 = A^*$  and  $Z_2 = \alpha I$ , where  $\alpha$  is a real parameter with the property  $\rho(P_{R(A)} - \alpha AA^*) < 1$ .

## 2. Nonstationary Methods

We can further modify the scheme (1.1) by varying B at each step. Such methods include gradient methods of finding approximate solutions. <u>Theorem 2.1.</u> Let  $A \in \ell_b(X,Y)$ ,  $P_{N(A)} \in \ell_b(X)$  and  $P_{R(A)} \in \ell_b(Y)$ . If  $B_k \in R\{P_{N(A)}c, P_{R(A)}\}$  for k = 0, 1, 2, ..., then the sequence  $x_k$  generated by

(2.1) 
$$x_{k+1} = x_k - B_k A x_k + B_k b$$
,  $k = 0, 1, 2, ...$ 

converges, for any  $x_0 \in X$ , to the approximate solution  $y^* = x^* + P_{N(A)} x_0$  of Ax = b, if  $||P_{N(A)}c - B_kA|| \le 1 - \alpha_k$ for k = 0, 1, 2, ..., where  $0 \le \alpha_k \le 1$  and  $\sum \alpha_k = \infty$ . Moreover, if  $x_0 \in N(A)^c$ , then  $y^* = x^*$ , i.e. we obtain the best approximate solution.

Proof. As in the proof of Theorem 1.1, we see that

$$x_{k+1} - x^* = \begin{bmatrix} k \\ \Pi \\ i=0 \end{bmatrix} (P_{N(A)}c - B_iA) = (x_0 - x^*) + P_{N(A)}x_0.$$

But

$$\left\| \begin{array}{c} k \\ \Pi \\ i=0 \end{array} \left( P_{N(A)^{c}} - B_{i}^{A} \right) \right\| \leq \begin{array}{c} k \\ \Pi \\ i=0 \end{array} \left( 1 - \alpha_{k}^{A} \right) \neq 0$$

if and only if the series  $\sum_{k=1}^{\infty} \alpha_{k}$  diverges, e.g. Knopp [32, p.92].

Specifying  $B_k$  in (2.1) again leads to various iterative schemes for computing the best approximate solution. For example, if X and Y are Hilbert spaces and we choose  $B_k = \alpha_k A^*$ , for some scalars  $\alpha_k$ , k = 0, 1, 2, ..., then we can rewrite (2.1) as

(2.2) 
$$x_{k+1} = x_k - \alpha_k A^* (Ax_k - b)$$
.

Since  $2A^*(Ax - b)$  is the gradient of the function  $||Ax - b||^2$ , we see that (2.2) defines a gradient method for minimizing ||Ax - b||, with step size given by  $\alpha_k/2$ .

If we let  $\Gamma_k = A^*(Ax_k - b)$  and  $\alpha_k = \|\Gamma_k\|^2 \|A\Gamma_k\|^{-2}$ , then (2.2) becomes the method of steepest descent. For a discussion on gradient methods, see e.g. Nashed [44, p.380] and Kammerer and Nashed [28].

#### 3. Examples

The iterative scheme (1.1) can be used to calculate an approximate solution of the equation (I.1.1) in abstract spaces. However, in many situations, it is actually used to calculate an approximate solution of an approximate equation (I.1.2) which is frequently a more manageable finite system of linear algebraic equations. The first case will now be demonstrated.

Example 3.1. Let us calculate the best least squares solution  $x^*$  of the inconsistent equation Ax(s) = (I - K)x(s) = b(s), where

$$Kx(s) = \frac{2}{\pi} \int_0^{\pi} (\sin s \sin \xi + \frac{1}{2}\cos s \cos \xi) x(\xi) d\xi$$

and b(s) = s. The operator K is chosen from Stakgold [56]. The problem will be solved in  $X = Y = L_2[0,\pi]$  using the iterative scheme (1.1).

We choose  $B = \alpha A^*$ , where  $0 < \alpha < 2 = \frac{2}{\|A^*A\|}$ . This guarantees here that  $\rho(P_{R(A)} - \alpha AA^*) < 1$  and that the scheme (1.1) converges to  $x^*$  for every choice of  $x_0(s)$  in  $N(A)^C = R(A^*)$ . For  $x_0(s) = 0$  one finds:

$$\begin{aligned} x_1(s) &= \alpha(s - 2\sin s + \frac{2}{\pi}\cos s) \\ &= \alpha_1^{(1)}s - \alpha_2^{(1)}\sin s + \alpha_3^{(1)}\cos s , \\ \end{aligned}$$
where  $\alpha_1^{(1)} &= \alpha$ ,  $\alpha_2^{(1)} &= 2\alpha$  and  $\alpha_3^{(1)} &= \frac{2\alpha}{\pi}$ .  
 $x_2(s) &= [(1 - \alpha)\alpha_1^{(1)} + \alpha]s - [(1 - \alpha)\alpha_2^{(1)} + 2\alpha]\sin s \\ &+ [(1 - \frac{\alpha}{4})\alpha_3^{(1)} - \frac{3\alpha}{\pi}\alpha_1^{(1)} + \frac{2}{\pi}\alpha]\cos s \\ &= \alpha_1^{(2)}s - \alpha_2^{(2)}\sin s + \alpha_3^{(2)}\cos s , \end{aligned}$ 
where  $\alpha_1^{(2)} &= (1 - \alpha)\alpha_1^{(1)} + \alpha$ ,  $\alpha_2^{(2)} &= (1 - \alpha)\alpha_2^{(1)} + 2\alpha$ , and  $\alpha_3^{(2)} &= (1 - \frac{\alpha}{4})\alpha_3^{(1)} - \frac{3\alpha}{\pi}\alpha_1^{(1)} + \frac{2}{\pi}\alpha$ . In general, if

$$x_{k}(s) = \alpha_{1}^{(k)}s - \alpha_{2}^{(k)}sin s + \alpha_{3}^{(k)}cos s$$

then

$$\begin{aligned} x_{k+1}(s) &= [(1-\alpha)\alpha_1^{(k)} + \alpha]s - [(1-\alpha)\alpha_2^{(k)} + 2\alpha]\sin s \\ &+ [(1-\frac{\alpha}{4})\alpha_3^{(k)} - \frac{3\alpha}{\pi}\alpha_1^{(k)} + \frac{2}{\pi}\alpha]\cos s \\ &= \alpha_1^{(k+1)}s - \alpha_2^{(k+1)}\sin s + \alpha_3^{(k+1)}\cos s . \end{aligned}$$

Since the iterative schemes

$$\alpha_{1}^{(k+1)} = (1 - \alpha)\alpha_{1}^{(k)} + \alpha$$

$$\alpha_{2}^{(k+1)} = (1 - \alpha)\alpha_{2}^{(k)} + 2\alpha$$

$$\alpha_{3}^{(k+1)} = (1 - \frac{\alpha}{4})\alpha_{3}^{(k)} - \frac{3\alpha}{\pi}\alpha_{1}^{(k)} + \frac{2}{\pi}\alpha_{3}^{(k)}$$

are convergent themselves with the solutions  $\alpha_1 = 1$ ,  $\alpha_2 = 2$  and  $\alpha_3 = -\frac{4}{\pi}$ , respectively, one concludes that

$$x^*(s) = s - 2\sin s - \frac{4}{\pi}\cos s$$

is the best least squares solution.

Using (1.10) with B defined by  $Bx(s) = x(s) - \frac{2}{\pi} \int_0^{\pi} \sin s \sin \xi x(\xi) d\xi$ , one can show that  $X_k$  converges to  $A^+$ , which is here

$$A^{+}x(s) = x(s) - \frac{2}{\pi} \int_{0}^{\pi} (\sin s \sin \xi - \cos s \cos \xi) x(\xi) d\xi$$

The best least squares solution of the equation introduced in Example 3.1 will be calculated in Example V.2.1 via Kantorovich's approximation theory. We conclude this section by demonstrating how the iterative scheme (1.1) can be applied to matrices. Example 3.2. We now calculate the best least squares solution, using the iterative scheme (1.1), of the inconsistent system

$$x_{1} + 3x_{3} = 1$$

$$-x_{1} + x_{2} = 1$$

$$x_{1} - x_{2} = 1$$

$$x_{2} + x_{3} = 0$$

Here

$$A = \begin{pmatrix} 1 & 0 & 3 \\ -1 & 1 & 0 \\ 1 & -1 & 0 \\ 0 & 1 & 1 \end{pmatrix}, \quad b = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 0 \end{pmatrix}$$

Specifying

$$B = \alpha A^{t}$$
, where  $\alpha = \frac{2}{trace A^{t}A}$ 

one obtains

$$B = \frac{1}{8} \begin{pmatrix} 1 & -1 & 1 & 0 \\ 0 & 1 & -1 & 1 \\ 3 & 0 & 0 & 1 \end{pmatrix}$$

It is easy to verify that for the above choice of B and  $\alpha$ ,  $\rho(P_{R(A^{t})} - BA) < 1$ , whenever rank A > 1. One can start iterating from  $x^{0} = 0$ , in which case the following numerical results are obtained:

27
k	x1k	x2 <sup>k</sup>	x3 <sup>k</sup>		
1	0.2500000	0.0000000	0.2500000		
2	0.2500000	0.0625000	0.3125000		
3	0.2656250	0.0625000	0.3281230		
4	0.2656260	0.0664063	0.3320313		
5	0.2666016	0.0664063	0.3330018		
6	0.2666016	0.0666504	0.3332520		
7	0.2666626	0.0666504	0.3333130		
8	0.2666626	0.0666657	0.3333282		
9	0.2666664	0.0666657	0.3333321		
10	0.2666664	0.0666666	0.3333330		
11	0.2666667	0.0666666	0.3333333		

The eleventh approximation  $x^{11}$  gives the best least squares solution correct to six decimal places:

 $x_1^* = \frac{4}{15}$ ,  $x_2^* = \frac{1}{15}$ ,  $x_3^* = \frac{1}{3}$ .

Example 3.3. In this example we apply scheme (1.1) to solve the problem

where X is a two-dimensional Banach space of scalars

$$\mathbf{x} = \begin{pmatrix} \xi_1 \\ \xi_2 \end{pmatrix}$$
,  $\mathbf{A} = \begin{pmatrix} 2 & 1 \\ 0 & 0 \end{pmatrix}$  and  $\mathbf{b} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$ .

If we specify the norm

$$\|\mathbf{x}\| = |\xi_1| + |\xi_2| + \max\{|\xi_1|, |\xi_2|, |\xi_1 + \xi_2|\}$$

(which is a TK norm with respect to the basis  $\begin{pmatrix} 1 \\ 0 \end{pmatrix}$ ,  $\begin{pmatrix} 0 \\ 1 \end{pmatrix}$ ), then  $E_{R(A)}$  is linear and equal to

$$\begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} \ .$$

Thus, in order to solve the problem, we must choose  $P_{R(A)} = E_{R(A)}$ . Since  $N(A)^{C}$  can be arbitrarily chosen, let it be  $N(A)^{C} = R(A)$ . Hence

$$P_{N(A)}^{c} = \begin{pmatrix} 1 & \frac{1}{2} \\ 0 & 0 \end{pmatrix}$$

Further we choose

B = 
$${}^{1}_{4}P_{N(A)}c^{P}_{R(A)} = \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}$$
 and  $x_{0} = 0$ .

The sequence  $x_1 = \frac{1}{4} \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$ ,  $x_2 = \frac{3}{8} \begin{pmatrix} 1 \\ 0 \end{pmatrix}$ , ... converging to

•

 $x^* = \frac{1}{2} \begin{pmatrix} 1 \\ 0 \end{pmatrix}$  is obtained.

# IV. KANTOROVICH'S THEORY FOR SINGULAR EQUATIONS

## 1. Condition for the Consistency of the Exact Equation

One of the most useful results in the formulation of the classical Kantorovich theory is a lemma which gives a condition for the consistency of the exact equation (I.1.1), e.g. [31, p.543]. This lemma will now be extended so that it also applies to singular equations.

Lemma 1.1. Let  $V \in l_b(X,Y)$ , E a closed subspace of X and F any subspace of Y containing V(E). If there exists, for every  $y \in F$ , an  $\hat{x} \in E$  such that

(1.1)  $\|V\hat{\mathbf{x}} - \mathbf{y}\| \le q\|\mathbf{y}\|$  and  $\|\hat{\mathbf{x}}\| \le \alpha \|\mathbf{y}\|$ 

where q < 1 and  $\alpha$  are constants, then the equation

(1.2) Vx = y

has, for every  $y \in F$ , a solution  $x \in E$  satisfying

(1.3) 
$$\| \mathbf{x} \| \le \frac{\alpha}{1-q} \| \mathbf{y} \|$$
.

<u>Proof.</u> Similarly to the proof in [31], we will construct an exact solution of (1.2) by recursion. Take an arbitrary  $y \in F$ . Set  $y_1 = y$ . By hypothesis, an  $\hat{x}_1 \in E$  exists such that

(1.4)  $\|V\hat{x}_1 - y_1\| \le q \|y_1\|, \|\hat{x}_1\| \le \alpha \|y_1\|.$ 

IV.1

Denote

(1.5) 
$$y_2 = y_1 - V\hat{x}_1$$
.

Clearly  $y_2 \in F$ , since  $y_1 \in F$  and F is a subspace containing V(E). We now apply the condition (1.1) to  $y_2$ . This implies the existence of  $\hat{x}_2 \in E$  such that

$$\| V \hat{x}_2 - y_2 \| \le q \| y_2 \| = q \| y_1 - V \hat{x}_1 \|, \quad by (1.5)$$
  
$$\le q^2 \| y_1 \|, \quad by (1.4).$$

Also  $\|\hat{x}_2\| \le \alpha \|y_2\| \le \alpha q \|y_1\|$ . Continuing this process, sequences  $\{y_k\}$  and  $\{\hat{x}_k\}$  are obtained such that

(1.6) 
$$y_{k+1} = y_k - V\hat{x}_k, \quad k = 1, 2, ...$$

and

(1.7) 
$$\|y_k\| \le q^{k-1} \|y_1\|$$
,  $\|\hat{x}_k\| \le \alpha q^{k-1} \|y_1\|$ .

By iteration, (1.5) and (1.6) give

(1.8) 
$$y_{k+1} = y_1 - V(\hat{x}_1 + \hat{x}_2 + \dots + \hat{x}_k)$$
,  $k = 1, 2, \dots$ 

Using the second inequality in (1.7), and recalling that  $y_1 = y$ , one obtains

$$\left\|\sum_{k=1}^{\infty} \hat{\mathbf{x}}_{k}\right\| \leq \sum_{k=1}^{\infty} \|\hat{\mathbf{x}}_{k}\| \leq \alpha \sum_{k=1}^{\infty} q^{k-1} \|\mathbf{y}\|$$

Since q < 1, the series  $\sum_{k=1}^{\infty} \hat{x}_k$  is convergent. Hence  $x \stackrel{\Delta}{=} \sum_{k=1}^{\infty} \hat{x}_k$ 

IV.1

belongs to E, since E is a closed subspace. Furthermore

$$\|\mathbf{x}\| \le \alpha \sum_{k=1}^{\infty} q^{k-1} \|\mathbf{y}\| = \frac{\alpha}{1-q} \|\mathbf{y}\|.$$

So taking limits in (1.8) gives

$$\lim_{k \to \infty} y_k = \lim_{k \to \infty} [y_1 - V(\hat{x}_1 + \hat{x}_2 + \dots + \hat{x}_k)]$$

=  $y_1 - Vx$ , by the continuity of V.

Also, since q < 1, we see from (1.7) that  $\lim_{k \to \infty} y_k = 0$ . Thus Vx = y. We have shown that  $x \in E$  is a solution of (1.2) and it satisfies (1.3).

Lemma 1.1 has been proved in [31] in the special case when E = X and F = Y. The above result will be used in the next section in the approximation theory. However, Lemma 1.1 is of an independent interest and in the remainder of this section we will show how, using the lemma, one can establish some new and some well-known estimates related to the equation (I.1.1)

<u>Proposition 1.1.</u> Let  $M \in \mathfrak{L}_{b}(X,Y)$  and let N(M) and  $\mathcal{R}(M)$ have topological complements. Denote by  $M^{+}$  the generalized inverse of M with respect to these topological complements. Consider

(1.9) A = M + N

such that  $A \in l_b(X,Y)$  and  $\|NM^+\| < 1$ . If

$$(1.10) R(N) \subset R(M) ,$$

then the equation

(1.11) 
$$Ax = y$$

has, for every  $y \in R(M)$ , a solution  $x^* \in N(M)^c$ . Also

(1.12) 
$$\|\mathbf{x}^*\| \leq \frac{\|\mathbf{M}^+\|}{1 - \|\mathbf{NM}^+\|} \|\mathbf{y}\|$$

and R(A) = R(M). In addition to the above assumptions, if (1.13)  $N(M)^{C} \cap N(A) = \{0\}$ ,

then

$$\|A^{+}\| \leq \frac{\|M^{+}\|}{1 - \|NM^{+}\|} ,$$

where  $A^+$  denotes the generalized inverse of A with respect to  $P_{N(A)}c = P_{N(M)}c$  and  $P_{R(A)} = P_{R(M)}$ .

<u>Proof.</u> We will show that the assumptions of Lemma 1.1 are satisfied with V = A,  $E = N(M)^{C}$ , F = R(M),  $q = \|NM^{+}\|$  and  $\alpha = \|M^{+}\|$ . Choose an arbitrary  $y \in R(M)$  and let  $\hat{x} = M^{+}y$ . Then

$$\|A\hat{x} - y\| = \|(M + N)\hat{x} - y\|$$
  
=  $\|(M + N)\hat{x} - MM^{+}y\|$ , since  $y \in \mathcal{R}(M)$   
 $\leq \|NM^{+}\|\|y\|$ ,

and  $\|\hat{x}\| = \|M^{+}y\| \leq \|M^{+}\|\|y\|$ . Thus (1.1) holds. Now we apply Lemma 1.1 to the equation (1.11) to conclude that for every  $y \in R(M)$ , the equation Ax = y has a solution  $x^{*} \in N(M)^{C}$ satisfying (1.12). This implies that

$$(1.15) R(M) \subseteq R(A) .$$

However,

(1.16) 
$$R(A) \subset R(M)$$
, by (1.9) and (1.10).

Now (1.15) and (1.16) imply

(1.17) 
$$R(A) = R(M)$$
.

Since  $N(M)^{c}$  is isomorphic to R(M) via M, it is also isomorphic via A, by (1.17) and the assumption (1.13). Thus, one can choose  $N(A)^{c} = N(M)^{c}$ . This determines  $A^{+}$  with respect to the topological complements  $N(M)^{c}$  and R(M). Now, since  $x^{*} \in N(A)^{c} = N(M)^{c}$ , we have  $x^{*} = A^{+}y$ , and (1.14) follows from (1.12).

<u>Corollary 1.1.</u> Let  $H \in l_b(X)$  with ||H|| < 1 and  $P \in l_b(X)$ such that  $P^2 = P$ . If  $R(H) \subset R(P)$ , then the equation (P+H)x = yhas a solution  $x^* \in R(P)$  for each  $y \in R(P)$ . Also  $||x^*|| \le \frac{1}{1 - ||H||} ||y||$ , R(P+H) = R(P) and  $||(P+H)^*|| \le \frac{1}{1 - ||H||}$ , where  $(P+H)^*$  denotes the generalized inverse with respect to  $P_{N(P+H)}c = P$  and  $P_{R(P+H)} = P$ . <u>Proof.</u> We will show that the hypotheses of Proposition 1.1 are satisfied with M = P and N = H. Clearly,  $P \in l_b(X)$  has topological complements  $N(P)^C = R(P)$  and  $R(P)^C = N(P)$ . So  $P^+$  (=P) is the generalized inverse of P with respect to these complements. Also  $P + H \in l_b(X)$ , since ||H|| < 1. Take an arbitrary  $y \in R(P)$ . Then  $\hat{x} \stackrel{\Delta}{=} P^+ y$  (see the proof of Proposition 1.1, where M = P) is equal to y, i.e.  $\hat{x} = y$ , since  $P = P^+$ . Therefore, the assumption  $||HP^+|| < 1$ , in Proposition 1.1, can be replaced by ||H|| < 1. Also, the assumption (1.13), which reads here

(1.18) 
$$N(P)^{C} \cap N(P+H) = \{0\}$$

is satisfied. If (1.18) were not true, there would exist an  $x \neq 0$  such that both

$$x \in N(P)^{c} = R(P)$$
 and  $x \in N(P + H)$ 

Hence (P + H)x = x + Hx = 0, which contradicts the assumption ||H|| < 1.

<u>Corollary 1.2.</u> (Ben-Israel [10].) Let H be an  $n \times n$  real matrix, ||H|| < 1 and L be a subspace of  $\mathbb{R}^{\mathbb{M}}$  such that  $\mathcal{R}(H) \subset L$ . Then

$$\| (P_{L} + H)^{+} \| \leq \frac{1}{1 - \| H \|}.$$

<u>Proof.</u> Specify  $P = P_{I}$  in Corollary 1.1.

The following classical result of Banach, e.g. [31], is also obtained.

Corollary 1.3. Let  $H \in \ell_b(X)$ . If ||H|| < 1, then  $(I + H)^{-1}$  exists and

$$\| (I + H)^{-1} \| \leq \frac{1}{1 - \|H\|}.$$

Proof. Specify P = I in Corollary 1.1.

<u>Corollary 1.4.</u> (Kantorovich and Akilov [31, p.172.]) Let  $M \in \ell_b(X)$  and suppose that  $M^{-1} \in \ell_b(X)$  exists. If  $N \in \ell_b(X)$  and  $\|NM^{-1}\| \le 1$ , then  $(M+N)^{-1}$  exists and

$$\| (M + N)^{-1} \| \leq \frac{\| M^{-1} \|}{1 - \| N M^{-1} \|}.$$

<u>Proof.</u> Apply Proposition 1.1 to the case when X = Y and  $M^{-1}$  exists.

## 2. Solvability of the Approximate Equation

Let  $\tilde{X}$  and  $\tilde{Y}$  be closed subspaces of the Banach spaces X and Y, respectively. Further, let  $\tilde{X}$  and  $\tilde{Y}$  be isomorphic via mappings  $J_0$  and  $H_0$  to the Banach spaces  $\overline{X}$  and  $\overline{Y}$ , respectively. Suppose also that J and H are linear extensions of  $J_0$  and  $H_0$  to all of X and Y, respectively. Such extensions always exist, for we may take  $J = J_0 P_{\widetilde{X}}$  and  $H = H_0 P_{\widetilde{Y}}$ . In many practical situations,  $\overline{X}$  and  $\overline{Y}$  are chosen to be finite.dimensional.

37

Consider again the equations

$$(I.1.1)$$
 Ax = b,

where A:  $X \rightarrow Y$ ,  $b \in Y$ , and

Y

 $(I.1.2) \qquad \overline{A}\overline{x} = \overline{b},$ 

where  $\overline{A}: \overline{X} \to \overline{Y}$ ,  $\overline{b} \in \overline{Y}$ . We shall refer to (I.1.1) as the "exact" equation and to (I.1.2) as its "approximate" equation. We assume that  $A \in \ell_b(X,Y)$ ,  $\overline{A} \in \ell_b(\overline{X},\overline{Y})$  and that the following decompositions are possible:

$$X = N(A) \oplus N(A)^{c}, \quad \overline{X} = N(\overline{A}) \oplus N(\overline{A})^{c}$$
$$= R(A) \oplus R(A)^{c}, \quad \text{and} \quad \overline{Y} = R(\overline{A}) \oplus R(\overline{A})^{c}$$

The symbol  $\oplus$  is here used to indicate that all eight complements are necessarily closed. Denote by  $A^+ \in \ell_b(Y,X)$  the generalized inverse of A relative to the continuous projections  $P_{N(A)}$  and  $P_{R(A)}$ , and by  $\overline{A}^+ \in \ell_b(\overline{Y}, \overline{X})$  the generalized inverse of  $\overline{A}$  relative to the continuous projections  $P_{N(\overline{A})}$  and  $P_{R(\overline{A})}$ . Let us denote by  $x^* = A^+b$  and  $\overline{x}^* = \overline{A}^+\overline{b}$ , the best approximate solutions of the equations (I.1.1) and (I.1.2) respectively.

In the sequel we will state results relating the best approximate solutions of the exact and approximate equations when some or all of the following conditions are satisfied:

(I) The operator A is represented as A = M + N, where M is bounded and  $X = N(M) \oplus N(M)^{C}$ ,  $N(A)^{C} \subset N(M)^{C}$  and  $Y = R(M) \oplus R(M)^{C}$ . (In this situation  $M^{+}$  denotes the generalized inverse of M, relative to the continuous projections  $P_{N(M)C}$  and  $P_{R(M)}$ .)

(II) 
$$J_0$$
 maps  $N(M)^{c} \cap \widetilde{X}$  into  $N(\overline{A})^{c} \cap \overline{X}$ .

- (III) H has the property  $\overline{A}^+\overline{b} = \overline{A}^+HP_{\mathcal{R}(A)}^{b}$ .
- (IV)  $\|AJ_0^{\mathfrak{X}} HA\widetilde{\mathfrak{X}}\| \leq \varepsilon \|\widetilde{\mathfrak{X}}\|$  for some constant  $\varepsilon \geq 0$  and all  $\widetilde{\mathfrak{X}} \in N(M)^{\mathbb{C}} \cap \widetilde{\mathfrak{X}}$ .
- (V) For every  $x \in N(A)^{c}$  there is a  $u \in N(M)^{c} \cap \widetilde{X}$  such that  $\|Mu - P_{R(M)}^{Nx}\| \leq \eta_{1} \|x\|$  for some constant  $\eta_{1} \geq 0$ .
- (VI) There exists a vector  $v \in N(M)^{C} \cap \widetilde{X}$  such that  $\|Mv - P_{R(M)}P_{R(A)}b\| \leq \eta_{2}\|P_{R(A)}b\|$  for some constant  $\eta_{2} \geq 0$ .

<u>Theorem 2.1.</u> (Conditions for the solvability of the approximate equation.) Let the conditions (I), (IV) and (V) be satisfied. In addition, suppose that

IV.2

(A) 
$$M^+: \widetilde{Y} \to \widetilde{X}$$
,

(B) 
$$H_0^{-1} H \mathcal{R}(A) \subseteq \mathcal{R}(A)$$
,

(C) 
$$J(N(M)^{c})$$
 is closed and

(D)  $R(\overline{A}) \subset HR(A)$ 

If

(2.1) 
$$q = [\varepsilon(1 + \eta_1 || M^+ ||) + \eta_1 || HA|| || M^+ ||] || A^+ H_0^{-1} || < 1,$$

then the equation

$$(2.2) \qquad \qquad \overline{A}\overline{x} = HP_{R(A)}b$$

has a solution  $\bar{x}^* \in J(N(M)^{C})$  for every  $b \in Y$ . Furthermore

(2.3) 
$$\|\bar{x}^*\| \leq \frac{\alpha}{1-q} \|HP_{R(A)}^b\|$$

where

(2.4) 
$$\alpha = \|J_0\| (1 + \eta_1 \|M^+\|) \|A^+ H_0^{-1}\|$$

<u>Proof.</u> It is sufficient to show that the conditions (1.1) of Lemma 1.1 are satisfied for the equation (2.2) with  $E = J(N(M)^{C})$ , F = HR(A),  $V = \overline{A}$ ,  $y = HP_{R(A)}b$  and  $\alpha$  and q as in (2.1) and (2.4). First, consider the equation  $Ax = H_0^{-1}HP_{R(A)}b$  and its solution  $x_0 = A^{+}H_0^{-1}HP_{R(A)}b$ . Denote  $z = Mx_0 - H_0^{-1}HP_{R(A)}b$ . Since  $x_0 \in N(A)^{C}$  and  $N(A)^{C} \subset N(M)^{C}$ , by condition (I), we find that

(2.5) 
$$x_0 = M^+ z + M^+ H_0^{-1} HP_{R(A)} b$$

and also

(2.6) 
$$Ax_0 = H_0^{-1} HP_{R(A)} b$$

by definition of  $x_0$  and condition (B). Therefore  $z = Mx_0 - Ax_0 = -Nx_0$ , since A = M + N.

Now, for  $x = -x_0^{-1}$ , condition (V) implies that there exists  $u \in N(M)^{C} \cap \widetilde{X}$  such that

(2.7) 
$$\| Mu - P_{R(M)} N(-x_0) \| = \| Mu - P_{R(M)} z \| \le \eta_1 \| x_0 \|,$$

for some  $\eta_1 \ge 0$ . Denote  $\tilde{x} = u + M^{+} H_0^{-1} HP_{\mathcal{R}(A)} b$ . Note that  $\tilde{x} \in N(M)^{\mathbb{C}} \cap \tilde{x}$ , by conditions (V) and (A). We will now show that  $J_0\tilde{x}$  is the required element  $\hat{x}$  of E in Lemma 1.1. First,

$$\|\overline{AJ}_{0}\widetilde{x} - HP_{\mathcal{R}(A)}b\| = \|\overline{AJ}_{0}\widetilde{x} - HH_{0}^{-1}HP_{\mathcal{R}(A)}b\|$$

$$= \|\overline{AJ}_{0}\widetilde{x} - HAx_{0}\|, \text{ by condition (B) and}$$

$$\det inition of x_{0}$$

$$\leq \|\overline{AJ}_{0}\widetilde{x} - HA\widetilde{x}\| + \|HA\widetilde{x} - HAx_{0}\|, \text{ by the}$$

$$\operatorname{triangular inequality}$$

$$\leq \epsilon \|\widetilde{x}\| + \|HA\|\|\widetilde{x} - x_{0}\|, \text{ by condition (IV)}$$

Since

$$\|\widetilde{\mathbf{x}} - \mathbf{x}_0\| = \|\mathbf{u} + \mathbf{M}^{\dagger} \mathbf{H}_0^{-1} \mathbf{HP}_{\mathcal{R}(\mathbf{A})} \mathbf{b} - \mathbf{M}^{\dagger} \mathbf{z} - \mathbf{M}^{\dagger} \mathbf{H}_0^{-1} \mathbf{HP}_{\mathcal{R}(\mathbf{A})} \mathbf{b}\|,$$
  
by definition of  $\widetilde{\mathbf{x}}$  and  $\mathbf{x}_0$ 

$$= \| u - M^{+} z \| = \| M^{+} M u - M^{+} P_{\mathcal{R}(M)} z \|$$

$$\leq \eta_{1} \| M^{+} \| \| x_{0} \| , \text{ by } (2.7)$$

$$\leq \eta_{1} \| M^{+} \| \| A^{+} H_{0}^{-1} \| \| H P_{\mathcal{R}(A)} b \| , \text{ by definition of } x_{0}$$

and

$$\begin{split} \|\widetilde{\mathbf{x}}\| &\leq \|\mathbf{x}_{0}\| + \|\widetilde{\mathbf{x}} - \mathbf{x}_{0}\| \\ &\leq (1 + \eta_{1} \|\mathbf{M}^{\dagger}\|) \|\mathbf{x}_{0}\|, \quad \text{by (2.8)} \\ &\leq (1 + \eta_{1} \|\mathbf{M}^{\dagger}\|) \|\mathbf{A}^{\dagger} \mathbf{H}_{0}^{-1}\| \|\mathbf{H}\mathbf{P}_{\mathcal{R}}(\mathbf{A})^{\mathbf{b}}\|, \end{split}$$

the above inequality gives

$$\|\overline{A}J_{0}\tilde{x} - HP_{R(A)}b\| \leq [\varepsilon(1 + \eta_{1}\|M^{+}\|)\|A^{+}H_{0}^{-1}\| + \eta_{1}\|HA\|\|M^{+}\|\|A^{+}H_{0}^{-1}\|]\|HP_{R(A)}b\|$$

$$(2.9) = q\|HP_{R(A)}b\|,$$

where q is the constant defined by (2.1). Also

(2.10) 
$$\|J_0^{\chi}\| \le \|J_0^{\chi}\| \le \alpha \|H^p_{\mathcal{R}(A)}^{b}\|,$$

where  $\alpha$  is defined by (2.4). The inequalities (2.9) and (2.10) correspond to the assumptions (1.1) of Lemma 1.1 with  $V = \overline{A}$ ,  $\hat{x} = J_0 \tilde{x}$  and  $y = HP_{R(A)}b$ . Conditions (C) and (D) guarantee that  $E = J(N(M)^c)$  be a closed subspace of X and F = HR(A)be a subspace of Y containing  $V(E) = \overline{A}(J(N(M)^c))$ . All conditions of Lemma 1.1 are now satisfied, and the conclusions of Theorem 2.1 follow.

<u>Corollary 2.1.</u> Let  $A \in l_b(X,Y)$  and suppose that the conditions (I), (IV), (V), (A), (B), (C) and (D) are satisfied. If q < 1 and  $\overline{A}$  satisfies the condition

(E) 
$$J(N(M)^{c}) \subset N(\overline{A})^{c}$$

then

$$\|\overline{A}^+\| \leq \frac{\alpha}{1-q} P_{R(\overline{A})} \|,$$

where q and  $\alpha$  are as in Theorem 2.1.

Proof. We need to show that

$$\|\overline{A}^{+} \overline{y}\| \leq \frac{\alpha}{1-q} \|P_{\mathcal{R}(\overline{A})}\| \|\overline{y}\|, \text{ for every } \overline{y} \in \overline{Y}.$$

Let  $\bar{x}^*$  denote the solution in  $J(N(M)^c)$  of the equation  $\overline{A}\bar{x} = P_{R(\overline{A})}\bar{y}$ . Such an  $\bar{x}^*$  exists, by Theorem 1.1. In fact, by condition (E),  $\bar{x}^* = \overline{A}^+ \bar{y}$ , i.e.  $\bar{x}^*$  is a unique best approximate solution in  $N(\overline{A})^c$ of the equation  $\overline{A}\bar{x} = \bar{y}$ . But  $P_{R(\overline{A})}\bar{y} = HP_{R(A)}b$  for some  $b \in Y$ , by condition (D). So  $\bar{x}^*$  is also a solution of the equation  $\overline{A}\bar{x} = HP_{R(A)}b$  and it satisfies

$$\|\bar{x}^*\| \leq \frac{\alpha}{1-q} \|HP_{\mathcal{R}(A)}b\|$$
, by (2.3).

Therefore,

$$\|\overline{A}^{+} \overline{y}\| \leq \frac{\alpha}{1-q} \|P_{R(A)}b\| = \frac{\alpha}{1-q} P_{R(\overline{A})} \overline{y}\| \leq \frac{\alpha}{1-q} P_{R(\overline{A})} \|\|\overline{y}\|.$$

If X = Y,  $\tilde{X} = \tilde{Y} = \overline{X} = \overline{Y}$ , M = I and  $A^{-1}$  exists then the conditions (A), (B), (C), (D), (I) and (II) are trivially satisfied, the requirement on H in (III) becomes  $\overline{A}^+ \overline{b} = \overline{A}^+ Hb$  while (IV), (V) and (VI) reduce to:

- (IV')  $\|\overline{Ax} HAx\| \le \varepsilon \|x\|$  for some constant  $\varepsilon \ge 0$  and all  $\overline{x} \in \overline{X}$ .
- (V') For every  $x \in X$  there is a  $u \in \widetilde{X}$  such that  $||u Nx|| \le \eta_1 ||x||$  for some constant  $\eta_1 \ge 0$ .
- (VI') There exists a vector  $v \in \widetilde{X}$  such that  $||v b|| \le n_2 ||b||$  for some constant  $n_2 \ge 0$ .

In the nonsingular case, Theorem 1.1 reduces to the following result of Kantorovich and Akilov [31, p.545].

<u>Corollary 2.2.</u> Let  $A \in {}^{k}_{b}(X)$  have an inverse and let the conditions (IV') and (V') be satisfied. If

$$q = [\varepsilon(1 + \eta_1) + \eta_1 \| HA\| ] \| A^{-1} \| < 1,$$

then the equation  $\overline{A}\overline{x} = \overline{b}$  has a solution  $\overline{x}^* \in \overline{X}$  for every  $\overline{b} \in \overline{X}$ . Also

$$\|\bar{\mathbf{x}}^*\| \leq \frac{\alpha}{1-q} \|\bar{\mathbf{b}}\|,$$

where  $\alpha = (1 + \eta_1) \| A^{-1} \|$ .

An estimate for the norm of the generalized inverse  $\overline{A}^+$  was obtained using Theorem 1.1. In the nonsingular case, Corollary 2.1 gives the following result of Kantorovich and Akilov [31, p.546].

<u>Corollary 2.3.</u> Let the hypotheses and the notation of Corollary 2.2 hold and let  $\overline{A}$  satisfy the condition:

(E') "The existence of a solution of the equation  $\overline{Ax} = \overline{b}$  for every  $\overline{b} \in \overline{X}$  implies its uniqueness".

Then

$$\|\overline{A}^{-1}\| \leq \frac{\alpha}{1-q}$$

<u>Proof.</u> Condition (E') and Corollary 2.2 imply the existence of  $\overline{A}^{-1}$ . The result now follows from Corollary 2.1 since the conditions (D) and (E) are satisfied when  $\overline{A}^{-1}$  and  $A^{-1}$  exist.

## 3. Error Estimates

The following theorem estimates the distance between the best approximate solution  $\bar{x}^*$  of the approximate equation and the best approximate solution  $x^*$  of the exact equation.

<u>Theorem 3.1.</u> (Estimate of the distance between best approximate solutions.) Consider the equation Ax = b and its approximate equation  $\overline{Ax} = \overline{b}$ . If the conditions (I) to (VI) are satisfied, then

$$\|\mathbf{x}^* - \mathbf{J}_0^{-1} \mathbf{\bar{x}}^*\| \le p \|\mathbf{x}^*\|,$$

where

(3.2) 
$$p = \varepsilon (1 + c) \| J_0^{-1} \| \| \overline{A}^+ \| + c (1 + \| J_0^{-1} \overline{A}^+ HA \|)$$

and

(3.3) 
$$c = \min\{1, (\eta_1 + \eta_2 \|A\|) \|M^{\dagger}\|\}$$

<u>Proof.</u> First we show that  $x^*$  can be approximated by an  $\tilde{X} \in N(M)^{\mathbb{C}} \cap \tilde{X}$  to the order of  $\eta_1 + \eta_2$ . We know, by conditions (V) and (VI), that there exist u and v in  $N(M)^{\mathbb{C}} \cap \tilde{X}$  such that (3.4)  $\|Mu - P_{\mathcal{R}(M)}Nx^*\| \leq \eta_1 \|x^*\|$ 

and

(3.5) 
$$\|Mv - P_{R(M)}P_{R(A)}b\| \le \eta_2 \|P_{R(A)}b\|$$

Denote  $\tilde{x} = M^+(Mv - Mu)$ . Clearly  $\tilde{x} \in N(M)^{\mathbb{C}} \cap \tilde{X}$ . We now show that  $x^*$  can be approximated by  $\tilde{x}$  to the order of  $n_1 + n_2$ .

$$\begin{aligned} \|x^* - \widetilde{x}\| &= \|M^{\dagger}Mx^* - \widetilde{x}\|, & \text{by condition (I)} \\ &= \|-M^{\dagger}Nx^* + M^{\dagger}(M+N)x^* - M^{\dagger}(Mv - Mu)\|, & \text{by definition of } \widetilde{x} \\ &= \|-M^{\dagger}Nx^* + M^{\dagger}P_{R(A)}b - M^{\dagger}(Mv - Mu)\|, & \text{since } A = M + N \text{ and } Ax^* = P_{R(A)}b \\ &\leq \|M^{\dagger}\|(\|Mu - P_{R(M)}Nx^*\| + \|Mv - P_{R(M)}P_{R(A)}b\|), & \text{since } M^{\dagger} = M^{\dagger}P_{R(M)} \\ &\leq \|M^{\dagger}\|(\eta_{1}\|x^*\| + \eta_{2}\|P_{R(A)}b\|), & \text{by (3.4) and (3.5)} \\ &\leq \|M^{\dagger}\|(\eta_{1} + \eta_{2}\|A\|)\|x^*\|, & \text{since } \|P_{R(A)}b\| = \|Ax^*\| \leq \|A\|\|x^*\| \end{aligned}$$

Hence we conclude that there exists an  $\ \widetilde{x} \in N(M)^{C} \cap \widetilde{X}$  such that

$$\|x^* - \widetilde{x}\| \le c \|x^*\|,$$

where  $c = \min\{1, (\eta_1 + \eta_2 \|A\|) \|M^*\|\}$ . (Note that  $c \le 1$ , since we can always choose  $\tilde{x} = 0$  in (3.6).)

Let us now prove (3.1). Denote  $\bar{x}_0 = \bar{A}^+ H A \tilde{x}$ . Then

$$(3.7) \|x^* - J_0^{-1}x^*\| \le \|x^* - \widetilde{x}\| + \|\widetilde{x} - J_0^{-1}\overline{x}_0\| + \|J_0^{-1}\overline{x}_0 - J_0^{-1}\overline{x}^*\| .$$

The first term on the right hand side  $||x^* - \tilde{x}||$  is estimated by (3.6). The two remaining terms will now be estimated. First

$$\begin{split} \widetilde{\mathbf{x}} - \mathbf{J}_0^{-1} \widetilde{\mathbf{x}}_0 &= \widetilde{\mathbf{x}} - \mathbf{J}_0^{-1} \overline{\mathbf{A}}^+ \mathbf{H} \overline{\mathbf{A}} \widetilde{\mathbf{x}} , \quad \text{by definition of } \widetilde{\mathbf{x}}_0 \\ &= \mathbf{J}_0^{-1} (\mathbf{J}_0 - \overline{\mathbf{A}}^+ \mathbf{H} \mathbf{A}) \widetilde{\mathbf{x}} , \quad \text{since } \widetilde{\mathbf{x}} \in \widetilde{\mathbf{X}} \end{split}$$

$$= J_0^{-1}\overline{A}^+ (\overline{A}J_0 - HA)\widetilde{x}, \text{ since } \overline{A}^+\overline{A}J_0\widetilde{x} = P_{N(\overline{A})}c^{-J_0}\widetilde{x}$$
$$= J_0\widetilde{x}, \text{ by condition (II)}$$

Hence

$$\begin{split} \|\widetilde{\mathbf{X}} - \mathbf{J}_0^{-1} \widetilde{\mathbf{X}}_0 \| &\leq \|\mathbf{J}_0^{-1} \overline{\mathbf{A}}^+\| \| (\overline{\mathbf{A}} \mathbf{J}_0 - \mathbf{HA}) \widetilde{\mathbf{X}} \| \\ &\leq \varepsilon \| \mathbf{J}_0^{-1} \overline{\mathbf{A}}^+\| \| \widetilde{\mathbf{X}} \| , \quad \text{by condition (IV)} \\ &\leq \varepsilon \| \mathbf{J}_0^{-1} \overline{\mathbf{A}}^+\| (\|\mathbf{x}^*\| + \|\mathbf{x}^* - \widetilde{\mathbf{X}}\|) , \quad \text{by the triangle inequality} \\ &\leq \varepsilon (1 + \mathbf{c}) \| \mathbf{J}_0^{-1} \overline{\mathbf{A}}^+\| \| \mathbf{x}^*\| , \quad \text{by (3.6).} \end{split}$$

The third term is estimated as follows:

$$\begin{split} \|J_0^{-1}\bar{x}_0 - J_0^{-1}\bar{x}^*\| &= \|J_0^{-1}(\overline{A}^*HA\tilde{x} - \overline{A}^*\overline{A}\bar{x}^*)\| , \text{ by definition of } \bar{x}_0 \\ &= \|J_0^{-1}(\overline{A}^*HA\tilde{x} - \overline{A}^*P_{\mathcal{R}(\overline{A})}\overline{b})\| \\ &= \|J_0^{-1}(\overline{A}^*HA\tilde{x} - \overline{A}^*HP_{\mathcal{R}(A)}b)\| , \text{ by condition (III)} \\ &\leq \|J_0^{-1}\overline{A}^*HA\|\|\tilde{x} - x^*\| , \text{ since } P_{\mathcal{R}(A)}b = Ax^* \\ &\leq c\|J_0^{-1}\overline{A}^*HA\|\|x^*\| , \text{ by (3.6).} \end{split}$$

After substituting the above estimates into (3.7) the conclusion follows.

<u>Remark 3.1.</u> It may happen that we could approximate  $x^*$  by some  $\tilde{x} \in N(M)^{C} \cap \tilde{x}$  directly. Then we no longer need the conditions (V) and (VI) and we can set  $c = \min\{1, \|x^* - \tilde{x}\|\}$  in (3.2).

<u>Remark 3.2.</u> If p < 1, we can write the estimate (3.1) as follows

$$\|x^* - J_0^{-1}\bar{x}^*\| \le \frac{p}{1-p}\|J_0^{-1}\bar{x}^*\|$$

This is true, since

$$\|x^* - J_0^{-1} \bar{x}^*\| \le p \|x^*\| \le p(\|J_0^{-1} \bar{x}^*\| + \|x^* - J_0^{-1} \bar{x}^*\|) .$$

<u>Remark 3.3.</u> If X and Y are Hilbert spaces and  $N(A)^{C} = R(A^{*})$ ,  $R(A)^{C} = N(A^{*})$ , then Theorem 3.1 reduces to a result obtained by Zlobec in [65]. However, the Hilbert space version of Theorem 3.1 is proved there under slightly different assumptions.

<u>Remark 3.4.</u> It may happen that one cannot satisfy condition (III) but that a constant  $n_3$ , such that

$$\|\overline{A}^{+}\overline{b} - \overline{A}^{+}HP_{\mathcal{R}(A)}b\| \leq \eta_{3}\|x^{*}\|,$$

is found. In this case the constant p in (3.1) is different. Now

$$\|J_0^{-1}\bar{x}_0 - J_0^{-1}\bar{x}^*\| \le (c\|J_0^{-1}\overline{A}^+HA\| + \eta_3\|J_0^{-1}\|)\|x^*\|$$

and hence

$$P = [\varepsilon(1 + c) \|\overline{A}^{+}\| + \eta_{3}] \|J_{0}^{-1}\| + c(1 + \|J_{0}^{-1}\overline{A}^{+}HA\|)$$

In the nonsingular case, we get the following result from [31, p.547]. (Recall that in our setting  $x^* \in N(A)^c$ . Also, as

in [31], we specify  $\overline{b} = Hb$ .)

<u>Corollary 3.1.</u> Let the conditions (IV'), (V') and (VI') be satisfied and let  $A^{-1}$  and  $\overline{A}^{-1}$  exist. Then

$$\|x^* - \bar{x}^*\| \le p \|x^*\|$$
,

where  $x^*$  is the solution of the exact equation (I.1.1),  $\bar{x}^*$  is the solution of the approximate equation (I.1.2) and

$$p = 2\varepsilon \|\overline{A}^{-1}\| + (n_1 + n_2 \|A\|) (1 + \|\overline{A}^{-1} HA\|).$$

<u>Proof.</u> Specify M = I, X = Y and  $\tilde{X} = \tilde{Y} = \overline{X} = \overline{Y}$  in Theorem 3.1.

## 4. Convergence Criteria

Our next result gives conditions for convergence of approximation schemes. Suppose that the exact equation Ax = b is approximated by a sequence of equations  $\overline{A}_n \overline{x} = \overline{b}_n$ , n = 1, 2, ..., rather than by a single equation. This determines a sequence of the spaces  $\widetilde{X}_n$ ,  $\widetilde{Y}_n$ ,  $\overline{X}_n$ ,  $\overline{Y}_n$ , the operators  $\overline{A}_n$ ,  $(J_0)_n$ ,  $(H_0)_n$ ,  $J_n$ ,  $H_n$  and the constants  $\varepsilon_n$ ,  $(n_1)_n$ ,  $(n_2)_n$ ,  $c_n$ ,  $q_n$ ,  $\alpha_n$ ,  $p_n$ , n = 1, 2, .... For the sake of notational simplicity these indices will generally

be omitted in the sequel. The following theorem gives conditions for the convergence of the sequence  $\bar{x}_n^*$ , the best approximate solution of  $\overline{A}_n \bar{x}_n = \overline{b}_n$ , n = 1, 2, ..., to  $x^*$ , the best approximate solution of the exact equation.

<u>Theorem 4.1.</u> (Convergence of the best approximate solutions.) Consider the equation Ax = b and a sequence of approximate equations  $\overline{Ax} = \overline{b}$ . Suppose that for each n = 1, 2, ... the conditions:

(i) (I) to (VI) and (A), (B), (C), (D) and (E)

- (ii)  $\sup_{n} \|J_0\| < \infty$ ,  $\sup_{n} \|H_0^{-1}\| < \infty$  and  $\sup_{n} \|P_{\mathcal{R}}(\overline{A})\| < \infty$
- (iii)  $\lim_{n \to \infty} \varepsilon \|J_0^{-1}\| = 0$ ,  $\lim_{n \to \infty} \eta_1 \|J_0^{-1}\| \|H\| = 0$ ,  $\lim_{n \to \infty} \eta_2 \|J_0^{-1}\| \|H\| = 0$

are satisfied. Then  $\lim_{n\to\infty} \eta_1 = 0$ ,  $\lim_{n\to\infty} \eta_2 = 0$  and the sequence of best approximate solutions of  $\overline{Ax} = \overline{b}$  converges to the best approximate solution  $x^*$  of Ax = b, i.e.

$$\lim_{n \to \infty} \| x^* - J_0^{-1} \bar{x}^* \| = 0.$$

More precisely,

$$\|x^* - J_0^{-1}\bar{x}^*\| \le \varepsilon c_1 \|J_0^{-1}\| + \eta_1 (c_2 + c_3 \|J_0^{-1}\|\|H\|) + \eta_2 (c_4 + c_5 \|J_0^{-1}\|\|H\|) \|x^*\|,$$

where  $c_1$  to  $c_5$  are some constants independent of the index n.

IV.4

<u>Proof.</u> Since  $H_0 H_0^{-1} = I$ , it follows that  $||H_0^{-1}|| ||H_0|| \ge 1$ , n = 1,2,.... Hence

$$(4.1) \qquad \inf_{n} \|H_0\| > 0$$

using the second assumption in (ii). Similarly, one concludes that

(4.2) 
$$\inf_{n} \|J_0^{-1}\| > 0$$

using the first assumption in (ii). Also  $\|H\| \ge \|H_0\|$ , since  $H_0 = H|_{\widetilde{Y}}$ . Therefore, by (4.1), (4.3)  $\inf \|H\| > 0$ .

From (4.2), (4.3) and condition (iii), we conclude that, in particular,

(4.4)  $\lim_{n\to\infty} \varepsilon = 0, \quad \lim_{n\to\infty} \eta_1 = 0, \quad \lim_{n\to\infty} \eta_2 = 0 \text{ and } \lim_{n\to\infty} \eta_1 \|H\| = 0.$ 

Recall the constant q introduced in Theorem 2.1:

$$q = [\varepsilon(1 + \eta_1 \| M^{\dagger} \|) + \eta_1 \| HA\| \| M^{\dagger} \|] \| A^{\dagger} H_0^{-1} \|.$$

For sufficiently large n, using (4.4), one has  $q < \frac{1}{2}$  and for such values of n Theorem 2.1 is applicable. But we can also apply Corollary 2.1 to obtain

$$\|\overline{A}^+\| \leq \frac{\alpha}{1-q} \|P_{\mathcal{R}(A)}\| \leq 2\alpha \|P_{\mathcal{R}(\overline{A})}\|.$$

Since  $\alpha = \|J_0\| (1 + \eta_1 \|M^{\dagger}\|) \|A^{\dagger}H_0^{-1}\|$ , one concludes, using the third assumption in (ii), that  $\|\overline{A}^{\dagger}\|$  is bounded independently of the index n, i.e.

$$(4.5) \qquad \sup_{n} \|\overline{A}^{+}\| = s < \infty$$

The desired estimate now follows for sufficiently large n:

$$\begin{aligned} \|x^* - J_0^{-1} \bar{x}^*\| &\leq [2\varepsilon \|J_0^{-1}\| \|\bar{A}^*\| + (\eta_1 + \eta_2 \|A\|) (1 + \|J_0^{-1} \bar{A}^* HA\|) \|M^*\|] \|x^*\|, \\ & \text{by (3.2) and (3.3)} \\ &\leq [\varepsilon c_1 \|J_0^{-1}\| + \eta_1 (c_2 + c_3 \|J_0^{-1}\| \|H\|) + \eta_2 (c_4 + c_5 \|J_0^{-1}\| \|H\|)] \|x^*\|, \end{aligned}$$

where  $c_1 = 2s$ ,  $c_2 = \|M^{\dagger}\|$ ,  $c_3 = s\|A\|\|M^{\dagger}\|$ ,  $c_4 = \|A\|\|M^{\dagger}\|$ ,  $c_5 = s\|A\|^2\|M^{\dagger}\|$  and s is defined by (4.5). The right-hand side in the above inequality tends to zero when  $n \rightarrow \infty$ , by (iii) and (4.4)

A corresponding result in the nonsingular case is given in [31, p.549] as follows.

<u>Corollary 4.1.</u> Consider the equation Ax = b and a sequence of approximate equations  $\overline{Ax} = \overline{b}$ . Assume that  $A^{-1}$  exists and that  $\overline{A}$  satisfies condition (E') for each n = 1, 2, ... Assume further that for each n = 1, 2, ... the conditions (IV'), (V'), (VI') are satisfied, and that

Then the approximate equations are consistent for sufficiently large n and the sequence of approximate solutions converges to the exact solution  $x^*$  of Ax = b, i.e.

$$\lim_{n\to\infty} \|x^* - \bar{x}^*\| = 0.$$

More precisely,

$$\| \mathbf{x}^{*} - \bar{\mathbf{x}}^{*} \| \leq \varepsilon c_{1} + \eta_{1} (c_{2} + c_{3} \| \mathbf{H} \|) + \eta_{2} (c_{4} + c_{5} \| \mathbf{H} \|),$$

where  $c_1$  to  $c_5$  are some constants indpendent of the index n.

<u>Proof.</u> Set X = Y,  $\overline{X} = \overline{Y} = \widetilde{X} = \widetilde{Y}$  and M = I in Theorem 2.1. Then  $J_0 = H_0 = I$  and  $||H|| = ||J|| \ge 1$ . Furthermore, from Corollary 1.3, we see that  $R(\overline{A}) = \overline{X}$ . So condition (E) is satisfied and  $||P_{R(\overline{A})}|| = 1$ . Conditions (i) to (iii) of Theorem 4.1 hold and the result follows.

<u>Remark 4.1.</u> Constants  $\varepsilon$ ,  $\eta_1$ ,  $\eta_2$  and  $c_1$  to  $c_5$  in Theorem 4.1 reduce, in the nonsingular case, to the corresponding constants in Corollary 4.1.

Let us recall that Corollary 2.1 gave us an estimate for  $\|\overline{A}^+\|$ in terms of  $\|A^+\|$  via constants  $\alpha$  and q. Our last result gives us the reverse estimate.

<u>Theorem 4.2.</u> (Estimate for the norm of the generalized inverse.) Let  $A \in l_b(X,Y)$  and  $\overline{A} \in l_b(\overline{X},\overline{Y})$ . If the conditions (I) to (V) are satisfied and

$$\mathbf{r} = \left[ \epsilon (1 + \eta_1 \| \mathbf{M}^+ \|) \| J_0^{-1} \| \| \overline{\mathbf{A}}^+ \| + \eta_1 \| \mathbf{M}^+ \| (1 + \| J_0^{-1} \overline{\mathbf{A}}^+ \mathbf{HA} \|) \right] < 1$$

then

$$\|A^{+}\| \leq (1 - r)^{-1} [\|J_{0}^{-1}\overline{A}^{+}H\| + \|M^{+}\| (1 + \varepsilon \|J_{0}^{-1}\|\|\overline{A}^{+}\| + \|J_{0}^{-1}\overline{A}^{+}HA\| )].$$

<u>Proof.</u> Take  $\hat{x} \in X$ ,  $\hat{x} \notin N(A)$ . Then  $x^* = P_{N(A)}c\hat{x}$  is clearly the best approximate solution of the equation

$$(4.6) Ax = A\hat{x}.$$

By condition (V), there exists a  $u \in N(M)^{C} \cap \widetilde{X}$  such that

(4.7) 
$$\| Mu - P_{\mathcal{R}(M)} N(-P_{N(A)} c^{\hat{x}}) \| \leq \eta_1 \| P_{N(A)} c^{\hat{x}} \|.$$

Now

$$\|x^* - u\| = \|P_{N(A)c}\hat{x} - u\| = \|M^* MP_{N(A)c}\hat{x} - M^* Mu\|$$
, by condition (I)

 $\leq (\|MP_{N(A)}^{c} \hat{x} + NP_{N(A)}^{c} \hat{x}\| + \|P_{R(M)}^{NP}_{N(A)}^{c} \hat{x} + Mu\|) \|M^{\dagger}\|$ 

54

$$\leq \left( \eta_1 + \frac{\|Ax^*\|}{\|P_{N(A)}^c \hat{x}\|} \right) \|P_{N(A)}^c \hat{x}\| \|M^*\|, \quad \text{by } (4.7).$$

Apply Remark 3.1 to the exact equation (4.6) and its approximate equation  $\overline{A}\overline{x} = HA\hat{x}$ , with

.

 $\bigcirc$ 

$$c = \min \{1, \|M^{+}\| \left(\eta_{1} + \frac{\|A\hat{x}\|}{\|P\|_{(A)}c\hat{x}\|}\right) \} \quad \text{ in } (3.2) .$$

Then

Now

$$\begin{split} \|P_{N(A)}c^{\hat{X}}\| &\leq \|P_{N(A)}c^{\hat{X}} - J_{0}^{-1}\overline{A}^{+}HA\hat{x}\| + \|J_{0}^{-1}\overline{A}^{+}H\|\|A\hat{x}\| , \\ & \text{by the triangle inequality} \\ &\leq r\|P_{N(A)}c^{\hat{X}}\| + [\|J_{0}^{-1}\overline{A}^{+}H\| + \|M^{+}\|(1 + \varepsilon\|J_{0}^{-1}\|\|\overline{A}^{+}\| + \|J_{0}^{-1}\overline{A}^{+}HA\|)]\|A\hat{x}\| , \\ & \text{by (4.8) and the definition of } r . \end{split}$$

Hence

0

$$\|A\hat{x}\| \geq \frac{1-r}{\|J_0^{-1}\overline{A}^+H\| + \|M^+\| (1+\epsilon\|J_0^{-1}\|\|\overline{A}^+\| + \|J_0^{-1}\overline{A}^+HA\|)} \|P_N(A)^c \hat{x}\|$$
  
=  $t\|P_{N(A)^c} \hat{x}\|$ ,

·

.

where t denotes the coefficient of  $\|P_{N(A)}c\hat{x}\|$ . Take an arbitrary  $0 \neq y \in R(A)$ . Then there exists an  $\hat{x} \in X$ ,  $\hat{x} \notin N(A)$  such that  $y = A\hat{x}$ . Hence  $A^{+}y = P_{N(A)}c\hat{x}$ . Furthermore, by the above inequality

$$\|A^{+}y\| = \|P_{N(A)}^{c}\hat{x}\| \le \frac{1}{t}\|A\hat{x}\| = \frac{1}{t}\|y\|$$

which gives the desired estimate for  $A^+$ . Note that here t > 0, since r < 1. If y = 0, the above inequality is trivially satisfied.

<u>Remark 4.2.</u> If  $A \in {}^{l}_{b}(X)$  and in addition X = Y,  $\overline{X} = \overline{Y} = \widetilde{X} = \widetilde{Y}$ , then the above result reduces to the bound for a left inverse of A given in [31, p.550].

## V. GALERKIN'S METHOD FOR BEST APPROXIMATE SOLUTIONS

### 1. Description of the Method

In this section we will use Kantorovich's theory to prove that a Galerkin type method, when applied to a certain kind of, possibly inconsistent, operator equation, produces the best approximate solution. This solution is obtained as the limit of a sequence of best approximate solutions of, possibly inconsistent, systems of linear algebraic equations. In the case of Hilbert space, another method is suggested by Nashed [44]. Unlike our approach he finds the best least squares solution by applying Galerkin's method, with a suitably chosen basis, to the <u>consistent</u> equations A\*Ax = A\*band  $Ax = P_{R(A)}b$ , rather than to Ax = b.

Consider an equation Ax = b in a separable Banach space X, where  $A \in l_b(X)$  and  $b \in X$ , not necessarily  $b \in R(A)$ , are given. We assume that A = I + N, where N is compact (which implies that N(A) is finite dimensional) and R(A) is closed. Further we assume that  $R(A) = N(A)^{C}$ . Denote by  $\{\phi_i: i = 1, ..., m\}$  a basis of N(A),  $m = \dim N(A)$ , and by  $\{\psi_i: i = 1, 2, ...\}$  a basis of R(A). It is assumed that R(A) has a countable basis. Then every  $x \in X$ can be written as

$$x = \sum_{i=1}^{m} c_{i}(x)\phi_{i} + \sum_{i=1}^{\infty} d_{i}(x)\psi_{i},$$

where  $c_i(x)$ , i = 1, ..., m and  $d_i(x)$ , i = 1, 2, ... are some

۷.1

coefficients which depend on x. The above situation occurs, for instance, in X = C[0,1] with A a Fredholm integral operator of the second kind with a continuous and symmetric kernel. The best approximate solution of such an equation can be calculated by Galerkin's method as follows: For sufficiently large n solve the system of n linear algebraic equations in n unknowns

(1.1) 
$$\sum_{j=1}^{n} d_{i} (A\psi_{j}) \xi_{j} = d_{i} (b) , \quad i = 1, ..., n.$$

We will show that, for sufficiently large n, the system (1.1) is consistent and that the sequence of solutions  $\bar{x} = (\xi_j)$  converges to the best approximate solution of Ax = b, with respect to  $P_{N(A)}c = P_{R(A)}$  and  $P_{R(A)}$ , when  $n \to \infty$ . (Note that in this situation both A and A<sup>+</sup> leave R(A) invariant.) In order to prove the consistency of (1.1), for large n, we will use a result from Krasnosel'skii et al. [33, p.212] which is stated here as the following lemma.

Lemma 1.1. Let  $T \in \ell_b(X)$  be compact and let  $\{P_n: n = 1, 2, ...\}$ be a sequence of projections in  $\ell_b(X)$ , where X is a Banach space. If  $P_n \neq I$  strongly, i.e., for every  $x \in X$ ,

$$\|P_x - x\| \to 0 \quad \text{as} \quad n \to \infty,$$

then  $\|(I - P_n)T\| \to 0$  as  $n \to \infty$ .

In our situation we specify

$$Y = X$$
,  $\tilde{Y} = \tilde{X} = \operatorname{span}\{\psi_1, \dots, \psi_n\}$  and  $P_{\tilde{Y}} = P_{\tilde{X}}$ ,

which is defined by

$$P_{\widetilde{X}} x = P_{\widetilde{X}} \left( \sum_{i=1}^{m} c_i(x) \phi_i + \sum_{i=1}^{\infty} d_i(x) \psi_i \right) = \sum_{i=1}^{n} d_i(x) \psi_i.$$

Further, J and H are defined by

$$Jx = Hx = \begin{pmatrix} d_1(x) \\ \vdots \\ d_n(x) \end{pmatrix}$$

while

$$J_0 = H_0 = J |_{\widetilde{X}}, \quad \overline{b} = Hb \quad and \quad \overline{A} = HAJ_0^{-1}$$

Note that  $P_{\widetilde{X}} = H_0^{-1}H$  and  $\overline{X} = \overline{Y}$  is the space of all n-tuples. The norm in  $\overline{X}$  and  $\overline{Y}$  is defined by

(1.2) 
$$\|\bar{\mathbf{x}}\| = \|\mathbf{J}\mathbf{x}\| = \sup_{k=1,...,n} \left\|\sum_{i=1}^{k} \mathbf{d}_{i}(\mathbf{x})\psi_{i}\right\|$$
.

We will first show that the system (1.1) is consistent for large n and then that all conditions of Theorem IV.4.1 are satisfied.

Matrix  $\overline{A} = (\overline{a}_{ij})$ ,  $\overline{a}_{ij} = d_i(A\psi_j)$ , i, j = 1, ..., n has an inverse if and only if  $P_{\widetilde{X}}A|_{\widetilde{X}}$  has an inverse. By Lemma 1.1, where X = R(A),  $I = I|_{R(A)}$ ,  $P_n = P_{\widetilde{X}}|_{R(A)}$  and  $T = N|_{R(A)}$ ,

$$\| \left( \mathbf{I} \Big|_{\mathcal{R}(\mathbf{A})} - \mathbf{P}_{\widetilde{\mathbf{X}}} \Big|_{\mathcal{R}(\mathbf{A})} \right)^{\mathbf{N}} \Big|_{\mathcal{R}(\mathbf{A})} \| = \| \left( \mathbf{P}_{\widetilde{\mathbf{X}}} - \mathbf{I} \right)^{\mathbf{P}}_{\mathcal{R}(\mathbf{A})} ^{\mathbf{N}} \Big|_{\mathcal{R}(\mathbf{A})} \| \neq 0$$

as  $n \to \infty$ . Since  $(A|_{\mathcal{R}(A)})^{-1}$  is bounded (by the assumptions on A), this further implies that

$$\| (P_{\widetilde{X}} - I) P_{\mathcal{R}(A)}^{N} |_{\mathcal{R}(A)} \| \| (A |_{\mathcal{R}(A)})^{-1} \| < 1$$

for sufficiently large n. Now, by specifying X = Y = R(A),

$$M = A \Big|_{\mathcal{R}(A)} \quad \text{and} \quad N = (P_{\widetilde{X}} - I) P_{\mathcal{R}(A)} N \Big|_{\mathcal{R}(A)}$$

in Proposition IV.1.1 we conclude that M+N is invertible, which is here

$$\begin{array}{l} A \Big|_{\mathcal{R}(A)} + (P_{\widetilde{X}} - I)P_{\mathcal{R}(A)}^{N}\Big|_{\mathcal{R}(A)} &= I \Big|_{\mathcal{R}(A)} + P_{\widetilde{X}}^{N}\Big|_{\mathcal{R}(A)} \\ &= P_{\widetilde{X}}^{A}\Big|_{\widetilde{X}} \text{, when restricted to } \widetilde{X} \text{.} \end{array}$$

Therefore  $\overline{A}$  is invertible, which implies that the system (1.1) is consistent for large n.

Let us now show that all assumptions of Theorem IV.4.1 are satisfied for sufficiently large n.

Condition I: Since M = I, this condition is obviously satisfied. Condition II: We know that  $\overline{A}$  is invertible, so this condition, for large n, reduces to  $J_0: \widetilde{X} \rightarrow \overline{X}$ , which is always satisfied.

Condition III: Since A is invertible, the condition becomes  
$$\overline{b} = HP_{R(A)}b$$
, which is satisfied by our construction

Condition IV: One can specify  $\varepsilon = 0$ , because  $\overline{A} = HAJ_0^{-1}$ . Condition V: For an arbitrary  $x \in N(A)^C$ , take  $u = P_{\widetilde{X}}Nx$ . Then

$$\|Mu - P_{\mathcal{R}(M)}Nx\| = \|(P_{\widetilde{X}} - P_{\mathcal{R}(A)})Nx\|, \text{ since}$$
$$M = I \text{ and } N(A)^{C} = \mathcal{R}(A)$$
$$\leq \|(P_{\widetilde{X}} - P_{\mathcal{R}(A)})N|_{\mathcal{R}(A)}\|\|x\|.$$

So, one can specify

Take

$$\eta_1 = \| (\mathbf{P}_{\widetilde{\mathbf{X}}} - \mathbf{P}_{\mathcal{R}(\mathbf{A})}) \mathbf{N} |_{\mathcal{R}(\mathbf{A})} \| .$$

Condition VI:

$$v = P_{\widetilde{X}} P_{\mathcal{R}(A)} b$$
. Then

 $\|Mv - P_{R(M)}P_{R(A)}b\| = \|(P_{\widetilde{X}} - P_{R(A)})P_{R(A)}b\|$ .

Therefore, one can choose

$$\eta_2 = \frac{1}{\|P_{R(A)}b\|} \|(P_{\widetilde{X}} - P_{R(A)})P_{R(A)}b\| .$$

If  $P_{\mathcal{R}(A)}b = 0$ , then set  $n_2 = 0$ .

Condition A:	Since	M = I	and	$\widetilde{X} = \widetilde{Y}$	, the c	onditi	on is satisfi	ied.
Condition B:	By our	constr	uctio	n of H	I and	<sup>н</sup> о,	$H_0^{-1}H = P_{\widetilde{X}}.$	
	Since	$\tilde{X} \subset R($	A), (	one com	ncludes	that	$H_0^{-1}HR(A) \subset \mathbb{R}$	R(A).
Condition C:	J(N(M))	$(z) = \overline{X}$	, whic	ch is d	losed.			

Condition D: Since 
$$R(A) = X$$
, this condition is satisfied by the construction of H.

Condition E: For large n,  $\overline{A}$  is invertible and  $N(\overline{A})^{c} = \overline{X}$ , so the condition is satisfied.

62

In order to prove conditions (ii) of Theorem IV.4.1 we proceed as follows: Define a linear mapping T from R(A) into the space of sequences  $Tx = (d_i(x))$ , i = 1, 2, ... such that  $\sum_{i=1}^{\infty} d_i(x)\psi_i$ is an element in R(A). It is shown in [37, p.135] that T is a linear bijection and that T is bounded and has the bounded inverse  $T^{-1}$ , if the norm in the sequence space is defined by

(1.3) 
$$\|Tx\| = \sup_{k} \left\| \sum_{i=1}^{k} d_{i}(x)\psi_{i} \right\|$$

Since  $\|J_0x\| \leq \|Tx\|$  for every  $x \in \tilde{X} \subset R(A)$ , where the norms are taken as in (1.2) and (1.3), respectively, one concludes that  $\|J_0\| \leq \|T\| < \infty$ , regardless of n. Hence  $\sup_n \|J_0\| < \infty$ . Space  $\overline{X}$  is homeomorphic with the subspace of the above sequence space consisting of all sequences with zero components from (n+1)-st on. Therefore  $H_0^{-1}\overline{x} = T^{-1}\overline{x}$  for all  $\overline{x} \in \overline{X}$ . Hence  $\|H_0^{-1}\| = \|T^{-1}|_{\overline{X}}\|$  $\leq \|T^{-1}\| < \infty$ , regardless of n, and one concludes that  $\sup_n \|H_0^{-1}\| < \infty$ . Since  $\overline{A}$  is invertible for large n,  $P_{R(\overline{A})} = I$  and thus  $\sup_n \|P_{R(\overline{A})}\| < \infty$ .

Finally, the conditions (iii) are satisfied, since  $\varepsilon = 0$ ,  $\eta_1 \rightarrow 0$  and  $\eta_2 \rightarrow 0$  as  $n \rightarrow \infty$ , by Lemma 1.1, while  $J_0 = H_0$ , and thus  $\sup_n \|J_0^{-1}\| = \sup_n \|H_0^{-1}\| < \infty$ , and  $\|Hx\| \le \|Tx\|$  for every

V.1

 $x \in R(A)$ , Hx = 0 for  $x \in R(A)^{C}$ , by the construction of H and T, which implies  $\sup \|H\| \le \|T\| < \infty$ , regardless of n.

63

All the conditions of Theorem IV.4.1 are satisfied and one concludes that  $\lim_{n\to\infty} \|x^* - J_0^{-1}x^*\| = 0$ , where  $x^*$  is the best approximate solution of Ax = b and  $\bar{x}^*$  is the exact solution (for large n) of the approximate equation (1.1).

The best approximate solution of Ax = b can also be calculated by solving systems of linear algebraic equations (1.3) in the case of a proper splitting A = M+N if, in addition to the proper splitting,  $M^*: \tilde{X} \to \tilde{X}$  for sufficiently large n. All the conditions of Theorem IV.4.1 are still satisfied. The only modification is that u and v in Conditions V and VI are taken as follows:  $u = M^* P_{\tilde{X}} Nx$  and  $v = M^* P_{\tilde{X}} P_{R(A)} b$ . Here  $\tilde{X}$  is still  $span\{\psi_1, \ldots, \psi_n\}$  in R(A). In fact, this requirement on  $\tilde{X}$  can be relaxed. One can choose  $\tilde{X} = span\{\tau_1, \ldots, \tau_n\}$ , where  $\{\tau_1, \ldots, \tau_n\}$  is an arbitrary set of linearly independent vectors in X, provided that  $P_{\tilde{X}} P_{R(A)} =$  $P_{R(A)} P_{\tilde{X}}$  for sufficiently large n and  $\tau_1, \ldots, \tau_n, \tau_{n+1}, \ldots$  is a basis of X. However, with this arbitrary construction of  $\tilde{X}$ , the system (1.1) may be inconsistent for sufficiently large n, in which case the best approximate solution  $\bar{x}^* = \bar{A}^* \bar{b}$  is obtained. Now one can show, using Lemma 1.1, that

 $J_0(N(A)^c \cap \widetilde{X}) = N(\overline{A})^c \cap \overline{X} \text{ and } H_0(N(A) \cap \widetilde{X}) = N(\overline{A}) \cap \overline{X}$ 

for sufficiently large n. These relations imply that the only

V.1
V.1

conditions which need verification, i.e. Conditions II and III are also satisfied.

A Galerkin method for calculating the best approximate solution of Ax = b can be formulated as follows:

- (i) Find a proper splitting A = M + N with N compact.
- (ii) Find a basis  $\{\tau_1, \tau_2, ...\}$  of X such that
- (1.4)  $M^{\dagger}: \widetilde{X} \to \widetilde{X} \text{ and } P_{\widetilde{X}}^{P} R(A) = P_{R(A)}^{P} P_{\widetilde{X}}^{P}$

for sufficiently large n, where  $\tilde{X} = \text{span}\{\tau_1, \ldots, \tau_n\}$ .

(iii) Calculate  $\overline{A} = HAJ_0^{-1}$  and  $\overline{b} = Hb$ . The elements of  $\overline{A} = (a_{ij})$  are determined by  $a_{ij} = e_i(A\tau_j)$ , where  $e_i(x)$  is the i-th coefficient of x in the expansion  $x = \sum_{i=1}^{\infty} e_i(x)\tau_i$ , while the elements of  $\overline{b} = (b_i)$  are determined by  $b_i = e_i(b)$ .

(iv) Calculate the best approximate solution of  $\overline{A}\overline{x} = \overline{b}$ , i.e.  $\overline{x}^* = \overline{A}^+\overline{b}$ .

If A is written as A = I + N, in which case we may not have a proper splitting, then the basis  $\{\tau_1, \tau_2, \ldots\}$  must be chosen as a basis of R(A). The conditions (1.4) are then redundant and  $\overline{A}$  is invertible for sufficiently large n. Example 2.1. The best least squares solution of the equation Ax = b from Example III.3.1 will now be calculated using Galerkin's method.

First, the operator A can be written as A = M + N, where

$$Mx(s) = x(s) - \left(x(s), \sqrt{\frac{2}{\pi}}\sin s\right)\sqrt{\frac{2}{\pi}}\sin s$$
$$Nx(s) = -\frac{1}{2}\left(x(s), \sqrt{\frac{2}{\pi}}\cos s\right)\sqrt{\frac{2}{\pi}}\cos s.$$

Here  $(\cdot, \cdot)$  denotes the inner product in  $L_2[0,\pi]$ . Since M is the orthogonal projection on  $(\operatorname{span}\{\sin s\})^1$  and  $[\operatorname{Nx}(s), \sin s] = 0$ for every  $x \in L_2[0,\pi]$ , one concludes that  $R(N) \subset (\operatorname{span}\{\sin s\})^1 = R(M)$ . Furthermore

$$\|N\| \leq \left(\int_0^{\pi} \int_0^{\pi} \left(\frac{1}{\pi} \cos t \, \cos \xi\right)^2 dt \, d\xi\right)^{\frac{1}{2}} = \frac{1}{2} < 1 ,$$

which implies that A = M + N is a proper splitting, by Corollary IV.1.1.

Second, we choose the following basis of X:

(2.1) 
$$\sqrt{\frac{2}{\pi}}\cos s$$
,  $\sqrt{\frac{2}{\pi}}\sin s$ ,  $\sqrt{\frac{2}{\pi}}\sin 3s$ ,  $\sqrt{\frac{2}{\pi}}\sin 5s$ , ...

The conditions (1.4) are now satisfied for every n.

Third, we calculate  $\overline{A}$  and  $\overline{b}$  for n = 1, 2, ....

$$a_{11} = \left(\sqrt{\frac{2}{\pi}}\cos s, A\left(\sqrt{\frac{2}{\pi}}\cos s\right)\right) = \frac{1}{2}, \text{ since } A(\cos s) = \frac{1}{2}\cos s$$
$$b_1 = \left(s, \sqrt{\frac{2}{\pi}}\cos s\right) = -2\sqrt{\frac{2}{\pi}}.$$

Thus  $\overline{A}\overline{x} = \overline{b}$  for n = 1 is given by  $\overline{x} = -2\sqrt{\frac{2}{\pi}}$ , which gives  $\overline{x}^* = -4\sqrt{\frac{2}{\pi}}$ . For n = 2, the system (1.1) is

$$\begin{pmatrix} \frac{1}{2} & 0 \\ 0 & 0 \end{pmatrix} \bar{\mathbf{x}} = \begin{pmatrix} -2\sqrt{\frac{2}{\pi}} \\ \pi\sqrt{\frac{2}{\pi}} \end{pmatrix}$$

and its best least squares solution is

$$\bar{\mathbf{x}}^{\star} = \begin{pmatrix} -4\sqrt{\frac{2}{\pi}} \\ 0 \end{pmatrix}$$

For n = 3, the system (1.1) becomes

$$\begin{pmatrix} \frac{1}{2} & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} \bar{\mathbf{x}} = \begin{pmatrix} -2\sqrt{\frac{2}{\pi}} \\ \pi\sqrt{\frac{2}{\pi}} \\ \frac{\pi}{3}\sqrt{\frac{2}{\pi}} \end{pmatrix}$$

with the best least squares solution

$$\bar{\mathbf{x}}^{\star} = \begin{pmatrix} -4\sqrt{\frac{2}{\pi}} \\ 0 \\ \frac{\pi}{3}\sqrt{\frac{2}{\pi}} \end{pmatrix}.$$

At the n-th step  $(n \ge 3)$  we obtain

V.2



and the best least squares solution is

$$\bar{x}^{*} = \sqrt{\frac{2}{\pi}} \begin{pmatrix} -4 \\ 0 \\ \frac{\pi}{3} \\ \frac{\pi}{5} \\ \vdots \\ \frac{\pi}{2n-3} \end{pmatrix}$$

Hence

$$J_0^{-1}\bar{x}^* = \frac{2}{\pi} \left[ -4\cos s + \frac{\pi}{3}\sin 3s + \frac{\pi}{5}\sin 5s + \ldots + \frac{\pi}{2n-3}\sin (2n-3) \right].$$

Since the coefficients  $(x^*(s), \tau_1)$ , i = 1, ..., n of the function  $x^*(s) = s - 2 \sin s - \frac{4}{\pi} \cos s$ , in the basis (2.1), are  $-4\sqrt{\frac{2}{\pi}}$ ,  $0, \frac{\pi}{3}\sqrt{\frac{2}{\pi}}$ ,  $\frac{\pi}{5}\sqrt{\frac{2}{\pi}}$ ,  $\dots, \frac{\pi}{2n-3}\sqrt{\frac{2}{\pi}}$  for every n, we conclude that  $J_0^{-1}\bar{x}^* \rightarrow x^*(s)$ , i.e.  $x^*(s)$  is the best least squares solution of Ax = b. The same result has been obtained in Example III.3.1 using iteration in  $L_2[0,\pi]$ .

#### REFERENCES

- N.N. Abdelmalek, "Computing the strict Chebyshev solution of overdetermined linear equations", <u>Mathematics of Computation</u>, <u>31</u> (1977), 974-983.
- [2] N.N Abdelmalek, "An efficient method for the discrete linear L<sub>1</sub> approximation problem", <u>Mathematics of Computation</u>, <u>29</u> (1975), 844-850.
- [3] N. Anderson, "A generalization of the method of averages for overdetermined linear systems", <u>Mathematics of Computation</u>, 29 (1975), 607-614.
- [4] P.M. Anselone, <u>Collectively Compact Approximation Theory</u>, Prentice-Hall, Englewood Cliffs, New Jersey, 1971.
- [5] P.M. Anselone and R.H. Moore, "Approximate solutions of integral and operator equations", Journal of Mathematical Analysis and <u>Applications</u>, 9 (1964), 268-277.
- [6] K.E. Atkinson, "The solution of non-unique linear integral equations", Numerische Mathematik, 10 (1967), 117-124.
- [7] K.E. Atkinson, <u>A Survey of Numerical Methods for the Solution</u> of Fredholm Integral Equations of the Second Kind, SIAM, Philadelphia, Pennsylvania, 1976.
- [8] A. Ben-Israel, "A note on an iterative method for generalized inversion of matrices", <u>Mathematics of Computation</u>, <u>20</u> (1966), 439-440.
- [9] A. Ben-Israel, "On direct sum decompositions of Hestenes algebras", Israel Journal of Mathematics, 2 (1964), 50-54.
- [10] A. Ben-Israel, "On error bounds for generalized inverses", <u>SIAM</u> Journal of Numerical Analysis, 3 (1966), 585-592.

- [11] A. Ben-Israel and Dan Cohen, "On iterative computation of generalized inverses and associated projections", <u>SIAM</u> Journal of Numerical Analysis, 3 (1966), 410-419.
- [12] A. Ben-Israel and T.N.E. Greville, <u>Generalized Inverses</u>, Theory and Applications, Wiley-Interscience, New York, 1974.
- [13] A. Berman and R.J. Plemmons, "Cones and iterative methods for best least squares solutions of linear systems", <u>SIAM</u> Journal of Numerical Analysis, 11 (1974), 145-154.
- [14] J. Blather, P.D. Morris and D.E. Wulbert, "Continuity of setvalued metric projection", <u>Mathematishe Annalen</u>, <u>178</u> (1968), 12-24.
- [15] Paul T. Boggs, "The convergence of the Ben-Israel iteration for nonlinear least squares problems", <u>Mathematics of Computation</u>, <u>30</u> (1976), 512-522.
- [16] J.H. Bramble and A.H. Shatz, "On the numerical solution of elliptic boundary value problems by least squares approximation of the data", <u>Numerical Solutions of Partial Differential</u> <u>Equations, II</u>, B.E. Hubbard, ed., Academic Press, New York (1971), 107-131.
- [17] J.H. Bramble and A.H. Shatz, "Least squares method for 2m-th order elliptic boundary value problems", <u>Mathematics of Compu-</u> tation, 25 (1971), 1-32.
- [18] J.H. Bramble and A.H. Shatz, "Rayleigh-Ritz-Galerkin methods for Dirichlet's problem using subspaces without boundary conditions", <u>Communications on Pure and Applied Mathematics</u>, <u>23</u> (1970), 653-676.
- [19] L.M. Delves and J. Walsh, <u>Numerical Solution of Integral Equations</u>, Clarendon Press, Oxford, 1974.

- [20] David Gay, "Modifying singular values: existence of solutions to systems of nonlinear equations having a possibly singular Jacobian matrix", Mathematics of Computation, 31 (1977), 962-973.
- [21] C.W. Groetsch, <u>Computational Theory of Generalized Inverses</u> of Bounded Linear Operators, Department of Mathematical Sciences, University of Cincinnati, Cincinnati, Ohio, 1975.
- [22] S.P. Gudder and M. Neumann, "Splittings and iterative methods for approximate solutions to singular operator equations in Hilbert spaces", Journal of Mathematical Analysis and Applications, 62 (1978), 272-294.
- [23] Apostolos Hadjidimos, "Accelerated overrelaxation method", Mathematics of Computation, 32 (1978), 149-157.
- [24] Michael P. Hanna, "Generalized overrelaxation and Gauss-Seidel convergence on Hilbert space", <u>Proceedings of the American Mathematical Society</u>, <u>35</u> (1972), 524-530.
- [25] E.N. Houstis and T.S. Papatheodorou, "A collocation method for Fredholm integral equations of the second kind", <u>Mathematics of</u> <u>Computation</u>, <u>32</u> (1978), 159-173.
- [26] Y. Ikebe, "The Galerkin method for the numerical solution of Fredholm integral equations of the second kind", Rept. CNA-S, University of Texas, Austin, Texas, (1970).
- [27] W.J. Kammerer and M.Z. Nashed, "On the convergence of the conjugate gradient method for singular linear operator equations", <u>SIAM</u> Journal of Numerical Analysis, 9 (1972), 165-181.
- [28] W.J. Kammerer and M.Z. Nashed, "Steepest descent for singular linear operators with nonclosed range", <u>Applicable Analysis</u>, 1 (1971), 143-159.

- [29] W.J. Kammerer and R.H. Plemmons, "Direct iterative methods for least squares solutions to singular operator equations", <u>Journal of Mathematical Analysis and Applications</u>, <u>49</u> (1975), 512-526.
- [30] L.V. Kantorovich, "Functional analysis and applied mathematics", Uspehi Matematičeskih Nauk, 3 (1948), 89-195 (Russian).
- [31] L.V. Kantorovich and G.P. Akilov, <u>Functional Analysis in Normed</u> <u>Spaces</u>, Translated from Russian by D.E. Brown, Pergamon Press, Oxford, 1964.
- [32] Konrad Knopp, <u>Infinite Sequences and Series</u>, Translated from German by F. Bagemihl, Dover Publications, New York, 1956.
- [33] M.A. Krasnosel'skii et al., <u>Approximate Solution of Operator</u> Equations, Walters-Noordhoff Publ., Groningen, Holland, 1972.
- [34] W.F. Langford, "The generalized inverse of an unbounded linear operator with unbounded constraints", <u>SIAM Journal of Mathematical</u> <u>Analysis</u> (to appear).
- [35] Charles Lawson and Richard Hanson, <u>Solving Least Squares Problems</u>, Prentice-Hall, Englewood Cliffs, New Jersey, 1974.
- [36] L. Marie Lawson, "Computational methods for generalized inverse matrices arising from proper splittings", <u>Linear Algebra and its</u> <u>Applications, 12</u> (1975), 111-126.
- [37] L.A. Liusternik and V.I. Sobolev, <u>Elements of Functional Analysis</u>, Frederic Ungar Publ. Comp., New York, 1961.
- [38] V. Lovass-Nagy and D.L. Powers, "On under- and over-determined initial value problems", <u>International Journal on Control</u>, <u>19</u> (1974), 653-656.
- [39] J.A. Meijerink and H.A. van der Vorst, "An iterative method for linear systems of which the coefficient matrix is a symmetric M-matrix", Mathematics of Computation, 31 (1977), 148-162.

- [40] R.H. Moore and M.Z. Nashed, "Approximations to generalized inverses of linear operators", <u>SIAM Journal of Applied Mathe-</u> matics, 27 (1974), 1-16.
- [41] F.J. Murray, "On complementary manifolds and projections in spaces L and L", <u>Transactions of the American Mathematical</u> <u>Society</u>, <u>41</u> (1937), 138-152.
- [42] N.I. Muskhelishvili, <u>Singular Integral Equations</u>, Noordhoff, Groningen, Holland, 1946.
- [43] M.Z. Nashed, "Generalized inverses, normal solvability, and iteration for singular operator equations", in <u>Nonlinear Func-</u> <u>tional Analysis and Applications</u> (L.B. Rall, editor), Academic Press, New York (1971), 311-359.
- [44] M.Z. Nashed, "Perturbations and approximations for generalized inverses and linear operator equations", in <u>Generalized Inverses</u> <u>and Applications</u> (M.Z. Nashed, editor), Academic Press, New York (1976), 325-396.
- [45] M.Z. Nashed and Grace Wahba, "Convergence rates of approximate least squares solutions of linear integral and operator equations of the first kind", Mathematics of Computation, 28 (1974), 69-80.
- [46] T.G. Newman and P.L. Odell, "On the concept of a p-q generalized inverse of a matrix", <u>SIAM Journal of Applied Mathematics</u>, <u>17</u> (1969), 520-525.
- [47] W.V. Petryshyn, "On generalized inverses and on the uniform convergence of  $(I \beta K)^n$  with application to iterative methods", Journal of Mathematical Analysis and Applications, 18 (1967), 417-439.
- [48] W.V. Petryshyn, "On the generalized overrelaxation method for operation equations", Proceedings of the American Mathematical Society, 14 (1963), 917-924.

- [49] W.V. Petryshyn, "On the extrapolated Jacobi or simultaneous displacements method in the solution of matrix and operator equations", Mathematics of Computations, 19 (1965), 37-56.
- [50] J.L. Phillips, "Collocation as a projection method for solving integral and other operator equations", Thesis, Purdue University (1969).
- [51] P.M. Prenter, "Collocation method for the numerical solution of integral equations", <u>SIAM Journal of Numerical Analysis</u>, <u>10</u> (1973), 570-581.
- [52] C.R. Rao and S.K. Mitra, <u>Generalized Inverse of Matrices and</u> its Applications, Wiley, New York, 1971.
- [53] W. Rudin, Functional Analysis, McGraw-Hill, New York, 1973.
- [54] S.M. Serbin, "Computational investigations of least-squares type methods for the approximate solutions of boundary value problems", <u>Mathematics of Computation</u>, <u>29</u> (1975), 777-793.
- [55] I. Singer, <u>Bases in Banach Spaces I</u>, Springer-Verlag, New York, 1970.
- .[56] I. Stakgold, <u>Boundary Value Problems of Mathematical Physics</u>, Volume I, Macmillan Series in Advanced Mathematics and Theoretical Physics, Macmillan, New York, 1969.
- [57] G.W. Stewart, "On the continuity of the generalized inverse", SIAM Journal of Applied Mathematics, 17 (1969), 33-45.
- [58] A.E. Taylor, <u>Introduction to Functional Analysis</u>, Wiley, New York, 1958.
- [59] K.S. Thomas, "On the approximate solution of operator equations", Numerische Mathematik, 23 (1975), 231-239.
- [60] R.S. Varga, <u>Matrix Iterative Analysis</u>, Prentice-Hall, Englewood Cliffs, New Jersey, 1962.

- [61] R.S. Varga, "Extensions of the successive overrelaxation theory with applications to finite element approximations", Department of Mathematics, Kent University, Kent, Ohio.
- [62] R.S. Varga, <u>Functional Analysis and Approximation Theory in</u> <u>Numerical Analysis</u>, SIAM, Regional Conference Series in Applied Mathematics, No.3, Chapter 6, 1971.
- [63] H. Wolkowicz and S. Zlobec, "Calculating the best approximate solution of an operator equation", <u>Mathematics of Computation</u>, October (1978), (to appear).
- [64] Gerhard Zielke, <u>Beiträge: Zur Theorie Und Berechnung Von</u> <u>Verallgemeinerten Inversen Matrizen</u>, Martin-Luther-Universität, Halle-Wittenberg, 1977.
- [65] S. Zlobec, "On computing the best least squares solutions in Hilbert spaces", <u>Rendiconti del Circolo Matematico di Palermo</u>, 26 (1977), 1-15.

#### ADDITIONAL REFERENCES

- [66] M.B. Gagua, "On the approximate solution of linear equations", Trudy Vycisl. Centra Akad. Nauk Gruzin. SSS R, 3 (1963), 27-47.
- [67] V.V. Ivanov, <u>The Theory of Approximate Methods and Their Applica-</u> tion to the Numerical Solution of Singular Integral Equations, Noordhoff International Publishing, Leyden, 1976.
- [68] J.J. Koliha, "On the iterative solution of linear operator equations with selfadjoint operators", <u>J. Australian Math. Soc.</u>, 13 (1972), 241-255.
- [69] J.J. Koliha, "Convergent and stable operators and their generalizations", J. Math. Anal. Appl., 43 (1973), 778-794.

- [70] J.J. Koliha, "Series representation of pseudoinverses and partial inverses of operators", <u>Bull. Amer. Math. Soc.</u>, <u>80</u> (1974), 325-328.
- [71] J.J. Koliha, "Power convergence and pseudoinverses of operators between Banach spaces", J. Math. Anal. Appl., 48 (1974) 446-469.
- [72] M.Z. Nashed and G.F. Votruba, "A unified approach to generalized inverses of linear operators: I. Algebraic topological and projectional properties", <u>Bull. Amer. Math. Soc.</u>, <u>80</u> (1974),
   [73]
- [73] M.Z. Nashed and G.F. Votruba, "A unified operator theory of generalized inverses", in <u>Generalized Inverses and Application</u>, M.Z. Nashed, ed., Academic Press, New York, 1976, 1-109.
- [74] L. Rakovshchik, "Approximate solutions of equations with normally solvable operators", U.S.S.R. Computational Math. and Math. Phys., 6 (1966), 3-11.
- [75] G.F. Votruba, "Generalized Inverses and Singular Equations in Functional Analysis", Doctoral Dissertation, University of Michigan, Ann Arbor, Michigan, 1963.

-74a

PART B

.

# CONVEX PROGRAMMING

## I. INTRODUCTION

75

1.

Consider the convex program

(P)  

$$f^{\circ}(x) \neq \min$$
s.t.  

$$f^{k}(x) \leq 0, k \in P = \{1, \dots, m\},$$

where  $f^k : X \rightarrow R$ ,  $k \in \{0\} \cup P$ , are continuous convex functions and X is a locally convex linear Hausdorff space. The convexity assumption represents the natural framework for the treatment of optimization problems. One can now develop an elegant theory that allows a wide range of applications. The optimality conditions and algorithms developed, using this theory, usually assume that Slater's condition or some other constraint qualification is satisfied. However, recently Ben-Israel, Ben-Tal and Zlobec [10], [11] have used the cones of directions of constancy of the constraints, to characterize optimality without assuming a constraint qualification. Algorithms for solving (P) were then given in [14]. In this part of the thesis we continue this approach and study the role played by the cones of directions of constancy when deriving optimality conditions and algorithms. In particular, we use the approach of Gould and Tolle [30] and study the relationship between the geometry of the feasible set and the analytic properties of the constraints. We will see that the cones of directions of constancy provide exactly

I.1

the missing analytic information needed in order to describe the feasible set and thus characterize optimality.

76

In Chapter II we summarize some basic results in convex analysis and optimization. We also introduce the set of 'badly behaved' constraints at x,  $p^{b}(x)$ , and present a closedness criterion for the sum of two, not necessarily convex, cones. In Chapter III we first give a lemma which shows the different relationships between the cones of directions of constancy, the tangent cone of the feasible set and the cone of subgradients. Using this lemma, we derive some old and some new characterizations of optimality. Chapter IV examines the notions of constraint qualifications and regularization. The most important results here are the regularization technique in Theorem IV.4.2 and the weakest constraint qualifications in Theorem IV.3.1. Finally, in Chapter V we present the Method of Reduction, which solves program (P) with faithfully convex constraints. The method includes an algorithm for finding the cone of directions of constancy of a faithfully convex function and also a generalization of the algorithm for finding the equality set, P, given by Abrams and Kerzner [3]. We conclude with several applications and examples.

I.1

# II. PRELIMINARIES

# 1. Notation

We will be dealing with real locally convex linear spaces X,Y,... and various mappings defined on them. We list below some of the abbreviations and symbols used and the sections where their meanings are explained.

tvs	a topological linear space
les	a locally convex (Hausdorff) space
l(X,Y)	the set of all linear operators from $X$ into $Y$
X'	the topological dual of $X$ , with the weak topology
$w^*$ -closed	weak* closed
R <sup>n</sup>	the real Euclidean n-space
R(B)	the range space of the operator B
N(B)	the null space of the operator B
κ <sup>⊥</sup>	the annihilator (in X') of a set K (in X)
span K	the vector subspace spanned by K
ĸ	topological closure
bdry K	the topological boundary of K
int K	the topological interior of K
cone	a set closed under nonnegative scalar multiplication
blunt cone	a set closed under positive scalar multiplication

II.1

cone K	the convex cone generated by K
conv K	the convex hull of K
ø	the empty set
φ•x	the value of $\phi(x)$ or $x(\varphi)$ , where $x \in X$ and $\varphi \in X'$
(P)	the convex program, II.2
(P <sub>r</sub> )	the regularized program, IV.4
Р	the indexing set of the constraints, II.2
S	the feasible set, II.2
P(x)	the binding constraints, II.2
P <sup>=</sup>	the equality set, II.2
$P^{<}(x) = P(x) \setminus P^{2}$	= , II.2
$P^{b}(x)$	the 'badly behaved' set, II.5
dist(x,K)	the distance from x to the set K, II.5
E <sub>k</sub> (x)	the directions of vanishing directional derivative, II.5
D <sup>&lt;,≤,=,&gt;</sup> f	the cone of directions of decrease, nonincrease,
	constancy and increase, II.3
$\nabla f(x;d)$	the directional derivative, II.4
$\partial f(x)$	the subdifferential, II.4
$\nabla f(x)$	the gradient, II.4
M*	the polar of M, II.6
F°(x)	the set of continuous, convex functions which achieve
	their minimum in Sat x II 7

•

$DF^{\circ}(x) = \{ \phi \in \partial f^{\circ}(x) \}$	): $\mathbf{f}^{\circ} \in \mathbf{F}^{\circ}(\mathbf{x})$ }, II.7
$C_{p(x)}(x)$	the linearizing cone, II.7
${}^{B}P(x)(x)$	the cone of subgradients, II.7
T(S,x)	the cone of tangents, II.8
K-T point	Kuhn-Tucker point, IV.1
K-T conditions	Kuhn-Tucker conditions, IV.1
regular point	IV.2
CQ	constraint qualification, IV.2
WCQ	weakest constraint qualification, IV.3
(PL)	the lexicographic problem, V.5
(PP)	the Pareto optimal problem, V.5
(PS)	the semi-infinite problem, V.5

.

## 2. The Convex Program

We consider the convex programming problem

(P)  

$$f^{\circ}(x) \neq \min$$
  
s.t.  
 $f^{k}(x) \leq 0$ ,  $k \in P = \{1, \dots, m\}$ ,

where  $f^k: X \to R$  are continuous convex functions, defined on the *lcs* X, for all  $k \in \{0\} \cup P$ . (Without loss of generality, we assume that none of the functions is constant.) Unless otherwise specified, we assume that the feasible set

$$S = \{x \in X: f^k(x) \le 0, \text{ for all } k \in P\}$$

is not empty. The set of binding constraints, at  $x \in S$ , is

$$P(x) = \{k \in P: f^k(x) = 0\}.$$

An important subset of P, independent of x, is the equality set

$$P^{=} = \{k \in P: f^{k}(x) = 0, \text{ for all } x \in S\}.$$

(See e.g. Zoutendijk [53], Rockafellar [43] and Abrams and Kerzner [3].) This is the set of indices k for which the constraint  $f^k$  vanishes on the entire feasible set. We then denote

$$P^{<}(\mathbf{x}) = P(\mathbf{x}) \setminus P^{=}$$
.

Note that unlike  $P^{=}$ ,  $P^{<}(x)$  depends on x.

#### 3. Cones of Directions and Faithfully Convex Functions

Following Ben-Israel, Ben-Tal and Zlobec [11], we define the relations

"relation" is "=", "<", "≤" or ">",

by

$$D_{f}^{"relation"}(x) = \{ d \in X: \text{ there exists } \overline{\alpha} > 0 \text{ with} \\ f(x + \alpha d) \text{ "relation" } f(x), \text{ for all } 0 < \alpha \leq \overline{\alpha} \}.$$

These are the cones of directions of constancy, descent, nonincrease and increase respectively. For simplicity of notation, we let

$$D_k^{"relation"}(x) = D_{fk}^{"relation"}(x)$$

and

$$D_{\Omega}^{"relation"}(x) = \bigcap_{k \in \Omega} D_{k}^{"relation"}(x), \text{ for } \Omega \subset \mathcal{P}.$$

<u>Remark 3.1.</u> For a function f in the class of faithfully convex functions, the cone  $D_f^{=}(x)$  is a subspace independent of x. Following Rockafellar [44], we say that a convex function f is <u>faithfully convex</u> if: f is affine on a line segment only if it is affine on the whole line containing that segment. If  $X = R^n$ , then Rockafellar has shown that f is faithfully convex if and only if it is of the form

(3.1) 
$$f(x) = h(Ax + b) + a \cdot x + \alpha$$

where  $A \in \mathbb{R}^{m \times n}$ ,  $b \in \mathbb{R}^{m}$ ,  $a \in \mathbb{R}^{n}$ ,  $\alpha \in \mathbb{R}$  and the function  $h:\mathbb{R}^{m} \to \mathbb{R}$ is strictly convex. It is easy to see that  $D_{f}^{=}(x) = N\begin{pmatrix} A \\ a^{t} \end{pmatrix}$ and is a subspace independent of x.

In the following lemma we collect some properties of the directions. We also show directly that the cone of directions of constancy of a continuous faithfully convex function on X, a *lcs*, is a subspace independent of  $x \in X$ .

Lemma 3.1. Suppose that  $f:X \rightarrow R$  is a convex function and  $x \in S$ . Then:

a)  $D_{f}^{\leq}(x)$  is a convex cone,  $D_{f}^{\leq}(x)$  is a convex blunt cone. b) conv  $D_{f}^{=}(x) \subset D_{f}^{\leq}(x)$ .

c) 
$$D_{P(x)}^{\leq}(x) = D_{P^{\leq}(x)}^{\leq}(x) \cap D_{p^{=}}^{=}(x)$$
  
=  $D_{P^{\leq}(x)}^{\leq}(x) \cap conv D_{p^{=}}^{=}(x)$ .

d)  $D_{p^{\pm}}^{\pm}(x) \cap D^{<}(x) \neq \emptyset$ .

e) If f is both faithfully convex and continuous, then  $D_{f}^{=}(x) = D_{f}^{=}$  is a subspace of X, independent of x. Proof. For (a)-(d), see e.g. [11], [12].

e) First, let us show that  $D_{f}^{=}(x)$  is a subspace. Suppose that  $d_{1}, d_{2} \in D_{f}^{=}(x)$  and let  $d = d_{1} + d_{2}$ . If  $\alpha \in \mathbb{R}$ , then  $f(x + \alpha d) = f(\frac{1}{2}(x + 2\alpha d_{1}) + \frac{1}{2}(x + 2\alpha d_{2}))$   $\leq \frac{1}{2}f(x + 2\alpha d_{1}) + \frac{1}{2}f(x + 2\alpha d_{2})$ , since f is convex = f(x), since  $d_{1}, d_{2} \in D_{f}^{=}(x)$  and f is faithfully convex.

Therefore f is bounded above on the whole line  $x + \alpha d$ ,  $\alpha \in \mathbb{R}$ , which implies that f is constant on this line (see e.g. Rockafellar [41, p.69]. Thus,  $d \in D_{\mathbf{f}}^{=}(\mathbf{x})$ . This shows that  $D_{\mathbf{f}}^{=}(\mathbf{x})$  is closed under addition. That  $D_{\mathbf{f}}^{=}(\mathbf{x})$  is closed under scalar multiplication is clear, from the definition of a faithfully convex function.

We have left to show that  $D_f^{=}(x) = D_f^{=}$ , i.e. it is independent of x. Suppose that x, y  $\in X$  and  $d \in D_f^{=}(x)$ . We will show that  $d \in D_f^{=}(y)$ .

Case (i): Suppose that  $f(y) \le f(x)$ . We will first show that

$$f(y + \alpha d) \leq f(x)$$
, for all  $\alpha \in \mathbb{R}$ .

Let  $\alpha \in \mathbb{R}$  and  $1 > t_k > 0$  with  $t_k \to 0$  as  $k \to \infty$ . Consider the directions  $z^k = \alpha d + t_k (x - y)$  and let  $\gamma_k = 1/t_k$ . Then

$$\begin{split} \mathbf{f}(\mathbf{y}) &\leq \mathbf{f}(\mathbf{x}) \\ &= \mathbf{f}(\mathbf{x} + \gamma_k \alpha d) , \quad \text{since } d \in \mathbb{D}_{\mathbf{f}}^{=}(\mathbf{x}) \\ &= \mathbf{f}(\mathbf{y} + \gamma_k \mathbf{z}^k) . \end{split}$$

By convexity of f and since  $\gamma_k^{} > 1$  , we conclude that

$$f(y + z^k) \leq f(x)$$

and thus, by continuity of f, we see that

$$f(y + \alpha d) = \lim_{k \to \infty} f(y + z^k) \le f(x)$$
.

This shows that f is bounded on the line  $y + \alpha d$ ,  $\alpha \in \mathbb{R}$ , and therefore, f is constant on this line, i.e.  $d \in D_{f}^{=}(y)$ .

Case (ii): Suppose that f(x) < f(y). By a similar argument to case (i), we see that

$$f(y + \alpha d) = \lim_{k \to \infty} f(y + \alpha d + t_k(x - y)) \le f(y)$$

for all  $\alpha \in \mathbb{R}$ , i.e.  $d \in D_{f}^{=}(y)$ .

<u>Remark 3.2.</u> For faithfully convex functions f, on  $\mathbb{R}^n$ , one can calculate  $D_f^=$  explicitly, see [49] and Section V.2. The class of faithfully convex functions is quite large. It includes all analytic convex, as well as strictly convex, functions. The algorithms in Chapter V will deal mainly with these functions. In general, however, the cone of directions of constancy may not be a subspace. In fact it may be neither convex nor closed. (See [10] for examples.)

### 4. Subdifferentiability

We have assumed that our functions are convex, but not necessarily differentiable. Nonsmooth, or nondifferentiable, functions occur quite often in convex analysis. Applications for these functions arise in approximation theory, e.g. Dem'yanov and Malozemov [23], duality theory, e.g. Rockafellar [42] and semi-infinite programming, e.g. Ben-Tal, Kerzner and Zlobec [13]. (See also Clarke [19] and Pshenichnyi [40].) For convex functions, it is possible to develop a complete calculus without assuming differentiability, e.g. Rockafellar [41], Pshenichnyi [40] and Holmes [34]. We now recall some concepts dealing with directional derivatives and subgradients of a convex function f, defined on the *les* X.

The <u>directional derivative</u> of f at x, in the direction d, is defined as

$$\nabla f(x;d) = \lim_{t \neq 0} \frac{f(x+td) - f(x)}{t}$$

Convex functions have the useful property that the directional derivatives exist universally, e.g. [41, Theorem 23.1].

A vector  $\phi \in X'$  is said to be a <u>subgradient</u> of a convex function f, at the point x, if

$$f(z) \ge f(x) + \phi \cdot (z - x)$$
, for all  $z \in X$ .

II.4

The set of all subgradients of f at x is then called the subdifferential of f at x and is denoted by  $\partial f(x)$ .

If the directional derivative of f at x is a continuous linear functional, i.e. if  $\nabla f(x; \cdot) = \phi \in X'$ , then

$$\phi \cdot d = \lim_{t \to 0^+} \frac{f(x + td) - f(x)}{t}$$

and  $\varphi$  is called the gradient of f at x and denoted  $\nabla f(x)$  . Note that in this case

$$\partial f(x) = \{ \nabla f(x) \}$$
.

We collect some useful properties in the following lemma. For more details and proofs, see e.g. [34], [41].

Lemma 4.1. Suppose that  $f:X \rightarrow R$  is convex. Then

a) ∇f(x;•) is a finite, sublinear functional on X, for all x ∈ X.
If, in addition, f is continuous at x, then:
b) ∇f(x;d) = max{φ•d:φ ∈ ∂f(x)}

and

c)  $\partial f(x)$  is a non-empty,  $w^*$ -compact convex subset of X'.

The next lemma presents some of the relations that exist between the subgradients and the directions introduced in Section 3. For the proofs see Ben-Tal and Ben-Israel [12].

Lemma 4.2. Suppose that 
$$f:X \rightarrow R$$
 is convex  
a)  $D_f^{\leq}(x) = \{d \in X: \nabla f(x;d) < 0\}$ .  
If  $\nabla f(x)$  exists, then:  
b)  $D_f^{\leq}(x) = \{d \in X: \nabla f(x) \cdot d < 0\}$   
and  
c) conv  $D_f^{\equiv}(x) = D_f^{\equiv}(x) \subset \{d \in X: \nabla f(x) \cdot d = 0\}$ .

5. The 'Badly Behaved' Constraints

For  $x \in S$ , let

$$P^{\mathbf{b}}(\mathbf{x}) \stackrel{\Delta}{=} \{\mathbf{k} \in P^{=}: (\mathbf{E}_{\mathbf{k}}(\mathbf{x}) \cap \mathbf{D}_{\mathbf{k}}^{>}(\mathbf{x}) \cap \mathbf{C}_{P(\mathbf{x})}^{-}(\mathbf{x})) \setminus \overline{\mathbf{D}_{p}^{=}}(\mathbf{x}) \neq \emptyset \},\$$

where

$$E_{k}(\mathbf{x}) = \{ \mathbf{d} \in \mathbf{X} : \nabla \mathbf{f}^{k}(\mathbf{x}; \mathbf{d}) = 0 \}$$

and

$$C_{P(x)}(x) = \{d \in X: \nabla f^k(x;d) \le 0, \text{ for all } k \in P(x)\}.$$

(See Section 7 below for further properties of the 'linearizing cone',  $C_{P(x)}(x)$ .) We call  $P^{b}(x)$  the <u>set of 'badly behaved'</u> <u>constraints</u> at  $x \in S$ , for program (P).

Then

The set  $P^b(x)$  is the set of constraints that creates problems in the Kuhn-Tucker theory. We can see that these are the constraints in  $P^{\overline{}}$ , whose analytic properties (given by the directional derivatives) do not fully describe the geometry of the feasible set (given by the feasible directions). It will be shown in IV.3 that

$$\mathcal{P}^{b}(\mathbf{x}) = \emptyset$$

is an essential condition for the Kuhn-Tucker theory to hold at x, independent of  $f^0$ . The set  $P^b(x)$  will also be used in the characterizations of optimality in III.4.

Abrams and Kerzner [3] have given an algorithm that finds the set  $P^{=}$ . (See V.3 for a modified version of their algorithm.) Once  $P^{=}$  is found, then, for any given index  $k_o \in P^{=}$ , we see that  $k_o \in P^{b}(x)$  if and only if the system

$$\nabla f^{k_{o}}(x;d) = 0$$
  

$$\nabla f^{k}(x;d) \leq 0, \quad \text{for all } k \in P(x) \setminus k_{o},$$
  

$$d \notin D_{k_{o}}^{=}(x) \cup \overline{D_{p}^{=}(x)}$$

is consistent. (Note that when  $D_{k_o}^{=}(x)$  is closed, then  $\overline{D_{p}^{=}(x)} \subset D_{k_o}^{=}(x)$ . This simplifies the above system and thus, the corresponding definition for the 'badly behaved' set.)

The set  $P^{b}(x)$  is not equal to  $P^{=}$  in general. In fact, if (5.1)  $E_{k}(x) = D_{k}^{=}(x)$ ,

(5.2)

then  $f^k$  is 'never badly behaved' at x, i.e.  $k \notin P^b(x)$  independent of the other constraints. This class of functions which are 'never badly behaved' at x includes all continuous linear functionals on X. Furthermore, if  $X = R^n$ ,  $\nabla f(x) \neq 0$  and f is a strictly convex function of one variable, considered as a function on  $R^n$ (i.e. if the <u>restriction</u> of f to  $R^1$  is strictly convex), then f is a nonlinear function which is 'never badly behaved' at x. (See Ben-Israel, Ben-Tal and Zlobec [10] for definitions and properties of functions whose restrictions are strictly convex.)

The class of functions which are 'never badly behaved' at x also includes the 'distance' functions defined below. We will see, in IV.4, that every program (P) can be 'regularized' by the addition of one such 'distance' function.

Lemma 5.1. Suppose that X is a Hilbert space, K is a nonempty closed convex cone in X,  $x \in S$  and  $k \in P$ . If, for  $y \in X$ ,

$$f^{K}(y) = dist(y - x, K)$$
$$\triangleq \min_{z \in K} \| (y - x) - z \|$$

then  $f^k$  is a convex function on X, which is 'never badly behaved' at x. Furthermore,

(5.3)  $\nabla f^k(x;d) = \begin{cases} 0 & \text{if } d \in K \\ positive & otherwise. \end{cases}$ 

<u>Proof.</u> First, let us show that the function  $f^k$  is convex. Suppose that  $y^1$ ,  $y^2 \in X$ . Since X is a Hilbert space and K is a nonempty, closed convex set, there exist unique points  $z^1$  and  $z^2$ , in K, which are closest to  $y^1 - x$  and  $y^2 - x$  resp., see e.g. Rudin [46, p.78]. Thus

$$f^{k}(y^{1}) = ||y^{1} - x - z^{1}||$$
,  $f^{k}(y^{2}) = ||y^{2} - x - z^{2}||$ .

Let  $0 \leq \lambda \leq 1$ . Then

$$\lambda z^{1} + (1 - \lambda) z^{2} \in K$$

and

$$f^{k}(\lambda y^{1} + (1 - \lambda)y^{2}) = dist(\lambda y^{1} + (1 - \lambda)y^{2} - x, K)$$

$$\leq \| (\lambda y^{1} + (1 - \lambda)y^{2} - x) - (\lambda z^{1} + (1 - \lambda)z^{2}) \|$$

$$\leq \lambda \| (y^{1} - x) - z^{1} \| + (1 - \lambda) \| (y^{2} - x) - z^{2} \|$$

$$= \lambda f^{k}(y^{1}) + (1 - \lambda) f^{k}(y^{2}) .$$

Therefore  $f^k$  is a convex function on X. Now let  $d \in X$ . Then

$$\nabla f^{k}(x;d) = \lim_{t \neq 0} \frac{f^{k}(x+td) - f^{k}(x)}{t}$$
$$= \lim_{t \neq 0} \frac{\text{dist}(td,K)}{t}$$

= dist(d,K), since K is a cone.

This yields (5.3) and further implies that (5.1) holds. Therefore  $f^k$  is 'never badly behaved' at x.

Example 5.1. Consider the program (P) with the single constraint in one variable,  $f^1(x) \le 0$ , where

$$f^{1}(x) = \begin{cases} x^{2} & \text{if } x \ge 0, \\ 0 & \text{otherwise.} \end{cases}$$

Then

$$P^{b}(x) = \begin{cases} \{1\} & \text{if } x = 0, \\ \emptyset & \text{otherwise.} \end{cases}$$

However, if

$$f^{1}(x) = \begin{cases} x^{2} + x & \text{if } x \ge 0, \\ 0 & \text{otherwise,} \end{cases}$$

then  $P^{b}(x) = \emptyset$  for all x, i.e.  $f^{1}$  is not 'badly behaved' at x, though  $1 \in P^{=}$ .

Example 5.2. Now consider the three functions

$$f^{1}(x) = \begin{cases} (x-1)^{2} & \text{if } x \ge 1, \\ 0 & \text{otherwise,} \end{cases}$$
$$f^{2}(x) = \begin{cases} x^{2} & \text{if } x \ge 0, \\ 0 & \text{otherwise,} \end{cases}$$

and

$$f^{3}(x) = \begin{cases} x^{2} + x & \text{if } x \ge 0, \\ 0 & \text{otherwise.} \end{cases}$$



If the program (P) has just the two constraints  $f^1$  and  $f^2$ , then

$$P^{b}(\mathbf{x}) = \begin{cases} \{2\} & \text{if } \mathbf{x} = 0, \\ \emptyset & \mathbf{x} \neq 0 \text{ and } \mathbf{x} \in S. \end{cases}$$

If, however, the program (P) has all three constraints, then

As mentioned above, we shall see that, (5.4) is essential for the Kuhn-Tucker theory to hold, independent of the choice of the objective function  $f^0$ .

# 6. Polar Sets and Closedness Criteria

In this section we collect some useful results on polar sets. These results can be found in e.g. Girsanov [28] and Holmes [34]. See also Borwein [16]. We also present a closedness criterion for the sum of two cones.

Recall that for  $M \subseteq X$  and X a *lcs*, the polar of M is

 $M^* = \{ \phi \in X' : \phi \cdot x \ge 0 \text{ for all } x \in M \}.$ 

 $M^*$  is then a  $w^*$ -closed convex cone in X'. However, if  $M \subset X'$ , then we define its polar to be

$$M^* = \{x \in X : \phi \cdot x \ge 0 \text{ for all } \phi \in M\}.$$

 $M^*$  is now a *w*-closed convex cone in X.

Lemma 6.1. Suppose that K and L are subsets of X and C is a subset of X'. Then:

a)  $K \subset L$  implies  $L^* \subset K^*$ .

b)  $K^* = (\overline{conv} K)^*$ ,  $C^* = (\overline{conv} C)^*$ ,  $K^{**} = \overline{cone} K$  and  $C^{**} = \overline{cone} C$ .

If, in addition, K and L are closed convex cones, then c)  $(K \cap L)^* = \overline{K^* + L^*}$ .

The following closedness criteria will be used in deriving the optimality conditions in Chapter III.

**II.6** 

Lemma 6.2. Suppose that L is a closed cone in X, C is a compact subset of X not containing the origin, and K is the cone generated by C, i.e.

$$\begin{array}{rcl} \mathsf{K} = & \cup & \lambda \mathsf{C} \\ & \lambda \ge 0 \end{array}$$

If

 $(6.1) C \cap (-L) = \emptyset,$ 

then

<u>Proof.</u> Suppose that the net  $k^n + l^n \rightarrow p$ , where  $k^n \in K$  and  $l^n \in L$ . We need to show that  $p \in K + L$ . For  $0 \neq \phi \in K$ , let

 $\|\phi\| \triangleq \inf\{t > 0: t^{-1}\phi \in C\}.$ 

Note that  $\|\phi\| < \infty$ , for all  $0 \neq \phi \in K$ , by the definition of K. (In fact,  $\|\cdot\|$  is the <u>Minkowski functional</u> of C and is a seminorm on X, when C is a balanced convex absorbing set.)

Case (i): Suppose that sup  $||k^n|| < \infty$ . Then

$$k^{n} = t_{n}c^{n}$$
, for some  $c^{n} \in C$ ,  $t_{n} > 0$  and  $\sup t_{n} < \infty$ .

Since C is a compact set, we can assume that  $k^n \rightarrow k \in K$ . Therefore

$$\lim l^n = \lim (p - k^n) = p - k$$

Let  $\ell \triangleq p - k$ . Then  $\ell \in L$ , since L is closed, and

$$p = k + \ell \in K + L.$$

Case (ii): Suppose that  $\sup \|k^n\| = \infty$ . We can assume that  $\|k^n\| \to \infty$ , or else we can extract a bounded subnet and use case (i). Let  $\alpha_n = 1/\|k^n\|$ . Then  $\alpha_n \to 0$  and  $\|\alpha_n k^n\| = 1$  for all n. As in case (i), we can assume that

$$\alpha_n k^n \neq k \in K$$
.

Furthermore, since  $\|\alpha_n k^n\| = 1$ , we see that

$$\alpha_n k^n \in C$$
, for all n.

This implies that  $k \in C$ . Now, since

$$k^{n} + \ell^{n} \rightarrow p \text{ and } \alpha_{n} \rightarrow 0$$
,

we see that

$$\alpha_n k^n + \alpha_n \ell^n \rightarrow \lim \alpha_n p = 0$$
,

i.e.

$$\alpha_n \ell^n \rightarrow \lim (-\alpha_n k^n) = -k$$
,

and  $-k \in L$ , since L is closed. Thus

$$(6.2) k \in C \cap (-L),$$

contradicting the hypothesis. This implies that  $\sup \|k^n\| < \infty$  and we can apply case (i).

Note that since X' with the  $w^*$ -topology is a lcs, the lemma holds when X is replaced with X'.

Remark 6.1. The condition

$$K \cap (-L) = \{0\}$$

alone, does not necessarily imply that

K+L is closed.

Halmos [33, p.28] has given an example of two subspaces, in a Hilbert space, with zero intersection and a non-closed sum.<sup> $\dagger$ </sup>

Closedness of a sum is related to closedness of the linear image of a set. Suppose that  $U \subset X$  and  $T:X \rightarrow Y$  is a bounded linear operator, where X and Y are Banach spaces. If T has closed range, then the linear image TU is closed if and only if U+N(T)is closed, see e.g. Atteia [5] and Holmes [35, p.142]. Therefore Lemma 6.2 implies that

### TU is closed

if, T has closed range, U is a cone generated by a compact set C which does not contain the origin and

 $C \cap N(T) = \emptyset$ .

<sup>+</sup> I would like to thank Professor Rockafellar for pointing out this example to me.

The above criteria is used to prove the existence of solutions in optimal control and spline approximation problems, see e.g. [35, Section 21], [21] and [27].

Remark 6.2. Lemma 6.2 holds if (6.1) is weakened to read

 $C \cap (-L) \cap (-bdry L) = \emptyset$ .

<u>Proof.</u> Following the proof of Lemma 6.2, we see that we fail to obtain a contradiction if

$$k \in C \cap (-int L)$$

in (6.2). Therefore, in this case, we need to show that  $p \in K+L$ . We accomplish this by showing that

 $p \in int (K + L)$ .

Since  $k \in -int L$ , we can find a convex neighborhood of the origin, *N*, such that

$$-k + n \subset L$$
.

(Recall that every *lcs* X has a convex local base, see e.g. [45].) It is therefore sufficient to show that the open set

 $p + \Pi \subset \overline{K + int L}$ .

Let
Then  $q \in \mathcal{N}$  and since  $\alpha_n \ell^n \to -k$ ,  $\alpha_n \to 0$  and  $\mathcal{N}$  is convex, we can assume that

$$\alpha_n \ell^n + \alpha_n q \in -k + \eta \subset int L$$
, for all n.

Therefore

$$l^{n} + q = \frac{1}{\alpha_{n}}(\alpha_{n}l^{n} + \alpha_{n}q) \in int L, \text{ for all } n,$$

since L is a cone. This implies that

$$k^{n} + \ell^{n} + q \rightarrow p + q = \overline{p},$$

where  $k^n \in K$  and  $l^n + q \in int L$ , i.e.  $\overline{p} \in \overline{K + int L}$ .

### 7. Some Special Cones

We now present some well-known definitions of cones used in mathematical programming, see e.g. Gould and Tolle [30] and Abadie [1]. However, the definitions are stated here in terms of subgradients.

By  $F^{\circ}(x)$ , we denote the cone of all continuous convex objective functions  $f^{\circ}$  with the property that x minimizes  $f^{\circ}$  over S. Then the cone

98

$$DF^{\circ}(x) \triangleq \{\phi \in X' : \phi \in \partial f^{\circ}(x), \text{ for some } f^{\circ} \in F^{\circ}(x)\}.$$

For every subset  $\Omega$  of P(x), the <u>linearizing cone</u> at  $x \in S$ , with respect to  $\Omega$ , is

$$C_{\Omega}(\mathbf{x}) = \{ \mathbf{d} \in \mathbf{X} : \phi \cdot \mathbf{d} \le 0 \text{ for all } \phi \in \partial \mathbf{f}^{k}(\mathbf{x}) \\ \text{and all } \mathbf{k} \in \Omega \}.$$

By Lemma 4.1(b), we see that

$$C_{\Omega}(x) = \{d \in X : \nabla f^{k}(x;d) \leq 0 \text{ for all } k \in \Omega\}.$$

This formulation corresponds with the definition of the linearizing cone given in Section 5.

The cone of subgradients at x is

$$B_{\Omega}(\mathbf{x}) = \{ \phi \in \mathbf{X}^{*} : \phi = \sum_{k \in \Omega} \lambda_{k} \phi^{k} \text{ for some } \lambda_{k} \ge 0 \text{ and} \\ \phi^{k} \in \partial \mathbf{f}^{k}(\mathbf{x}) \}.$$

This cone is convex. It is also  $w^*$ -closed, if  $0 \notin conv \cup \partial \mathbf{f}^k(\mathbf{x})$ , ke $\Omega$ by Lemma 6.2. We now set

$$B_{\emptyset}(x) = \{0\}.$$

The linearizing cone and the cone of subgradients have the following dual property.

II.7

Lemma 7.1. Suppose that  $x \in S$  and  $\Omega \subseteq P(x)$ . Then

$$\overline{B_{\Omega}(x)} = -C_{\Omega}^{*}(x).$$

Proof. Since

(7.1) 
$$B_{k}^{k}(x) = \partial f^{k}(x)^{**}, \text{ by Lemma 6.1(b)}$$
$$= -C_{k}^{*}(x), \text{ by definition,}$$

we conclude that

$$-C_{\Omega}^{*}(x) = -\frac{\Sigma C_{k}^{*}(x)}{k \epsilon \Omega}, \text{ by Lemma 6.1(c)}$$
$$= \overline{B_{\Omega}(x)}, \text{ by (7.1)}.$$

Gould and Tolle [31] used Farkas' Lemma to prove the above result, for differentiable functions on  $\mathbb{R}^n$ . Note that in the differentiable case,  $B_{\Omega}(x)$  is closed. This may fail in the nondifferentiable case. (Consider  $B_1(0)$ , where  $f^1$  is the support function of the set  $\{x \in \mathbb{R}^2 : \|x - (0,1)\| \le 1\}$ .<sup>†</sup>) The closure of  $B_{P(x)}(x)$ will play an essential role in the characterization of regularity in IV.3.

<sup>†</sup> I would like to thank Professor J. Borwein for pointing this out to me.

# 8. The Cone of Tangents

For  $x \in M$ , where M is an arbitrary set in X, the <u>cone of</u> tangents to M at x is defined by

$$T(M,x) = \{ d \in X : d = \lim \lambda_k (x^k - x), \text{ where } x^k \in M, \\ \lambda_k \ge 0 \text{ and } x^k \to x. \}$$

This cone is closed and it is convex if M is. In fact. when M is convex, it is exactly the  $\overline{cone}$  (M - x), the <u>support cone</u> of M at x. For further properties, see e.g. Guignard [32] and Holmes [34]. For a relationship of the cone of tangents with derivatives in mathematical programming, see e.g. Massam [38] and Massam and Zlobec [39].

The cone of tangents is used in optimization theory to describe the geometry of the feasible set. For example, one gets the following characterization of optimality.

<u>Theorem 8.1.</u> [34, p.30]  $x \in S$  is optimal for (P) if and only if

$$\partial \mathbf{f}^{\circ}(\mathbf{x}) \cap \mathbf{T}^{*}(\mathbf{S},\mathbf{x}) \neq \emptyset.$$

This result will be the starting point for our characterizations of optimality. Note that the characterization is in terms of the feasible set, rather than the constraints.

### III. CHARACTERIZATIONS OF OPTIMALITY

### 1. Introduction

Optimality conditions of the Kuhn-Tucker type usually use the analytic properties of the constraints, given by the cone of subgradients,  $B_{p(x)}(x)$ , or the linearizing cone,  $C_{P(x)}(x)$ . However, these cones may not provide all the required information needed to characterize optimality. As seen in Theorem II.8.1, the cone of tangents of the feasible set does provide enough information to characterize optimality. However, this cone is hard to use computationally. In [10], it was shown that the cones of directions of constancy can be used to characterize optimality. We will see how these cones provide the missing information. Using Theorem II.8.1, we first present optimality conditions of the type given by Guignard [32] and Gould and Tolle [30]. (See Theorem 3.1.) From this result we deduce several new, as well as known, optimality conditions.

### 2. A Basic Lemma

The following lemma presents several relationships between the types of cones mentioned in Section 1. This lemma is of key importance in the proofs in the rest of this chapter and in Chapter IV.

<u>Lemma 2.1.</u> Suppose that  $x \in S$ ,  $\Omega$  satisfies  $P^b(x) \subset \Omega \subset P^=$ and either *conv*  $D_{\Omega}^{=}(x)$  is closed or  $\Omega = P^{=}$ .

Then:

a) 
$$T(S,x) = \overline{D_{P(x)}^{\leq}(x)}$$
.  
b)  $\overline{conv} D_{\Omega}^{=}(x) \cap C_{P(x)}(x) = \overline{D_{\Omega}^{=}(x)} \cap C_{P(x)}(x) = \overline{D_{p}^{=}(x)} \cap C_{P(x)}(x)$ .  
c)  $T(S,x) = \overline{conv} D_{\Omega}^{=}(x) \cap C_{P(x)}(x)$ .  
d)  $conv \left\{ \bigcup_{k \in P^{\leq}(x)} \partial f^{k}(x) \right\} \cap \left( D_{\Omega}^{=}(x) \right)^{*} = \emptyset$ .  
e)  $T^{*}(S,x) = \overline{(D_{\Omega}^{=}(x))^{*}} + C_{P(x)}^{*}(x) = \overline{(D_{\Omega}^{=}(x))^{*}} - B_{P(x)}(x)$ .  
f)  $T^{*}(S,x) = (D_{p}^{=}(x))^{*} + C_{P(x)}^{*}(x) = (D_{p}^{=}(x))^{*} - B_{P(x)}(x)$ .  
(Recall that X' is given the w\*-topology and thus is a lcs.)

<u>Proof.</u> Since the point  $x \in S$  is fixed throughout, we will omit it in this proof when the intended meaning is clear, e.g.  $D_p^{\leq}$  will denote  $D_{P(x)}^{\leq}(x)$ ,  $P^{\leq}$  will denote  $P^{\leq}(x)$ , etc. a) Suppose that  $d \in T(S,x)$ , with associated nets  $\{x^k\}$  and  $\lambda_k \ge 0$ . Let  $d^k = x^k - x$ . Since  $x^k \in S$ , we see that  $x + d^k \in S$ , i.e.  $d^k \in D_p^{\le}$ . Furthermore, since  $D_p^{\le}$  is a cone,  $\lambda_k d^k \in D_p^{\le}$ . But

$$d = \lim \lambda_k (x^k - x) = \lim \lambda_k d^k,$$

which implies that  $d \in \overline{D_p^{\leq}}$ . Thus

$$T(S,x) \subset \overline{D_p^{\leq}}.$$

Conversely, suppose that  $d \in D_p^{\leq}$ . Then, there exists  $\overline{\alpha} > 0$  such that

$$x + \alpha d \in S$$
, for all  $0 < \alpha \leq \overline{\alpha}$ .

Let

$$x^{k} = x + \frac{\overline{\alpha}}{k} d$$
, for  $k \ge 1$ .

Then  $x^k \in S$  and, since X is a *tvs*,  $x^k \rightarrow x$  as  $k \rightarrow \infty$ . By choosing  $\lambda_k = \frac{k}{\overline{\alpha}}$ , we see that

$$\lambda_k(x^k - x) = d$$
, for all  $k \ge 1$ .

Therefore  $d \in T(S,x)$ , i.e. we have shown that  $D_p^{\leq} \subset T(S,x)$ . The result now follows since T(S,x) is closed.

b) (i) First, let us show that

(2.1) 
$$\overline{conv} \ D_{\Omega}^{=} \cap C_{p} \subset \overline{conv} \ D_{\Omega}^{=} \cap C_{p}.$$

By hypothesis and the fact that  $\overline{conv} D^{=}_{p} \subset C$  and  $C_{p} = p^{=}_{p} p^{=}_{p}$ , it is sufficient to show that

(2.2) 
$$\overline{conv} \stackrel{=}{D} \stackrel{=}{p} \stackrel{C}{p} \stackrel{conv}{p} \stackrel{=}{p} \stackrel{C}{p} \stackrel{c$$

ς,

But this follows since

(2.3) 
$$\hat{d} \in D_{p^{=}}^{\overline{p}} \cap D^{\leq} \subset conv D_{p^{=}}^{\overline{p}} \cap int C_{p^{\leq}} \neq \emptyset,$$

by Lemma II.3.1(d).

Let us define

(2.4) 
$$d_{\lambda} \stackrel{\Delta}{=} \lambda \hat{d} + (1 - \lambda)d,$$

for scalars  $\lambda$  and vectors d.

(ii) Next, let us show that

(2.5) 
$$\overline{D_{\Omega}^{-} \cap C_{p}} \subset \overline{D_{p}^{-}} \cap C_{p}$$

Suppose that

$$d \in (D_{\Omega}^{=} \cap C_{p}) \setminus (\overline{D_{p}^{=}} \cap C_{p})$$
.

We will find a set  $I \subset P^{=}$  and feasible directions  $d_{\lambda} \in D_{p}^{\leq}$ , which are directions of decrease for  $f^{k}$ ,  $k \in I$ . This will contradict the definition of  $P^{=}$ .

By the assumption, we can find a nonempty set  $\mbox{I} \subset \ensuremath{\mathcal{P}}^{=} \backslash \Omega$  such that

$$d \in C_p = \bigcap (-\partial f^k)^*, d \in D_p^= \setminus I$$
 but  
 $k \in P$   
 $d \notin D_k^= \cup D_p^=, \text{ for each } k \in I.$ 

Recall that when  $k_{o} \in P^{=}$ , then  $f^{k_{o}}$  is 'badly behaved' at x if the system

$$\nabla f^{k_{\circ}}(x;d) = 0$$

$$\nabla f^{k}(x;d) \leq 0, \ k \in P(x) \setminus k_{\circ}$$

$$d \notin D^{=}_{k_{\circ}} \cup \overline{D^{=}_{p^{=}}}$$

is consistent. Therefore, since

4

$$I \subset P^{=} \setminus \Omega \subset P^{=} \setminus P^{b},$$

we see that

$$\nabla f^{k}(x;d) \leq 0$$
, for all  $k \in I$ ,

i.e.

$$(2.6) d \in D_{I}^{\leq} \cap D_{p=\backslash I}^{\equiv}.$$

Let  $\hat{d}$  and  $d_{\lambda}$  be defined as in (2.3) and (2.4) respectively. Then, by (2.3), (2.4) and (2.6),

(2.7) 
$$d_{\lambda} \in D_{I}^{<}$$
, for all  $0 \leq \lambda < 1$ .

Furthermore, Lemma II.3.1(b) implies that

(2.8) 
$$d_{\lambda} \in D_{p=\backslash I}^{\leq}$$
, for all  $0 < \lambda < 1$ .

Now, by continuity and (2.3), there exists  $0 < \beta < 1$  such that

(2.9) 
$$d_{\lambda} \in D_{p}^{\leq}$$
, for all  $\beta \leq \lambda < 1$ .

From (2.7), (2.8) and (2.9), we conclude that

$$d_{\lambda} \in D_{p}^{\leq} \cap D_{I}^{\leq}$$
, for all  $\beta \leq \lambda < 1$ ,

contradiction. Thus we have shown that

$$\mathbf{D}_{\Omega}^{=} \cap \mathbf{C}_{p} \subset \overline{\mathbf{D}_{p}^{=}} \cap \mathbf{C}_{p}.$$

The inclusion (2.5) follows, since both  $C_p$  and  $\overline{D_p^{=}}$  are closed.

(iii) By a similar argument, in particular employing LemmaII.3.1(b) again, we see that

(2.10) 
$$\operatorname{conv} \operatorname{D}_{\Omega}^{\overline{c}} \cap \operatorname{C}_{p} \subset \operatorname{D}_{\Omega}^{\overline{c}} \cap \operatorname{C}_{p}^{\overline{c}}$$

The desired result now follows from (2.1), (2.5), (2.10) and

$$\overline{\mathbb{D}_{p}^{=}} \cap \mathbb{C}_{p} \subset \overline{\mathbb{D}_{\Omega}^{=}} \cap \mathbb{C}_{p} \subset \overline{conv} \mathbb{D}_{\Omega}^{=} \cap \mathbb{C}_{p}.$$

c) By (a), (b) and (2.1), it is sufficient to show that

$$\overline{D_p^{\leq}} = \overline{conv \ D_p^{=}} \cap C_p.$$

That

(2.11) 
$$\overline{D_p^{\leq}} \subset \overline{conv} \ \overline{D_p^{=}} \cap C_p$$

is clear from the definitions and Lemma II.3.1(b). To prove the converse, we first show that

(2.12) 
$$\operatorname{conv} \operatorname{D}_{p}^{=} \cap \operatorname{C}_{p} \subset \operatorname{D}_{p}^{\leq}.$$

Suppose that we are given

d 
$$\epsilon$$
 conv  $D_p^{=} \cap C_p$ 

and the neighbourhood of the origin, n. We need to show that

(2.13) 
$$D_p^{\leq} \cap (n+d) \neq \emptyset.$$

Let  $\hat{d}$  and  $d_{\lambda}^{}$  be defined as in (2.3) and (2.4) respectively. Then

$$d_{\lambda} \cdot \phi = \lambda \hat{d} \cdot \phi + (1 - \lambda) d \cdot \phi < 0,$$

for all  $0 < \lambda \le 1$  and all  $\phi \in \bigcup_{k \in \mathcal{P}^{\le}} \partial f^k$ . Therefore

$$d_{\lambda} \in D_{p}^{\leq} \cap conv \ D_{p}^{=} \subset D_{p}^{\leq}, \text{ for all } 0 < \lambda \leq 1.$$

Furthermore,  $d_{\lambda} \in \mathcal{N} + d$  for sufficiently small  $\lambda$ . This proves (2.13) and thus (2.12).

The desired result now follows since  $\overline{D_p^{\leq}}$  is closed.

d) Let

(2.14) 
$$C \stackrel{\Delta}{=} conv \left\{ \begin{array}{c} \cup & \partial f^k \\ k \in P^{\leq} \end{array} \right\}$$
.

Suppose that the intersection is not empty. Then  $P^{\leq} \neq \emptyset$  and there exists

$$\phi \in C \cap (\bar{D}_{\Omega}^{\bar{-}})^*$$

where  $\phi = \sum_{k \in P} \langle \lambda_k \phi^k$  for some  $\lambda_k \ge 0$ ,  $\Sigma \lambda_k = 1$  and  $\phi^k \in \partial f^k$ . By Lemma II.6.1(a) and (b),

$$D_{\Omega}^{=} \subset \{\phi\}^{*} \text{ and } - C^{*} \subset - \{\phi\}^{*}.$$

Therefore

$$D_{\Omega}^{=} \cap - C^{*} \subset \{\phi\}^{\perp}.$$

Let  $\hat{d}$  be as in (2.3). Then

$$\hat{d} \in D_p^{=} \cap D_p^{\leq} \subset D_{\Omega}^{=} \cap - C^* \subset \{\phi\}^{\perp},$$

i.e.  $\hat{d} \cdot \phi = 0$ . But

$$\hat{\mathbf{d}} \cdot \boldsymbol{\phi} = \hat{\mathbf{d}} \cdot \sum_{k \in \mathcal{P}^{\leq}} \lambda_{k} \boldsymbol{\phi}^{k} < 0,$$

e) The result follows from (c) and dLemmas II.6.1(c) and II.7.1.

f) First, note that

(2.15) 
$$D_p^{=*} + C_p^* = D_p^{=*} + C_p^* <,$$

since  $conv D_p^{=} \subset C_p^{=}$ . Now, let C be as in (2.14). Then C is  $w^*$ -compact, since the subdifferentials  $\partial f^k$  are  $w^*$ -compact and  $p^<$  is finite. Furthermore, since  $0 \notin C$  by (d), we get that  $B_{p^<}$  is closed by Lemma II.6.2 and thus

$$-C_{p<}^{*} = B_{p<} = cone C$$
, by Lemma II.7.1.

The result now follows from (d), (e), (2.15) and Lemma II.6.2.

# 3. Gould and Tolle Optimality Criteria

We will be interested in optimality criteria of the following type:

 $x \in S$  is optimal if and only if the system

$$(G - T) \begin{cases} \phi^{\circ} + \sum_{k \in P(x)} \lambda_{k} \phi^{k} \in G \\ \phi^{\circ} \in \partial f^{\circ}(x) \\ \phi^{k} \in \partial f^{k}(x), \lambda_{k} \geq 0 \end{cases}$$

### is consistent,

where G is a nonempty cone in X'. Gould and Tolle have considered such optimality criteria for differentiable, not necessarily convex, functionals on  $R^n$ . They have shown that if

(3.1)  $T^*(S,x) = C^*_{p(x)}(x) + G,$ 

then the (G - T) conditions are necessary for  $x \in S$  to be optimal. (Note that the condition (3.1) depends only on the constraints and not on the objective function  $f^{\circ}$ .) One obvious candidate for G is

$$T^{*}(S,x) \setminus C^{*}_{P(x)}(x) \cup \{0\}.$$

Moreover, if  $G = \{0\}$  satisfies (3.1), then we get the well-known Kuhn-Tucker conditions, e.g. [36], [37]. In our setting, i.e. for the convex program (P), we can say: <u>Theorem 3.1.</u> Suppose that  $x \in S$  and  $G \subseteq X'$ . Then the statement:

"x is optimal for (P) if and only if the system

(3.2) 
$$\begin{cases} \phi^{\circ} + \sum_{k \in P(x)} \lambda_{k} \phi^{k} \in G \\ k \in P(x) \end{cases}$$
$$\phi^{\circ} \in \partial f^{\circ}(x) \\ \phi^{k} \in \partial f^{k}(x), \lambda_{k} \geq 0 \end{cases}$$

is consistent",

holds, for any fixed objective function  $f^{\circ}$ , if and only if

G satisfies (3.1).

<u>Proof.</u> We need to show that we can choose G in (3.2) if and only if

(3.3) 
$$T^*(S,x) = -B_{P(x)}(x) + G$$
,

Sufficiency: Suppose that G satisfies (3.3). By Theorem II.8.1, we know that x is optimal if and only if  $\partial f^{\circ}(x) \cap T^{*}(S,x) \neq \emptyset$ . By (3.3), this implies that x is optimal if and only if  $\partial f^{\circ}(x) \cap (-B_{P(x)}(x) + G) \neq \emptyset$ , i.e. if and only if (3.2) is consistent. Necessity: We need to show that (3.3) holds. Suppose that  $\phi \in T^*(S,x)$  and  $f^\circ$  is defined by the linear functional  $\phi(\cdot)$ on X. Then  $\phi \in \partial f^\circ(x) \cap T^*(S,x)$  and we can conclude that x is optimal for (P), i.e.  $\phi = f^\circ \in F^\circ(x)$ . Therefore, by the conditions (3.2), we see that  $\phi \in -B_{P(x)} + G$ . Thus

$$T^*(S,x) \subseteq -B_{\mathcal{P}(x)}(x) + G.$$

Conversely, let  $\phi \in -B_{\mathcal{P}(x)}(x) + G$ . Then we can find  $\lambda \geq 0$ and  $\phi^k \in \partial f^k(x)$  such that

(3.4) 
$$\phi + \sum_{k \in P(\mathbf{x})} \lambda_k \phi^k \in G.$$

Again we let  $f^{\circ}$  be the linear functional  $\phi$ . Then  $\phi = f^{\circ} \in F^{\circ}(x)$ , by (3.2). Since  $\partial f^{\circ}(x) = \{\phi\}$ , Theorem II.8.1 implies that  $\phi \in T^{*}(S,x)$ . Thus

$$-B_{P(x)}(x) + G \subset T^{*}(S,x).$$

The condition (3.1) is frequently referred to as a necessary and sufficient constraint qualification, or a weakest constraint qualification.

# 4. Some Choices for the Cone G

By specifying G in (3.3), we get necessary and sufficient conditions for optimality which hold without a constraint qualification. For example,

<u>Theorem 4.1.</u> The point  $x \in S$  is optimal for (P) if and only if the system

$$\begin{cases} \phi^{\circ} + \sum_{k \in P(x)} \lambda_{k} \phi^{k} \in (D_{P(x)}^{\leq}(x))^{*} \\ \phi^{\circ} \in \partial f^{\circ}(x) \\ \phi^{k} \in \partial f^{k}(x), \quad \lambda_{k} \geq 0 \end{cases}$$

is consistent.

<u>Proof.</u> By Theorem 3.1, we need only show that  $G = \left(D_{P(x)}^{\leq}(x)\right)^{*}$  satisfies (3.3). Now

$$T^{*}(S,x) = \left(D_{P(x)}^{\leq}(x)\right)^{*}, \text{ by Lemmas 2.1(a) and II.6.1(a)}$$
$$= -B_{P(x)}(x) + \left(D_{P(x)}^{\leq}(x)\right)^{*}, \text{ since } D_{P(x)}^{\leq}(x) \subset C_{P(x)}(x) \text{ and }$$
$$-B_{P(x)}(x) \subset C_{P(x)}^{*}(x). \square$$

Other, possibly more useful candidates for G are given in the next theorem.

<u>Theorem 4.2.</u> Suppose that  $x \in S$ , the set  $\Omega$  satisfies  $p^{b}(x) \subset \Omega \subset p^{=}$  and both *conv*  $D_{\Omega}^{=}(x)$  and  $-B_{P(x)}(x) + (D_{\Omega}^{=}(x))^{*}$  are closed. Then, x is optimal for (P) if and only if the system

(4.1) 
$$\begin{cases} \phi^{\circ} + \sum_{k \in \mathcal{P}(x)} \lambda_{k} \phi^{k} \in \left(D_{\Omega}^{=}(x)\right)^{2} \\ \phi^{\circ} \in \partial f^{\circ}(x) \\ \phi^{k} \in \partial f^{k}(x), \quad \lambda_{k} \geq 0 \end{cases}$$

is consistent.

<u>Proof.</u> The result follows immediately from Theorem 3.1 and Lemma 2.1(e).

We have assumed that  $conv D_{\Omega}^{=}$  and  $-B_{P(x)}(x) + (D_{\Omega}^{=}(x))^{*}$  are closed. (This can be considered as a kind of constraint qualification.) The sets are closed, for example, when the constraints are faithfully convex and differentiable. For then both cones in the sum are polyhedral. The following two examples show that the closure assumptions are necessary.

Example 4.1. Consider the program

 $f^{\circ}(x) \rightarrow \min$ s.t.  $f^{k}(x) \le 0, \quad k \in P = \{1, 2, 3\},\$ 

where  $x = (x_1) \in \mathbb{R}^3$ ,  $f^1(x) = x_1$ ,  $f^2(x) = -x_1$ ,  $f^3(x) = (dist(x,K))^2$ (see (II.5.2)) and K is the self-polar, 'ice-cream' cone

$$K \stackrel{\Delta}{=} \{x \in R^3 : x_1 \ge 0, x_1 x_2 \ge x_3^2\}.$$

Note that now

$$\nabla f^{3}(0,d) = \lim_{t \neq 0} \frac{\min \|td - z\|^{2}}{\frac{z \in K}{t}}$$
$$= 0, \text{ for all } d \in \mathbb{R}^{3}.$$

Let  $\overline{x} = 0$ . Then  $\overline{x} \in S$ ,  $P^{=} = P$  while  $P^{D}(\overline{x}) = \{3\}$ . Furthermore,

$$C_{P(0)}^{*}(0) = \operatorname{span} \left\{ \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \right\}$$
 and  $\left( D_{pb(0)}^{*}(0) \right)^{*} = K.$ 

Let us show that

$$C_{P(0)}^{*}(0) + (D_{p}^{-b}(0)(0))^{*}$$
 is not closed.

Choose

$$k^{i} = \begin{pmatrix} i \\ \frac{1}{i} \\ 1 \end{pmatrix} \in K \text{ and } k^{i} = \begin{pmatrix} -i \\ 0 \\ 0 \end{pmatrix} \in C^{*}_{P(0)}(0), \quad i = 1, 2, \dots$$

Then

$$k^{i} + k^{i} \rightarrow \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} \notin C^{*}_{P(0)}(0) + (D^{=}_{P^{b}(0)}(0))^{*}.$$

Example 4.2. Consider the program

$$f^{\circ}(x) \rightarrow \min$$
  
s.t.  
 $f^{k}(x) \leq 0, k \in P = \{1,2\},\$ 

where  $x = (x_i) \in \mathbb{R}^2$ ,

$$f^{1}(x) = \begin{cases} (x_{1}^{2} + x_{2}^{2} - 1)^{2} \text{ if } x_{1}^{2} + x_{2}^{2} - 1 \ge 0\\ 0 & \text{otherwise,} \end{cases}$$

 $f^{2}(x) = dist(x - \overline{x}, K), K = \{x \in R^{2} : x_{1} \ge 0, x_{2} \ge 0\}$  and  $\overline{x} = (1,0)^{t}$ . Then  $S = \{\overline{x}\}, P^{\overline{x}} = P$  while  $P^{b}(\overline{x}) = \{1\}$ . Let  $\Omega = \{1\}$ . Then

$$D_{\Omega}^{=}(\overline{x}) = \{ d \in R^{2} : d_{1} < 0 \} \cup \{ 0 \}$$

is not closed. Furthermore

$$T^*(S,\overline{x}) = R^2$$
 and  $C^*_{P(\overline{x})}(\overline{x}) = K$ .

Therefore

$$(D_{\Omega}^{=}(\overline{x}))^{*} + C_{P(\overline{x})}^{*}(\overline{x}) = \{x \in \mathbb{R}^{2} : x_{2} \ge 0\}$$
  
 $\neq T^{*}(S,x).$ 

This implies that (3.1) fails and that we cannot choose  $\Omega = \{1\}$  in Theorem 4.2.

<u>Remark 4.1.</u> When the sum  $C_{P(x)}^{*}(x) + D_{\Omega}^{**}(x)$  is not closed, we can, however, get the following asymptotic conditions: If  $P^{b}(x) \subset \Omega \subset P^{=}$  and *conv*  $D_{\Omega}^{=}$  is closed, then  $x \in S$  is optimal if and only if there exists nets

$$\{b^n\} \subset -B_{\mathcal{P}(\mathbf{x})}(\mathbf{x}) \text{ and } \{d^n\} \subset \{D_{\Omega}^{=}(\mathbf{x})\}^*$$

such that

$$\lim(b^n + d^n) \in \partial f^{\circ}(x).$$

For more details on asymptotic conditions see e.g. Zlobec [50], [51].

Using the fact that the sum  $-B_{P(x)}(x) + D_{p}^{*}(x)$  is always closed (see Lemma 2.1(f)), we get the following theorem.

<u>Theorem 4.2.</u> Let  $x \in S$ . Then x is optimal for (P) if and only if the system

$$\begin{pmatrix}
\phi^{\circ} + \sum_{k \in P(x)} \lambda_{k} \phi^{k} \in (D_{p}^{=}(x))^{*} \\
\phi^{\circ} \in \partial f^{\circ}(x) \\
\phi^{k} \in \partial f^{k}(x), \quad \lambda_{k} \geq 0
\end{pmatrix}$$

is consistent.

<u>Proof.</u> The result follows from Theorem 3.1 and Lemma 2.1(f).  $\Box$ 

This result is equivalent to the following characterization of optimality given in [12]. (See also [3] and [11] for the differentiable case.)

<u>Corollary 4.1.</u> Let  $x \in S$ . Then x is optimal for (P) if and only if the system

$$\begin{cases} \phi^{\circ} + \sum_{k \in P^{\leq}(x)} \lambda_{k} \phi^{k} \in \left(D_{p}^{=}(x)\right)^{2} \\ \phi^{\circ} \in \partial f^{\circ}(x) \\ \phi^{k} \in \partial f^{k}(x), \lambda_{k} \geq 0 \end{cases}$$

is consistent.

Proof. The result follows from (2.15) and Lemma II.7.1.

The above optimality criteria holds without a constraint qualification. However, when Theorem 4.2 is applicable, it may prove useful to choose  $G = \left(D_{p}^{-b}b_{(x)}(x)\right)^{*}$ , since this is a smaller G. One would therefore, have simpler necessary conditions to check. This question will be examined in more detail in Section IV.5.

The case when  $G = \{0\}$  deserves special attention because then, the system (3.2) reduces to the Kuhn-Tucker conditions. We will study this case in the next chapter.

# IV. CONSTRAINT QUALIFICATIONS AND REGULARIZATION TECHNIQUES

120

# 1. Kuhn-Tucker Points

A point  $x \in S$  is called a <u>Kuhn-Tucker (K-T) point</u> for (P) if the Kuhn-Tucker (K-T) conditions are satisfied at x, i.e. if the system

(1.1) 
$$\begin{cases} \phi^{\circ} + \sum_{k \in P(x)} \lambda_{k} \phi^{k} = 0 \\ k \in P(x) \end{cases}$$
$$\phi^{\circ} \in \partial f^{\circ}(x) \\ \phi^{k} \in \partial f^{k}(x), \quad \lambda_{k} \ge 0 \end{cases}$$

is consistent. It is well-known that if  $x \in S$  is a *K-T* point, then x solves program (P). However, the converse does not always hold.

### 2. Regular Points and Slater's Condition

We call  $x \in S$  a <u>regular point</u> (Lagrange regular point), if the *K-T* conditions (1.1) hold for every  $f^{\circ} \in F^{\circ}(x)$ , i.e. we can choose  $G = \{0\}$  in (III.3.2), see e.g. [31]. A <u>constraint qualifi-</u> <u>cation (CQ)</u> is then a condition on the set of constraints which guarantees that x is a regular point.

Slater's condition, i.e. the requirement that

"there exists 
$$\hat{x} \in X$$
 with  $f^k(\hat{x}) < 0$ , for all  $k \in P$  "

is a well-known CQ which guarantees that each point  $x \in S$  is a regular point. Slater's condition is equivalent to the fact that

$$P^{-} = \emptyset,$$

see e.g. [11]. (This follows from Lemma II.3.1(d).) Note that Slater's condition is not equivalent to

int 
$$S \neq \emptyset$$
,

see e.g. [11]. However, we can say the following.

<u>Theorem 2.1.</u> Suppose that, when  $P^{-} \neq \emptyset$  there exists  $k \in P^{-}$  such that  $f^{k}$  is faithfully convex. Then

(i) Slater's condition holds

if and only if

(ii) int  $S \neq \emptyset$ .

(In particular, the equivalence holds when all the constraints are faithfully convex.)

<u>Proof.</u> That (i) implies (ii) follows by continuity. It is now sufficient to show that (ii) implies (2.1). Suppose that (2.1) fails, i.e. there exists  $k \in P^{=}$  and, by hypothesis,  $f^{k}$  is faithfully convex. Then,  $D_k^{\bar{e}}(x) = D_k^{\bar{e}}$  is a proper subspace of X independent of  $x \in X$ , see Lemma II.3.1(e). This implies that  $S \subset x + D_k^{\bar{e}}$ , which implies that *int*  $S = \emptyset$ .

For a discussion on constraint qualifications, see e.g. [4], [7], [8].

### 3. Weakest Constraint Qualifications

A weakest constraint qualification (WCQ) is a constraint qualification, that holds if and only if x is a regular point. In other words, it is a condition that holds if and only if x is a K-T point for all  $f \in F(x)$ . Gould and Tolle [30], [31] have shown that, in their setting (see e.g. Section III.3)

$$T^{*}(S,x) = C^{*}_{P(x)}(x)$$
 is a WCQ.

By Theorem III.3.1, we see that in our more general setting,

(3.1) 
$$T^*(S,x) = -B_{P(x)}(x)$$
 is a WCQ.

Note that this requires  $B_{P(x)}(x)$  be closed. In this section, we present two different WCQ's.

<u>Theorem 3.1.</u> Suppose that  $x \in S$ ,  $P^b(x) \subset \Omega \subset P^{=}$  and either conv  $D_{\Omega}^{=}(x)$  is closed or  $\Omega = P^{=}$ . Then:

a) 
$$P^{D}(x) = \emptyset$$
 and  $B_{P(x)}(x)$  is closed is a WCQ.  
b)  $C_{P(x)}(x) \subseteq conv D_{\Omega}^{=}(x)$  and  $B_{P(x)}(x)$  is closed is a WCQ.

<u>Proof.</u> Suppose that (a) holds. Then, Lemma III.2.1(c) implies that

$$\Gamma(S,x) = C_{P(x)}(x)$$
 and  $B_{P(x)}(x)$  is closed.

By (3.1) and Lemma II.6.1(a), the above is a WCQ. This implies that x is a regular point.

Conversely, suppose that  $p^{b}(x) \neq \emptyset$ . Recall that  $k \in p^{b}(x)$  if  $k \in p^{\overline{b}}$  and there exists

$$d \in \left( D_{k}^{>}(x) \cap E_{k}(x) \cap C_{p(x)}(x) \right) \setminus \overline{D_{p}^{=}(x)}.$$

But, this implies that

d  $\notin \overline{D_{p}^{=}(x)} \cap C_{P(x)}(x)$ = T(S,x), by Lemma III.2.1 (b) and (c).

Therefore,

$$T(S,x) \neq C_{P(x)}(x),$$

which, as above, implies that x is not a regular point. That  $B_{P(x)}(x)$  is not closed implies that x is not a regular point, follows from (3.1). This proves (a). To prove (b) note that, as in the proof of (a), we need only show that:  $C_{P(x)}(x) \subset \overline{conv} D_{\Omega}^{=}(x)$  if and only if  $T(S,x) = C_{P(x)}(x)$ . But this follows directly from Lemma III.2.1(b) and (c).

<u>Remark 3.1.</u> Suppose that  $B_{P(x)}(x)$  is closed and we can find  $\hat{x} \in S$  and  $\Omega \subset P$  such that  $f^k(\hat{x}) < 0$ , for all  $k \in P \setminus \Omega$ , and  $f^k$  is 'never badly behaved', for all  $k \in \Omega$ , i.e.

$$E_k(x) = D_k^{=}(x)$$
, for all  $x \in S$  and  $k \in \Omega$ .

Then, since  $P^{b}(x) \subset P^{=}$ , this implies that  $P^{b}(x) = \emptyset$ , for all  $x \in S$ , i.e.  $x \in S$  is a regular point. Thus we see that, when checking if Slater's condition holds, we need not worry about the functions which are 'never badly behaved'. In particular, we can ignore all linear functionals.

<u>Remark 3.2.</u> The condition given in (b) may be easier to check computationally than the one in (a). For example, when  $X = R^n$ , the constraints  $f^k$ ,  $k \in P(x)$ , are differentiable and the constraints  $f^k$ ,  $k \in \Omega$ , are faithfully convex, then  $C_{P(x)}(x)$  is a polyhedral cone while  $D_{\Omega}^{=}(x) = D_{\Omega}^{=}$  is a subspace, independent of x. Furthermore,  $D_{\Omega}^{=}$  can be calculated explicitly, see Section V.2.

<u>Remark 3.3.</u> Suppose that S contains two distinct points. Then Slater's condition is a WCQ with respect to the Fritz John optimality conditions. IV.3

$$\begin{cases} \sum_{k \in P(x) \cup \{0\}} \lambda_k \phi^k = 0 \\ \\ \phi^k \in \partial f^k(x), \quad \lambda_k \ge 0, \quad \sum_{k \in P(x) \cup \{0\}} \lambda_k = 1 \end{cases}$$

is consistent. Necessity always holds. We need to show that, if S contains two distinct points (Note that when  $S = \{x\}$ , then x is optimal for any  $f^{\circ}$  chosen, the Fritz John optimality conditions hold, but Slater's condition fails.) then the Fritz John conditions are sufficient for optimality, independent of the objective function  $f^{\circ}$ , if and only if Slater's condition holds. This follows from the fact that the system

$$\begin{cases} \sum_{k \in P(x)} \lambda_k \phi^k = 0 \\ k \in P(x) \end{cases}$$
$$\phi^k \in \partial f^k(x), \ \lambda_k \ge 0, \ \sum_{k \in P(x)} \lambda_k = 1 \\ k \in P(x) \end{cases}$$

is consistent if and only if Slater's condition fails (i.e. if and only if  $P^{=} \neq \emptyset$ ), which in turn follows from Motzkin's Theorem of the alternative, see e.g. [52].

# 4. Regularization

Gould and Tolle have posed the question: "Can the program (P) be regularized by the addition of a finite number of constraints?" Augunwamba [6] has considered the nonconvex, differentiable case and has shown that one can always regularize with the addition of an infinite number of constraints. He has also given necessary and sufficient conditions to insure the number of constraints added may be finite. In this section, we show that one can always regularize (P) at x, by the addition of one (possibly nondifferentiable) constraint. Furthermore, in the case of faithfully convex constraints, we can regularize (P) by the addition of a finite number of linear constraints. (In the following theorem, we assume that  $B_{P(x)}(x)$  is closed.)

<u>Theorem 4.1.</u> Suppose that  $\overline{x} \in S$ , X is a Hilbert space  $P^{b}(\overline{x}) \subset \Omega \subset P^{=}$  and either *conv*  $D_{\Omega}^{=}(\overline{x})$  is closed or  $\Omega = P^{=}$ . Consider program (P) with the additional constraint

 $f^{m+1}(x) \stackrel{\Delta}{=} dist ((x - \overline{x}), \overline{conv} D_{\Omega}^{=}(\overline{x})).$ 

Then x is a regular point.

<u>Proof.</u> By Lemma II.5.1,  $f^{m+1}$  is not 'badly behaved' at  $\overline{x}$  and therefore,  $P^{b}(\overline{x})$  is not increased by the addition of  $f^{m+1}$ . Now, by Theorem 3.1, we need only show that

$$\begin{array}{ll} C & (\overline{\mathbf{x}}) & \subset & \overline{conv} & D_{\Omega}^{-}(\overline{\mathbf{x}}) \\ \mathcal{P}(\overline{\mathbf{x}}) \cup \{m+1\} \end{array}$$

But

$$C_{m+1}(\overline{x}) = \{ d \in X : \nabla f^{m+1}(\overline{x}; d) \le 0 \}$$
$$= \overline{conv} D_{\Omega}(\overline{x}), \text{ by (II.5.3).}$$

Note that the feasible set remains unchanged after the addition of  $f^{m+1}$ . For, let  $\overline{S}$  denote the feasible set after the addition. Then

$$x \in \overline{S} \Leftrightarrow x \in S$$
 and  $x - \overline{x} \in \overline{conv} D_{\overline{\Omega}}^{\overline{c}}(\overline{x})$   
 $\Leftrightarrow x \in S$ , since  $\Omega \subseteq P^{\overline{c}}$  and  $D_{\overline{\Omega}}^{\overline{c}}(\overline{x}) \subseteq \overline{conv} D_{\overline{\Omega}}^{\overline{c}}(\overline{x})$ 

We have, therefore, regularized the point  $\overline{x}$ , by the addition of a 'redundant' constraint.

<u>Theorem 4.2.</u> Let  $\overline{x} \in S$  and  $f^k$ ,  $k \in P^=$ , be faithfully convex. Suppose that  $B : Y \to X$  is the linear operator satisfying

$$D_p^{=} = \mathcal{R}(B),$$

where Y is a lcs. Consider the program, in the variable  $y \in Y$ ,

$$(P_r) \qquad \qquad f^{\circ}(\overline{x} + By) \neq \min \\ s.t. \\ f^{k}(\overline{x} + By) \leq 0, \ k \in P \setminus P^{=}.$$

Then Slater's condition is satisfied for  $(P_r)$ , and y = 0 is a feasible point of  $(P_r)$ . Moreover, if  $y^*$  solves  $(P_r)$ , then  $\overline{x} + By^*$  solves (P).

<u>Proof.</u> The result follows from the characterization of optimality in [3], [11] (see also Corollary III.5.1) and the fact that  $P^{=} = \emptyset$ if and only if Slater's condition holds.

In the next chapter we will see how to calculate  $D_{p}^{=}$  and how  $p^{=}$ to apply the above theorem to find a feasible point and solve program (P). Note that, after the substitution, (P<sub>r</sub>) has fewer constraints (and as we shall see, fewer variables) than (P).

### 5. Strongest and Weakest Optimality Conditions

In Chapter III, we presented several optimality criteria of the type:

(5.1)  $\begin{cases} x \in S \text{ is optimal if and only if the system} \\ \varphi^{\circ} + \sum_{k \in P(x)} \lambda_k \varphi^k \in G \\ \varphi^{\circ} \in \partial f^{\circ}(x) \\ \lambda_k \ge 0, \ \varphi^k \in \partial f^k(x), \ k \in P(x), \\ \text{ is consistent,} \end{cases}$ 

where G is a closed convex cone satisfying

(5.2) 
$$T^*(S,x) = -B_{P(x)}(x) + G.$$

We have seen that we can always choose

(5.3) 
$$G = (D_{p=}^{=}(x))^{*} \text{ or } (D_{p(x)}^{\leq}(x))^{*}.$$

Furthermore, when

$$p^{b}(x) \subset \Omega \subset p^{=}$$

and the sets

conv 
$$D_{\Omega}^{=}(x)$$
 and  $-B_{P(x)}(x) + (D_{\Omega}^{=}(x))$ 

are closed, we can choose

(5.4) 
$$G = \left(D_{\Omega}^{*}(x)\right)^{*}.$$

Among the choices of G in (5.3) and (5.4), clearly

$$G = \left( D^{=}_{p^{b}(x)}(x) \right)^{*}$$

is the smallest.

Gould and Tolle [30] have posed the following question: "When does there exist a 'smallest' G for (5.1) ?" By a 'smallest' G we shall mean a nonempty, closed convex cone that satisfied (5.2) but which contains no proper convex subset which also satisfied (5.2). It is of interest to use the 'smallest' G, because then algorithms which use (5.1) have fewer necessary conditions to check, i.e. we have a 'tighter' theory. From Lemma III.2.1(a), we know that

$$T^{*}(S,x) = (D_{p(x)}^{\leq}(x))^{*}.$$

Therefore  $G = (D_{P(x)}^{\leq}(x))^{*}$  will always satisfy (5.2) and this G will then be the 'largest' possible. On the other hand, we have seen that, when  $B_{P(x)}(x)$  is closed,  $G = \{0\}$  satisfies (5.2) if and only if  $P^{b}(x) = \emptyset$ . In this case,  $G = \{0\}$  is clearly a 'smallest' G. However, when  $P^{b}(x) \neq \emptyset$ , the cone

$$G = \left( D^{=}_{p^{b}(x)} (x) \right)^{*}$$

need not be a 'smallest' G. In fact a 'smallest' G need not exist (See Example 5.1 below.). The following theorem gives conditions for finding a 'smallest' G.

<u>Theorem 5.1.</u> Suppose that  $x \in S$ . If there exists a closed convex cone H such that

$$C_{\mathcal{P}(\mathbf{x})} \cap H = T(S,\mathbf{x})$$

and

$$-B_{P(x)}(x) + H^*$$
 is closed,

then

 $G = H^*$ 

satisfies (5.2). Furthermore, if there exists a largest, by inclusion,

such H, then  $H^*$  is a 'smallest' G for (5.2).

Corollary 5.1. If 
$$P^{b}(x) \subset \Omega \subset P^{=}$$
 and  
 $C_{P(x)}(x) = X$ ,

then

$$\left(D_{p=}^{=}(x)\right)^{*} = \left(D_{\Omega}^{=}(x)\right)^{*} = \left(D_{pb}^{=}(x)\right)^{*},$$

and, furthermore, this cone is the unique and so 'smallest' G satisfying (5.2).

Corollary 5.2. If there exists a halfspace H such that

 $C_{\mathcal{P}(x)}(x) \cap H = T(S,x)$ 

and

$$-B_{P(x)}(x) + H^*$$
 is closed,

then

$$G = \begin{cases} \{0\} & \text{if } P^{b}(x) = \emptyset \\ H^{*} & \text{otherwise} \end{cases}$$

is a 'smallest' G.

Example 5.1. A smallest G need not exist. Consider the program (P) with the two constraints

$$f^{1}(x) = x_{1}^{2} + x_{3}^{2} \le 0$$
  
$$f^{2}(x) = x_{1} \le 0.$$

Then  $S = \{x = (x_1) \in \mathbb{R}^3 : x_1 = x_3 = 0\}, \ \overline{x} = (0,0,0)^t \in S, \ P^= = \{1,2\}, D_{\overline{p}(\overline{x})}^{\leq}(\overline{x}) = D_{\overline{p}}^{=}(\overline{x}) = S, \ T^*(S,x) = \left(D_{\overline{p}(\overline{x})}^{\leq}(\overline{x})\right)^* = \{x \in \mathbb{R}^3 : x_2 = 0\} \text{ and } C_{\overline{p}(\overline{x})}(\overline{x}) = \{x \in \mathbb{R}^3 : x_1 \leq 0\}.$  We can now set  $G = cone\{(1,0,\varepsilon)^t, (1,0,-\varepsilon)^t\}$ 

where  $\varepsilon > 0$ , since

$$C_{P(\overline{x})}^{*}(\overline{x}) + G = \{x \in \mathbb{R}^{3} : x_{2} = x_{3} = 0, x_{1} \le 0\} + G$$
  
=  $T^{*}(S, \overline{x}).$ 

However, we cannot set  $\varepsilon = 0$  in G.

Example 5.2. Consider the program (P) with the single constraint

$$f^{1}(x) = x_{1}^{2} + x_{3}^{2} \le 0.$$
  
Then  $\overline{x} = (0,0,0)^{t} \in S, P^{=} = P^{b}(\overline{x}) = \{1\}, C_{P(\overline{x})}(\overline{x}) = R^{3}$  and

$$D_{p}^{=}(\overline{x}) = \{x \in \mathbb{R}^{3} : x_{1} = x_{3} = 0\}. \text{ By Corollary 5.1,}$$
$$G = (D_{p}^{=}(\overline{x}))^{*} = (D_{p}^{=}(\overline{x})(\overline{x}))^{*} = \{x \in \mathbb{R}^{3} : x_{2} = 0\}$$

is the unique and so the 'smallest' G.
#### V. THE METHOD OF REDUCTION

#### 1. Introduction

The purpose of this chapter is to introduce a numerical method for solving finite dimensional convex programs (P), regardless of whether Slater's condition is satisfied or not. Ben-Tal and Zlobec [14] have presented a feasible directions algorithm that solves program (P) without assuming a constraint qualification. They find feasible directions d, by solving the system

$$d \cdot \nabla f^{k}(x) < 0, \ k \in \{0\} \cup \Omega$$
$$d \in D_{k}^{=}(x), \ k \in P(x) \setminus \Omega,$$

where  $\Omega$  is some subset (possibly empty) of P(x). If no solution exists, for any  $\Omega \subset P(x)$ , then x is an optimal point for (P). Abrams and Kerzner [3] have shown that one need only consider the single system,  $\Omega = P^{=}$ . They have also presented an algorithm that finds  $P^{=}$ . (Note that it may be computationally better to use other subsets  $\Omega$ , rather than just  $P^{=}$ , since this allows more feasible directions to choose from.) Zoutendijk [53] has suggested that, in the absence of Slater's condition, one should solve the perturbed program

 $(P_{\varepsilon}) \qquad f^{\circ}(x) \neq \min \\ s.t. \\ f^{k}(x) \leq \varepsilon, \ k \in \Omega \\ f^{k}(x) \leq 0, \ k \in P \setminus \Omega, \\ where \ \varepsilon > 0 \ and \ \Omega = P \ or \ P^{=}.$ 

V.1

If the feasible set is compact, this is a 'stable' perturbation of (P), i.e. if  $x(\varepsilon)$  denotes a solution of  $(P_{\varepsilon})$ , then every cluster point of the net  $\{x^*(\varepsilon)\}_{\varepsilon \neq 0}$  is a solution of (P). Moreover,  $f^{\circ}(x^*(\varepsilon)) \neq f^{\circ}(x^*)$  as  $\varepsilon \neq 0$ , where  $x^*$  is a solution of (P), see e.g. G. Wolkowicz [48].

In this chapter, we combine the approaches in [3] and [14] with the regularization technique in Theorem IV.4.2, to formulate the method, which we call the "Method of Reduction". The method first finds a feasible point  $x^{\circ}$  and, in the process of finding  $x^{\circ}$ , it reduces program (P) to an equivalent program, in fewer variables and fewer constraints, for which Slater's condition is satisfied. The solution of (P) is now calculated by any method that works when Slater's condition is satisfied.

An integral part of the algorithm is finding the cone  $D_{p^{=}}^{=}(x)$ . This is done in Sections 2 and 3. The method of reduction is then presented in Section 4. Applications and examples follow in Sections 5 and 6.

# 2. Calculating the Cone of Directions of Constancy+

Recall that the cone of directions of constancy  $D_{f}^{=}$  of a faithfully convex function  $f : R^{n} \rightarrow R$  is a subspace of  $R^{n}$  independent of the choice of x, see Lemma II.3.1(e). We now formulate an algorithm that finds  $D_{f}^{=} \cap R(A_{0})$ , where f is a faithfully convex function,  $A_{0}$  is any specified  $n \times p$  matrix and  $R(A_{0})$  denotes the range space of  $A_{0}$ . Calculation of the intersection  $D_{f}^{=} \cap R(A_{0})$ is useful in the situation when the intersection of two or more cones of directions of constancy is needed. If  $A_{0} = I$ , the identity matrix, then the algorithm calculates the cone of directions of constancy of f.

The algorithm is based on the fact that  $D_f$  lies in the orthogonal complement of  $\phi$ , for any  $\phi \in \partial f(x)$ . By repeatedly considering the restriction of f to this orthogonal complement, we calculate  $D_f^{=}$ .

First we present the following two lemmas.

Lemma 2.1. Suppose that  $0 \neq d \in \mathbb{R}^k$  and  $i_0$  is the smallest positive integer such that the  $i_0$ -th component of d is nonzero, i.e.  $d_i \neq 0$ . Let

V.2

<sup>†</sup> This algorithm has been published, in the case of differentiable faithfully convex functions, in [49].



Then R(A) = N(d), where N(d) denotes the null space of d.



<u>Lemma 2.2.</u> Let  $\phi \in \partial f(x)$ , where  $f : \mathbb{R}^n \to \mathbb{R}$  is a faithfully convex function and  $x \in R^n$ . Then

$$D_{f}^{=} \subset N(\phi).$$

<u>Proof.</u> Let  $d \in D_{f}^{=}$ . Then  $\nabla f(x;d) = 0$  and Lemma II.4.1(b) implies that  $\phi \cdot d \leq 0$ . But  $-d \in D_{f}^{-}$ , since  $D_{f}^{-}$  is a subspace when f is faithfully convex. Thus  $\phi \cdot d = 0$ .

Let  $E_k = \{e^i : i = 1, ..., k\}$  denote the set of unit vectors in  $R^k$  and  $A_0 \in R^{nxp}$  be given.

# Algorithm A:

Initialization: Set  $P_0 = A_0$  and i = 1. <u>i-th step  $(1 \le i \le p)$ </u>: Find a point x in the set of p - i + 2 vectors  $\{0\} \cup E_{p-i+1}$  such that

(2.1)  $\phi P_{i-1} \neq 0$ , for some  $\phi \in \partial f(P_{i-1}x)$ .

Case (i): If such an x exists and i < p, then using Lemma 2.1, determine  $A_i \in R^{(p-i+1)x(p-i)}$  such that

(2.2) 
$$R(A_i) = N(\phi P_{i-1})$$

Set P = P A and proceed to step i + 1.

Case (ii): If such an x exists but i = p, then STOP. <u>Conclusion</u>:  $D_{f}^{=} \cap R(A_{0}) = \{0\}$ . Case (iii): If such an x does not exist, then STOP. <u>Conclusion</u>:  $D_{f}^{=} \cap R(A_{0}) = R(P_{i-1})$ .

<u>Theorem 2.1.</u> Suppose that  $f : \mathbb{R}^n \to \mathbb{R}$  is a faithfully convex function and  $A_0$  is some given nxp matrix. Then the above algorithm finds  $D_f^{=} \cap \mathcal{R}(A_0)$  in at most p - s + 1 steps, where  $s = \dim(D_f^{=} \cap \mathcal{R}(A_0))$ .



<u>Proof.</u> Let  $x^{i}$  denote the point x which satisfies (2.1) at the i-th step and for  $i \ge 0$  let  $f_{i} = f \circ P_{i}$  denote the composite function formed by applying first  $P_{i}$  and then f. By the linearity of  $P_{i}$ ,  $f_{i}$  is a faithfully convex function and so  $D_{f_{i}}^{=}$  is a fixed subspace of  $R^{p-i}$ . Furthermore,  $\partial f_{i}(x) = \partial f(P_{i}x)P_{i}$ .

Now suppose that case (i) has occurred, i.e.  $x^{i} \in \{0\} \cup \underset{p-i+1}{\mathbb{P}}$ ,  $\phi \in \partial f(\underset{i-1}{\mathbb{P}} x^{i})$ ,  $\phi \underset{i-1}{\mathbb{P}} \neq 0$  and i < p. Let us show that

(2.3) 
$$D_{f}^{\dagger} \cap R(A_{0}) = P_{i}D_{f_{i}}^{\dagger}.$$

First, let us show that

(2.4) 
$$D_{f}^{=} \cap R(A_{0}) = A_{0}D_{f_{0}}^{=}.$$

Suppose that  $d \in D_{f_0}^{=}$ . This means that  $f_0(\alpha d) = f_0(0)$  for all  $\alpha \in \mathbb{R}$ . By definition of  $f_0$  and the linearity of  $A_0$ , this gives  $f(\alpha A_0 d) = f(0)$  for all  $\alpha \in \mathbb{R}$ , i.e.  $A_0 d \in D_f^{=}$ . Furthermore, since  $A_0 d \in \mathbb{R}(A_0)$ ,  $A_0 d \in D_f^{=} \cap \mathbb{R}(A_0)$ .

Conversely, suppose that  $d \in D_{f}^{=} \cap R(A_{0})$ . Then there exists a  $\overline{d} \in R^{p}$  such that  $d = A_{0}\overline{d}$  and  $f(\alpha A_{0}\overline{d}) = f(0)$  for all  $\alpha \in R$ . Again, by definition of  $f_{0}$  and the linearity of  $A_{0}$ , we get that  $f_{0}(\alpha \overline{d}) = f_{0}(0)$  for all  $\alpha \in R$ , i.e.  $\overline{d} \in D_{f_{0}}^{=}$  where  $d = A_{0}\overline{d}$ . This proves (2.4).

(2.5) 
$$D_{f_{i-1}}^{=} = A_i D_{f_i}^{=}$$
 for  $i \ge 1$ .

Suppose that  $d \in D_{f_i}^{=}$ . This means that  $f_i(\alpha d) = f_i(0)$  for all  $\alpha \in \mathbb{R}$ . Since  $f_i = f_{i-1} \circ A_i$ , we get that  $f_{i-1}(\alpha A_i d) = f_{i-1}(0)$  for all  $\alpha \in \mathbb{R}$ , i.e.  $A_i d \in D_{f_{i-1}}^{=}$ .

Conversely, suppose that  $d \in D_{f_{i-1}}^{=}$ , i.e.  $f_{i-1}(\alpha d) = f_{i-1}(0)$ for all  $\alpha \in \mathbb{R}$ . By Lemma 2.2,  $D_{f_{i-1}}^{=} \subset \partial f_{i-1}^{\perp}(x^i) \subset N(\phi P_{i-1})$  and  $N(\phi P_{i-1}) = R(A_i)$  by (2.2). Therefore there exists a  $\overline{d} \in \mathbb{R}^{p-i}$ such that  $d = A_i \overline{d}$ . So  $f_i(\alpha \overline{d}) = f_{i-1}(\alpha A_i \overline{d}) = f_{i-1}(\alpha d) = f_{i-1}(0) = f_i(0)$ for all  $\alpha \in \mathbb{R}$ , i.e.  $\overline{d} \in D_{f_i}^{=}$  and  $d = A_i \overline{d}$ . This proves (2.5).

By repeated substitution of (2.5) into (2.4), one gets that  $D_{f}^{=} \cap R(A_{0}) = A_{0}D_{f_{0}}^{=} = A_{0}A_{1}D_{f_{1}}^{=} = \dots = P_{i}D_{f_{i}}^{=}$ , which proves (2.3).

Now suppose that case (ii) has occurred, i.e.  $x^{1} \in \{0\} \cup E_{p-i+1}$ ,  $\phi P_{i-1} \neq 0$  but i = p. Since  $f_{p-1} : R \neq R$  is faithfully convex, we get that  $D_{f_{p-1}}^{=} = \{0\}$ . But, by (2.3), the (p-1) -st step implies that  $D_{f}^{=} \cap R(A_{0}) = P_{p-1}D_{f_{p-1}}^{=}$ . Substituting for  $D_{f_{p-1}}^{=}$  yields the desired result that  $D_{f}^{=} \cap R(A_{0}) = \{0\}$ .

Finally, suppose that case (iii) has occurred, i.e.  $\partial f_{i-1}(y) = \{0\}$  for all  $y \in \{0\} \cup E_{p-i+1}$ . Then, by the convexity of  $f_{i-1}$ , the conplete set  $E_{p-i+1}$  lies in  $D_{f_{i-1}}^{=}$ . But  $D_{f_{i-1}}^{=}$  is a subspace  $f_{i-1}$ 

V.2

of  $\mathbb{R}^{p-i+1}$  and so we conclude that  $D_{f_{i-1}}^{=} = \mathbb{R}^{p-i+1}$ . Substituting into (2.3) yields  $D_{f}^{=} \cap \mathbb{R}(\mathbb{A}_{0}) = \mathbb{R}(\mathbb{P}_{i-1})$ .

The algorithm will be illustrated by two examples.

Example 2.1. Consider the function

$$\mathbf{f}(\mathbf{x}) = -\left(4 + (x_1 + x_2)^2\right)^{\frac{1}{2}} + x_1 + x_2 + x_3^2.$$

This function is convex and analytic and so is faithfully convex. Let us determine its cone of directions of constancy  $D_f^{=}$ .

Initialization: 
$$P_0 = A_0 = I_{3x3}$$
 and  $i = 1$ .  
Step 1: Since  $\nabla f(x) = \left(1 - \frac{x_1 + x_2}{(4 + (x_1 + x_2)^2)^{\frac{1}{2}}}, 1 - \frac{x_1 + x_2}{(4 + (x_1 + x_2)^2)^{\frac{1}{2}}}, 2x_3\right)$ ,

we see that  $0 \in \{0\} \cup E_3$  and  $\nabla f(P_0 0)P_0 = \nabla f(0) = (1, 1, 0) \neq 0$ . Furthermore, since i = 1 , we are in case (i). UsingLemma 2.1, we find that

$$P_{1} = A_{0}A_{1} = A_{1} = \begin{bmatrix} 1 & 0 \\ -1 & 0 \\ 0 & -1 \end{bmatrix}$$

Step 2: For 
$$\mathbf{x} \in \mathbb{R}^2$$
,  $\nabla f(\mathbb{P}_1 \mathbf{x}) \mathbb{P}_1 = \nabla f\left(\begin{pmatrix} x_1 \\ -x_1 \\ -x_2 \end{pmatrix}\right) \begin{bmatrix} 1 & 0 \\ -1 & 0 \\ 0 & -1 \end{bmatrix} = (0, 2x_2).$ 

Therefore,  $\nabla f(P_1e^2)P_1 = (0,2) \neq 0$ , where  $e^2 \in E_2$ . Furthermore, since i = 2 < p we are in case (i) again and so we find that

$$A_{2} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \text{ and } P_{2} = P_{1}A_{2} = \begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix}.$$

Step 3: The finite point set  $\{0\} \cup E_{p-i+1}$  is  $\{0,1\}$  and

$$\nabla f(P_2^0)P_2 = \nabla f(P_2^1)P_2 = (1,1,0) \begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix} = 0.$$
 Therefore, we are in case

(iii) and

V.2

$$D_{f}^{-} = R(P_{2})$$
$$= \left\{ \begin{pmatrix} d \\ -d \\ 0 \end{pmatrix} \in R^{3} : d \in R \right\}.$$

Initialization: 
$$P_0 = A_0 = \begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix}$$
 and  $i = 1$ .

<u>Step 1</u>: Since p = 1 and  $\nabla g(x) = (-1, -1, 2x_3)$  we see that  $\{0\} \cup E_1 = \{0, 1\}$  and that  $\nabla g(P_0 0)P_0 = \nabla g(P_0 1)P_0 = 0$ . Therefore, we are in case (iii) and

$$D_{\mathbf{f}}^{\mathbf{z}} \cap D_{\mathbf{g}}^{\mathbf{z}} = \mathcal{R}(P_{0}) = \left\{ \begin{pmatrix} \mathbf{d} \\ -\mathbf{d} \\ \mathbf{0} \end{pmatrix} \in \mathbb{R}^{3} : \mathbf{d} \in \mathbb{R} \right\}.$$

A computer program for the above algorithm appears in the appendix.

3. Calculating the Sets  $P^{=}$  and  $D^{=}_{p^{=}}$ 

In [3], an algorithm for calculating  $P^{=}$  is given, for the program (P). We now present a modified version of this algorithm, in the case that the constraints  $f^{k}$ ,  $k \in P^{=}$ , are faithfully convex. In actual fact, the algorithm finds

$$\overline{P}^{=}$$
 and  $\overline{D}^{=}_{\overline{P}} \cap R(A_{0})$ ,

V.3

V.3

$$\overline{P}^{=} = \{ k \in P : f^{k}(x) = 0 \text{ for all } x \in S \cap (\overline{x} + R(A_{0})) \}$$

and  $A_0$  is any specified  $n \times n_0$  matrix. (Recall that  $D_k^{=}$  is independent of x, when the function  $f^k$  is faithfully convex.) If  $A_0$  is specified to be the identity, then  $P^{=}$  and  $D_{P^{=}}^{=}$  are found. (The generalization, to find  $\overline{P}^{=}$  and  $D_{\overline{P}^{=}}^{=}$ , will be needed in the sequel.)

The algorithm is a (finite) iterative method. We start with  $P_0^{=} = \emptyset$  and find the sets  $P_{i+1}^{=} = P_i^{=} \cup J_i$  at each iteration. (The sets  $J_i$  are defined below.) The algorithm terminates when  $P_i^{=} = \overline{P}^{=}$  is reached. The difference between this algorithm and the one in [3], is that, at each iteration, we discard the constraints  $f^k$ ,  $k \in J_i$ , and, by a substitution technique, we then consider the remaining constraints as being restricted to the subspace  $D_{J_i}^{=}$ . In addition, when finding the set  $J_i$ , we first check if  $\phi = 0$  is in the subdifferential of any of the (remaining) binding constraints. (Recall that, if  $0 \in \partial f(x)$  and f is convex, then f achieves a global minimum at x.) The algorithm is demonstrated in Example 3.1 below.

Algorithm B:

<u>Initialization</u>: Let  $\overline{x} \in S$ ,  $P_0 = P(\overline{x})$ ,  $P_0^{=} = \emptyset$ ,  $P_0 = A_0$ and i = 0.

i-th step 
$$(0 \le i \le t)$$
: Find  $k \in P_i$  such that  
 $0 \in \partial f^k(\overline{x}) P_i$ .

Case (i): If such a  $\,k\,$  exists, use Algorithm A to find the  $n_{i}\,\times\,n_{i+1}\,$  matrix A satisfying

(3.1) 
$$R(A_{i+1}) = \bigcap_{k \in J_i} D^{=}_{k \circ P_i},$$

where

(3.2) 
$$J_{i} = \{k \in P_{i} : 0 \in \partial f^{k}(\overline{x})P_{i}\}.$$

Then set

(3.3) 
$$\begin{cases} P_{i+1} = P_{i}/J_{i} \\ P_{i+1} = P_{i}A_{i+1} \\ P_{i+1}^{=} = P_{i}^{=} \cup J_{i} \end{cases}$$

and proceed to step i + 1.

Case (ii): If such a k does not exist but the system



Flowchart to find  $\overline{P}$  and  $D_{\overline{P}} \cap R(A_0)$ 

(3.4) 
$$\begin{cases} \sum_{k \in P_{i}} \lambda_{k} \phi^{k} P_{i} = 0 \\ \sum_{k \in P_{i}} \lambda_{k} = 1, \lambda_{k} \ge 0, \phi^{k} \in \partial f^{k}(\overline{x}) \\ k \in P_{i} \end{cases}$$

is consistent, then find  $A_{i+1}$ ,  $P_{i+1}$ ,  $P_{i+1}$ , and  $P_{i+1}^{=}$  satisfying (3.1) and (3.3), where

$$J_{i} = \{k \in P_{i} : \lambda_{k} \neq 0 \text{ in } (3.4)\}.$$

Now proceed to step i + 1.

Case (iii): If such a k does not exist, but the system (3.4) is inconsistent, then STOP.

Conclusion:

(3.5) 
$$\begin{cases} \overline{P}^{=} = P_{i}^{=} \\ D_{\overline{P}^{=}}^{=} \cap R(A_{0}) = R(P_{i}) \\ \overline{P}^{=} & 0 \end{cases}$$

Before proving the convergence of the algorithm, let us first prove the following rather technical lemma.

Lemma 3.1. Denote

 $f_{i}^{k}(y) \stackrel{\Delta}{=} f^{k}(\overline{x} + P_{i}y)$  and  $S_{i} \stackrel{\Delta}{=} \{x \in R^{n_{i}} : f_{i}^{k}(x) \leq 0 \text{ for all } k \in P_{i}\}.$ 

Then

(3.6) 
$$f_{i}^{k}(y) = f_{i-1}^{k}(A_{i}y).$$

(3.7) 
$$D_{k}^{=} = D_{k}^{=}, \text{ for all } k \in P^{=},$$
$$f_{i}^{k} f^{k} \circ P_{i}$$

(3.8) 
$$D_{p_{i}}^{=} \cap R(A_{0}) = R(P_{i}).$$

$$(3.9) S_i \subseteq \mathcal{R}(A_{i+1}).$$

(3.10) 
$$S \cap (\overline{x} + R(A_0)) \subset \overline{x} + R(P_i).$$

Proof. Since

$$f_{i}^{k}(y) = f_{i}^{k}(x + P_{i}y) = f_{i-1}^{k}(x + P_{i-1}A_{i}y) = f_{i-1}^{k}(A_{i}y),$$

relation (3.6) is proved.

Now, when  $f^k$  is faithfully convex, there exists a strictly convex function g, a matrix B, vectors a and b and a constant c such that  $f^k(x) = g(Ax + b) + a^t x + c$ , with  $D_{f^k}^{=} = N {A \choose a^t}$ , (see Remark II.3.1). Therefore

$$f_{i}^{k}(y) = f^{k}(\overline{x} + P_{i}y)$$
$$= g(A(\overline{x} + P_{i}y) + b) + a^{t}(\overline{x} + P_{i}y) + c$$
$$= g(AP_{i}y + A\overline{x} + b) + a^{t}P_{i}y + a^{t}\overline{x} + c$$

with

$$D_{f_{i}}^{=} = N \begin{pmatrix} AP_{i} \\ a^{t}P_{i} \end{pmatrix},$$

V.3

$$f^{k} \circ P_{i}(y) = g(AP_{i}y + b) + a^{t}P_{i}y + c.$$

This proves (3.7).

Let us prove (3.8) by finite induction on i. The result holds for i = 0, since  $P_0^= = \emptyset$  and  $P_0 = A_0^-$ . So, let us assume that  $i \ge 1$  and that

(3.11) 
$$D^{=} \cap R(A_{0}) = R(P_{i-1}).$$
  
 $i-1$ 

(Note that we will consider  $D_k^{=}$  as a subset of  $R^n$  and as a subset of  $R^i$  depending on the context, i.e. depending on whether we are considering the function  $f^k$  or  $f_i^k$ .) First, suppose that  $d \in R(P_i)$ , i.e.

d = 
$$A_0 A_1 \dots A_i \overline{y}$$
, for some  $\overline{y} \in \mathbb{R}^n$ .

This implies that  $d \in R(P_{i-1}) = D^{=} \cap R(A_{0})$ , by (3.11). Now,  $p_{i-1}^{=}$ to show that  $d \in D^{=} \cap R(A_{0})$ , it is sufficient to show that  $p_{i}^{=}$   $d \in D^{=} P_{i-1}^{=} = D_{J_{i-1}}^{=}$ , by (3.3). So, let  $k \in J_{i-1}$  and  $\alpha \in R$ .  $p_{i}^{=} \setminus P_{i-1}^{=} = J_{i-1}^{=}$ 

Then

$$f^{k}(\overline{x} + \alpha d) = f^{k}(\overline{x} + \alpha A_{0} \cdots A_{\underline{i}}\overline{y})$$
$$= f^{k}_{\underline{i}-1}(\alpha A_{\underline{i}}\overline{y}), \quad by \quad (3.6)$$

= 
$$f_{i-1}^{k}(0)$$
, since  $k \in J_{i-1}$  and  $R(A_i) \subset D_k^{=}$ ,  
 $f_{i-1}^{k}$  by (3.1) and (3.7)

=  $f^{k}(\bar{x})$ , by (3.6).

This implies that  $d \in D_{j-1}^{=}$ . Thus we have shown that  $R(P_{i}) \subset D_{i-1}^{=}$ .  $P_{i-1}^{=}$ 

Conversely suppose that .

$$d \in D^{=} \cap R(A_{0})$$

Since  $P_i^{\overline{i}} \supset P_{i-1}^{\overline{i}}$ , (3.11) implies that  $d = A_0 A_1 \dots A_{i-1} \overline{y}$ , for some  $\overline{y} \in \mathbb{R}^{n-1}$ . To show that  $D_{p_i^{\overline{i}}}^{\overline{i}} \cap \mathbb{R}(A_0) \subset \mathbb{R}(P_i)$ , it is now sufficient  $P_i^{\overline{i}}$ .

to show that

$$\overline{y} = A_{i}\overline{z}$$
, for some  $\overline{z} \in R^{n}$ .

Suppose that  $k \in P^{=} \setminus P^{=} = J$  and  $\alpha \in \mathbb{R}$ . Then

$$f_{i-1}^{k}(0) = f^{k}(\overline{x}), \text{ by } (3.6)$$

$$= f^{k}(\overline{x} + \alpha d), \text{ since } d \in D_{p_{i}}^{=} \text{ and } k \in P_{i}^{=}$$

$$= f^{k}(\overline{x} + A_{0} \cdots A_{i-1}(\alpha \overline{y}))$$

$$= f_{i-1}^{k}(\alpha \overline{y}), \text{ by } (3.6).$$

This implies that  $\overline{y} \in D_{j-1}^{=} = R(A_j)$ , by (3.1). Thus,  $\overline{y} = A_i \overline{z}$ 

for some  $\overline{z} \in \mathbb{R}^{n}$ . This completes the proof of (3.8).

To prove (3.9), we consider two separate cases.

Case (a): Suppose that  $0 \in \partial f_i^k(0)$ , for some  $k \in P_i$ . (Note that  $\partial f_i^k(0) = \partial f^k(\overline{x})P_i$ .) By (3.2),  $0 \in \partial f_i^k(0)$  for all  $k \in J_i$ . Therefore, y = 0 is a global minimum for the convex functions  $f_i^k$ ,  $k \in J_i$ . Now, suppose that  $\overline{y} \in S_i$ , i.e.  $f_i^k(\overline{y}) \leq 0$  for all  $k \in P_i$ . Then  $f_i^k(\overline{y}) = 0$  for all  $k \in J_i$ , since y = 0 is a global minimum for these functions and  $f_i^k(0) = 0$ . Since  $S_i$  is convex and  $0 \in S_i$ , we conclude that  $\alpha \overline{y} \in S_i$  for all  $0 \leq \alpha \leq 1$ . This further implies that  $f_i^k(\alpha \overline{y}) = 0$ , for all  $k \in J_i$  and  $0 \leq \alpha \leq 1$ , i.e.  $\overline{y} \in D_{J_i}^{=} = R(A_{i+1})$ . This proves (3.9), in case (a).

Case (b): Suppose that  $0 \notin \partial f_i^k(0)$  for all  $k \in P_i$ . Also, assume that the system (3.4) is consistent, i.e. there exist  $\lambda_k > 0$ such that

(3.12) 
$$\sum_{k \in J_{i}}^{\Sigma} \lambda_{k} \phi^{k} = 0, \ \phi^{k} \in \partial f_{i}^{k}(0).$$

(Note that if no such  $\lambda_k$ 's exist, then the algorithm stops and (3.9) does not require proof.) As in case (a), we need only show that

if 
$$\overline{y} \in S_i$$
 and  $k \in J_i$ , then  $f_i^k(\overline{y}) = 0$ .

Suppose not. Then there exists  $y \in S_i$  and  $k \in J_i$  such that  $f_i^k(y) \leq 0$ , for all  $k \in P_i$ , and  $f_i^{k \circ}(y) < 0$ . This implies that

$$\phi^{k} \cdot y \leq 0$$
, for all  $k \in J_{i}$ , and  $\phi^{k} \cdot y < 0$ ,

153

for all  $\phi^k \in \partial f_i^k(0)$ ,  $k \in J_i$  and  $\phi^{k_o} \in \partial f_i^{k_o}(0)$ . By Motzkin's theorem of the alternative [52], this contradicts (3.12). Therefore, (3.9) is proved.

Let us now prove (3.10), by finite induction on i. The result holds for i = 0, since  $P_0 = A_0$ . So, let us assume that  $i \ge 0$  and

$$S \cap (\overline{x} + R(A_0)) \subset \overline{x} + R(P_{i-1}).$$

Let  $x \in S \cap (\overline{x} + R(A_0))$ . Then the above implies that  $x = \overline{x} + P_{i-1}\overline{y}$ for some  $\overline{y} \in R^{n_{i-1}}$ . Thus

$$f_{i-1}^{k}(\overline{y}) = f^{k}(x), \text{ by } (3.6)$$
$$\leq 0, \text{ since } x \in S.$$

Therefore,  $\overline{y} \in S_{i-1}$ . Now, by (3.9),  $\overline{y} = A_i \overline{z}$  for some  $\overline{z} \in R^{n-1}$ . Substituting for  $\overline{y}$  in the expression for x, implies that  $x = \overline{x} + A_0 \dots A_i \overline{z}$ , which proves (3.10). (Note that the sets  $\overline{x} + R(P_i)$  are decreasing linear manifolds containing the set  $S \cap (\overline{x} + R(A_0))$ . The algorithm essentially stops when  $\overline{x} + R(P_i)$  is the smallest such linear manifold.)

We are now ready to prove the convergence of the algorithm. Recall that

$$\overline{P} = \{k \in P : f^k(x) = 0 \text{ for all } x \in S \cap (\overline{x} + R(A_0))\}.$$

<u>Theorem 3.1.</u> Suppose that  $\overline{x} \in S$ ,  $A_0$  is an arbitrary  $n \times n_0$  matrix and the constraints  $f^k$ ,  $k \in \overline{P}^{-}$  are faithfully convex. Then the above algorithm finds

$$\overline{P}$$
 and  $D_{\overline{P}} \cap R(A_0)$ 

in at most t = min{card $P(\overline{x})$ ,  $n_0 + 1 - dim(S \cap (\overline{x} + R(A_0)))$ } steps.

<u>Proof.</u> We need to prove that (3.5) holds when case (iii) occurs. So, suppose that  $0 \notin \partial f_i^k(0)$ , for all  $k \in P_i$ , and the system (3.4) is inconsistent. This implies that the system

$$\sum_{k \in P_{i}}^{\Sigma} y^{k} = 0, y^{k} \in cone \ \partial f^{k}(\overline{x}) P_{i},$$

is inconsistent. We now conclude, by the Dubovitski-Milyutin theorem of the alternative [52] and by Lemma II.7.1, that

$$\bigcap_{k \in \mathcal{P}_{i}}^{n} \{ y \in \mathbb{R}^{n_{i}} : \phi^{k_{i}} \varphi < 0, \text{ for all } \phi^{k_{i}} \in \partial_{f}^{k_{i}}(\overline{x}) \} \neq \emptyset.$$

This yields  $\hat{y} \in \mathbb{R}^n$  such that

(3.13)  $\phi^k P_i \cdot \hat{y} < 0$ , for all  $k \in P_i$  and  $\phi^k \in \partial f^k(\overline{x})$ .

Let

(3.14) 
$$\mathbf{x}(\alpha) = \mathbf{x} + \mathbf{P}_{\mathbf{i}} \alpha \mathbf{\hat{y}}.$$

Then, (3.13) and (3.14) imply that

(3.15) 
$$\begin{cases} f^{k}(x(\alpha)) < 0, & \text{for all } k \in P \setminus P(\overline{x}), \\ f^{k}(x(\alpha)) < 0, & \text{for all } k \in P_{i}, \end{cases}$$

for all  $0 \le \alpha \le \overline{\alpha}$ , for some  $\overline{\alpha} > 0$ . Furthermore, if  $0 \le \alpha \le \overline{\alpha}$ , then

(3.16) 
$$f^{k}(x(\alpha)) = f^{k}(\overline{x} + P_{i}\alpha\hat{y}), \text{ by (3.14)}$$
  
=  $f^{k}(\overline{x}), \text{ for all } k \in P_{i}^{=}, \text{ by (3.8)}$   
= 0, for all  $k \in P_{i}^{=}, \text{ since } P_{i}^{=} \subset P(\overline{x}).$ 

Therefore, (3.15) and (3.16) imply that  $x(\alpha) \in S \cap (\overline{x} + R(A_0))$ and, moreover,

 $\overline{P}^{=} \subset (P(\overline{x}) \setminus P_{i}) = P_{i}^{=}.$ 

Since  $D_{i}^{=} \cap R(A_{0}) = R(P_{i})$  by (3.8), to prove (3.5) we have only  $P_{i}^{=}$ 

left to show that

$$(3.17) P_{i}^{\overline{e}} \subset \overline{P}^{\overline{e}}.$$

Let us prove this by finite induction on i. Now, (3.17) holds for i = 0, since  $P_0^{=} = \emptyset$ . Therefore, let us assume that  $i \ge 1$  and

 $P_{i-1}^{=} \subset \overline{P}^{=}$ . Since  $P_{i}^{=} = P_{i-1}^{=} \cup J_{i-1}$ , by iteration, it is sufficient to show that  $J_{i-1} \subset \overline{P}^{=}$ . Suppose not. Then, there exists  $x \in S \cap (\overline{x} + R(A_0))$  and  $k_{\circ} \in J_{i-1}$  such that

(3.18) 
$$f^{k}(x) \leq 0$$
 for all  $k \in P$  and  $f^{k \circ}(x) < 0$ .

But  $x = \overline{x} + A_0 \dots A_i y$  for some  $y \in \mathbb{R}^{n_i}$ , by (3.10), and

$$f^{k_{\circ}}(x) = f^{k_{\circ}}(\overline{x} + A_{0} \dots A_{i}y)$$
  
=  $f^{k_{\circ}}_{i-1}(A_{i}y)$ , by (3.6)  
= 0, since  $D^{=}_{J_{i-1}} = R(A_{i})$ , by (3.1) and (3.7).

This contradicts (3.18).

Example 3.1. Suppose  $S \subseteq R^5$  is defined by the constraints  $f^{1}(x) = e^{x_1} + x_2^2$   $-1 \le 0$   $f^{2}(x) = x_1^2 + x_2^2 + e^{-x_3}$   $-1 \le 0$   $f^{3}(x) = x_1 + x_2^2 + e^{-x_3}$   $-1 \le 0$   $f^{4}(x) = e^{-x_2} - 1 \le 0$   $f^{5}(x) = (x_1 - 1)^2 + x_2^2 - 1 \le 0$   $f^{6}(x) = x_1 + e^{-x_4} - 1 \le 0$  $f^{7}(x) = x_2 + e^{-x_5} - 1 \le 0.$ 

Let us find  $P^{=}$  and  $D^{=}_{p^{=}}$ .

Initialization: Let  $\overline{\mathbf{x}} = (0,0,1,\frac{1}{2}\sqrt{2},\frac{1}{2}\sqrt{2})$  be the chosen feasible point. Then  $P_0 = A_0 = I_{5\times5}$ ,  $P_0 = P(\overline{\mathbf{x}}) = \{1,3,4,5\}$  and  $P_0^{=} = \emptyset$ . The corresponding gradients are  $\nabla \mathbf{f}^1(\overline{\mathbf{x}}) = (1,0,0,0,0)$  $\nabla \mathbf{f}^3(\overline{\mathbf{x}}) = (1,0,0,\sqrt{2},\sqrt{2})$  $\nabla \mathbf{f}^4(\overline{\mathbf{x}}) = (0,-1,0,0,0)$  $\nabla \mathbf{f}^5(\overline{\mathbf{x}}) = (-2,0,0,0,0)$ . <u>Step 0:</u> Since  $\nabla \mathbf{f}^k(\overline{\mathbf{x}}) A_0 = \nabla \mathbf{f}^k(\overline{\mathbf{x}}) \neq 0$  for all  $k \in P_0$ , we solve the system given by (3.4), i.e.

А

V.3

$$\underbrace{\text{Step 1:}}_{1} \text{ Since } \nabla f^{4}(\overline{x}) P_{1} = 0 \text{ while } \nabla f^{3}(\overline{x}) P_{1} \neq 0 \text{ we get that}$$

$$J_{1} = \{4\}, P_{2} = \{3\}, P_{2}^{\pm} = \{1, 4, 5\}, A_{2} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \text{ with } R(A_{2}) = D_{1}^{\pm} A_{2} = P_{1}^{A} A_{2} = A_{1}^{\pm}.$$

$$\underbrace{\text{Step 2:}}_{1} \text{ Since } P_{2} = \{3\} \text{ and } \nabla f^{3}(\overline{x}) P_{2} \neq 0 \text{ case (iii)}$$
occurs. STOP.

Conclusion:

$$P^{=} = P_{2}^{=} = \{1, 4, 5\}$$

and

$$D_{p}^{=} = R(P_{2}) = \{ \begin{pmatrix} 0 \\ 0 \\ 0 \\ d_{3} \\ d_{4} \\ d_{5} \end{pmatrix} \in R^{5} : d_{3}, d_{4}, d_{5} \in R \}.$$

<u>Remark 3.1.</u> Using the substitution technique and checking whether  $\nabla f^{k}(\bar{x})P_{i} = 0$ , reduces the number of computations required to find  $P^{=}$  and  $D_{p^{=}}^{=}$ .

### 4. The Method of Reduction

We now collect the machinery presented in the previous two sections and formulate the method of reduction. This algorithm first finds a feasible point and then solves the general convex program (P) with faithfully convex constraints. No constraint qualification need be assumed. Let us denote by

### S(P),

a method that solves program (P) under the assumption that Slater's condition is satisfied (e.g. Zoutendijk's feasible directions method [53], Robinson's method [52] or Powell's method [52]). The method of reduction finds the regularized program (P<sub>r</sub>) of Theorem IV.4.2, in the process of finding a feasible point. It then solves (P<sub>r</sub>), using  $S(P_r)$ . Furthermore, if Slater's condition was not satisfied for the original program (P), the regularized program (P<sub>r</sub>) will always have fewer constraints and fewer variables than (P).

# Algorithm C:

V.4



160

Flowchart for the Method of Reduction

V.4

161

and

(4.2) 
$$T_i = T_{i-1}/R_i$$
.

Now consider the program

$$(R_{i}) \qquad s.t. \qquad \begin{array}{c} \Sigma \quad f^{k}(\overline{x^{i}} + P_{i-1}y) \rightarrow \min \\ k \in T_{i} \\ s.t. \\ f^{k}(\overline{x^{i}} + P_{i-1}y) \leq 0, \ k \in R_{i}, \ y \in R^{i-1}. \end{array}$$

Using the feasible point 0 and Theorem IV.4.2, regularize this program, i.e. find  $R_{i}^{=}$  (Note that  $R_{i}^{=}$  is the equality set for  $(R_{i})$ .) and the  $n_{i-1} \times n_{i}$  matrix  $A_{i}$  satisfying  $D_{i}^{=} = R(A_{i})$  (Use  $R_{i}^{=}$  algorithm B). Now set

$$(4.3) \qquad \qquad \hat{R}_{i} = R_{i} \setminus R_{i}^{=}$$

and

$$(4.4) P_i = P_{i-1}A_i$$

to get the 'reduced' program

$$(\hat{R}_{i}) \qquad \begin{array}{c} \Sigma \quad f^{k}(\overline{x}^{i} + P_{i}y) \neq \min \\ k \in T_{i} \\ \text{s.t.} \\ f^{k}(\overline{x}^{i} + P_{i}y) \leq 0, \ k \in \hat{R}_{i} \quad \text{and} \quad y \in \mathbb{R}^{n}. \end{array}$$

(Note that Slater's condition is satisfied and 0 is a feasible point, by Theorem IV.4.2.)

Case (i): Suppose that  $T_i = \emptyset$ . Set  $\overline{x} = \overline{x}^i$  and  $P = P_i$ and solve the 'reduced' program,

162

(R) 
$$f^{\circ}(\overline{x} + Py) \neq \min$$
  
s.t.  
 $f^{k}(\overline{x} + Py) \leq 0, k \in \hat{R}_{i}$  and  $y \in \mathbb{R}^{i}$ ,

using the initial feasible point y = 0 and S(R).

<u>Conclusion:</u> If  $y^*$  is a solution of R, then  $x^* = \overline{x} + Py^*$  is a solution of the original program (P).

Case (ii): Suppose that  $T_i \neq \emptyset$ . Then, set  $z_0^i = 0$  and  $z_{j+1}^i = S(\hat{R}_i) z_j^i$ , j = 0, 1, ...

i.e.  $z_{j+1}^{1}$  is the point obtained after one iteration of  $S(\hat{R}_{i})$ , applied to the point  $z_{i}^{i}$ .

Case (ii) (a): Suppose that after j iterations of  $S(\hat{R}_i)$ , we find  $k \in T_i$  such that  $f^k(\overline{x}^i + P_i z_j^i) \le 0$ . Then set

(4.5) 
$$\begin{cases} z^{i} = z^{i}_{j} \\ \overline{x^{i+1}} = \overline{x^{i}} + P_{j} z^{i} \end{cases}$$

and proceed to step i + 1.

Case (ii) (b): Suppose that after j iterations of  $S(\hat{R}_i)$ , we have not found  $k \in T_i$  such that  $f^k(\overline{x^i} + P_i z_j^i) \leq 0$  but  $z_j^i$  solves the program  $(\hat{R}_i)$ .

Before proving convergence, we prove the following lemma.

Lemma 4.1.

(4.6) 
$$\stackrel{i}{\cap} _{D}^{=} = R(P_{i}),$$
  
 $j=1 R_{j}^{=}$ 

(4.7) 
$$f^k(\overline{x^i}) \leq 0$$
, for every  $k \in P \setminus T_i$ .

(4.8) 
$$S \subset \overline{x}^i + \mathcal{R}(P_{i-1}).$$

$$(4.9) R_i^{=} \subset P^{=}.$$

<u>Proof.</u> Let us prove (4.6) by finite induction on i. The result holds for i = 1, by (4.4). So, let us assume that  $i \ge 2$  and that

(4.10) 
$$\begin{array}{c} i-1 \\ \cap D^{=} = R(P_{i-1}), \\ j=1 R_{i}^{=} \end{array}$$

(Recall that we consider  $D^{=}_{R_{j}^{=}}$  as a subset of  $R^{n}$  and as a subset  $R_{j}^{n}$  of  $R^{j-1}$ , depending on the context.) Suppose that  $d \in R(P_{i})$ , i.e.  $d = P_{i} y$  for some  $y \in R^{n}$ . Then  $d \in R(P_{i-1}) = \bigcap_{j=1}^{i-1} D^{=}_{j}$ . But, if  $k \in R_{i}^{=}$  and  $\alpha \in R$ , then

$$f^{k}(\overline{x^{i}} + \alpha d) = f^{k}(\overline{x^{i}} + P_{i-1}(\alpha A_{i}y))$$
$$= f^{k}(\overline{x^{i}}), \text{ since } D^{=}_{R_{i}} = R(A_{i}).$$

Therefore,  $d \in D^{\overline{}}$ . This shows that  $\mathcal{R}_{i}^{\overline{}}$ 

$$R(P_{i}) \subset \bigcap_{j=1}^{i} \mathbb{R}_{j}^{=}$$

Conversely, suppose that  $d \in \bigcap_{j=1}^{1} D^{=}$ . By (4.10), d = P y for j=1  $\mathcal{R}_{j}^{=}$ 

some  $y \in \mathbb{R}^{n}$ . To show  $\bigcap_{D}^{=} \subset \mathcal{R}(\mathbb{P}_{j})$  it is now sufficient to show j=1  $\mathbb{R}_{j}^{=}$ 

that  $y = A_i z$  for some  $z \in R^i$ . Suppose that  $k \in R_i^=$  and  $\alpha \in R$ . Then

$$f^{k}(\overline{x^{i}}) = f^{k}(\overline{x^{i}} + \alpha d), \text{ since } d \in D^{=}$$

$$R^{=}_{i}$$

$$= f^{k}(\overline{x^{i}} + P_{i-1}\alpha y).$$

This implies that  $y \in D^{=}$ . Now, since  $D^{=} = R(A_{i})$ , we see that  $R^{=}_{i}$ 

 $y = A_i z$  for some  $z \in \mathbb{R}^{n}$ . This proves (4.6).

Let us also prove (4.7) by finite induction on i. The result holds for i = 1 by the initialization and the definition of  $T_1$ . So, let us assume that  $i \ge 2$  and

(4.11) 
$$f^k(\overline{x}^{i-1}) \leq 0$$
, for every  $k \in P \setminus T_{i-1}$ 

Suppose that  $k \in R_i$ . Then

$$f^{k}(\vec{x^{i}}) = f^{k}(\vec{x^{i-1}} + P_{i-1}z^{i-1})$$
  
 $\leq 0$ , by (4.1) and the fact that  $z^{i-1}$  is a  
feasible point of the problem  $\hat{R}_{i-1}$ .

On the other hand, suppose that  $k \in \bigcup_{j=1}^{i-1} R_j^{=} = (P \setminus T_i) \setminus R_i$ . Then

$$f^{k}(\vec{x^{i}}) = f^{k}(\vec{x^{i-1}} + P_{i-1}z^{i-1})$$
$$= f^{k}(\vec{x^{i-1}}), \text{ by } (4.6)$$
$$\leq 0, \text{ by } (4.11).$$

This proves (4.7).

As above, let us prove (4.8) by finite induction on i. The result holds for i = 1, since  $P_0 = I_{n \times n}$ . So, let us assume that  $i \ge 2$  and

(4.12) 
$$S \subseteq \overline{x^{i-1}} + R(P_{i-2}).$$

Suppose that  $x \in S$ . Then,  $x = \overline{x^{i-1}} + P_{i-2}y$  for some  $y \in \mathbb{R}^{n-2}$ , by (4.12). So,

> $f^{k}(x) = f^{k}(\overline{x^{i-1}} + P_{i-2}y)$  $\leq 0, \text{ for every } k \in P, \text{ since } x \in S.$

Therefore, y is a feasible point for the program  $(R_{i-1})$ . Since

 $A_{i-1}z^{i-1}$  is a feasible point for  $(R_{i-1})$  and  $D_{i-1}^{=} = R(A_{i-1})$ , we get that  $y = A_{i-1}z^{i-1} + A_{i-1}z$  for some  $z \in R^{n_{i-1}}$ . Substituting for y gives

$$x = \overline{x^{i-1}} + P_{i-2}(A_{i-1}z^{i-1} + A_{i-1}z)$$
  
=  $\overline{x^{i}} + P_{i-1}z$ ,

i.e.  $x \in \overline{x}^{i} + R(P_{i-1})$ . This proves (4.8).

Suppose that (4.9) fails to hold. Then, there exists  $k_o \in R_i^{=}$ such that  $k_o \in P^{=}$ . This implies that there exists an  $\hat{x} \in R^n$  such that  $f^k(\hat{x}) \leq 0$ , for  $k \in P$ , and  $f^{k_o}(\hat{x}) < 0$ . But, by (4.8), we see that  $\hat{x} = \overline{x^{i-1}} + P_{i-1}y$  for some  $y \in R^{n_i-1}$ . Therefore,  $f^k(\overline{x^{i-1}} + P_{i-1}y) \leq 0$ , for  $k \in P$ , but  $f^{k_o}(\overline{x^{i-1}} + P_{i-1}y) < 0$ . This contradicts the assumption that  $k_o \in R_i^{=}$ .

Let us now prove convergence of the algorithm.

<u>Theorem 4.1.</u> Assume that (P) and  $f^k$ ,  $k \in P$ , are such that S(P) is a convergent method when Slater's condition is satisfied. Furthermore, suppose that the constraints  $f^k$ ,  $k \in P^=$ , are faithfully convex. Then, the method of reduction first finds a feasible point  $\overline{x}$  and then solves program (P) by solving the 'reduced' program (R). (If  $y^*$  solves (R), then  $x^* = \overline{x} + Py^*$  solves (P).)

<u>Proof.</u> From (4.7), we get that y = 0 is a feasible point for  $(R_i)$ . Therefore, from Theorem IV.4.2, we get that Slater's condition is satisfied and y = 0 is a feasible point for the regularized program  $(\hat{R}_i)$ . Now, let us treat each of the cases separately.

Case (i): From (4.6), (4.7), (4.9) and the fact that Slater's condition is satisfied for (R), we see that program (R) is actually the program obtained in Theorem IV.4.2, after regularizing our original program (P). We can thus solve (R) using S(R). That  $x^* = \overline{x}^i + P_i y^*$  solves (P), if  $y^*$  solves (R), follows from Theorem IV.4.2.

Case (ii) (a): In this case we just proceed to step i + 1.

Case (ii) (b): Suppose that  $S \neq \emptyset$  and  $x \in S$ . Then,  $x = \overline{x^{i}} + P_{i-1}^{z}$  for some  $z \in \mathbb{R}^{n_{i-1}}$ , by (4.8). Moreover,

 $f^{k}(\overline{x}^{i} + P_{i-1}^{z}) \leq 0, \quad k \in \mathcal{P}.$ 

Therefore, z is a feasible point for the program  $(R_i)$  and  $\sum_{k \in T_i} f^k (\overline{x}^{-i} + P_{i-1} z) \leq 0$ . But then,  $z = A_i y$  for some  $y \in R^{n_i}$ . Substituting for z, we get that  $f^k (\overline{x}^i + P_i y) \leq 0$ ,  $k \in P$ , and  $\sum_{i} f^k (\overline{x}^i + P_i y) \leq 0$ . Since this implies that y is a feasible  $k \in T_i$ point for the program  $(\hat{R}_i)$ , we have contradicted the optimality of  $z_i^i$ .

<u>Remark 4.1.</u> We have assumed that S(P) solves programs of type (P), when Slater's condition is satisfied. However, it may happen that the objective function and the feasible set (for one of the regularized programs  $(\hat{R}_i)$ ) may have a common direction of recession. Such programs are called degenerate, see e.g. Abrams [2]. The infimum for degenerate programs may not be achieved or may be achieved on an unbounded set. Abrams [2] has shown how to reduce such a program to a nondegenerate program in a finite number of steps. Another possible way of handling this situation is to find a K > 0, large enough, and add the constraint

 $f^{m+1}(x) = ||x||^2 - K \le 0.$ 

Such a K can be found, if a solution  $x^*$  exists, for our original program (P), and if we can approximate  $||x^*||$ . Adding the constraint  $f^{m+1}$  and choosing  $\overline{x}^1$  such that  $||\overline{x}^1|| < K$ , will ensure that the programs  $(\hat{R}_i)$  and the final regularized program (R) are nondegenerate.

### 5. Applications

The method of reduction is in particular applicable to convex programming problems for which Slater's condition is not satisfied. As mentioned in [14], one class of problems for which Slater's condition is never satisfied is the class of multicriteria problems. This includes the lexicographic problem and the Pareto optimal problem. Following [15], let us first define the lexicographic problem as follows: Suppose that  $f^1$ , ...,  $f^m$  is an ordered set of objectives. The corresponding <u>lexicographic problem (PL)</u> consists in choosing decisions successively subject to  $x \in R^n$ . The set of all optimal solutions of (PL) can then be obtained by solving the following sequence of programming problems:

> (PL<sub>1</sub>)  $a_1 = \min\{f^1(x) : x \in R^n\}.$

Determine  

$$(PL_2)$$
  
 $a_2 = \min\{f^2(x) : f^1(x) - a_1 \le 0\}.$ 

(PL<sub>m</sub>) Solve  $({}^{\text{PL}}_{\text{m}}) = \min\{f^{\text{m}}(x) : f^{\text{k}}(x) - a_{\text{k}} \le 0, \text{ k} = 1, \dots, m - 1\}.$ 

Note that Slater's condition is never satisfied for the programs (PL<sub>i</sub>). In fact,

$$\{1, \ldots, j - 1\} \subset PL_{j}^{=}, \quad j = 2, \ldots, m,$$

V.5
where  $PL_{j}^{=}$  is the equality set of  $(PL_{j})$ . Now suppose that the functions  $f^{k}$ , k = 1, ..., m - 1, are faithfully convex. Then, if we use the method of reduction, solving (PL) reduces to solving m unconstrained problems. After finding the matrix A such that  $R(A) = \bigcap_{i=1}^{j-1} D_i^{=}$  we see that, to solve  $\hat{i}=1$ (PL<sub>i</sub>), we need only find

$$\min_{y} f^{j}(x^{j-1} + Ay),$$

where  $x^{j-1}$  is a solution of  $(PL_{j-1})$ . The lexicographic problem is treated in greater detail by Ben-Tal and Zlobec [15]. They provide two different methods for finding solutions.

Let us now define the Pareto optimum problem:

(PP)  $\begin{cases} Find \ \overline{x} \in R^n & \text{such that there is no other point } x \\ satisfying \\ f^k(x) \leq f^k(\overline{x}), \ k \in P, \\ \text{with at least one strict inequality.} \end{cases}$ 

The point  $\overline{x}$  is then called Pareto optimal or efficient, see e.g. [9]. Note that any lexicographic solution of  $f^k$ ,  $k \in P$ , in any order, is a Pareto optimal solution. In fact, one can say even more.

<u>Theorem 5.1.</u> Suppose that  $\Omega_1, \Omega_2, \ldots, \Omega_r$  is any disjoint partition of the index set P and

$$g^{k}(\cdot) \stackrel{\Delta}{=} \sum_{i \in \Omega_{k}} f^{i}(\cdot), \text{ for } k = 1, \dots, r.$$

If  $x^*$  is the lexicographic solution with the ordered set of objectives  $g^1, \ldots, g^r$ , then  $x^*$  is a Pareto optimal solution of (PP).

<u>Proof.</u> Suppose not. Then there exists  $x \neq x^*$ ,  $k_{\circ}$  and  $\ell$ such that  $k_{\circ} \in \Omega_{\ell}$ ,  $f^{k_{\circ}}(x) < f^{k_{\circ}}(x^*)$  and  $f^{k}(x) \leq f^{k}(x^*)$  for all  $k \in P$ . But then  $g^{k}(x) \leq g^{k}(x^*)$  for k = 1, ..., r and  $g^{\ell}(x) < g^{\ell}(x^*)$ , which contradicts the fact that  $x^*$  solves the lexicographic problem. Thus  $x^*$  is a Pareto optimal point.

Using this theorem, we can find efficient points by solving a finite number of unconstrained optimization problems. However, not all the efficient points of (PP) can be found in this way. Charnes and Cooper [18] have given the following characterization of efficiency.

Theorem 5.2.  $\overline{x}$  is a Pareto optimum for (PP) if and only if  $\overline{x}$  solves the program

$$(P_{\overline{x}}) \qquad \begin{array}{c} \Sigma \ f^{k}(x) \rightarrow \min \\ k \in P \\ \text{s.t.} \\ f^{k}(x) - f^{k}(\overline{x}) \leq 0, \ k \in P. \end{array}$$

This result leads to the following characterization of efficiency in terms of the equality set of  $(P_{\overline{x}})$ .

Theorem 5.3.  $\overline{x}$  is a Pareto optimum for (PP) if and only if

$$P = P = \frac{1}{x},$$

where  $P_{\overline{x}}^{=}$  is the equality set of  $(P_{\overline{x}})$ .

<u>Proof.</u> The proof is immediate from the definition of Pareto optimum.

The above characterization suggests the following algorithm, when the functions  $f^k$ ,  $k \in P$ , are faithfully convex.

## Algorithm D:

Initialization: Let  $P = I_{n \times n}$  and  $\overline{x} = \overline{y} = 0$ .

<u>i-th step  $(1 \le i \le card P)$ </u>: Use the feasible point  $\overline{x}$  and  $S(P_{\overline{x}})$  to perform one iteration for  $(P_{\overline{x}})$ . Redefine  $\overline{x}$  to be the new point obtained, i.e.  $\overline{x} \stackrel{\Delta}{=} S(P_{\overline{x}})(\overline{x})$ . Continue until  $P_{\overline{x}}^{\overline{=}} \neq \emptyset$ . Now, using Theorem 2.1, find the matrix A such that  $R(A) = D_{\overline{x}}^{\overline{=}}$  and let

$$P = P \setminus P_{\overline{x}}^{=}, f^{k}(\cdot) = f^{k}(\overline{x} + A \cdot)$$
  
 $\overline{y} = \overline{y} + P\overline{x}, P = PA \text{ and } \overline{x} = 0$ 



Flowchart to find a Pareto optimum



Case (i): If  $P \neq \emptyset$ . Proceed to step i + 1.

Case (ii): If  $P = \emptyset$ , STOP.

<u>Conclusion:</u>  $\overline{y}$  is a Pareto optimum.

Let us now consider the semi-infinite programming problem,

(PS) 
$$g^{\circ}(x) \rightarrow \min$$
  
s.t.  
 $g^{k}(x,t) \leq 0, t \in T^{k}, k \in P = \{1,...,m\}, x \in \mathbb{R}^{n},$ 

where  $g^{\circ}$  is convex,  $g^{k}(x;t)$  is convex in x for each  $t \in T^{k}$ and continuous in t for each x and  $T^{k}$  is compact in  $R^{\ell}$ . Ben-Tal, Kerzner and Zlobec [13] have presented a characterization of optimality for (PS) which does not require a constraint qualification. Another way of treating program (PS) is by considering the convex functions

$$f^{k}(x) \stackrel{\Delta}{=} \sup_{t \in T^{k}} g^{k}(x,t), \quad k \in \mathcal{P}.$$

This reduces (PS) to the form (P). We can now apply the results in the previous chapters. A third way of treating (PS) is by discretization, i.e. let  $T_i^k$ , i = 1, 2, ... be finite subsets of  $T^k$  such that  $\lim T_i^k$  is dense in  $T^k$  for each  $k \in P$ . Now solve the sequence of programs (of type (P))

$$(P_{i}) \qquad s.t. \qquad g^{\circ}(x) \rightarrow \min \\ s.t. \qquad g^{k}(x,t) \leq 0, t \in T_{i}^{k}, k \in P.$$

Let  $X_i^*$  denote the set of optimal solutions  $(P_i)$  and let  $x_i^* \in X_i^*$ . Then every cluster point of the sequence  $\{x_i^*\}$  is an optimal point of (PS). Conversely, every optimal point of (PS) is the limit of a subsequence of  $\{x_i^*\}$  for some  $x_i^* \in X_i^*$ . Let us outline the proof: First, using the fact that a compact set-valued map is upper semi-continuous (u.s.c.) if and only if it's graph is closed (see e.g. Debreu [22]) we can show that the maps  $\Omega^k : T^k \to \{x \in R^n : g^k(x,t) \le 0\}$  are u.s.c. for each  $k \in P$ . This then implies that

 $\bigcap_{t \in T_{i}^{k}} \Omega^{k}(t) = \bigcap_{t \in T^{k}} \Omega^{k}(t), \quad k \in P,$ 

and therefore  $S_i \rightarrow S$  where  $S_i$  is the feasible set of  $(P_i)$  and S is the feasible set of (PS). The result now follows from Fiacco [26, Theorem 2.1].

### 6. Examples

To illustrate the method of reduction we consider the following three examples, which were solved using the computer program given in the appendix.

# Example 6.1. Consider the program

(P) s.t. 
$$f'(x) = x_1 + x_2 + x_3 \neq \min$$
  
 $f^1(x) = x_1^2 + x_2^2 - 2 \le 0$   
 $f^2(x) = (x_1^{-2})^2 + (x_2^{-2})^2 - 2 \le 0$   
 $f^3(x) = e^{-x_3} - 1 \le 0$ .

For the initial estimate, let us choose

$$x^{\circ} = (0,0,0)^{t}$$
.

Since

$$f^{1}(x^{\circ}) < 0, f^{3}(x^{\circ}) = 0$$
 while  $f^{2}(x^{\circ}) > 0,$ 

the algorithm begins by considering the program

$$(P_1) \qquad f^2(x) \neq \min_{\substack{f \\ f^1(x) \leq 0 \\ f^3(x) \leq 0.}}$$

Slater's condition is satisfied and  $x^{\circ}$  is a feasible starting point for  $(P_1)$ . Applying Zoutendijk's method, yields the solution

$$\bar{x} = (1, 1, 5.65)^{t}$$
.

Since  $f^2(\overline{x}) = 0$ , we can eliminate the last constraint from the objective function. We now consider the original program (P) with the feasible point  $\overline{x}$ . The algorithm now finds that

$$P^{=} = \{1, 2\}$$
 and  $D^{=} = R(P) = R\begin{pmatrix} 0\\ 0\\ 1 \end{pmatrix}$ 

After substituting P, we get the equivalent reduced program with one variable y and one constraint:

$$f^{\circ}(\overline{x} + Py) = \overline{x} + y \neq \min$$
  
s.t.  
$$f^{3}(\overline{x} + Py) = e^{-(5.65+y)} - 1 \le 0$$
  
$$y \in \mathbb{R}.$$

Zoutendijk's method yields the solution

and thus

$$x^* = \overline{x} + Py^* = (1,1,0)^t$$

is the solution of our original program (P). The above problem, when solved by the computer program in the appendix, gave the solution correct to 8 decimal places.

Example 6.2. Consider the program

(P) s.t. 
$$f'(x) = x_1 - x_2 + (x_3 - 1)^2 + (x_4 - 2)^2 + (x_5 - 2)^2 + \min$$
  
 $f^1(x) = e^{x_1} + x_2^2 - 1 \le 0$   
 $f^2(x) = x_1^2 + x_2^2 + e^{-x_3} - 1 \le 0$   
 $f^3(x) = x_1 + x_2^2 - 1 \le 0$   
 $f^4(x) = e^{-x_2} - 1 \le 0$   
 $f^5(x) = (x_1 - 1)^2 + x_2^2 - 1 \le 0$   
 $f^6(x) = x_1 + e^{-x_4} - 1 \le 0$   
 $f^7(x) = x_2 + e^{-x_5} - 1 \le 0$ 

Starting with the initial (not feasible) point

 $x^{\circ} = (2, 2, 2, 2, 2)^{t}$ ,

the algorithm finds the feasible point

$$\overline{\mathbf{x}} = (0, 0, 1.95, .8095, .5871)^{\mathsf{T}},$$

and

$$P^{=} = \{1, 4, 5\}; \quad D_{P^{=}}^{=} = R(P) = R\left(\begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}\right)$$

Using Zoutendijk's method, applied to the reduced program

$$f^{\circ}(\overline{x} + Py) = (.95 + y_{1})^{2} + (y_{2} - 1.1905)^{2} + (y_{3} - 1.4129)^{2} + min$$
s.t.  $f^{2}(\overline{x} + Py) = e^{-1.95 + y_{1}}$   $-1 \le 0$ 
 $f^{3}(\overline{x} + Py) = (.8095 + y_{2})^{2} + (.5871 + y_{3})^{2} - 1 \le 0$ 
 $f^{6}(\overline{x} + Py)$   $e^{-(.8095 + y_{2})}$   $-1 \le 0$ 
 $f^{7}(\overline{x} + Py)$   $e^{-(.5871 + y_{3})} - 1 \le 0$ 

we get that

$$x^* = (0,0,1,.707,.707)^t$$
; and  $f(x^*) = 3.343146$ 

is the solution of program (P). The computer program found the solution correct to 7 decimal places. (The above problem was also solved in [14] using the feasible direction methods MELP1 and MELP2.)

Example 6.3. Our last example is the program

$$f^{\circ}(x) = 10\sqrt{2}x_{1} + (12 - 2\sqrt{2})x_{2} + 10x_{3} - 2x_{4} - (12 + 2\sqrt{2})x_{5}$$
$$+ 2x_{1}^{2} + 2x_{2}^{2} + 3x_{3}^{2} + 3x_{4}^{2} + 2x_{5}^{2} + 2\sqrt{2}x_{1}x_{3} - 2\sqrt{2}x_{1}x_{4}$$
$$- 4x_{2}x_{5} + 2x_{3}x_{4} + \min$$

(P) s.t.  

$$f^{1}(x) = 2e^{\frac{1}{2}(\sqrt{2}x_{1}-x_{3}+x_{4})} + x_{2}^{2} + x_{5}^{2} + 2x_{2}x_{5} - 2 \le 0$$

$$f^{2}(x) = 2x_{1}^{2} + 2x_{2}^{2} + x_{3}^{2} + x_{4}^{2} + 2x_{5}^{2} - 2\sqrt{2}x_{1}x_{3} + 2\sqrt{2}x_{1}x_{4}$$

$$+ 4x_{2}x_{5} - 2x_{3}x_{4} + 4e^{\frac{1}{2}(x_{2}-x_{3}-x_{4}-x_{5})} - 4 \le 0$$

$$f^{3}(x) = 2\sqrt{2}x_{1} - 2x_{3} + 2x_{4} + 2x_{1}^{2} + x_{2}^{2} + 2x_{3}^{2} + 2x_{4}^{2} + x_{5}^{2}$$
$$+ 2\sqrt{2}x_{1}x_{3} - 2\sqrt{2}x_{1}x_{4} + 2x_{2}x_{3} + 2x_{2}x_{4} - 2x_{2}x_{5}$$
$$- 2x_{3}x_{5} - 2x_{4}x_{5} - 4 \le 0$$

$$f^{4}(x) = x^{-\frac{1}{2}(\sqrt{2}x_{2}+\sqrt{2}x_{5})} - 1 \le 0$$

$$f^{5}(x) = -4\sqrt{2}x_{1} + 4x_{3} - 4x_{4} + 2x_{1}^{2} + 2x_{2}^{2} + x_{3}^{2} + x_{4}^{2} + 2x_{5}^{2}$$

$$- 2\sqrt{2}x_{1}x_{3} + 2\sqrt{2}x_{1}x_{4} + 4x_{2}x_{5} - 2x_{3}x_{4} \le 0$$

$$f^{6}(x) = \sqrt{2}x_{1} - x_{3} + x_{4} + 2e^{\frac{1}{2}(\sqrt{2}x_{1}+x_{3}-x_{4})} - 2 \le 0$$

$$f^{7}(x) = x_{2} + x_{5} + \sqrt{2}e^{\frac{1}{2}(x_{2}+x_{3}+x_{4}-x_{5})} - \sqrt{2} \le 0.$$

Using the initial (not feasible) starting point

 $\mathbf{x}^{\circ} = (0,0,2,1,1)^{t},$ 

the computer program in the appendix finds the feasible point

$$\overline{x}$$
 = (-.0547, -1.4604, 1.1583, 1.2357, 1.4605)<sup>t</sup>,

and

,

$$P^{=} = \{1,4,5\}; D_{p^{=}}^{=} = R(P) = R\left( \begin{bmatrix} .707 - .707 & 0 \\ 0 & 0 & -1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \right).$$

It then finds the solution

$$x^{*} = (-.5, -.855, -.206, .501, .845)^{t}; f^{\circ}(x^{*}) = -22.627.$$

The output for this program is given in the appendix.

<u>Remark 6.1.</u> Example 6.3 was constructed from Example 6.2 by substituting the unitary matrix

$$T = \frac{1}{2} \begin{bmatrix} \sqrt{2} & 0 & -1 & 1 & 0 \\ 0 & \sqrt{2} & 0 & 0 & \sqrt{2} \\ 0 & -1 & 1 & 1 & 1 \\ -\sqrt{2} & 0 & -1 & 1 & 0 \\ 0 & -1 & -1 & -1 & 1 \end{bmatrix}$$

i.e. we replaced the functions  $f^{k}(x)$  in Example 6.2 by the functions  $f^{k}(Tx)$ , for all  $k \in P$ . The objective function  $f^{\circ}(x)$  was replaced by  $4f^{\circ}(Tx) - 36$ . (Multiplying by 4 and subtracting 36 were done in order to eliminate fractions and constants.) Example 6.3 is there-fore the equivalent problem obtained, from Example 6.2, after a rotation of the axis. Since an exact solution of Example 6.2 is (see[14])

$$x^* = (0,0,1,\sqrt{2}/2,\sqrt{2}/2); f^{\circ}(x^*) = 9 - 4\sqrt{2},$$

we see that an exact solution of Example 6.3 is

$$T^{t}x^{*} = \frac{1}{2}(-1, -1 - \sqrt{2}/2, 1 - \sqrt{2}, 1, 1 + \sqrt{2}/2)^{t}.$$

The exact value of the objective function is

$$4(9 - 4\sqrt{2}) - 36$$

Furthermore,  $p = p^{-1}$  for Example 6.3 is  $T^{t}P$ , where P is the matrix  $p^{-1}$  in Example 6.2. Example 6.3 was constructed to illustrate how the

computer program works when the cones of directions of constancy of the constraints are difficult to calculate by inspection.

<u>Remark 6.2.</u> The above three examples illustrate that the computer program in the appendix can be used to solve convex programs, with faithfully convex constraints, independent of Slater's condition. We have not as yet compared this algorithm with any other existing techniques. Nor have we studied the stability of the algorithm, though we have had to account for round-off error in several instances. The author hopes to study these questions in the near future. METHED OF PEDUCTION STLVES CONVEX PROGRAMS (P). USES ZOUTENDIJK'S METHOD WITH SIMPLEX METHOD AND FIBONNACCI SEARCH. TAKES JAMMING INTE ACCOUNT. THE USER NEEDS TO ADD: ANALYTIC REPRESENTATIONS OF ALL THE FUNCTIONS AND GRADIENTS THIS IS ADDED BETWEEN THE LINES \*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\* THE USER ADDS FOR DATA: THE DIMENSION OF THE X VARIABLE - N THE NUMBER OF CONSTRAINTS - PONSTR INITIAL VECTOR ESTIMATE - XBAR IMPLICIT REAL#8(A-H,C-Z) INTEGER COLS,PCMSTP COMMON P(10,10), E(10,10), GRD(10), GRADS(10,10), PO(10), XEAR(10) COMMON F0, SMETH, FNK(10), REDUC, EPS1, EPS2, F0SA VE COMMON EPS3.EPS4.EPS5.EPS6.EPS7.EPS3.EPS9.EPS10 COMMON N.KP.PCNSTR.JI(10).COLS.ITERS APPROXIMATIONS OF 0 DUE TO ROUNDOFF: EPS1 - CONSTRAINT SATISFIED EPS2 - BINDING (ANTIZIG.) GRADIENT IN SUBR. CONE GRADIENT IN SUBR. PEQUAL EPS3 -EPS4 ----EPS5 -CONSISTENCY OF SYSTEM IN PEQUAL EPS6 - BASIS ELEMENTS IN SIMPLEX EPS7 - NOPM OF GRADIENT OF OBJ. F EPS8 - SCLUTION FOUND EPS9 - ENSURE PIVOT>O IN SIMPLEX DBJ. FUNCTION EPS1=EPS2=EPS3=EPS5=EPS7=EPS9=1.D-5 EPS4=1.0-4 EP\$8=1.0-6 EPS6=1.D-8 IT = 0READ, N, PCNSTR KP=N READ, (XBAR(I), I=1,N) PEINT 1003 1003 FORMAT('1', 'OUTPUT FCR EXAMPLE 6.3') **PEINT 1004** FCPMAT( +++ + INITIAL ESTIMATE +) 1004 PFINT 1005, (XBAR(1), I=1,N) 1005 FORMAT(' ' ,10F10.4) DO 1010 K=1. PONSTR 1010 PO(K) = 007 1020 I=1.N 07 1020 J=1.N E(I,J)=0P(I,J)=01020 IF(I.E0.J)E(I.J)=P(I.J)=1 DD 1030 K=1.PC4STR CALL GRDF(K, XBAP) D-1030 I=1.4 GRADS(I,K)=GED(I) 1030

APPENDIX

```
SMETH=1
ITIPS=0
1FE45=0
        CALL FNEVAL (XBAR. IFEAS)
        FOSAVE =F0+1
CALL SMTHD1
CALL SMTHD1
CALL SMTHD
IF (REDUC.EQ.1)G0 TO 1045
IF (SMETH.EQ.0)G0 TO 1050
 1040
        IT = IT + 1
        SMETHED
DT 1047 K=1 +PCNSTR
 IF (PO(K).CO.O.AND.FNK(K).GE.EPS1 )SMETH=1
1047 IF (PO(K).TO.).AND.FNK(K).LT.EPS1 )PO(K)=1
IF(IT.LT.PCNSTR)GD TC 1045
        PRINT, FEASIBLE SET IS EMPTY!
        STOP
 1045 CENTINUE
        IF (SMETH.FR.1)G0 TO 1051
PF INT 1044
 1044 FORMAT ( !- ! , 'FEASIBLE POINT FOUND ! )
 PFINT 1046. (XBAR(I).I=1.N)
1046 FURMAT(* *
                                                          ,10F10.4)
        CALL PEQUAL
 PFINT 1048, N,KP
1048 FORMAT('-',II,'X',II,'MATRIX FOR THE CONE OF CONSTANCY')
        DO 1049 I=1.N
 1049 PRINT 1054.(P(I,J),J=1.KP)
1054 F(RMAT(' ',10F12.4)
 GO TO 1040
1051 CALL PEQUAL
        GO TO 1040
PEINT 1053
 1050 PRINT
 1050 PRIMAT('-','SOLUTION FOUND')

PRIMT 1052 (X34F(I),I=1,M)

1052 FIRMAT('','X IS ',10F10.4)

PRIMT,'VALUE OF CBJECTIVE FUNCTION IS
                                                               ',F0
        PPINT 1055
 1055 FORMAT( ! 1 !)
        STOP
        END
        DOUBLE PRECISION FUNCTION FNO(X)
00000
        THE STATEMENT FUNCTION
EVALUATES THE OBJECTIVE FUNCTION
THE USER MUST ADD THE LINE CONTAINING FNO
        IMPLICIT REAL#8(A-H.C-Z)
        DIMENSION X(10)
        PEAL #4 SOFT
С
        ** ****
c
        FN0=10* SORT(2.)*X(1) +(12.-2.* SORT(2.))*X(2)+10.*X(3)-2*X(4)-(
       1+2*5QRT(2.))*X(5)-4*X(2)*X(5)+2*X(3)*X(4)+2*X(2)**2+3*X(3)**2+
       1 3#X(4)**2+2#X(5)*#2 +2#X(1)**2+2*SQRT(2.)*X(1)*X(3)-2*
       1SORT(2.)*X(1)*X(4)
С
        等于 泽东 法审查法 化二氯化二溴苯苯 苯汞腙 人名卡尔 法发送者 医胆液检测 医法拉尔 医法法法 医法法法
```

C 杂义 2 法承担 我家本让学校大别 大笑 林浩岩 长路尔语斯 或字文学 连 法本本学学校 李大 辛辛 辛沙语香港 RETURN END SU BROUTINE GEDEVO(X) 0000 SUBROUTINE EVALUATES THE GRADIENT OF FO THE USER MUST ADD THE VECTOR OF PARTIAL DERIVATIVES IMPLICIT REALIR(A-H.C-Z) INTEGEP COLS.PCNSTP COLS.PONSTP CLMMON P(10,10),F(10,10),GRD(10),GRADS(10,10),FO(10),XBAP(10) CDMMON F0,SMFTH,FNK(10),REDUC,EPS1,EPS2,FOSAVF CDMMON F0,SMFTH,FNK(10),REDUC,EPS1,EPS2,FOSAVF CUMMIN N, KP, PCHSTR, JI(10), COLS, ITERS DIMENSION X(10) FEAL 4 SOST ç \* 冰学说 医眼球球子的变形 化化合合物 化水溶液 人名巴克 医克克克德德德德德德氏液体 化活动计算器 使气力变形 GPり(1)=10本SORT(2.)+4本X(1)+2本SORT(2.)\*X(3)-2本SORT(2.)\*X(4) GFD(2)=4+X(2)-4+X(5)+(12-2\*SOFT(2.)) GFD(3)=6+X(3)+2+X(4)+10 +2+SOFT(2.)\*X(1) GRD(4)=6\*X(4)+2\*X(3)-2-2\*SORT(2.)\*X(1) GRD(5)=4+X(3)-4+X(2)-2\*SORT(2.) -12 С C RETURM END SUBROUTINE GROF(K,X) 0000 THE SUBFOUTINE FINDS THE GRADIENT OF THE CONSTRAINTS. THE USER MUST ADD THE VECTORS OF PARTIAL DERIVATIVES IMPLICIT FEALSB(A-H,C-Z) INTEGER COLS.PONSTE COMMON P(10,10),E(10,10),GRD(10),GRADS(10,10),P0(10),XBAR(10) COMMON F0,SMETH,FNK(10),REDUC.EPS1,EPS2,F0S4VE CAMMON FPS3,FPS4,EPS5,EPS6,EPS7,EPS8,EPS9,EPS10 CEMMON N.KF,PONSTR,JI(10),COLS,ITERS DIMENSION X(10) PEAL#4 SOFT EXP(T) = DEXP(T)G. TO(1,2.3.4.5.6.7.8.9.10.11.12.13.14.15.16.17.18.19.20.21.22. 1 23.24.25.25.27.28.29.30).K ĉ 化学校 法法学院 建冷水管 化化化合物 化化合合体 化化合合体 化化合合合合体 化化化合合体 化化化合合体 黑头月 建口 医海道 的复数 法用户结 化环化锌 法通行 铁铁 把字 表示 化分子 网络大学 化分子 化分子 化分子 化分子 化分子 化分子 化分子 化分子分子 化分子分子 化分子分子 1 CONTINUE GRD(1)= - 范XP(•5年(SORT(2•)\*X(1)-X(3)+X(4)))\*SCRT(2•) \*X(5) GFD(2)=2\*X(2)+2 GPD(3)=EXP(.5\*(SORT(2.)\*X(1)-X(3)+X(4))) \*(-1) GRD(4)=EXP(.5\*(SORT(2.)\*X(1)-X(3)+X(4))) GRD(5)=2 \*X(2) +2\*X(5) RE TU-2N CONTINUE 2 GFD(1)=44X(1)-2\*SORT(2.)\*X(3)-2\*SORT(2.)\*X(4) GFD(2)=4\*X(2)+4\*X(5) +2\*EXF(.5\*(X(2)-X(3)-X(4) 1 - X(5)))

	GPD(3)=2*X(3)-2*SQRT(2.)*X(1)-2*X(4) -2*EXF(.5*(X(2)-X(3)-X(4)
	. → X (5))) GPD(4)=2#X(4)+2#SQRT(2•)#X(1)→2#X(3) → 2#EXF(•5#(X(2)→X(3)→X(4
	- ~X(5))) - CED(5)-4キャ(5)+4キャ(2)
	L +X(5)))
7	RETURN
2	GED(1)=2*SQRT(2·)+4*X(1)+2*SQRT(2·)*X(3)-2*SQRT(2·)*X(4)
	GRD(2)=2*X(2)+2*X(3)+2*X(4)-2*X(5)
	GPD(3)==2+4米X(3)+2※5QPT(2。)*X(1)+2*X(2)=2*X(5) GPD(4)=2+4和X(4)=2※5QPT(2。)*X(1)+2*X(2)=2*X(5)
	GPD(5)=2*X(5)-2*X(2)-2*X(3)-2*X(4)
<i>I</i> .	RETURI
4	GPD(1)=0
	GRD(2) = - EXP(5*(SQRT(2.)*X(2)+SQRT(2.)*X(5)))/SQRT(2.)
	GRD(4)=0
	GRD(5)=-EXP(5*(SQRT(2.)*X(2)+SQRT(2.)*X(5)))/SQRT(2.)
5	
,	GFD(1)=4*X(1)-2*SQFT(2.)*X(3)+2*SQRT(2.)*X(4)-4*SQRT(2.)
	GRD(2)=4*X(2)+4*X(5) CFD(3)=2*X(3)=2*SODT(2))*X(1)+4=2*X(4)
	GFD(4)=2+X(4)+2+SGFJ(2+)+X(1)+4=2+X(4) GFD(4)=2+X(4)+2+SGRT(2+)+X(1)+2+X(3)+4
	GFD(5) = 4 * X(5) + 4 * X(2)
6	CONTINUE
	GED(1)=SQRT(2.)+SQRT(2.)*EXP(.5*(SQRT(2.)*X(1)+X(3)-X(4)))
	GRD(2)=0 GRD(3)=+1 + EXP(-5*(SOPT(2.)*X(1)+X(3)-X(4)))
	GFD(4) = 1 - EXP(.5*(SORT(2.)*X(1)*X(3)-X(4)))
	GFD(5)=0
7	CENTINUE
	GRD(1)=0
	GRD(2) = EXP(.5*(X(2)+X(3)+X(4)-X(5)))/SQRT(2.) +1 $GRD(3) = EXP(.5*(X(2)+X(3)+X(4)-X(5)))/SQRT(2.)$
	GFD(4) = EXP(.5*(X(2)+X(3)+X(4)-X(5)))/SQRT(2.)
	GRD(5)=1-EXP(.5*(X(2)+X(3)+X(4)-X(5)))/SQRT(2.) RETURN
9	CONTINUE
0	
,	RETURN
10	
11	CONTINUE
• •	PETURN
12	FETURN
13	CONTINUE
14	FETURN
• •	FETURN
15	
15	CONTINUE
	RETURN

 $\bigcirc$ 

1	[7	CONTINUE
1	13	PETUPN CONTINUE
	10	
		FF TURN
2	20	CONTINUE Return
	21	CONTINUE
2	22	CONTINUE
	23	CONTINUE
	24	EFITURN CENTINUE
	- · 	RETURN
6	20	RETURN STATES
	26	CONTINUE PETURN
3	27	
1	28	CUNTINUE
2	29	CONTINUE
	30	RETURN CONTINUE
c		
č		· · · · · · · · · · · · · · · · · · ·
		END
		SURROUTINE ENEVAL(X.IFEAS)
ŝ		THIS SUBSCHITTE EVALUATES THE EUNCTIONS AT THE BOINT Y
č		THE USEP MUST SUPPLY THE ANALYTIC REPRESENTATIONS OF THE FUNCTION
C		IMPLICIT PEAL#8(A-H,C-Z)
		INTEGER CTUS.PONSTR CEMMIN P(10.10).F(10.10).GRD(10).GRADS(10.10).F0(10).XBAP(10)
		CCMMON F0, SMETH, FNK (10), REDUC, EPS1, EPS2, F0SAVE
		CUMMER PASSIERS (1930, 2930, 2937, 2930, 2939, 2937, 2
		DIMENSION X(1)) REAL¥4 SORT
		EXP(T)=DEXP(T) IE(SMETH+E0,1+CP+LEEAS+ME+0)60 TO 3001
		$F_0 \neq F_{NO}(X)$
3(	001	IF = 1
		F0=0 00 80 K=1.PCHSTR
		IF(P0(K)+UT+D)GC TA 80 IF(IFCAS+F0+D+AFD+P0(K)+N5+D)GC TA 80
		IF (IFEAS .EQ.1.AND.PO(K).NE.1.AND.PO(K).NE.2) (C TO BO
		- いい - いい - いいは つどり つびりつじりつびりつびりつ ダリチレッチエリチどりチ どうチャワリ チレリチ イリチン・チャ
_	1	1,50),K
00	1	1 • 5 3 ) • K 本书 水水水 水水 水 水水水水 水水 水水 水水 水水 水水 水水 水水 水水

,

•

.

.

31	CONTINUE E=2482XP(+8 1 X(2)*#2+>	5% (SOR1 X (5) **2	(2.)9X(1 2-2	1)-x(3)+x(4)	))+2*	X(2)4X(5)+
32	FNK(K)=F G0 T0 70 CONTINUE F=44EXP(*9	54 (X(2)	-x (3)-x(	(4)-X(5)))-2	#SQRT(2.)#X(1	)#X(3)+2#506T(2.
	1 *X(1)*X(4 1 X(4)**2+7 F <sup>K</sup> K(K)=F	4)-2*X(5)>	(3)*X(4)+ **2-4	+4*X(2)*X(5)	+2*X(1)**2+2*	X(2) **2+X(3) **2+
33	CENTINUE F=2*SORT(2 1 X(1)*X(4)	2•)*X() )+2+X(2	L)-2*X(3) 2)%X(3)+2	)+2*X(4)+2*S 2*X(2)*X(4)-	]RT(2.)*X(1)* 2*X(2)+X(5)-2	X(3)-2*SQRT(2.)* *X(3)*X(5)-2*X(4
34	1 *X(5)+2*) FNK(K)=F GC TC 70 CUNTINUE	X(1)**/	2+X(2)**2	2+2*X(3)**2+	2*X(4)**2+X(5	)**2-4
75	F= EXP( FNK(K)=F G( T() 70 CONTINUE	(	54 (SQRT (3	2•) #X ( 2 ) + SOR	T(2•)*X(5)))-	1
و. ر	$F = -4\pi \text{ SORT}(1 \times (1) \times (4))$ $1 \times (1) \times (4)$ $1 \times (2 + X)$	(2•)*X )−2*X(] **2+2*)	(1)+4*X(3 3)*X(4)+4 ((5)**2	3)-4*X(4)-2* 4*X(2)*X(5)+	SQRT(2.)*>(1) 2*X(1)**2+2*X	*X(3)+2*SORT(2) (2)**2+X(3)
36	FMK(K)=F GO TO 70 CONTINUE F=2*SXP(•*	5* (SQF)	「(2•)*X()	1)+	x(3)-x(4)))+s	QRT(2.)*X(1)
37	I - X(3)+ FNK(K)=F GO TO 70 CONTINUE	X(4)•	* 2			
38	F=SQRT(2.) FNK(K)=F GD TO 70 CENTINUE	)*2XP()	5*(X(2)1	+ X ( 3) + X ( 4 ) - X	(5)})+X(2)+X(	5)-SQRT(2.)
39	FNK(K)=F G: TC 70 CCNTINUE FNK(K)=F		,	1		
40	GO TO 70 CONTINUE FNK(K)=F					
41	CONTINUE					
42	CONTINUE FNK(K)=F GD TD 70					
43	CONTINUE FNK(K)=F GO TO 70					
44	CONTINUE ENK(K)=E					
45	CONTINUE FNK(K)=F					
4.5	CONTINUE					

 $\bigcirc$ 

47	FINK(K)=F GO TO 70 7 CHATIMUE FNK(K)=F
4	GUNTE 70 B. CENTIBUE FNK(K)=F
49	G( TC 70 * 9 CONTINUE ENK(K)=E
50	GF TO 70 D CUNTINUE FNK(K)=F GU TR 70
c C 7(	477 米ホオ オネチ キャオ オキャガオ 本キ キカニ オキ オオオ キキ キャキ オキ
7:	IF(F+LT+EPS1 +AND+F+GT+-EPS2 )PO(K)=1 IF(F+LT+-EPS2 )PO(K)=2 5 IF(IFEAS+EQ+1)GD TD 80
80	FU=FD+F D CONTINUE IFEAS=IF RETUEN END
~	SUBROUTINE CONE
č	WE FIND THE RANGE SPACE OF P INTERSECT THE CONE OF CONSTANCY. The result is put back into p
c	IMPLICIT REAL#8(A-H,C-Z) INTEGER COLS,PONSTP CL MMON P(10,10),E(10,10),GRD(10),GRADS(10,10),F0(10),XBAR(10) COMMON F0,SMETH,FNK(10),REDUC,EPS1,EPS2,F0SAVE CLMMON EPS3,EPS4,EPS5,EPS6,EPS7,EPS3,EPS9,EPS10 COMMON N,KP,PONSTR,J1(10),COLS,ITERS REAL#8 B(10,10),Y(10),A(10,10),X(10),TEMP(10)
	IF JI(K)=0 WE SKIP THE K-TH FUNCTION
1 0 ( 1 1 ( 1 1 ±	DO 250 NCNSTR=1,PCNSTR IF(JI(NCNSTR).EQ.0)GC TÙ 240 D J=0 DO 110 I=1,N O X(I)=0 GG TO 140 5 DY 120 I=1,N
120 143	X(I)=0 DD 120 K=1.KP X(I)=X(I)+P(I.K)ME(K.J) CALL GRDF(NCNSTP.X) DC 150 I=1.KP X(I)=0
151	D: 150 K=1.N Y(I)=Y(I)+GRD(K)>P(K.I) D0 180 I=1.KP

IF (DAES(Y(I)).GT.EPS3 )GC TO 190 IF (J.EO.KP)GC TO 240 180 J = J + 1GC TC 115 CONTINUE IF (KP.E0.1) GC TO 220 190 KPM1=KP-1 DC 170 K1=1,KP DC 170 K2=1,KPM1 K2PL=K2+1  $4(K1 \cdot KS) = 0$ IF (K1.EQ.K2.AND.K1.LT.I)A(K1.K2)=-1 IF (K1-1.E0.K2.AND.K1.GT.I)A(K1,K2)=-1 170 IF (K1.EQ.I.AND.K2.GE.I)A(K1.K2)=Y(K2PL)/Y(I) DC 200 I1=1.M DO 200 I2=1.KPM1 D(11.12)=0 DD 200 K=1.KP 200 B(I1,I2) = P(I1,K) \* A(K,I2) + B(I1,I2)210 11=1 .N Ð DO 210 I2=1,KPM1 P(I1.I2)=8(I1.I2) DO 214 K=1.FCNSTF 210 IF(PO(K).EQ.-1)G0 TO 214 D. 212 I=1,KPV1 TF. VP(I)=0 DC 212 J=1.KP TEMP(I) = TEMP(I) + GRADS(J,K) \* A(J,I)212 D' 213 I=1,KPM1 GRADS(I.K)=TEMP(I) 213 CONTINUE 214 KD=KDM1 GE TO 100 CONTINUE. 220 DU 230 I=1.N 230 P(I,1)=0 240 CONTINUE CONTINUE 250 RETURN END SUPRCUTINE SMPLEX(A.C., ARTIF1, ROWS, BASIS, COLS) 000 THE SUBROUTINE SOLVES PROBLEMS AX=B BY THE SIMPLEX METHOD IMPLICIT REAL\*8(A-H,0-Z) INTEGER TOLS INTEGER PCNSTR.COLS.APTIF1.COLSM1.ROWS.PIV(T(2).BASIS(11) CCMM3N P(10.10).E(10.10).GRD(10).GRADS(10.10).F0(10).XBAR(10) COMMON FO, SMETH, FNK(10), REDUC, EPS1, EPS2, FOSAVE COMMON FPS3, EPS4, EPS5, EPS6, EPS7, EPS8, EPS9, EPS10 COMMON N, KP, PCNSTR, JI(10), TOLS, ITERS FE4L\*8 A(11,22), C(21) ABS(T)=DAES(T) COLSM1=COLS-1 NOITER=0 DC 406 I=1.ROWS K=ARTIF1+I BASIS(I)=K

D 405 J=1, REWS 405 A(J,K)=0406 A(I,K)=1ZJMIN=0407 00 420 J=1.COLSM1 IF(C(J).E0.-1000.)G0 T0 420 2J=0DC 410 I=1.ROWS ZJ = ZJ + A(I,J) + C(BASIS(I))410 IF (ZJ -C(J).GE.ZJMIN) GG TO 420 ZJAIN=ZJ -C(J) PIVOT(2) = J420 OF NT INUE IF (ZJMIN .EQ.0) RETURN THETA0=-1 DF 430 I=1.RCWS IF(A(I.PIVOT(2)).LE.0)G0 T0 430 T=A(I,CDLS)/A(I,PIVOT(2))IF (T.LT.O.DR.(T.GT.THETAO.AND.THETAO.GE.O))GC TO 430 PIVT(1) = ITHETAD=T CONT INUE 430 IF (THETA0.GE.J)GC TO 432 PE INT. 'ERROR1' ST IP 432 @ASIS(PIVOT(1))=PIVOT(2) D: 435 I=1, RIWS IF(I.EQ.PIVOT(1))GO TO 435 A(I,COLS)=A(I,COLS)-THETAU\*A(I,PIVOT(2)) IF(A(I,CELS).GT.0.AND.A(I.COLS).LT.EPS9 )A(I.CELS)=0 435 CONTINUE A(PIVT(1), CTLS)=THETAO IF (A (PIVCT(1),COLS).GT.0.AND.A(PIVOT(1),COLS).LT.EPS9 ) 1 = A(PIVCT(1),COLS)=0DP 455 J=1.CCLSM1 IF(C(J).EQ.-1000)GD TO 455 IF(J.EQ.PIVOT(2))GO TO 455 DO 450 I=1. POWS IF(I.EQ.PIVOT(1))G0 T0 450 A(I,J) = A(I,J) - A(FIVOT(1),J) + A(I,PIVOT(2)) / A(FIVOT(1),PIVOT(2))IF (A(I,J).GT.0.AND.A(I,J).LT.EPS9 )A(I,J)=0 450 CONTINUE A(PIVOT(1), J) = A(PIVOT(1), J)/A(PIVOT(1), PIVOT(2))455 CONTINUE . D7 460 I=1.ROWS A(I,PIVOT(2))=0 460 A(PIVOT(1), PIVOT(2))=1 NOITER=NCITER +1 IF (NOITEP.LE.50)GD TC 407 PRINT, 'EPPOR2' STOP END SUBROUTINE PEQUAL C C SUBROUTINE FINDS THE EQUALITY SET C IMPLICIT REAL\*8(A-H.O-Z) INTEGER PCNSTR, ROWS, COLS, ARTIF1, BASIS(11), ILAMDA(10)

بحديرهم منصد عنا

and decrease on the second

COMMON P(10.10).E(10.10).GFD(10).GRADS(10.10).FO(10).XBAP(10) COMMON F0.SMETH.FNK(10).PEDUC.EPS1.EPS2.FOSAVE COMMON FPS3.UPS4.EPSE.EPS6.EPS7.EPS8.EPS9.EPS10 CEMMON N.KP.PCNSTR.JI(10).CDLS.ITERS RFAL48 A(11.22).C(21) ADS(1)-DAS(1). ABS(T)=DABS(T) 500 DO 510 K=1.PCNSTR JI(K)=0510 D. 350 K=1, PONSTR IF (PO(K) .NE .1) GD TO 550 512 I=1,KP DE IF (DARS(GRADS(I+K)) .GT.EPS4 )GD TO 550 512 JI(K)=1PO(K)=-1 IF (K.EQ. PONSTR) GD TO 630  $K_{1} = K + 1$ DC 520 I=K1,PCNSTR IF (PO(K) .NE.1)GP TO 520 00 00 515 J=1.KP IF (DABS(GPADS(J.K)).GT.EPS4 )G0 T0 520 515 JI(K)=1PO(K)=-1 CONTINUE GC TO.630 520 CONTINUE 550 ARTIF1=0 ROWS=KP+1 DD 570 K=1, PCNSTR IF(PO(K).NE.1)GO TO 570 AFTIF1=ARTIF1+1 ILAMDA(AFTIF1)=K 00 560 I=1,KP A(1,ARTIF1)=GRADS(I,K) 560 A(ROWS, APTIF1)=1 570 CONT INUE IF (ARTIF1.LT.1)PETURN CCLS=ARTIF1+ROWS+1 D0 580 I=1.KP A(I.COLS)=0 A(ROWS,COLS)=1 580 00 500 I=1.APTIF1 C(I)=0 IP1=ARTIF1+1 590 IMINUS=CCLS-1 DO 500 I=IPL, IMINUS C(I) = -1000600 CALL SMPLEX (A, C, AFTIF1, ROWS, BASIS, COLS) SUM=0 DC 510 I=1, ROWS SUM=SUM+A(I,CCLS)\*C(EASIS(I)) 610 IF (SUM .LT .- EPS5) FETURN 520 I=1, ROWS 0: IF (ABS(A(I.COLS)).I T.EPS6 )GO TO 620 PO(ILAMDA(BASIS(I)))=-1 JI(ILAMDA(BASIS(I)))=1 CONTINUE 620 630 CALL CONF END

000000 THIS SUBROUTINE PERFORMS ITERATIONS ASSUMING SLATER'S CONDITION THE CONSTRAINTS ARE THOSE FOR WHICH PO(K).GT.0 THE OBJECTIVE FUNCTION IS THE SUM OF FO AND ANY CONSTRAINTS ZI UTENDIJK'S METHOD IS USED IMPLICIT REAL\*8(A-H,0-Z) INTEGER ROWS, PCNSTR, COLS, ARTIF1, BASIS(11) COMMON P(10,10), E(10,10), GRD(10), GRADS(10,10), P0(10), XBAR(10) COMMON P(5, SMETH, FNK(10), PEDUC, EPS1, EPS2, FDSAVE COMMON EPS3, EPS4, EPS5, EPS6, EPS7, EPS8, EPS9, EPS10 CUMMON N, KP, PCNSTR, JI(10), COLS, ITERS DIMENSION GROFO(10) DIMENSION D(10), DKP(10) DIMENSION XY(10),XZ(10) DIMENSION TEMP1(10) REAL#8 A(11,22),C(21) ABS(T)=DABS(T) ARTIF1=2\*KP+1 ITERS=0 FOWS=0 700 DSUM=0 ITERS=ITERS+1 REDUC=0 DO 730 K=1. FONSTP IF (PO(K) .NE .1)GP TP 730 ROWS=ROWS+1 DO 720 J=1.KP J2=2\*J J2M=2\*J-1 A(ROWS+J2M)=GRADS(J+K) 720 A(RBWS,J2) = -A(RCWS,J2M)730 CUNTINUE ROWS=ROWS+1 DU 7301 J=1,ARTIF1 7301 A(RCW5,J)=0 IF(SMETH.E0.1)GD TO 732 FUSAVE=FNO(XBAR) CALL GRDEVO(XEAR) DD 7302 I=1,KP TEMP1(I)=0 D0 7302 J=1,N TEMP1(I)=TEMP1(I)+GPD(J)\*P(J,I) 7302 00 731 I=1.KP 12M=2\*I-1 A(FOWS, 12M) =TEMP1(1) GO TO 736 IFEAS=0 731 732 CALL FNEVAL (XBAR . IFEAS) FOSAVE=FO DC 735 K=1.PCNSTR IF(PO(K).NE.0)GC TO 735 DC: 733 I=1.KP 12M=2\*1-1 A(ROWS, I2M) = A(ROWS, I2M) + GRADS(I,K) 733 735 736 CONT INUE CONTINUE D) 737 I=1.KP 12=2\*1

SUBRIUTINE SATHD

```
DEUM=DSUM+ ABS(A(HOWS.IZM))
      A(R)WS, 12) = OKP(I) = -A(RCWS, 12M)
737
      IF (DSUM.LT. EPS7) FETUEN
      IF (RCWS.GT.1)GD TO 739
DT 739 I=1,KP
         739 I=1,KP
      DKP(I)=DKP(I)/DSUM
733
      GO TO
             845
739
      COLS=POWS+ARTIF1+2
      DO 740 I=1.ROWS
      A(I.COLS)=0
740
      A(I,AOTIF1)=1
      ROWS=ROWS+1
      I2=ARTIF1-1
      D3 750 J=1,12
      C(J) = 0
750
      A(ROWS,J)=1
      A(POWS,ARTIF1)=0
A(ReWS,COLS)=1
DI 760 J=4RTIF1,COLS
760
      C(J) = 0
      C(ARTIF1)=1
      CALL SMPLEX (A.C.ARTIF1.ROWS.BASIS.COLS)
D0 790 I=1.RCWS
790
      IF (BASIS(I).E0.AFTIF1.AND.ABS(A(I.COLS)).GT.EPS6)GO TO 810
      RETURN
810
      CENTINUE
      DB 320 I=1.KP
      DKD(I)=0
820
      DC 840 1=1.80WS
      IF(545IS(I).GT.2*KP)G0 TC 840
IDIV=845IS(I)/2
      IF (BASIS(I).EQ.IDIV#2)GD TO 830
      IDIV=IDIV+1
      DKP(IDIV)=DKP(IDIV)+A(I,COLS)
      Gr Th 840
      DKP(IDIV)=DKP(IDIV)-A(I.COLS)
830
840
      CONT INUE
      DC 850 I=1.N
845
      D(I) = 0
      D<sup>6</sup> 850 J=1,KP
C(I)=P(I,J)*OKP(J)+D(I)
850
      SUM=20
      AL 2H4=20
      DU 360 I=1,N
XEAR(I)=X3AR(I)+ALPHA*D(I)-
860
      DO 330 ITEP=1.30
      SIGN=-1
      IFEAS=1
      CALL FNEVAL (XBAR+IFEAS)
IF (IFEAS.E0.1.AND.ITER.EQ.1)GD TO 885
IF (IFEAS.F0.1)SIGN=1
      ALPHA=ALPHAZ2
      SUM=SUM+SIGN#ALPHA
      D0 380 I=1.N
XEAR(I)=XEAR(I)+ALPHA*D(I)*SIGN
880
      CONTINUE
885
      IFEAS=1
      CALL FNEVAL (XBAP. IFEAS)
IF (IFEAS.ED.1)GD TO 900
```

12M=12-1

890	ALPHA=ALPHA/2 SUM=SUM-ALPHA DD B90 I=1.N XEAR(I)=XBAR(I)-ALPHA*D(I) IFEAS=1 CALL FNEVAL(XEAR.IFEAS)
900 0 0 0 0	CONTINUE . Now find the minimum of Fo in the one direction
ou c	FIBONACCI SEARCH
960	TOL=10000 NN=1 MG=1 M1=1 N2=N0+N1 IF(N2.GE.TOL)GD TO 970 NN=NN+1 N0=01 N1=N2 GD TL 960
97)	Y=SUM+NO/N2 Z=SUM+Y DC 975 I=1.N
975	XY(I)=XEAF(I)-Y*D(I) XZ(I)=XEAF(I)-Z*D(I) IFEAS=0
	CALL FNEVAL (XY.IFEAS) FY==F0 IFEAS=0 CALL FNEVAL (XZ.IFEAS) FZ==F0 DC 990 I=1.NN DIST=SUM=40 IF(FY.LT.FZ)XN=Z IF(FY.GE.FZ)XN=Y IF(FY.LT.FZ)GT TF 980 SUM=Z Z=Y FZ=FY Y=A0+SUM-Y DD 077 I=1.N
977	XY(J)=XBAR(J)-Y+D(J) IFEAS=0 CALL FNEVAL(XY+IFEAS) FY=-F0
980	A0 = Y Y=Z FY = F Z Z= A0 + SUM - Z
982	D: 982 J=1.N XZ(J)=XBAP(J)-Z*D(J) IFEAS=0 CALL FNEVAL(XZ,IFEAS) FZ=#F0

۰.

90 CONTINUE CONTINUE 995 DC 996 I=1.N XEAR(I) = XBAR(I) - XM + D(I)995 IFEAS=1 CALL FNEVAL (XEAR . IFEAS) IFEAS=0 CALL FNEVAL (XBAP, IFEAS) IF (SMETH.FO.0)FU=FNO(XBAR) D2 998 K=1,FCNSTE IF(PO(K).E0.-1)GF TO 998 CALL GROF(K.XEAF) DO 997 I=1.KP GRADS(I+K)=0 D0 997 J=1+N GFADS(I+K)=GFADS(I+K)+GRD(J)\*P(J+I) 997 . CUNTINUE 993 IF( DABS(FO-FOSAVE).LT.EPS8)G0 T0 913 IF(ITERS.GT.5))RETURN IF(SMCTH.EQ.0)GO TO 918 ENTRY SMTHD1 ARTIF1=2\*KP+1 REDUC=0 918 SMETH=0 D0 950 K=1, PCNSTF IF(P3(K).LT.0)G0 T0 950 IF(FNK(K).LE. -EPS2)G0 TC 920 IF(FNK(K).GE.EPS1.AND.P0(K).E0.0)GJ T0 940 IF(P0(K).E0.0)REDUC=1 P0(K).E0. PO(K)=1GO TO 950 IF (P0(K) .EQ.0)PEDUC=1 920 PU(K)=2 GU TU 950 940 IF(PO(K).E0.0)SMETH=1 CONTINUE 95) IF (REDUC.E0.0.AND. DARS(FO-FOSAVE).GT.EPS8 )63 TO 700 RE TUPN END

.

\$DAT4

### OUTPUT FOR EXAMPLE 5.3

INITIAL ESTIMA 0.0000	0.0000	2.0000	1.0000	1.0000
FEASIBLE POINT -0.0547	F FOUND -1.4604	1.1583	1.2357	1.4605
5X3MATRIX FOR 0.7071 -0.0000 1.0000 0.0000 0.0000 0.0000	THE CONE -0.707 0.000 0.000 1.000 0.009	CF CONSTAN 1 0.0 0 -1.0 0 0.0 0 0.0 0 0.0 0 1.0	CY 000 000 000 000 000	• •

.

SOLUTION FOUND X IS -0.4999 -0.8549 -0.2059 0.5011 0.8549 VALUE OF OBJECTIVE FUNCTION IS -22.62742138784544

197

Card in the second

•

#### REFERENCES

- J. Abadie, "On the Kuhn-Tucker theorem", <u>Nonlinear Programming</u>,
   J. Abadie, ed., North-Holland Publishing Co., Amsterdam, 1967.
- [2] R.A. Abrams, "Projections of convex programs with unattained infima", SIAM Journal of Control, 13 (1975), 706-718.
- [3] R.A. Abrams and L. Kerzner, "A simplified test for optimality", Journal of Optimization Theory and Applications (forthcoming).
- [4] K. Arrow, L. Hurwicz and H. Uzawa, "Constraint qualifications in maximization problems", <u>Naval Research Logistics Quarterly</u>, <u>8</u> (1961), 175-191.
- [5] M. Atteia, "Fonctions spline definies sur un ensemble convexe", Numerische Mathematik, 12 (1968), 192-210.
- [6] C.C. Augunwamba, "Optimality condition: constraint regularization", Mathematical Programming, 13 (1977), 38-48.
- [7] M.S. Bazaraa, J.J. Goode and C.M. Shetty, "Constraint qualifications revisited", Management Science, 18 (1972), 567-573.
- [8] M.S. Bazaraa, C.M. Shetty, J.J. Goode and M.Z. Nashed, "Nonlinear programming without differentiability in Banach spaces: necessary and sufficient constraint qualification", <u>Applicable</u> Analysis, 5 (1976), 165-173.
- [9] A. Ben-Israel, A. Ben-Tal and A. Charnes, "Necessary and sufficient conditions for a Pareto optimum in convex programming", Econometrica, 45 (1977), 811-820.
- [10] A. Ben-Tal, A. Ben-Israel and S. Zlobec, "Characterization of optimality in convex programming without a constraint qualification", Journal of Optimization Theory and Applications, 20 (1976), 417-437.

- [11] A. Ben-Israel, A. Ben-Tal and S. Zlobec, "Optimality conditions in convex programming", The IX International Symposium on Mathematical Programming, Budapest, Hungary, August, 1976.
- [12] A. Ben-Tal and A. Ben-Israel, "Characterizations of optimality in convex programming: the nondifferentiable case", Report (1976).
- [13] A. Ben-Tal, L. Kerzner and S. Zlobec, "Optimality conditions for convex semi-infinite programming problems", Report (1976).
- [14] A. Ben-Tal and S. Zlobec, "A new class of feasible direction methods", Report (1977).
- [15] A. Ben-Tal and S. Zlobec, "Convex programming and the lexicographic multicriteria problem", <u>Math. Operation-forsch. Statist.</u> Set. Optimization, 8 (1977), 61-73.
- [16] J. Borwein, "Proper efficient points for maximizations with respect to cones", <u>SIAM Journal of Control and Optimization</u>, <u>15</u> (1977), 57-63.
- [17] F. Burns, M. Fiedler and E. Haynsworth, "Polyhedral cones and positive operators", <u>Linear Algebra and its Applications</u>, <u>8</u> (1974), 547-559.
- [18] A. Charnes and W.W. Cooper, <u>Management Models and Industrial</u> <u>Applications of Linear Programming</u>, Vol.1, John Wiley & Sons, New York, 1961.
- [19] F.H. Clarke, "Nondifferentiable functions in optimization", Applied Math. Notes, 2, Dec., (1976).
- [20] B.D. Craven and J.J. Koliha, "Generalizations of Farkas' Theorem", <u>SIAM Journal of Math. Analysis</u>, <u>8</u>, Nov., (1977), 983-997.
- [21] J.W. Daniel and L.L. Schumaker, "On the closedness of the linear image of a set, with applications to generalized spline functions", Applicable Analysis, 4 (1974), 191-205.

- [22] G. Debreu, <u>Theory of Value</u>. An Axiomatic Study of Economic Equilibrium, John Wiley & Sons, 1959.
- [23] V.F. Dem'yanov and V.N. Malozemov, <u>Introduction to Minimax</u>, Translated from Russian by D. Louvish, John Wiley & Sons, New York, 1974.
- [24] J. Dieudonné, "Sur la séparation des ensembles convexes", Mathematische Annalen., 163 (1966), 1-3.
- [25] A. Ya Dubovitskii and A.A. Milyutin, "Extremum problems in the presence of restrictions", <u>Zh. vychish. Mat. mat. Fiz.</u>, <u>5</u> (1965), 395-453.
- [26] A.V. Fiacco, "Convergence properties of local solutions of sequences of mathematical problems in general spaces", <u>Journal</u> of Optimization Theory and Applications, 13 (1974), 1-12.
- [27] S.D. Fisher and J.W. Jerome, <u>Minimum Norm Extremals in Function</u> <u>Spaces</u>, Lecture Notes in Mathematics, 479, Springer-Verlag, New York, 1975.
- [28] I.V. Girsanov, <u>Lectures on Mathematical Theory of Extremum</u> <u>Problems</u>, Lecture Notes in Economics and Mathematical Systems, No.67, Springer-Verlag (1972).
- [29] J.L. Goffin, "On convergence rates of subgradient optimization methods", Mathematical Programming, 13 (1977), 329-347.
- [30] F.J. Gould and J.W. Tolle, "Geometry of optimality conditions and constraint qualifications", <u>Mathematical Programming</u>, <u>2</u> (1972), 1-18.
- [31] F.J. Gould and J.W. Tolle, "A necessary and sufficient qualification for constrained optimization", <u>SIAM Journal of Applied</u> Mathematics, 20, 2 (1971), 164-171.

- [32] M. Guignard, "Generalized Kuhn-Tucker conditions for mathematical programming problems in a Banach space", <u>SIAM Journal of</u> Control, 7 (1969), 232-241.
- [33] P.R. Halmos, <u>Introduction to Hilbert space</u>, Chelsea Publishing Co., New York, 1957.
- [34] R.B. Holmes, <u>A Course on Optimization and Best Approximation</u>, Lecture Notes in Mathematics, No.257, Springer-Verlag (1972).
- [35] R.B. Holmes, <u>Geometric Functional Analysis and its Applications</u>, Graduate Tests in Mathematics, 24, Springer-Verlag (1975).
- [36] H.W. Kuhn, "Nonlinear programming: a historic view", <u>SIAM-</u> AMS Proceedings, 9 (1969), 1-26.
- [37] O.L. Mangasarian, Nonlinear Programming, McGraw-Hill (1969).
- [38] H. Massam, "Mathematical programming with cones", Ph.D. Thesis, McGill University, 1977.
- [39] H. Massam and S. Zlobec, "Various definitions of the derivative in mathematical programming", <u>Mathematical Programming</u>, 7 (1974), 144-161.
- [40] B.N. Pshenichnyi, <u>Necessary Conditions for an Extremum</u>, Marcel Dekker, INC., New York (1971).
- [41] R.T. Rockafellar, Convex Analysis, Princeton University Press (1972).
- [42] R.T. Rockafellar, <u>Conjugate Duality and Optimization</u>, SIAM, Regional Conference Series No.16.
- [43] R.T. Rockafellar, "Ordinary convex programs without a duality gap", Journal of Optimization Theory and Applications, 7, 3 (1971) 143-148.
- [44] R.T. Rockafellar, "Some convex programs whose duals are linearly constrained", <u>Nonlinear Programming</u>, J.B. Rosen, O.L. Mangasarian and K. Ritter, eds., Academic Press, New York, 1970, 293-322.

- [45] W. Rudin, Functional Analysis, McGraw-Hill, 1973.
- [46] W. Rudin, Real and Complex Analysis, McGraw-Hill, 1966.
- [47] M. Slater, "Lagrange multipliers revisited. A contribution to nonlinear programming", <u>Cowles Commission Paper, Math.</u>, <u>403</u> (1950).
- [48] G. Wolkowicz, "Some aspects of stability in nonlinear programming", M.Sc. Thesis, McGill University, 1978.
- [49] H. Wolkowicz, "Calculating the cone of directions of constancy", Journal of Optimization Theory and Applications, <u>25</u> (1978), 451-457.
- [50] S. Zlobec, "Asymptotic Kuhn-Tucker conditions for mathematical programming problems in a Banach space", <u>SIAM Journal of Control</u>, 8, 4 (1970), 505-512.
- [51] S. Zlobec, "Extensions of Asymptotic Kuhn-Tucker conditions in mathematical programming", <u>SIAM Journal of Applied Mathematics</u>, 21, 3 (1971), 448-460.
- [52] S. Zlobec, <u>Mathematical Programming</u>, Lecture Notes, McGill University, 1977.
- [53] G. Zoutendijk, <u>Methods of Feasible Directions, A Study in Linear</u> and Non-linear Programming, Elsevier, New York, 1960.