LINEAR LEAST SQUARES AND QUADRATIC PROGRAMMING

• •

BY

GENE H. GOLUB MICHAEL A. SAUNDERS

-...

TECHNICAL REPORT NO. CS 134 MAY 1969

COMPUTER SCIENCE DEPARTMENT School of Humanities and Sciences STANFORD UNIVERSITY



Linear Least Squares and Quadratic Programming *

Gene H. Golub

Michael A. Saunders

Reproduction in whole or in part is permitted for any purpose of the United States Government.

^{*}Invited paper to be presented at the NATO Advanced Study Institute on Integer and Nonlinear Programming, Ile de Bendor, France, June 8-22, 1969. This research was supported, in part, by the National Science Foundation and the Atomic Energy Commission project at Stanford University.

Abstract

Several algorithms are presented for solving linear least squares problems; the basic tool is orthogonalization techniques. A highly accurate algorithm is presented for solving least squares problems with linear inequality constraints. A method is also given for finding the least squares solution when there is a quadratic constraint on the solution.

Ĺ

 $\left| \right|$

0. Introduction

One of the most common problems in any computation center is that of finding linear least squares solutions. These problems arise in a variety of areas and in a variety of contexts. For instance, the data may be arriving sequentially from a source and there may be some constraint on the solution. Linear least squares problems are particularly difficult to solve because they frequently involve large quantities of data, and they are ill-conditioned by their very nature.

In this paper, we shall present several numerical algorithms for solving linear least squares problems in a highly accurate manner. In addition, we shall give an algorithm for solving linear least squares problem with linear inequality constraints.

1. Linear least sauares

Let A be a given mxn real matrix of rank r and b a given vector. We wish to determine $\mathbf{\hat{x}}$ such that

$$\sum_{i=1}^{m} (b_{i} - \sum_{i=1}^{n} a_{ij} x_{j})^{2} = \min.$$

or using matrix notation

$$\left\| \underline{b} - A \underline{x} \right\|_{2} = \min.$$
 (1.1)

If m>n and r<n, then there is no unique solution. Under these conditions, we require amongst those vectors \underline{x} which satisfy (1.1) that

$$\left\| \hat{\mathbf{x}} \right\|_{2} = \min$$
.

For r = n, \hat{x} satisfies the normal equations

$$A^{T}A\hat{\mathbf{x}} = A^{T}b \quad . \tag{1.2}$$

Unfortunately, the matrix $A^{T}A$ is frequently ill-conditioned and influenced greatly by roundoff errors. The following example illustrates this well. Suppose

which is clearly of rank 4 . Then .

1

$$\mathbf{T}_{\mathbf{A}} = \begin{bmatrix} \mathbf{1} + \boldsymbol{\varepsilon}^2 & 1 & 1 & 1 \\ & 2 & & \\ 1 & 1 & 1 & \mathbf{1} + \boldsymbol{\varepsilon}^2 \\ \mathbf{1} & \mathbf{1} + \boldsymbol{\varepsilon} & \mathbf{1} + \boldsymbol{\varepsilon}^2 & \mathbf{1}^2 \end{bmatrix}$$

and the eigenvalues of $A^{T}A$ are $4+\epsilon^{2}$, ϵ^{2} , ϵ^{2} , ϵ^{2} , ϵ^{2} . Assume that the elements of $A^{T}A$ are computed using double-precision arithmetic, and then rounded to single precision accuracy. Now let η be the largest number on the computer such that $fl(1+\eta) = 1$ where $fl(\ldots)$ indicates the floating point computation. Then if $\epsilon < \sqrt{\eta}$,

a matrix of rank one, and consequently, no matter how accurate the linear equation solver it will be impossible to solve the normal equations (1.2).

LONGLEY[1967] has given examples in which the solution of the normal equations leads to almost no digits of accuracy of the least squares problem.

2. A matrix decomposition

Now $\|\mathbf{y}\|_2 = (\mathbf{y}^T \mathbf{y})^{1/2}$ so that $\|\mathbf{Q}\mathbf{y}\|_2 = \|\mathbf{y}\|_2$ when Q is an orthogonal matrix, <u>viz</u>., $\mathbf{Q}^T \mathbf{Q} = \mathbf{I}$. Thus

$$\|\mathbf{b}-\mathbf{A}\mathbf{x}\|_{2} = \|\mathbf{c}-\mathbf{Q}\mathbf{A}\mathbf{x}\|_{2}$$

where c = Qb and Q is an orthogonal matrix. We choose Q so that

$$QA = R = \begin{bmatrix} \tilde{R} \\ \vdots \\ 0 \end{bmatrix}_{(m-n) \times n}$$
(2.1)

where \widetilde{R} is an upper triangular matrix (V). $_{\rm Let}$

then

$$\begin{aligned} \left\| \underbrace{\mathbf{b}}_{2}^{\mathbf{h}} - \mathbf{A} \mathbf{x} \right\|_{2}^{2} &= (\mathbf{c}_{1}^{\mathbf{h}} - \mathbf{r}_{12}^{\mathbf{h}} \mathbf{x}_{2}^{\mathbf{h}} - \mathbf{r}_{12}^{\mathbf{h}} \mathbf{x}_{2}^{\mathbf{h}} - \mathbf{r}_{1n}^{\mathbf{h}} \mathbf{x}_{n}^{\mathbf{h}} \right)^{2} \\ &+ (\mathbf{c}_{2}^{\mathbf{h}} - \mathbf{r}_{22}^{\mathbf{h}} \mathbf{x}_{2}^{\mathbf{h}} - \mathbf{r}_{2n}^{\mathbf{h}} \mathbf{x}_{n}^{\mathbf{h}} \right)^{2} \\ &+ \dots + (\mathbf{c}_{n}^{\mathbf{h}} - \mathbf{r}_{nn}^{\mathbf{h}} \mathbf{x}_{n}^{\mathbf{h}} \right)^{2} \\ &+ \mathbf{c}_{n+1}^{2} + \mathbf{c}_{n+2}^{2} + \dots + \mathbf{c}_{m}^{2} \end{aligned}$$

Thus $\left\| \underbrace{\mathbf{b}}_{\mathbf{x}} - \mathbf{A} \underbrace{\mathbf{x}}_{\mathbf{x}} \right\|_{2}^{2}$ is minimized when

$$r_{11}\hat{x}_{1} + r_{12}\hat{x}_{2} + ... + r_{1n}\hat{x}_{n} = c_{1}$$
$$r_{22}\hat{x}_{2} + ... + r_{2n}\hat{x}_{n} = c_{2}$$
$$\vdots$$
$$r_{nn}\hat{x}_{n} = c_{n}$$

i.e., $R\hat{x} = \tilde{c}$, where

$$\tilde{c}^{T} = (c_{1}, c_{2}, \dots, c_{n})$$
,

and

$$|\hat{\mathbf{b}} - A\hat{\mathbf{x}}||_2^2 = c_{n+1}^2 + c_{n+2}^2 + \dots + c_m^2$$
 (2.2)

Then

$$R^{T}R = [\widetilde{R} \circ]^{T}[\widetilde{R} \circ] = \widetilde{R}^{T}\widetilde{R}$$
$$= [QA]^{T}[QA] = A^{T}A , \qquad (2.3)$$

and thus $\mathbb{R}^T\mathbb{R}$ is simply the Cholesky decomposition of A^TA .

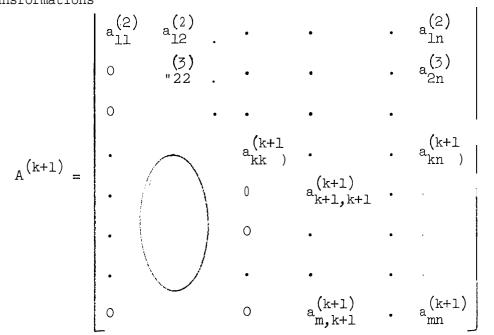
There are a number of ways to achieve the decomposition of (2.1); e.g. one could apply a sequence of plane rotations to annihilate the elements below the diagonal of A . A very effective method to realize the decomposition (2.1) is via Householder transformations. A matrix P is said to be a <u>Householder transformation</u> if

$$P = I - 2uu^T$$
, $u^T u = 1$.

Note that 1) $P = P^{T}$ and 2) $PP^{T} = I - 2uu^{T} - 2uu^{T} + 4uu^{T}uu^{T} = I$ so that P is a symmetric, orthogonal transformation. Let $A^{(1)} = A$ and let $A^{(2)}$, $A^{(3)}$,..., $A^{(n+1)}$ be defined as follows:

$$A^{(k+1)} = P^{(k)}A^{(k)}$$
 (k = 1,2,...,n)

where $P^{(k)} = I - 2w^{(k)}w^{(k)^{T}}$, $w^{(k)^{T}}w^{(k)} = I$. The matrix $P^{(k)}$ is chosen so that $a_{k+1,k}^{(k+1)} = a_{k+2,k}^{(k+1)} = \cdots = f_{n,k}^{(k+1)} = 0$. Thus after k transformations



Note that $|a_{kk}^{(k+1)}| = (\sum_{i=k}^{m} (a_{ik}^{(k)})^2)^{1/2}$ since $P^{(k)}$ is an orthogonal transformation. The details of the computation are given in BUSINGER and GOLUB [1965] and GOLUB [1965]. The Householder transformations have been used in a highly effective manner by KALFON et al. [1968] in the implementation of the projection gradient method.

Clearly

 $R = A^{(n+1)}$

and

$$Q = P^{(n)}P^{(n-1)} \dots P^{(1)}$$

although one need not compute Q explicitly. The number of multiplications required to produce R is roughly $mn^2-(n^3/3)$ whereas approximately $mn^2/2$ multiplications are required to form the normal equations (1.2).

3. The practical procedure

It is known that the Cholesky method for solving systems of equations is numerically stable even if no interchanges of rows and columns are performed. Since we are in effect performing a Cholesky decomposition of $A^{T}A$ no interchanges of the columns of A are needed in most situations. However, numerical experiments have indicated that the accuracy is slightly improved by the interchange strategies outlined below, and consequently, in order to ensure the utmost accuracy one should choose the columns of A by some strategy. In what follows, we shall refer to the matrix $A^{(k)}$ even if some of the columns have been interchanged.

One possibility is to choose at the k^{th} stage the columns of $A^{(k)}$ which will maximize $|a_{kk}^{(k+1)}|$. This is equivalent to searching for the maximum diagonal element in the Cholesky decomposition of $A^{T}A$. Let

$$s_{j}^{(k)} = \sum_{j=k}^{m} (a_{i,j}^{(k)})^{2}$$
 for $j = k, k+1, ..., n$.

-Then since $|a_{kk}^{(k+1)}| = (s_k^{(k)})^{1/2}$, one should choose that column for which (k) is maximized. After $A^{(k+1)}$ has been computed, one can compute $s_j^{(k+1)}$ as follows: $s_j^{(k+1)} = s_j^{(k)} - (a_{k,j}^{(k+1)})^2$ (j = k+1,...,n)

since the orthogonal transformations leave the column lengths invariant. Naturally, the $s_j^{(k)}$'s must be interchanged if the columns of $A^{(k)}$ are interchanged.

The above strategy is useful in determining the rank of a matrix. If the rank of A is r and the arithmetic is performed exactly, then after r transformations

$$A^{(r+1)} = \begin{bmatrix} \tilde{R}_{rxr} & S_{(n-r)xr} \\ 0 & N \end{bmatrix} ,$$

and

L

$$s_{j}^{(r+1)} = 0$$
 for $j = r+1,...,n$

which implies N = 0 . In most situations, however, where rounded arithmetic is used $\|N\|$ = ϵ . It is not easy to determine bounds on ϵ when the rank of A is unknown.

The strategy described above is most appropriate when one has a sequence of vectors b_1, b_2, \ldots, b_p for which one desires a least squares estimate. In many problems, there is but one vector b and one wishes to express it in as few columns of A as possible. Or more precisely, one wishes to determine the k indices such that

$$\sum_{i=1}^{n} (b_i - \sum_{\upsilon=1}^{k} a_{ij} \hat{x}_{j})^2 = \min.$$

We cannot solve this problem, but we shall show how to choose index k when the first k-l indices are given so that the sum of squares of residuals is maximally reduced. This is the stage-wise regression problem.

Let $c_{\tilde{k}}^{(1)} = b$ and $c_{\tilde{k}+1}^{(k+1)} = P_{c_{\tilde{k}}^{(k)}}^{(k)} \cdot Now \tilde{R}_{\tilde{k}}^{(k)} \hat{x}_{\tilde{k}}^{(k-1)} = \tilde{c}_{\tilde{k}}^{(k)}$ where $\hat{x}^{(k-1)}$ is the least squares estimate based on (k-1) columns of A and $\tilde{c}_{\tilde{k}}^{(k)} = (c_{1}^{(k)}, c_{2}^{(k)}, \dots, c_{k-1}^{(k)})$. Thus by (2.2)

$$\begin{aligned} \|\underline{c}^{(k+1)} - \widetilde{R}^{(k+1)} \underline{\hat{x}}^{(k)}\|_{2}^{2} &= \sum_{j=k+1}^{m} (c_{j}^{(k+1)})^{2} \\ &= \sum_{j=k}^{m} (\underline{c}_{j}^{(k+1)})^{2} - (c_{k}^{(k+1)})^{2} \\ &= \sum_{j=k}^{m} (c_{j}^{(k)})^{2} - (c_{k}^{(k+1)})^{2} \end{aligned}$$

since length is preserved under an orthogonal transformation. Consequently, we wish to choose that column of ${\rm A}^{(k)}$ which will maximize $|c_k^{(k+1)}|$. Let

$$t_j^{(k)} = \left(\sum_{i=k}^m a_{ij}^{(k)} c_i^{(k)}\right)$$
 for $j = k+1, \dots, n$.

Then since $|c_k^{(k+1)}| = |(\sum_{i=k}^m a_{ik}^{(k)} c_i^{(k)})/s_k^{(k)}|$, one should choose that column of $A^{(k)}$ for which $(t_j^{(k)})^2/s_j^{(k)}$ is maximized. After $P^{(j)}$ is applied to $A^{(k)}$, one can adjust $t_j^{(k)}$ as follows:

$$t_{j}^{(k+1)} = t_{j}^{(k)} - a_{kj}^{(k+1)}c_{k}^{(k+1)}$$

In many statistical applications, if $(t_j^{(k)})^2/s_j^{(k)}$ is sufficiently small, then no further transformations are performed.

4. Statistical calculations

In many statistical calculations, it is necessary to compute certain -auxiliary information associated with A^TA . These can readily be obtained from the orthogonal decomposition. Thus

$$\det(A^{T}A) = (r_{11} \times r_{22} \times \ldots \times r_{nn})^{2}$$

Since

L

$$A^{T}A = \widetilde{R}^{T}\widetilde{R}$$
, $(A^{T}A)^{-1} = \widetilde{R}^{-1}\widetilde{R}^{-T}$

The inverse of R can be readily obtained since \tilde{R} is an upper triangular matrix. It is possible to calculate $(A^{T}A)^{-1}$ directly from R . Let

$$(A^{T}A)^{-1} = x = (x_{1}, x_{2}, \dots, x_{n}).$$

Then from the relationship

 $\tilde{R} X = \tilde{R}^{-T}$

and by noting that $\{\tilde{R}^{-T}\}_{ii} = 1/r_{ii}$, it is possible to compute $\underset{n}{x}, \underset{n}{x}, \underset{n-1}{x}, \ldots, \underset{n+1}{x}$. The number of operations is roughly the same as in the first method but more accurate bounds may be established for this method provided all inner products are accumulated to double precision.

In some applications, the original set of observations are augmented by an additional set of observations. In this case, it is not necessary to begin the calculation from the beginning again if the method of orthogonalization is used. Let \tilde{R}_1, \tilde{c}_1 correspond to the original data after it has been reduced by orthogonal transformations and let A_2, \tilde{b}_2 correspond to the additional observations. Then the up-dated least squares solution can be obtained directly from

$$A = \begin{bmatrix} A_2 \\ \cdots \\ \widetilde{R}_1 \end{bmatrix}, \quad b = \begin{bmatrix} b_2 \\ \cdots \\ \widetilde{c}_1 \end{bmatrix}$$

This follows immediately from the fact that the product of two orthogonal transformations is an orthogonal transformation.

The above observation has another implication. One of the arguments frequently advanced for using normal equations is that only n(n+1)/2 memory locations are required. By partitioning the matrix A by rows, however, then similarly only n(n+1)/2 locations are needed when the method of orthogonalization is used.

In certain statistical applications, it is desirable to <u>remove</u> a row of the matrix A after the least squares solution has been obtained. This can be done in a very simple manner. Consider the matrix

$$A = \begin{bmatrix} \tilde{R} \\ \dots \\ i & \alpha \\ \tilde{R} \end{bmatrix} \text{ and } d = \begin{bmatrix} \tilde{c} \\ \dots \\ i & \beta \end{bmatrix}$$

where α is the row of A which one wishes to remove, β is the corresponding element of b , and i = $\sqrt{-1}$. Note that

8

$$S^{T}S = \widetilde{R}^{T}\widetilde{R} - \alpha^{T}\alpha = A^{T}A - \alpha^{T}\alpha$$

Let

L

L

-

We choose $\cos \theta$ so that $\{s^{(2)}\}_{n+1,1} = 0$. Thus

$$\{s^{(2)}\}_{1,1} = \sqrt{(r_{11}^2 - \alpha_1^2)}$$

$$\{s^{(2)}\}_{1,j} = (r_{11}r_{1j} - \alpha_1\alpha_j) / \sqrt{(r_{11}^2 - \alpha_1^2)} \quad j = 2,3,\ldots,n$$

$$\{s^{(2)}\}_{n+1,j} = i(\alpha_1r_{1j} - \alpha_jr_{11}) / \sqrt{(r_{11}^2 - \alpha_1^2)} \quad j = 2,3,\ldots,n$$

Note no complex arithmetic is really necessary. The process is continued as follows:

Let

Then

$$S^{(k+1)} = Z_{k,n+1} S^{(k)}$$
, $k = 1,2,...,n$,

and $\cos \Theta_k$ is determined so that $(S^{(k+1)})_{k,n+1} = 0$. Thus roughly $3n^2$ multiplications and divisions and n square roots are required to form the new \tilde{R} .

Suppose it is desirable to add an additional variable so that the matrix A is augmented by a vector g (say). The first n columns of $\tilde{R}^{(n)}$ are unchanged. Now one computes

$$h = P^{(n)} \dots P^{(2)} P^{(1)} g$$

From h one can compute $P^{(n+1)}$ and apply it to $P^{(n)} \dots P^{(1)}b$. This technique is also useful when an auxiliary serial storage (e.g. magnetic tape) is used.

It is also possible to drop one of the variables in a simple fashion after \tilde{R} has been computed. For example, suppose we wish to drop variable 1 , then

$$\widetilde{\mathbf{R}} = \begin{bmatrix} \mathbf{r}_{12} & \cdots & \mathbf{r}_{\ln} \\ \mathbf{r}_{22} & \cdots & \cdots \\ \vdots & & & \\ \vdots & & & \mathbf{r}_{nn} \end{bmatrix}_{nx(n-1)}$$

By using plane rotations, similar to those given by (4.1), it is possible to reduce R to the triangular form again.

5. Gram-Schmidt orthogonalization

In \$2, it was shown that it is possible to write

$$QA = R \qquad (5.1)$$

The matrix Q is constructed as a' product of Householder transformations.

From (5.1), we see that

$$A = Q^{T}R = PS$$

where $P^{T}P = I_{n}$, $S : \nabla$. Each row of S and each column of P is uniquely determined up to a scalar factor of modulus one. In order to avoid computing square roots, we modify the algorithms so that S is an upper triangular matrix with ones on the diagonal. Thus $P^{T}P = D$, a diagonal matrix. The calculation of P and S may be calculated in two ways.

a) Classical Gram-Schmidt Algorithm (CGSA)

The elements of S are computed one column at a time. Let

$$\mathbf{A}^{(k)} = [\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_{k-1}, \mathbf{a}_k, \dots, \mathbf{a}_n]$$

and assume

$$\sum_{i=1}^{T} p_{ij} = \delta_{ij} d_{ij} , \quad 1 \leq i,j \leq k-1 .$$

At step \boldsymbol{k} , we compute

$$s_{ik} = (p_i^T a_k / d_i) , \quad 1 \le i \le k-1$$
$$p_k = a_k - \sum_{i=1}^{k-1} s_{ik} p_i , \quad d_k = ||p_k||_2^2$$

b) Modified Gram-Schmidt Algorithm (MGSA)

Here the elements of S are computed one row at a time. We define

$$A^{(k)} (\underline{p}_1, \underline{p}_2, \dots, \underline{p}_{k_1}, \underline{a}_{k_k}^{(k)}, \dots, \underline{a}_{n}^{(k)})$$

and assume

$$\sum_{i=1}^{T} p_{i} = \delta_{ij} d_{i} , \qquad q_{i=1}^{T} a_{\ell}^{(k)} = 0 , \qquad l \leq i, j \leq k-1 , \qquad k \leq \ell \leq n .$$

At step k , we take $p = \frac{a^{(k)}}{k}$, and compute

$$\mathbf{d}_{\mathbf{k}} = \|\mathbf{p}_{\mathbf{k}}\|_{2}^{2}, \mathbf{s}_{\mathbf{k}\boldsymbol{\ell}} = (\mathbf{p}_{\mathbf{k}}^{\mathrm{T}} \mathbf{a}_{\boldsymbol{\ell}}^{(\mathbf{k})})/\mathbf{d}_{\mathbf{k}}, \mathbf{a}_{\boldsymbol{\ell}}^{(\mathbf{k}+1)} = \mathbf{a}_{\boldsymbol{\ell}}^{(\mathbf{k})} \mathbf{s}_{\mathbf{k}\boldsymbol{\ell}} \mathbf{p}_{\mathbf{k}}, \mathbf{k}+1 \leq \boldsymbol{\ell} < n.$$

In both procedures, $s_{kk} = 1$. The two procedures in the absence of roundoff errors, produce the same decomposition. However, they have completely different numerical properties when n > 2. If A is at all "ill-conditioned", then using the CGSA, the computed columns of P will soon lose their orthogonality. Consequently, one should never use the CGSA without reorthogonalization, which greatly increases the amount of computation. Reorthogonalization is never needed when using the MGSA. A careful roundoff analysis is given by BJORK [1967]. RICE [1966] has shown experimentally that the MGSA produces excellent results.

The MGSA has the advantages that it is relatively easy to program, and experimentally (cf. JORDAN [1968]), it seems to be slightly more accurate than the Householder procedure. However, it requires roughly $mn^2/2$ operations which is slightly more than that necessary in the Householder procedure. Furthermore, it is not as simple as the Householder procedure to add observations.

6. Sensitivity of the solution

We consider first the inherent sensitivity of the solution of the least squares problem. For this purpose it is convenient to introduce the condition number $\kappa(A)$ of a non-square matrix A. This is defined by

$$\kappa(\mathbf{A}) = \sigma_1 / \sigma_n, \sigma_1 = \max_{\mathbf{x} \neq 0} ||\mathbf{A}\mathbf{x}||_2 / ||\mathbf{x}||_2, \sigma_n = \min_{\mathbf{x} \neq 0} ||\mathbf{A}\mathbf{x}||_2 / ||\mathbf{x}||_2$$

so that σ_1^2 and σ_n^2 are the greatest and the least eigenvalues of $A^T A$. From its definition it is clear that $\kappa(A)$ is invariant with respect to unitary transformations. If \tilde{R} is defined as in (2.1) then

$$\sigma_{1}(\tilde{R}) = \sigma_{1}(A) , \sigma_{n}(\tilde{R}) = \sigma_{n}(A) , \kappa(\tilde{R}) = \kappa(A) ,$$

while

$$\sigma_1(\tilde{\mathbf{R}}) = \|\tilde{\mathbf{R}}\|_2$$
 and $\sigma_n(\tilde{\mathbf{R}}) = 1/\|\tilde{\mathbf{R}}^{-1}\|_2$.

The commonest method of solving least squares problems is via the normal equations

$$A^{T}Ax = A^{T}b \qquad (6.1)$$

The matrix $\boldsymbol{A}^{T}\!\boldsymbol{A}$ is square and we have

$$\kappa (A^{T}A) = \kappa^{2}(A)$$
.

This means that if A has a condition number of the order of $2^{t/2}$ then $A^{T}A$ has a condition number of order 2^{t} and it will not be possible using t-digit arithmetic to solve (6.1). The method of orthogonal transformations replaces the least squares problem by the solution of the equations $\tilde{R}x = \tilde{g}$ and $\kappa(\tilde{R}) = \kappa(A)$. It would therefore seem to have substantial advantages since we avoid working with a matrix with condition number $\kappa^{2}(A)$.

We now show that this last remark is an oversimplification. To this end, we compare the solution of the original system [A : b] with that of a perturbed system. It is convenient to assume that

$$\sigma_1 = \|A\|_2 = \|b\|_2 = 1;$$

this is not in any sense a restriction since we can make $\|A\|_2$ and $\|b\|_2$ of order unity merely by scaling by an appropriate power of two. We now have

$$\kappa(\mathbf{A}) = \kappa(\tilde{\mathbf{R}}) = \|\tilde{\mathbf{R}}^{-1}\|_2 = 1/a_n$$

Consider the perturbed system

$$(\mathbf{A} + \boldsymbol{\varepsilon}\mathbf{E}; \mathbf{b} + \boldsymbol{\varepsilon}\mathbf{e})$$
 , $\|\mathbf{E}\|_2 = \|\mathbf{e}\|_2 = 1$,

where $\boldsymbol{\epsilon}$ is to be arbitrarily small. The solution $\boldsymbol{\bar{x}}$ of the perturbed system satisfies the equation

$$(A + \varepsilon E)^{T}(A + \varepsilon E)\bar{x} = (A + \varepsilon E)^{T}(b + \varepsilon e)$$
(6.2)

If $\hat{\mathbf{x}}$ is the exact solution of the original system and Q is the exact orthogonal transformation corresponding to A we have

$$QA = \begin{bmatrix} \tilde{R} \\ \vdots \\ 0 \end{bmatrix}, \quad Q(A + \varepsilon E) = \begin{bmatrix} \tilde{R} + \varepsilon F \\ \vdots \\ \varepsilon G \end{bmatrix}, \quad Qe = \begin{bmatrix} f \\ \vdots \\ g \\ \vdots \end{bmatrix}$$

and

$$\mathbf{r} = \mathbf{b} - \mathbf{A} \mathbf{\hat{x}} , \mathbf{A}^{\mathbf{T}} \mathbf{r}_{\mathbf{x}} = \mathbf{\Theta}$$

Equation (6.2) therefore becomes

$$(\mathbf{A} + \varepsilon \mathbf{E})^{\mathrm{T}}(\mathbf{A} + \varepsilon \mathbf{E}) = (\mathbf{A}^{\mathrm{T}} + \varepsilon \mathbf{E}^{\mathrm{T}})(\mathbf{A}\mathbf{x} + \mathbf{x} + \varepsilon \mathbf{e})$$

giving

$$\begin{bmatrix} \mathbf{R} + \varepsilon \mathbf{F} \\ \vdots \\ \varepsilon \mathbf{G} \end{bmatrix}^{\mathrm{T}} \begin{bmatrix} \widetilde{\mathbf{R}} + \varepsilon \mathbf{F} \\ \vdots \\ \varepsilon \mathbf{G} \end{bmatrix}^{\widetilde{\mathbf{x}}} = \begin{bmatrix} \mathbf{R} + \varepsilon \mathbf{F} \\ \vdots \\ \varepsilon \mathbf{G} \end{bmatrix}^{\mathrm{T}} \left(\begin{bmatrix} \widetilde{\mathbf{R}} \\ \vdots \\ \mathbf{O} \end{bmatrix}^{\mathrm{x}} + \varepsilon \begin{bmatrix} \mathbf{f} \\ \vdots \\ \mathbf{g} \end{bmatrix} \right) + \varepsilon \mathbf{E}^{\mathrm{T}} \mathbf{r} \quad .$$

Neglecting ε where advantageous,

$$(\widetilde{\mathbf{R}} + \varepsilon \mathbf{F})^{\mathrm{T}} (\widetilde{\mathbf{R}} + \varepsilon \mathbf{F}) \widetilde{\mathbf{x}} = (\widetilde{\mathbf{R}} + \varepsilon \mathbf{F})^{\mathrm{T}} \widetilde{\mathbf{R}} \widetilde{\mathbf{x}} + \varepsilon (\widetilde{\mathbf{R}} + \varepsilon \mathbf{F})^{\mathrm{T}} \mathbf{f} + \varepsilon \mathbf{E}^{\mathrm{T}} \mathbf{x} + \mathbf{0} (\varepsilon^{2})$$

$$\widetilde{\mathbf{x}} = (\widetilde{\mathbf{R}} + \varepsilon \mathbf{F})^{-1} \widetilde{\mathbf{R}} \widetilde{\mathbf{x}} + \varepsilon (\widetilde{\mathbf{R}} + \varepsilon \mathbf{F})^{-1} \mathbf{f} + \varepsilon (\widetilde{\mathbf{R}}^{\mathrm{T}} \widetilde{\mathbf{R}})^{-1} \mathbf{E}^{\mathrm{T}} \mathbf{r} + \mathbf{0} (\varepsilon^{2})$$

$$= \widehat{\mathbf{x}} - \varepsilon \widetilde{\mathbf{R}}^{-1} \mathbf{F} \widetilde{\mathbf{x}} + \varepsilon \mathbf{R}^{-1} \mathbf{f} + \varepsilon \mathbf{R}^{-1} \mathbf{f} + \varepsilon (\widetilde{\mathbf{R}}^{\mathrm{T}} \widetilde{\mathbf{R}})^{-1} \mathbf{E}^{\mathrm{T}} \mathbf{r} + \mathbf{0} (\varepsilon^{2})$$

giving

$$\begin{split} \|\bar{\mathbf{x}} - \hat{\mathbf{x}}\|_{2} &\leq \varepsilon \|\tilde{\mathbf{R}}^{-1}\|_{2} \|\mathbf{F}\|_{2} \|\hat{\mathbf{x}}\|_{2} + \varepsilon \|\tilde{\mathbf{R}}^{-1}\|_{2} \|\mathbf{f}\|_{2} + \varepsilon \|\tilde{\mathbf{R}}^{-1}\|_{2}^{2} \|\mathbf{E}\|_{2} \|\mathbf{f}\|_{2} + o(\varepsilon^{2}) \\ &\leq \varepsilon \kappa(\mathbf{A}) \|\hat{\mathbf{x}}\|_{2} + \varepsilon \kappa(\mathbf{A}) + \varepsilon \kappa^{2}(\mathbf{A}) \|\mathbf{f}\|_{2} + o(\varepsilon^{2}) . \end{split}$$

We observe that the bounds include a term $\epsilon \kappa^2(A) \|\mathbf{r}\|_2$. It is easy to verify by means of a 3 x 2 matrix A that this bound is realistic and that an error of this order of magnitude does indeed result from almost any such perturbation E of A. We conclude that although the use of the orthogonal transformation avoids some of the ill effects inherent in the use of the normal equations the value κ (A) is still relevant to some extent.

When the equations are compatible $\|\mathbf{r}\|_2 = 0$ and the term in $\kappa^2(\mathbf{A})$ d&appears. In the non-singular linear equation case r is always null and hence it is always $\kappa(\mathbf{A})$ rather than κ (\mathbf{A}) 'which-is relevant.

Since the sensitivity of the solution depends on the condition number, it is frequently desirable to replace the original unknowns x by a new vector of unknowns $D^{-1}x$ where D is a-diagonal matrix with-non-zero diagonal elements. Thus we wish to find \hat{y} for which

 $\left\| \underbrace{\mathbf{b}}_{\widetilde{\mathbf{c}}} - C \underbrace{\mathbf{\hat{y}}}_{\widetilde{\mathbf{c}}} \right\|_{2} = \min.$

where C = AD and $\hat{y} = D^{-1}\hat{x}$. Let \hat{y} be the set of all $n \times n$ diagonal matrices with non-zero diagonal elements. We wish to choose ${\tt D}$ so that

$$\kappa(A\hat{D}) \leq \kappa(AD)$$
 for all D_{eq}

Let
$$\tilde{D}_{e3}n$$
 and $\{\tilde{D}\}_{ii} = 1/||a_i||_2$. VAN, DER SIUIS [1968] has shown that $\kappa(\tilde{AD}) \leq \sqrt{n} \kappa(\tilde{AD})$.

Therefore in the absence of other information, it would appear that it is best to precondition the matrix A so that all columns of the matrix A have equal length. In practice, one adjusts the exponents of the stored elements of A so that the mantissa of the floating point representation is not changed,

Iterative refinement for least squares problems 7.

The iterative refinement method may be used for improving the solution to linear least squares problems. Let

$$\alpha \rho = b - A \hat{x}$$
, $\alpha > 0$

so that

. .

5

$$\alpha \mathbf{A}^{\mathrm{T}} \boldsymbol{\rho} = \mathbf{A}^{\mathrm{T}} \mathbf{b} - \mathbf{A}^{\mathrm{T}} \mathbf{A} \mathbf{\hat{x}} = \mathbf{\Theta}$$

When $\alpha = 1$, the vector ρ is simply the residual vector r . Thus

$$\begin{bmatrix} \alpha \mathbf{I} & \mathbf{A} \\ \hline \mathbf{A}^{\mathrm{T}} & \mathbf{O} \end{bmatrix} \begin{bmatrix} \rho \\ \vdots \\ \hat{\mathbf{x}} \\ \neg \end{bmatrix} = \begin{bmatrix} b \\ \vdots \\ \mathbf{O} \\ \neg \end{bmatrix} , \qquad (7.1)$$

or

$$Cy = g$$

One of the standard methods for solving linear equations may now be used to solve (7.1). However, this is quite wasteful of memory space since the dimension of the system to be solved is (m+n) . We may simplify this problem somewhat by noting with the aid of (2.3) that

$$\begin{bmatrix} \alpha \mathbf{I} & \mathbf{A} \\ \hline \mathbf{A}^{\mathrm{T}} & \mathbf{0} \end{bmatrix} = \begin{bmatrix} \sqrt{\alpha} \mathbf{I} & \mathbf{0} \\ \hline \frac{1}{\sqrt{\alpha}} \mathbf{A}^{\mathrm{T}} & \frac{1}{\sqrt{\alpha}} \tilde{\mathbf{R}}^{\mathrm{T}} \end{bmatrix} \begin{bmatrix} \sqrt{\alpha} \mathbf{I} & \frac{1}{\sqrt{\alpha}} \mathbf{A} \\ \hline \mathbf{0} & -\frac{1}{\sqrt{\alpha}} \tilde{\mathbf{R}} \end{bmatrix} = \mathbf{IU} \cdot (7.2)$$

Once an approximate solution to Cy = g has been obtained, it is frequently possible to improve the accuracy of the approximate solution. Let \bar{y} be an approximate solution, and let $v = g - C\bar{y}$. Then if $y = \bar{y} + \delta_{z}$, δ satisfies the equation

$$C_{\tilde{v}}^{\delta} = v \qquad (7.3)$$

Equation (7.3) can be solved approximately from the decomposition (7.2). Of course, it is not possible to solve precisely for 8 so that the process may be repeated.

We are now in a position to use the <u>iterative refinement method</u> (cf. MOLER [1967], WILKINSON [1967]) for solving linear equations. Thus one might proceed as follows:

1) Solve for $x^{(\circ)}$ using one of the orthgonalization procedures outlined in § 2 or 5. \tilde{R} must be saved but it is not necessary to retain Q. Then

$$\rho^{(o)} = \frac{1}{\alpha} (b - Ax^{(o)})$$

2) The vector $y^{(s+1)}$ is determined from the relationship

$$y^{(s+1)} = y^{(s)} + \delta^{(s)}$$

where

$$C\delta^{(s)} = g^{-Cy}(s) \equiv v^{(s)}$$
(7.4)

۰ ۴

This calculation is simplified by solving

$$L_{z}(s) = v(s)$$
$$U_{z}(s) = (z)s$$

The vector $v^{(s)}$ must be calculated using double precision accuracy and then rounding to single precision.

3) Terminate the iteration when $\|\delta^{(s)}\| / \|y^{(s)}\|$ is less than a prescribed number.

Note that the computed residual vector is an approximation to the residual vector when the exact solution $\hat{\mathbf{x}}$ is known. This may differ from the residual vectorcomputed from the approximate solution to the least squares problem.

There are three sources of error in the process: (1) computation of the vector $v_{x}^{(s)}$, (2) solution of the system of equations for the correction vector $\delta_{x}^{(s)}$, and (3) addition of the correction vector to the approximation $y^{(s)}$. It is absolutely necessary to compute the components of the vector $v_{x}^{(s)}$ using double precision inner products and then to round to single precision accuracy. The convergence of the iterative refinement process has been discussed in detail by MOLER [1967]. Generally speaking, for a large class of matrices for $k \ge k_0$ all components of $y^{(s)}$ are the correctly rounded single precision approximations to the components of \bigvee . There are exceptions to this, however, (cf. KAHAN [1966]). Experimentally, it has been observed, in most instances, that if $\|\delta_{x}^{(0)}\|_{\infty} / \|y_{x}^{(0)}\|_{\infty} \le 2^{-p}$ where

$$\|\underbrace{\mathbb{N}}_{\infty} = \max |\mathbf{y}_{i}| \\ 1 \le i \le n$$

then $k_0 \ge [t/p]$. We shall return to the subject of iterative refinement when we discuss the solution of linear least squares problem with linear constraints.

A variant of the above procedure has been analyzed by BJÖRCK [1967b], [1968], and he has also given an ALGOL procedure. This has proved to be a very effective method for obtaining highly accurate solutions to linear : least squares problems.

8. Least squares problems with constraints

Frequently, one wishes to determine $\hat{\mathbf{x}}$ so that $\|\underline{\mathbf{b}}-\mathbf{A}\hat{\mathbf{x}}\|_2$ is minimized subject to the condition that $G\hat{\mathbf{x}} = \mathbf{h}$ where G is a pxn matrix of rank p. One can, of course, eliminate $\hat{\mathbf{p}}$ of the columns of A by Gaussian elimination after a pxp non-singular submatrix of G has been determined and then solve

17

the resulting normal equations. This, unfortunately, would not be a numerically stable scheme since no row interchanges between A and G would be permitted.

If one uses Lagrange multipliers, then one must solve the $(n+p)\chi(n+p)$ system of equations.

$$\begin{bmatrix} A^{T}A & G^{T} \\ \hline G & O \end{bmatrix} \begin{bmatrix} \hat{x} \\ \vdots \\ \lambda \\ \vdots \end{bmatrix} = \begin{bmatrix} A^{T}b \\ \vdots \\ h \\ a \end{bmatrix}$$
(8.1)

where λ is the vector of Lagrange multipliers. Since $\hat{\mathbf{x}} = (\mathbf{A}^{\mathrm{T}} \mathbf{A})^{-1} \mathbf{A}^{\mathrm{T}} \mathbf{b} - (\mathbf{A}^{\mathrm{T}} \mathbf{A})^{-1} \mathbf{G}^{\mathrm{T}} \lambda$,

$$G(A^{T}A)^{-1} G^{T} \lambda = G_{Z-g}$$

where

$$z = (\underline{A}^{T}A)^{-1} A^{T} b$$

Note z is the least squares solution of the original problem without constraints and one would frequently wish to compare this vector with the final solution $\hat{\mathbf{x}}$. The vector z , of course, should be computed by the orthogonalization procedures discussed earlier.

Since $A^T A = \tilde{R}^T \tilde{R}$, $G(A^T A)^{-1} G^T = W^T W$ where $W = \tilde{R}^{-T} G^T$. After W is computed, it should be reduced to a pxp upper triangular matrix K by orthogonalization. The matrix equation

$$K^{T}K\lambda = Hz-g$$

should be solved by the obvious method. Finally, one computes

$$\hat{\mathbf{x}} = \mathbf{z} - (\mathbf{A}^{\mathrm{T}} \mathbf{A})^{-1} \mathbf{G} \boldsymbol{\lambda}$$

where $(A^{\rm T} A)^{-1} G \! \lambda$ can be easily computed by using \widetilde{R}^{-1} .

It is also possible to use the techniques described in §7. Again, let $r = b-A\hat{x}$ so that from (8.1)

$$\begin{bmatrix} I & A & O \\ \hline A^{T} & O & G^{T} \\ \hline O & G & O \end{bmatrix} \begin{bmatrix} r \\ \vdots & \vdots \\ \hat{x} \\ \vdots \\ \lambda \end{bmatrix} = \begin{bmatrix} b \\ \vdots \\ O \\ a \vdots \\ b \\ \vdots \\ h \end{bmatrix}$$
(8.2)

Dz=g

Note D is an $(m+n+p)\times(m+n+p)$ matrix. We may simplify the solution of (8.2), however, by noting that

$$\begin{bmatrix} \mathbf{I} & \mathbf{A} & \mathbf{0} \\ \mathbf{A}^{\mathrm{T}} & \mathbf{0} & \mathbf{G}^{\mathrm{T}} \\ \mathbf{0} & \mathbf{G} & \mathbf{0} \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{0} & \mathbf{0} \\ \mathbf{A}^{\mathrm{T}} & -\mathbf{\widetilde{R}}^{\mathrm{T}} & \mathbf{0} \\ \mathbf{0} & \mathbf{B}^{\mathrm{T}} & \mathbf{S}^{\mathrm{T}} \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{\widetilde{R}} & -\mathbf{B} \\ \mathbf{0} & \mathbf{0} & \mathbf{S} \end{bmatrix}$$
(8.3)

where $B = (\tilde{GR}^1)^T = PS$ and $P^TP = I$ with $S : \nabla$. The decomposition (8.3) can be used very effectively in conjunction with the method of iterative refinement. BJÖRCK and GOLUB [1967] have given a variant of the above procedure which requires Q and P.

9. Linear least squares solutions with inequality constraints

Again let A,G be given real matrices of orders $m\chi n$, $p\chi n$, with $m\geq n$, and let b , h be given real vectors of orders m , p . For any vector x we define

$$r = b-Ax$$

and we wish to determine an x such that

$$r_{UN}^{T} = min.'$$

· subject to

$$G_{x} \geq h$$

Our problem can therefore be stated as follows: find r , x , w such that

r + Ax = b G x - w = h $w \ge Q$ $r^{T}r = min.$

or

These problems can be solved by quadratic programming but we present an algorithm in this section which leads to a much smaller ^{system} of equations and highly accurate results.

If we define

$$f(\mathbf{r}, \mathbf{w}, \mathbf{x}, \mathbf{y}, \mathbf{z}) = \frac{1}{2} \stackrel{T}{\underset{\sim}{x}} \mathbf{r} - \underbrace{\mathbf{y}}^{\mathrm{T}}(\mathbf{r} + A\mathbf{x} - b) - \underbrace{\mathbf{z}}^{\mathrm{T}}(G\mathbf{x} - \mathbf{w} - b)$$

where we require without loss of generality that $z \ge Q$, then an equivalent problem is to determine r, w, x, y, z such that

 $w, z \ge \Theta$ f = min.

Equating to zero the partial derivatives of f with respect to x, x, y, z respectively, we get

 $r - y = \Theta$ $-A^{T}y - G^{T}z = \Theta$ $r + Ax - b = \Theta$ Gx - w - h = 8

Further, let the elements of w, z be w_i, z_i (i = 1,2,...,p). Then

Now if $w_i > 0$ in the optimal solution, the constraint $w_i \ge 0$ is not -binding and we have

$$\frac{\partial f}{\partial w_i} = 0$$
,

i.e.,

$$w_{i} > 0 => z_{i} = 0$$

Since $\textbf{z}_{i} \geq 0$, this further means that

$$z_i > 0 \Rightarrow w_i = 0$$
.

(For otherwise, $\mathbf{z}_i > 0 \Rightarrow \mathbf{w}_i > 0 \Rightarrow \mathbf{z}_i = 0$ which is a contradiction.) Accordingly, our problem has become one of finding a solution of the system

$$\begin{array}{ccc} r + Ax & = b \\ \sim & \sim & \sim \\ \end{array} \tag{9.1}$$

$$A^{\mathrm{T}}r + G^{\mathrm{T}}z = \Theta$$
(9.2)

 $G_{X} - w = h$ (9.3)

such that

$$z \geq \Theta \quad , \quad w \geq \Theta \quad , \quad z = 0$$

We now determine an orthogonal matrix ${\tt Q}$ and an upper-triangular matrix ${\tt R}$ such that

$$A = QR$$
,

where R is nxn and non-singular if rank(A) = n. Then

$$A^{T}A = R^{T}Q^{T}QR = R^{T}R$$
.

Letting $B = (GR^{-1})^T$ and eliminating **r** from (9.1) and (9.2) it is easily verified that

$$x = \hat{x} + R^{-1}Bz$$
, (9.4)

where

$$\mathbf{\hat{x}} = (\mathbf{R}^{\mathrm{T}}\mathbf{R})^{-1} \mathbf{A}^{\mathrm{T}}\mathbf{b}$$

is the unconstrained least squares solution (i.e., the solution of (9.1) and (9.2) with z = 0). \hat{x} is found by the methods of §7.

We now determine if $\hat{\mathbf{x}}$ satisfies the original inequalities: if we define $q = G\hat{\mathbf{x}} - \hat{\mathbf{h}}$ and find that $q \geq \underline{\Theta}$ then the constraints are satisfied and $\hat{\mathbf{x}}$ solves the problem.

"Otherwise, we substitute (9.4) in (9.3) and obtain

$$G(\hat{x} + R^{-1}Bz) - w = h$$

or

$$B^{T}Bz + q = w$$

a solution, or with an indication that none exists.

 $\begin{array}{c} \underline{z} \geq \Theta \\ \underline{z} \geq \Theta \\ \underline{z} \end{array}, \quad \underline{w} \geq \Theta \\ \underline{z} \quad \underline{z}^{\mathrm{T}} \underline{w} = 0 \end{array}$

Thus we find that z,w solve the <u>linear complementarity problem</u> (LCP) defined by (9.5). This is a fundamental mathematical programming problem and several algorithms have been developed for finding solutions (e.g. see . LEMKE [1968], COTTLE [1968], COTTLE and DANTZIG [1968]). The matrix $M = B^{T}B$ is positive semi-definite, and this is one of the cases when, for example, the principal pivoting method in COTTLE [1968] guarantees termination with

Once z has been found it would be a simple matter to substitute into (9.1), (9.2) and find \mathbf{r}, \mathbf{x} from

$$\left. \begin{array}{c} \mathbf{r} + \mathbf{A} \mathbf{x} = \mathbf{b} \\ \mathbf{a}^{\mathrm{T}} \mathbf{r} &= -\mathbf{G}^{\mathrm{T}} \mathbf{z} \end{array} \right\}$$
(9.6)

In practice, however, if we are concerned with the accuracy of our estimate of x we use the solution of the LCP (9.5) only to determine which elements of \tilde{w} are exactly zero. These are the w_i which are non-basic in the solution of (9.5). (There is certainly at least one such w_i , for otherwise we would have z = 0, $w \ge 0$, which is the case checked for earlier in determining whether or not \hat{x} solved the problem.)

We now delete from (9.3) those constraints for which w_i is basic, obtaining an l_{XN} system of equations

$$\tilde{\mathbf{G}}\mathbf{x} = \tilde{\mathbf{h}}$$

where $l \leq \ell \leq p$.

If \tilde{z} is the vector z with the corresponding elements deleted, the remaining step is to solve the system

$$r + Ax = b$$

$$A^{T}r + \tilde{G}^{T}\tilde{z} = 0 \qquad (9.7)$$

$$\tilde{G}x = h$$

where we are now working with original data and can therefore expect a more accurate solution than could be obtained from (9.6). We can now apply the methods of §8 to this system of equations.

The standard methods for solving the linear complementarity problem employ the elements of \underline{w} as the initial set of basic variables, with all elements of \underline{z} initially non-basic. In general, it is probable that only a small proportion of the inequalities in the original problem will be constraining the system, which means that only a small proportion of the w_1 will be non-zero. Hence it might be expected in general that only a small number of iterations (relative to p) should be required to bring some of the z_i into the basis and reach a feasible solution.

In our particular form of the problem, since the matrix $M = B^T B$ has its largest elements on the diagonal, accuracy can be conserved, to within the limits of the error in forming M, by interchanging rows whenever a column of M is brought into the basis in such a way that the diagonal elements of M become diagonal elements of the basis matrix. This is easily done if the LU decomposition of the basis is calculated each iteration as in the treatment of the simplex method by BARTELS [1968] and BARTELS and GOLUB [1969].

Note that $B = (GR^{-1})^T$ can be determined column by column via repeated back-substitution on the system

$$R^{T}B = G^{T}$$

subject

The algorithm presented here can be used for any quadratic programming problem when a positive definite quadratic form is given. Suppose we wish to determine an x such that

$$x^{T}Cx + d^{T}x = \min.$$
to
$$Gx \ge h$$

$$Gx = h$$

$$Gx = h$$

$$Gx = h$$

$$Gx = h$$

23

Since C is positive definite, we may write

 $C = R^{T}R$

where $R(\nabla)$ is the Cholesky factor of C. Such a decomposition can easily be computed. If we now define 'b = $-\frac{1}{2}R^{-T}d$ (and calculate b from $R^{T}b = -\frac{1}{2}d$) we find that

 $\|\mathbf{b} - \mathbf{R}\mathbf{x}\|_{2}^{2} = \mathbf{b}_{\mathbf{x}}^{\mathrm{T}}\mathbf{b}_{\mathbf{x}} - 2\mathbf{b}_{\mathbf{x}}^{\mathrm{T}}\mathbf{R}\mathbf{x}_{\mathbf{x}} + \mathbf{x}^{\mathrm{T}}\mathbf{R}^{\mathrm{T}}\mathbf{R}\mathbf{x}_{\mathbf{x}}$ $= \mathbf{b}_{\mathbf{x}}^{\mathrm{T}}\mathbf{b}_{\mathbf{x}} + \mathbf{d}_{\mathbf{x}}^{\mathrm{T}}\mathbf{x}_{\mathbf{x}} + \mathbf{x}_{\mathbf{x}}^{\mathrm{T}}\mathbf{C}\mathbf{x}_{\mathbf{x}}$

and consequently if we determine an x such that

$$\left\| \underbrace{\mathbf{b}}_{\mathbf{x}} - \operatorname{Rx}_{\mathbf{x}} \right\|_{2} = \min.$$

subject to $Gx \ge h$

then x will satisfy (9.8) as required.

10. Singular systems

If the rank of A is less than n and if column interchanges are performed to maximize the diagonal elements of R , then

$$A^{(r+1)} = \begin{bmatrix} \frac{\tilde{R}_{rxr} & S_{(n-r)xr} \\ 0 & 0 \end{bmatrix}$$

when rank(A) = r . A sequence of Householder transformations may now be applied on the right of $A^{(r+1)}$ so that the elements of $S_{(n-r)xr}$ become annihilated. Thus dropping subscripts and superscripts, we have

$$QAZ = T = \left[\begin{array}{c} \widetilde{T} & O \\ \hline \mathbf{0} & \mathbf{0} \end{array} \right]$$

where T is an **rxr** upper triangular matrix. Now

$$\|\mathbf{b} - \mathbf{A}\mathbf{x}\|_{2} = \|\mathbf{b} - \mathbf{Q}^{T} \mathbf{T} \mathbf{Z}^{T}\mathbf{x}\|_{2}$$
$$= \|\mathbf{c} - \mathbf{T}\mathbf{y}\|_{2}$$

where c = Qb and $y = Z^T x$. Since T is of rank r, there is no unique solution so "that we impose the condition that $\|\hat{x}\|_2 = \min$. But $\|y\|_2 = \|x\|_2$ since T is orthogonal and $\|y\|_2 = \min$. when

$$\forall_{r+1} = y_{r+2} = \dots = y_m = 0$$

Thus

$$\hat{\mathbf{x}} = \mathbf{Z} \begin{bmatrix} \widetilde{\mathbf{T}}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{Q} \mathbf{b}$$

This solution has been given by FADEEV, et. al. [1968] and HANSON and LAWSON [1968]. The problem still remains how to numerically determine the rank which will be discussed in §12.

11. Singular value decomposition

$$A = U\Sigma V^{\mathrm{T}}$$
(11.1)

where

 $UU^{T} = I_{m}$, $VV^{T} = I_{n}$

and

L

1

The matrix U consists of the orthonormalized eigenvectors of $A\!A^T$, and the matrix V consists of the orthonormalized eigenvectors of $A^+\!A$. The

diagonal elements of Σ are the non-negative square roots of the eigenvalues of $A^{\rm T}A$; they are called <u>singular values</u> or <u>principal values</u> of A . We assume

$$\sigma_1 \ge \sigma_2 > \ldots > \sigma_n \ge 0$$
 .

Thus if rank(A) = r , $\sigma_{r+1} = \sigma_{r+2} = \dots = \sigma_n = 0$. The decomposition (11.1) is called the singular value decomposition (SVD).

Let

$$\tilde{A} = \begin{bmatrix} 0 & A \\ A^T & 0 \end{bmatrix}$$
(11.2)

It can be shown that the non-zero eigenvalues of \tilde{A} always occur in \pm pairs, viz.

$$\lambda_{j}(\tilde{A}) = + \sigma_{j}(A) \quad (j = 1, 2, ..., r) .$$
 (11.3)

12. Applications of the SVD

The singular value decomposition plays an important role in a number of least squares problems, and we will illustrate this with some examples. Throughout this discussion, we use the Euclidean or Frobenius norm of a matrix, viz.

$$\|A\| = (\Sigma |a_{ij}|^2)^{1/2}$$

A) Let u_n be the set of all nxn orthogonal matrices. For an arbitrary nxn real matrix A , determine $Q \in U_n$ such that

$$\|A-Q\| \leq \|A-X\|$$
 for any $X \in \mathbf{u}_n$.

It has been shown by FAN and HOFFMAN [1955] that if

$$A = U\Sigma V^{T}$$
 , then $Q = UV^{T}$

B) An important generalization of problem A occurs in factor analysis. For arbitrary nxn real matrices A and B , determine $Q \in U_n$ such that

$$\|A-BQ\| \leq \|A-BX\|$$
 for any $X \in U_n$.

It has been shown by GREEN [1952] and by SCHÖNEMANN [1966] that if

$$B^{T}A = U\Sigma V^{T}$$
 , then $Q = W^{T}$

C) Let $\mathfrak{M}_{m,n}^{(k)}$ be the set of all mxn matrices of rank k. Assume $A\mathfrak{M}_{m,n}^{(r)}$. Determine $B\mathfrak{M}_{m,n}^{(k)}$ $(k \leq r)$ such that $\||A-B\|| \leq \||A-X\||$ for all $X\mathfrak{M}_{m,n}^{(k)}$

It has been shown by ECKART and YOUNG [1936] that if

$$A = U\Sigma V^{T} , \quad \text{then } B = U\Omega_{k} V^{T}$$
 (12.1)

where

Note that

$$\|A-B\| = \|\Sigma - \Omega_{k}\| = (\sigma_{k+1}^{2} + \ldots + \sigma_{r}^{2})^{1/2} \cdot (12.3)$$

D) An nxm matrix X is said to be the pseudo-inverse of an mxn matrix A if X satisfies the following four properties:

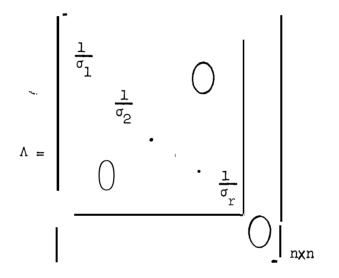
i) AXA = A,
ii) XAX = X,
iii)
$$(AX)^{T} = AX$$
,
iv) $(XA)^{T} = xA$.

We denote the pseudo-inverse by A^+ . We wish to determine A^+ numerically. It can be shown (cf. PENROSE [1955]) that A^+ can always be determined and is unique. It is easy to verify that

$$A + = V \Lambda U^{T}$$
(12.4)

where

.



In recent years there have been a number of algorithms proposed for computing the pseudo-inverse of a matrix. These algorithms usually depend upon a knowledge of the rank of the matrix or upon some suitable chosen parameter. For example in the latter case, if one uses (12.4) to compute the pseudo-inverse, then after one has computed the singular value decomposition numerically it is necessary to determine which of the singular values are zero by testing against some tolerance.

Alternatively, suppose we know that the given matrix A can be represented as

$A = B + \delta B$

where δB is a matrix of perturbations and

 $\|\delta B\| \leq \eta$.

Now, we wish to construct a matrix \hat{B} such that

$$|A - \hat{B}|| \leq \eta$$

and

rank $(\hat{B}) = \min$.

This can be accomplished with the aid of the solution to problem (C).

Let

 $B_{k} = U\Omega_{k}V^{T}$

where Ω_k is defined as in (12.2). Then using (12.3),

if

$$(\sigma_{p+1}^2 + \sigma_{p+2}^2 + \cdots + \sigma_n^2)^{1/2} \leq \eta$$

an.d

$$(\sigma_{p}^{2} + \sigma_{p+1}^{2} + ... + \sigma_{n}^{2})^{1/2} > \eta$$
 .

Since rank(\$) = p by construction,

 $\hat{B} = B_{D}$

$$\hat{B}^{+} = V\Omega_{p}^{+}U^{T}$$

Thus, we take \hat{B}^{\dagger} as our approximation to A^{\dagger} .

• E) Let A be a given matrix, and b be a known vector. Determine \hat{x} so that amongst all x for which $\|\hat{b}-Ax\|_2 = \min$, $\|\hat{x}\|_2 = \min$. It is easy to verify that

13. Calculation of the SVD

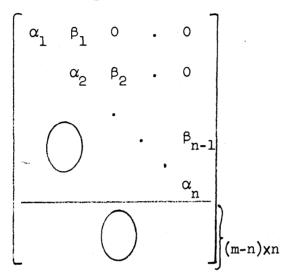
It was shown by GOLUB and KAHAN[1965] that it is possible to construct a sequence of orthogonal matrices

$$\left\{P^{(k)}\right\}_{k=1}^{n}, \left\{Q^{(k)}\right\}_{k=1}^{n-1}$$

via Householder transformation so that

$$P^{(n)}P^{(n-1)} = P^{T}AQ = J$$

and J is an mxn bi-diagonal matrix of the form



The singular values of $\ensuremath{\mathsf{J}}$ are the same as those of $\ensuremath{\mathsf{A}}$. Thus if the singular value decomposition of

$$J = X\Sigma Y^{T}$$

then

$$A = PX\Sigma Y^{T}Q^{T}$$

so that

U = PX, V = QT.

GOLUB[1968] has given an algorithm for computing the SVD of J ; the algorithm is based on the highly'effective QR algorithm of FRANCIS [1961,1962] for computing the eigenvalues.

It is not necessary to compute the complete SVD when a vector **b** is given. Since $\hat{\mathbf{x}} = V\Sigma^{\dagger}U^{T}b$, it is only necessary to compute V, Σ and $U^{T}b$; note, this has a strong "flavor of principal component analysis. An ALGOL" procedure for the SVD has been given by GOLUB and REINSCH[1969].

30

14. Quadratic constraints

We wish to determine $\hat{\mathbf{x}}$ so that

$$\|\mathbf{b}-\mathbf{A}\mathbf{\hat{x}}\|_2 = \min.$$

when

$$\|\hat{\mathbf{x}}\|_2 = \alpha$$

Such problems occur in a number of situations, e.g. in the numerical solution of integral equations of the first kind (cf. PHILLIPS [1962]), and in the solution of non-linear least squares problems (cf. MARQUARDT [1963]). Using Lagrange multipliers, we are led to the equation

$$(\mathbf{A}^{\mathrm{T}}\mathbf{A} - \lambda \mathbf{*}\mathbf{I})\mathbf{\hat{x}} = \mathbf{A}^{\mathrm{T}}\mathbf{b}$$

where the real constant $\lambda \star$ is determined as the smallest root of

$$\alpha^{2} - b^{T} A (A^{T} A - \lambda I)^{-2} A^{T} b = 0 . \qquad (14.1)$$

Using the decomposition $A = U\Sigma V^T$ and $c = U^T b$, equation (14.1) becomes

$$\alpha^2 - c^T \Sigma (\Sigma^2 - \lambda I)^{-2} \Sigma c = 0$$

A combination of bisection and Newton iteration may be used to determine λ^* . It is easily shown that $\lambda^* < \sigma_{\min}^2$ (cf. FORSYTHE and GOLUB[1965]).

It is also possible to determine λ^* as a solution to an eigenvalue • problem using a technique given by FORSYTHE and GOLUB [1965]. Consider the identity

$$\det \begin{bmatrix} X & Y \\ Z & W \end{bmatrix} = \det(X) \det(W - ZX^{-1}Y)$$

which is valid for any partitioned matrix with X and W square and $det(X) \neq 0$. Thus (14.1) is equivalent to the determinantal equation

$$\det \begin{bmatrix} (A^{T}A - \lambda I)^{2} & A^{T}b \\ & \ddots \\ & b^{T}A & \alpha^{2} \end{bmatrix} = 0$$

Now there exists a vector p and a number q such that

L

$$(\mathbf{A}^{\mathrm{T}}\mathbf{A}-\lambda\mathbf{I})^{2}\mathbf{p} + \mathbf{A}^{\mathrm{T}}\mathbf{b}\mathbf{q} = \Theta$$
, $\mathbf{b}^{\mathrm{T}}\mathbf{A}\mathbf{p} + \alpha^{2}\mathbf{q} = 0$.

A simple elimination shows that $\lambda {\boldsymbol \star}$ -must satisfy the determinantal equation

$$det[(A^{T}A - \lambda I)^{2} - \alpha^{-2} A^{T}bb^{T}A] = 0$$
 (14.2)

It is possible to transform (14.2) into a 2nx2n ordinary eigenvalue problem.

Once $\lambda \star$ is determined, the solution \hat{x} can be computed from the SVD of A . Thus,

$$\hat{\mathbf{x}} = \mathbf{V}(\boldsymbol{\Sigma} - \boldsymbol{\lambda} * \boldsymbol{\Sigma}^{-1})^{-1} \mathbf{c} \cdot \boldsymbol{\Sigma}^{-1}$$

References

- Bartels, R. H. [1968], "A numerical investigation of the simplex method", Technical Report No. CS 104, Computer Science Dept., Stanford University, California.
- Bartels, R. H. and G. H. Golub [1969], "The simplex method of linear programming using LU decomposition", <u>Comm</u>. ACM 12, 5, pp. 266-268.
- Björck, A. [1967a], "Solving linear least squares problems by Gram-Schmidt orthogonalization", <u>BIT</u>, 7, pp. 1-21.
- Björck, A. [1967b], "Iterative refinement of linear least squares solution I", BIT, 7, pp. 257-278.
- Björck, A.[1968], "Iterative refinement of linear least squares solutions II", BIT, 8, pp. 8-30.
- Björck, Å. and G. H. Golub [1967], "Iterative refinement of linear least squares solutions by Householder-transformation", BIT, 7, pp. 322-337.
- Businger, P. and G. H. Golub [1965], "Linear least squares solutions by Householder transformations", Num. Math., 7, pp. 269-276.
- Cottle, R. W. [1968], "The principal pivoting method of quadratic programming", <u>Mathematics of the Decision Sciences</u>, <u>Part 1</u>, G. B. Dantzig and A. F. Veinott, eds., pp. 144-162.
- Cottle, R. W. and G. B. Dantzig [1968], "Complementary pivot theory of mathematical programming", in <u>Mathematics of the Decision Sciences</u>, <u>Part 1</u>, G. B. Dantzig and A. F. Veinott, eds., pp. 115-136.
- Eckart, C. and G. Young [1936], "The approximation of one matrix by another of lower rank", <u>Psychometrika</u>, 1, pp. 211-218.
- Fadeev, D. K., V. N. Kublanovskaya, and V. N. Fadeeva[1968], "Sur les systemes linearires algebriques de matrices rectangulaires et malconditionnees", in <u>Programmation en Mathematriques Numeriques</u>, Editions du Centre National de la Recherche Scientifique, Paris VII.
- Fan, K. and A. Hoffman [1955], "Some metric inequalities in the space of matrices'", Proc. Amer. Math. Soc., 6, pp. 111-116.
- Forsythe, G. E. and G. H. Golub [1965], "On the stationary values of a second-degree polynomial on the unit sphere", <u>J</u>. SIAM, 13, pp. 1050-1068.
- Forsythe, G. E. and C. Moler [1967], Computer Solution of Linear Algebraic Systems, Prentice-Hall, Englewood Cliffs, New Jersey.
- Francis, J. [1961,1962], "The QR transformation. A unitary analogue to the LR transformation'", Comput. J., 4, pp. 265-271.
- Golub, G. H. [1965], "Numerical methods for solving linear least squares problems'", Num. Math., 7, pp. 206-216.

Golub, G. H. and W. Kahan [1965], "Calculating the singular values and pseudo-inverse of a matrix", <u>J. SIAM</u>, <u>Numer</u>. Anal. Ser. <u>B</u>, 2, pp. 205-224.

-

- Golub, G. H. and J. Wilkinson [1966], "Iterative refinement of least squares solution", Num. Math., 9, pp. 139-148.
- Golub, G. H. [1968], "Least squares, "singular values and matrix approximations", Aplikace Matematiky, 13, pp. 44-51.
- Golub, G. H. and C. Reinsch [1969], "Singular value decomposition and least squares solution", Technical Report No. CS 131, Computer Science Dept., Stanford University, California.
- Green, B. [1952], "The orthogonal approximation of an oblique structure in factor analysis", <u>Psychometrika</u>, 17, pp. 429-440.
- Hanson, R. and C. Lawson [1968], "Extensions and applications of the Householder algorithm for solving linear least squares problems", Jet Propulsion Laboratory.
- Householder, A.-S. [1958], "Unitary triangularization of a nonsymmetric matrix", <u>J. Assoc.</u> Comput. Mach., 5, pg. 339-342.
- Jordan, T. [1968], "Experiments on error growth associated with some linear least-squares procedures", Math. Comp., 22, pp. 579-588.
- Kahan, W. [1966], "Numerical linear algebra", Canad. Math. Bull., 9, pp. 757-801.
- Kalfon, P., G. Ribiere, and J. Sogno [1968], "Méthode du gradient projete utilisant la triangularisation unitaire". Publication no. FT/11.3.8/AI Centre National de la Recherche Scientifique Institut Blaise Pascal.
- Lanczos, C. [1961], Linear Differential Operators, Van Nostrand, London, Chap. 3.
- Lemke, C. E. [1968], "On complementary pivot theory", Mathematics of the Decision Sciences, Part 1, G. B. Dantzig and A. F. Veinott, eds., pp. 95-114.
- Longley, J. [1967], "An appraisal of least squares problems for the electronic computer from the point of view of the user", JASA, 62, pp. 819-841.
- Marquardt, W. [1963], "An algorithm for least-squares estimation of nonlinear parameters", J. SIAM, 11, pp. 431-441.
- Moler, C. B. [1967], "Iterative refinement in floating point", <u>J</u>. <u>Assoc.</u> <u>Mompute h . </u>, 14, pp. 316-321.
- Penrose, R. [1955], "A generalized inverse for matrices", Proc. Cambridge Philos. Soc., 51, pp. 406-413.

Phillips, P. [1962], "A technique for the numerical solution of certain integral equations of the first kind", <u>J. Assoc. Comput. Mach.</u>, 9, pp. 84-97.

L

- Rice, J. [1966], "Experiments on Gram-Schmidt Orthogonalization", <u>Math.</u> <u>Comp.</u>, 20, pp. 325-328.
- Schönemann, P. [1966], "A generalized solution of the orthogonal procrustes problem", Psychometrika, 31, pp. 1-10.
- Wilkinson, J. H. [1963], Rounding Errors in Algebraic Processes, Prentice-Hall, Englewood Cliffs, New Jersey.
- Wilkinson, J. H. [1967], "The solution of ill-conditional linear equations", Mathematical Methods for Digital Computers, Vol. II, A. Ralston, Ph.D. and H. Wilf, Ph.D., eds., John Wiley, New York, pp. 65-93.