

Numerical Methods for Solving Linear Least Squares Problems*

By
G. GOLUB

Abstract. A common problem in a Computer Laboratory is that of finding linear least squares solutions. These problems arise in a variety of areas and in a variety of contexts. 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 consider stable numerical methods for handling these problems. Our basic tool is a matrix decomposition based on orthogonal Householder transformations.

1. Introduction

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

$$\|\mathbf{b} - A\hat{\mathbf{x}}\| = \min. \tag{1.1}$$

where $\|\dots\|$ indicates the euclidean norm. If $m \geq n$ and $r < n$ then there is no unique solution. Under these conditions, we require simultaneously to (1.1) that

$$\|\hat{\mathbf{x}}\| = \min. \tag{1.2}$$

Condition (1.2) is a very natural one for many statistical and numerical problems.

If $m \geq n$ and $r = n$, then it is well known (cf. [4]) that $\hat{\mathbf{x}}$ satisfies the equation

$$A^T A \mathbf{x} = A^T \mathbf{b}. \tag{1.3}$$

Unfortunately, the matrix $A^T A$ is frequently ill-conditioned [6] and influenced greatly by roundoff errors. The following example of LÄUCHLI [3] illustrates this well. Suppose

$$A = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ \varepsilon & 0 & 0 & 0 & 0 \\ 0 & \varepsilon & 0 & 0 & 0 \\ 0 & 0 & \varepsilon & 0 & 0 \\ 0 & 0 & 0 & \varepsilon & 0 \\ 0 & 0 & 0 & 0 & \varepsilon \end{bmatrix},$$

* Reproduction in Whole or in Part is permitted for any Purpose of the United States government. This report was supported in part by Office of Naval Research Contract Nonr-225(37) (NR 044-11) at Stanford University.

then

$$A^T A = \begin{bmatrix} 1 + \varepsilon^2 & 1 & 1 & 1 & 1 \\ 1 & 1 + \varepsilon^2 & 1 & 1 & 1 \\ 1 & 1 & 1 + \varepsilon^2 & 1 & 1 \\ 1 & 1 & 1 & 1 + \varepsilon^2 & 1 \\ 1 & 1 & 1 & 1 & 1 + \varepsilon^2 \end{bmatrix}. \tag{1.4}$$

Clearly for $\varepsilon \neq 0$, the rank of $A^T A$ is five since the eigenvalues of $A^T A$ are $5 + \varepsilon^2, \varepsilon^2, \varepsilon^2, \varepsilon^2, \varepsilon^2$.

Let us 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 $fl(1.0 + \eta) \equiv 1.0$ where $fl(\dots)$ indicates the floating point computation. Then if $\varepsilon < \frac{\sqrt{\eta}}{2}$, the rank of the computed representation of (1.4) will be one. Consequently, no matter how accurate the linear equation solver, it is impossible to solve the normal equations (1.3).

In [2], HOUSEHOLDER stressed the use of orthogonal transformations for solving linear least squares problems. In this paper, we shall exploit these transformations and show their use in a variety of least squares problems.

2. A Matrix Decomposition

Throughout this section, we shall assume $m \geq n = r$.

Since the euclidean norm of a vector is unitarily invariant,

$$\|b - Ax\| = \|c - QAx\|$$

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

$$QA = R = \begin{pmatrix} \tilde{R} \\ \dots \\ 0 \end{pmatrix}_{(m-n) \times n} \tag{2.1}$$

where \tilde{R} is an upper triangular matrix. Clearly,

$$\hat{x} = \tilde{R}^{-1} c$$

where \tilde{c} is the first n components of c and consequently,

$$\|b - A\hat{x}\| = \left(\sum_{j=m+1}^n c_j^2 \right)^{\frac{1}{2}}.$$

Since \tilde{R} is an upper triangular matrix and $\tilde{R}^T \tilde{R} = A^T A$, $\tilde{R}^T \tilde{R}$ is simply the Choleski decomposition of $A^T A$.

There are a number of ways to achieve the decomposition (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 [2]. Let $A = A^{(1)}$, and let $A^{(2)}, A^{(3)}, \dots, A^{(n+1)}$ be defined as follows:

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

$P^{(k)}$ is a symmetric, orthogonal matrix of the form

$$P^{(k)} = I - 2\mathbf{w}^{(k)}\mathbf{w}^{(k)T}$$

for suitable $\mathbf{w}^{(k)}$ such that $\mathbf{w}^{(k)T}\mathbf{w}^{(k)} = 1$. A derivation of $P^{(k)}$ is given in [9].

In order to simplify the calculations, we redefine $P^{(k)}$ as follows:

$$P^{(k)} = I - \beta_k \mathbf{u}^{(k)}\mathbf{u}^{(k)T}$$

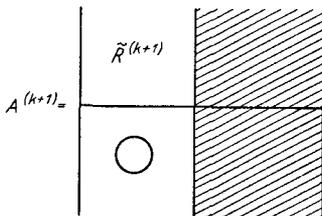
where

$$\begin{aligned} \sigma_k &= \left(\sum_{i=k}^m (a_{i,k}^{(k)})^2 \right)^{\frac{1}{2}}, \\ \beta_k &= [\sigma_k(\sigma_k + |a_{k,k}^{(k)}|)]^{-1}, \\ \mathbf{u}_i^{(k)} &= 0 && \text{for } i < k, \\ \mathbf{u}_k^{(k)} &= \text{sgn}(a_{k,k}^{(k)}) (\sigma_k + |a_{k,k}^{(k)}|), \\ \mathbf{u}_i^{(k)} &= a_{i,k}^{(k)} && \text{for } i > k. \end{aligned}$$

Thus

$$A^{(k+1)} = A^{(k)} - \mathbf{u}^{(k)}(\beta_k \mathbf{u}^{(k)T} A^{(k)}).$$

After $P^{(k)}$ has been applied to $A^{(k)}$, $A^{(k+1)}$ appears as follows:



where $\tilde{R}^{(k+1)}$ is a $k \times k$ upper triangular matrix which is unchanged by subsequent transformations. Now $a_{k,k}^{(k+1)} = -(\text{sgn } a_{k,k}^{(k)})\sigma_k$ so that the rank of A is less than n if $\sigma_k = 0$. Clearly,

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

and

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

although one need not compute Q explicitly.

3. The Practical Procedure

WILKINSON [10] has shown that the Choleski decomposition is stable for a positive definite matrix even if no interchanges of rows and columns are performed. Since we are in effect performing a Choleski 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 $|a_{k,k}^{(k+1)}|$. This is equivalent to searching for the maximum diagonal element in the Choleski decomposition of $A^T A$. Let

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

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

$$s_j^{(k+1)} = s_j^{(k)} - (a_{k,j}^{(k+1)})^2 \quad (j = k+1, \dots, m)$$

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. Although it is possible to compute σ_k directly from the $s_j^{(k)}$'s, it is best to compute σ_k at each stage using double precision inner products to ensure maximal accuracy.

The strategy described above is most appropriate when one has a sequence of vectors $\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_p$ for which one desires a least squares estimate. In many problems, there is one vector \mathbf{b} and one wishes to express it in as few columns of A as possible. This is the stagewise multiple regression problem. We cannot solve this problem, but we shall show how one can choose that column of $A^{(k)}$ for which the sum of squares of residuals is maximally reduced at the k^{th} stage.

Let $\mathbf{c}^{(1)} = \mathbf{b}$ and $\mathbf{c}^{(k+1)} = P^{(k)} \mathbf{c}^{(k)}$. Now $\tilde{R}^{(k)} \hat{\mathbf{x}}^{(k-1)} = \tilde{\mathbf{c}}^{(k)}$ where $\hat{\mathbf{x}}^{(k-1)}$ is the least squares estimate based on $(k-1)$ columns of A , and $\tilde{\mathbf{c}}^{(k)}$ is the first $(k-1)$ elements of $\mathbf{c}^{(k)}$, and consequently

$$\|\mathbf{c}^{(k)} - \tilde{R}^{(k)} \hat{\mathbf{x}}^{(k-1)}\| = \left(\sum_{j=k}^m (c_j^{(k)})^2 \right)^{\frac{1}{2}}.$$

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

$$t_j^{(k)} = \sum_{i=k}^m a_{i,j}^{(k)} c_i^{(k)} \quad \text{for } j = k, k+1, \dots, m.$$

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

$$t_j^{(k+1)} = t_j^{(k)} - a_{k,j}^{(k+1)} c_k^{(k+1)}.$$

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

Once the solution to the equations has been obtained then it is possible to obtain an improved solution by a simple iterative technique. This technique, however, requires that the orthogonal transformations be saved during their application. The best method for storing the transformation is to store the elements of $\mathbf{u}^{(k)}$ below the diagonal of the k^{th} column of $A^{(k+1)}$.

Let $\bar{\mathbf{x}}$ be the initial solution obtained, and let $\hat{\mathbf{x}} = \bar{\mathbf{x}} + \mathbf{e}$. Then

$$\|\mathbf{b} - A\hat{\mathbf{x}}\| = \|\mathbf{r} - A\mathbf{e}\|$$

where

$$\mathbf{r} = \mathbf{b} - A\bar{\mathbf{x}}, \quad \text{the residual vector.}$$

Thus the correction vector \mathbf{e} is itself the solution to a linear least squares problem. Once A has been decomposed then it is a fairly simple matter to compute \mathbf{r} and solve for \mathbf{e} . Since \mathbf{e} critically depends upon the residual vector, the components of \mathbf{r} should be computed using double precision inner products and then rounded to single precision accuracy. Naturally, one should continue to iterate as long as improved estimates of $\hat{\mathbf{x}}$ are obtained.

The above iteration technique will converge only if the initial approximation to $\hat{\mathbf{x}}$ is sufficiently accurate. Let

$$\mathbf{x}^{(q+1)} = \mathbf{x}^{(q)} + \mathbf{e}^{(q)} \quad (q=0, 1, \dots)$$

with $\mathbf{x}^{(0)} = \mathbf{0}$. Then one should iterate only if $\|\mathbf{e}^{(1)}\|/\|\mathbf{x}^{(1)}\| \leq c$ where $c < \frac{1}{2}$, i.e. "at least one bit of the initial solution is correct"; otherwise there is little likelihood that the iterative method will converge. Since convergence tends to be linear, one should terminate the procedure as soon as

$$\|\mathbf{e}^{(k+1)}\| > c\|\mathbf{e}^{(k)}\| \quad \text{or} \quad \|\mathbf{e}^{(k)}\| < \eta\|\mathbf{x}^{(1)}\|$$

where η is the maximum positive number such that $f!(1+\eta) \equiv 1$.

4. A Numerical Example

In Table 1, we give the results of an extensive calculation. The matrix consists of the first 5 columns of the inverse of the 6×6 Hilbert matrix. The calculations were performed in single precision arithmetic. The columns were chosen so that the diagonal elements were maximized at each stage. The iteration procedure was terminated as soon as $\|\mathbf{e}^{(k+1)}\| > 0.25\|\mathbf{e}^{(k)}\|$. Three iterations were performed but since $\|\mathbf{e}^{(2)}\| > 0.25\|\mathbf{e}^{(1)}\|$, $\mathbf{x}^{(2)}$ was taken to be the correct solution.

In Table 2, we show the results of using double precision inner products on the same problem. Note that the first iterate in Table 1 is approximately as accurate as the first iterate in Table 2. The double precision inner product routine converged to a solution for which all figures were accurate. The normal equations were formed using double precision inner products but even with a very accurate linear equation solver described by MCKEEMAN [5] no solution could be obtained.

5. An Iterative Scheme

For many problems, even with the use of orthogonal transformations it may be impossible to obtain an accurate solution. Or, the rank of A may truly be less than n . In this section, we give an algorithm for finding the least squares solution even if $A^T A$ is singular.

In [7], RILEY suggested the following algorithm for solving linear least squares problems for $r=n$. Let $\mathbf{x}^{(0)}$ be an arbitrary vector, then solve

$$(A^T A + \alpha I)\mathbf{x}^{(q+1)} = A^T \mathbf{b} + \alpha \mathbf{x}^{(q)}. \quad (5.1)$$

Table 1

		A					B				
		5	4	3	2	1	5	4	3	2	1
	Pivot										
	R										
	Iteration										
1	X	3.60000000 ₁₀ +01	-6.30000000 ₁₀ +02	3.36000000 ₁₀ +03	-7.56000000 ₁₀ +03	7.56000000 ₁₀ +03	4.63000000 ₁₀ +02				
	R	-6.30000000 ₁₀ +02	1.47000000 ₁₀ +04	-8.82000000 ₁₀ +04	2.11680000 ₁₀ +05	-2.20500000 ₁₀ +05	-1.38600000 ₁₀ +04				
	E	3.36000000 ₁₀ +03	-8.82000000 ₁₀ +04	5.64480000 ₁₀ +05	-1.41120000 ₁₀ +06	1.51200000 ₁₀ +06	9.70200000 ₁₀ +04				
	X	-7.56000000 ₁₀ +03	2.11680000 ₁₀ +05	-1.41120000 ₁₀ +06	3.62880000 ₁₀ +06	-3.96900000 ₁₀ +06	-2.58720000 ₁₀ +05				
	R	7.56000000 ₁₀ +03	-2.20500000 ₁₀ +05	1.51200000 ₁₀ +06	-3.96900000 ₁₀ +06	4.41000000 ₁₀ +06	2.91060000 ₁₀ +05				
	E	-2.77200000 ₁₀ +03	8.31600000 ₁₀ +04	-5.82120000 ₁₀ +05	1.55232000 ₁₀ +06	-1.74636000 ₁₀ +06	-1.16424000 ₁₀ +05				
	X	-6.3706872199 ₁₀ +06	5.7760446368 ₁₀ +06	-2.2224495814 ₁₀ +06	3.2875491420 ₁₀ +05	-1.1522407952 ₁₀ +04					
	R	-6.7135626821 ₁₀ +04	6.0308471363 ₁₀ +04	6.0308471363 ₁₀ +04	-1.6102552917 ₁₀ +04	9.4910780925 ₁₀ +02					
	E										
	X										
	R										
	E										
2	X	9.9999912362 ₁₀ -01	4.9999972493 ₁₀ -01	3.3333322021 ₁₀ -01	2.4999995198 ₁₀ -01	1.9999998332 ₁₀ -01					
	R	1.4528632164 ₁₀ -06	-4.7683715820 ₁₀ -07	0.00000000 ₁₀ +00	0.00000000 ₁₀ +00	0.00000000 ₁₀ +00	0.00000000 ₁₀ +00				
	E	5.8284221383 ₁₀ -07	1.7292551586 ₁₀ -07	6.938181279 ₁₀ -08	2.9087566395 ₁₀ -08	1.0038052093 ₁₀ -08					
	X	9.9999970646 ₁₀ -01	4.9999989786 ₁₀ -01	3.3333328959 ₁₀ -01	2.4999998107 ₁₀ -01	1.9999999336 ₁₀ -01					
	R	3.2037496567 ₁₀ -07	4.7683715820 ₁₀ -07	3.8146972656 ₁₀ -06	-3.8146972656 ₁₀ -06	3.8146972656 ₁₀ -06	0.00000000 ₁₀ +00				
	E	-3.7819874906 ₁₀ -07	-1.1487774268 ₁₀ -07	-4.6569725113 ₁₀ -08	-1.9660767363 ₁₀ -08	-6.8227426601 ₁₀ -09					

Table 2

		Pivot						
		5	4	3	2	1		
		<i>R</i>						
		-6.3706872199 ₁₀ +06	5.7760446368 ₁₀ +06	-2.2224495814 ₁₀ +06	3.2875491420 ₁₀ +05	-1.1522407952 ₁₀ +04		
			-6.7135626821 ₁₀ +04	6.0308471363 ₁₀ +04	-1.6102552917 ₁₀ +04	9.4910780925 ₁₀ +02		
				-1.4425503500 ₁₀ +03	9.4707042643 ₁₀ +02	-1.0982644637 ₁₀ +02		
					4.6175131642 ₁₀ +01	-1.4949038077 ₁₀ +01		
						2.0095559061 ₁₀ +00		
		<i>Iteration</i>						
1	<i>X</i>	9.9999912347 ₁₀ -01	4.9999972435 ₁₀ -01	3.3333321985 ₁₀ -01	2.4999995180 ₁₀ -01	1.9999998325 ₁₀ -01		
	<i>R</i>	1.4480865502 ₁₀ -06	-8.8621163741 ₁₀ -08	9.3205017038 ₁₀ -08	-2.6189081836 ₁₀ -06	2.1314917831 ₁₀ -07	3.2270327211 ₁₀ -07	
	<i>E</i>	8.7653116638 ₁₀ -07	2.7564789022 ₁₀ -07	1.1348434536 ₁₀ -07	4.8201409291 ₁₀ -08	1.6752167940 ₁₀ -08		
2	<i>X</i>	1.0000000000 ₁₀ +00	5.0000000000 ₁₀ -01	3.3333333333 ₁₀ -01	2.5000000000 ₁₀ -01	2.0000000000 ₁₀ -01		
	<i>R</i>	-7.5378920883 ₁₀ -09	2.1391315386 ₁₀ -07	-1.4423858374 ₁₀ -06	3.7434801925 ₁₀ -06	-4.1254679672 ₁₀ -06	1.6236008378 ₁₀ -06	
	<i>E</i>	9.4773720382 ₁₀ -18	2.6923536909 ₁₀ -18	-6.0632875784 ₁₀ -13	4.25556290006 ₁₀ -19	-7.2759561775 ₁₀ -13		
3	<i>X</i>	1.0000000000 ₁₀ +00	5.0000000000 ₁₀ -01	3.3333333333 ₁₀ -01	2.5000000000 ₁₀ -01	2.0000000000 ₁₀ -01		
	<i>R</i>	-7.5378920883 ₁₀ -09	2.1391315386 ₁₀ -07	-1.4423858374 ₁₀ -06	3.7434801925 ₁₀ -06	-4.1254679672 ₁₀ -06	1.6236008378 ₁₀ -06	
	<i>E</i>	9.4773720382 ₁₀ -18	2.6923536909 ₁₀ -18	-6.0632875784 ₁₀ -13	4.25556290006 ₁₀ -19	-7.2759561775 ₁₀ -13		

The sequence $\mathbf{x}^{(q)}$ converges to $\hat{\mathbf{x}}$ if $\alpha > 0$ since the spectral radius of $\alpha(A^T A + \alpha I)^{-1}$ is less than 1. Again we may implement this algorithm more effectively by the use of orthogonal transformations.

First, let us note that (5.1) is equivalent to the following:

$$\mathbf{r}^{(q)} = \mathbf{b} - A\mathbf{x}^{(q)}, \tag{5.2a}$$

$$(A^T A + \alpha I)\mathbf{e}^{(q)} = A^T \mathbf{r}^{(q)}, \tag{5.2b}$$

$$\mathbf{x}^{(q+1)} = \mathbf{x}^{(q)} + \mathbf{e}^{(q)}. \tag{5.2c}$$

The vector $\mathbf{e}^{(q)}$ is itself the solution of a linear least squares problem since $\mathbf{e}^{(q)}$ minimize $\|\mathbf{d}^{(q)} - C\mathbf{e}^{(q)}\|$ where

$$C = \begin{pmatrix} A \\ \dots \\ \sqrt{\alpha} I \end{pmatrix}, \quad \mathbf{d}^{(q)} = \begin{pmatrix} \mathbf{r}^{(q)} \\ \dots \\ \mathbf{0} \end{pmatrix}.$$

Thus the numerical procedure should go as follows. Decompose C by the methods described in Section 2 so that

$$PC = S = \begin{pmatrix} \tilde{S} \\ \dots \\ O \end{pmatrix}$$

where $P^T P = I$ and \tilde{S} is an upper triangular matrix. Then let $\mathbf{x}^{(0)} = \mathbf{0}$,

$$\begin{aligned} \tilde{S}\mathbf{e}^{(q)} &= \tilde{\mathbf{f}}^{(q)}, \\ \mathbf{x}^{(q+1)} &= \mathbf{x}^{(q)} + \mathbf{e}^{(q)} \end{aligned}$$

and $\tilde{\mathbf{f}}^{(q)}$ is the vector whose components are the first n components of $P\mathbf{d}^{(q)}$. We choose $\mathbf{x}^{(0)} = \mathbf{0}$ since otherwise there is no assurance that $\mathbf{x}^{(q)}$ will converge to $\hat{\mathbf{x}}$.

Now going back to the original process (5.1),

$$\mathbf{x}^{(q+1)} = G\mathbf{x}^{(q)} + \mathbf{h} \tag{5.3}$$

where

$$G = \alpha(A^T A + \alpha I)^{-1} \quad \text{and} \quad \mathbf{h} = (A^T A + \alpha I)^{-1} A^T \mathbf{b}.$$

Thus

$$\mathbf{x}^{(q+1)} = (G^q + G^{q-1} + \dots + I)\mathbf{h}. \tag{5.4}$$

It is well known (cf. [6]) that A may be written as

$$A = U \Sigma V^T$$

where Σ is an $m \times n$ matrix with the singular values σ_j on the diagonal and zeros elsewhere, and U and V are the matrices of eigenvectors of AA^T and $A^T A$, respectively. Then

$$A^T \mathbf{b} = V \Sigma^T U^T \mathbf{b} = \beta_1 \sigma_1 \mathbf{v}_1 + \beta_2 \sigma_2 \mathbf{v}_2 + \dots + \beta_r \sigma_r \mathbf{v}_r,$$

where $\boldsymbol{\beta} = U^T \mathbf{b}$, and r is the rank of A . Then from (5.4) we see that

$$\mathbf{x}^{(q)} = \gamma_1^{(q)} \mathbf{v}_1 + \dots + \gamma_r^{(q)} \mathbf{v}_r,$$

where

$$\gamma_j^{(q)} = \left[1 - \left(\frac{\alpha}{\alpha + \sigma_j^2} \right)^q \right] \frac{\beta_j}{\sigma_j} \quad (j = 1, 2, \dots, r).$$

Thus as $q \rightarrow \infty$

$$\mathbf{x}^{(q)} \rightarrow \frac{\beta_1}{\sigma_1} \mathbf{v}_1 + \dots + \frac{\beta_r}{\sigma_r} \mathbf{v}_r = \hat{\mathbf{x}}.$$

The choice of α will greatly affect the rate of convergence of the iterative method, and thus one must choose α with great care. If α is too small then the equations will remain ill-conditioned. If δ is a lower bound of the smallest non-zero singular value, then α should be chosen so that

$$\frac{\alpha}{\alpha + \delta^2} < 0.1, \quad \text{say.}$$

This means at each stage, there will be at least one more place of accuracy in the solution. There are a number of methods for accelerating the convergence of (5.4) (cf. [I]).

It is easy to see that

$$\mathbf{e}^{(q+1)} = G\mathbf{e}^{(q)} = \alpha(A^T A + \alpha I)^{-1} \mathbf{e}^{(q)}.$$

Since $\mathbf{e}^{(q)}$ lies in the space spanned by $\mathbf{v}_1, \dots, \mathbf{v}_r$, it follows immediately that

$$\|\mathbf{e}^{(q+1)}\| \leq \frac{\alpha}{\alpha + \sigma_r^2} \|\mathbf{e}^{(q)}\| < \|\mathbf{e}^{(q)}\|.$$

Thus a good termination procedure is to stop iterating as soon as $\|\mathbf{e}^{(q)}\|$ increases or does not change.

6. Statistical Calculations

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

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

Since

$$A^T A = \tilde{R}^T \tilde{R}, \quad (A^T A)^{-1} = \tilde{R}^{-1} \tilde{R}^{-T}.$$

The inverse of \tilde{R} can be readily obtained since \tilde{R} is an upper triangular matrix. WAUGH and DWYER [8] have noted that it is possible to calculate $(A^T A)^{-1}$ directly from \tilde{R} . Let

$$(A^T A)^{-1} = X = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n).$$

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 $\mathbf{x}_n, \mathbf{x}_{n-1}, \dots, \mathbf{x}_1$. The number of operations are 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 statistical 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, \mathbf{b}_2 correspond to the additional observations. Then the up-dated least squares solution can be obtained directly from

$$A = \begin{pmatrix} \tilde{R}_1 \\ \dots \\ A_2 \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} \tilde{c}_1 \\ \dots \\ \mathbf{b}_2 \end{pmatrix}.$$

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.

7. Least Squares Problems with Constraints

Frequently, one wishes to determine $\hat{\mathbf{x}}$ so that $\|\mathbf{b} - A\hat{\mathbf{x}}\|$ is minimized subject to the condition that $H\hat{\mathbf{x}} = \mathbf{g}$ where H is a $p \times n$ matrix of rank p . One can, of course, eliminate p of the columns of A by Gaussian elimination after a $p \times p$ submatrix of H has been determined and then solve the resulting normal equations. This, unfortunately, would not be a numerically stable scheme since no row interchanges between A and H would be permitted.

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

$$\left(\begin{array}{c|c} A^T A & H^T \\ \hline H & 0 \end{array} \right) \begin{pmatrix} \hat{\mathbf{x}} \\ \boldsymbol{\lambda} \end{pmatrix} = \begin{pmatrix} A^T \mathbf{b} \\ \mathbf{g} \end{pmatrix}$$

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

$$H(A^T A)^{-1} H^T \boldsymbol{\lambda} = H\mathbf{z} - \mathbf{g}$$

where

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

Note \mathbf{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 \mathbf{z} , of course, should be computed by the orthogonalization procedures discussed earlier.

Since $A^T A = \tilde{R}^T \tilde{R}$, $H(A^T A)^{-1} H^T = W^T W$ where $W = \tilde{R}^{-T} H^T$. After W is computed, it should be reduced to a $p \times p$ upper triangular matrix K by orthogonalization which is the Choleski decomposition of $W^T W$. The matrix equation

$$K^T K \boldsymbol{\lambda} = H\mathbf{z} - \mathbf{g}$$

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

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

where $(A^T A)^{-1} H \boldsymbol{\lambda}$ can be easily computed by using \tilde{R}^{-1} .

Acknowledgements. The author is very pleased to acknowledge the programming efforts of Mr. PETER BUSINGER and Mr. ALAN MERTEN. He also wishes to thank Professor THOMAS ROBERTSON for his critical remarks.

References

[1] GOLUB, G. H., and R. S. VARGA: Chebyshev Semi-Iterative Methods, Successive Over-Relaxation Iterative Methods, and Second Order Richardson Iterative Method. *Numer. Math.* **3**, 147–168 (1961).

[2] HOUSEHOLDER, A. S.: Unitary Triangularization of a Nonsymmetric Matrix. *J. Assoc. Comput. Mach.* **5**, 339–342 (1958).

[3] LÄUCHLI, P.: Jordan-Elimination und Ausgleichung nach kleinsten Quadraten. *Numer. Math.* **3**, 226–240 (1961).

[4] LINNIK, Y.: Method of Least Squares and Principles of the Theory of Observations. Translated from Russian by R. C. ELANDT. New York: Pergamon Press 1961.

[5] MCKEEMAN, W. M.: Crout with Equilibration and Iteration. Algorithm 135. *Comm. Assoc. Comput. Mach.* **5**, 553–555 (1962).

[6] OSBORNE, E. E.: On Least Squares Solutions of Linear Equations. *J. Assoc. Comput. Mach.* **8**, 628–636 (1961).

[7] RILEY, J. D.: Solving Systems of Linear Equations with a Positive Definite, Symmetric, but Possibly Ill-Conditioned Matrix. *Math. Tables Aids Comput.* **9**, 96–101 (1956).

[8] WAUGH, F. V., and P. S. DWYER: Compact Computation of the Inverse of a Matrix. *Ann. Math. Stat.* **16**, 259–271 (1945).

[9] WILKINSON, J. H.: HOUSEHOLDERS Method for the Solution of the Algebraic Eigenproblem. *Comput. J.* **3**, 23–27 (1960).

[10] — Error Analysis of Direct Methods of Matrix Inversion. *J. Assoc. Comput. Mach.* **8**, 281–330 (1961).

Stanford University
 Computation Center
 Stanford, California 94305 (USA)

(Received September 24, 1964)