

Gene Howard Golub

Numerical methods for solving linear least squares problems

*Aplikace matematiky*, Vol. 10 (1965), No. 3, 213–216

Persistent URL: <http://dml.cz/dmlcz/102951>

## Terms of use:

© Institute of Mathematics AS CR, 1965

Institute of Mathematics of the Czech Academy of Sciences provides access to digitized documents strictly for personal use. Each copy of any part of this document must contain these *Terms of use*.



This document has been digitized, optimized for electronic delivery and stamped with digital signature within the project *DML-CZ: The Czech Digital Mathematics Library* <http://dml.cz>

NUMERICAL METHODS FOR SOLVING LINEAR LEAST  
SQUARES PROBLEMS<sup>1)</sup>

GENE H. GOLUB

(to topic d)

One of the problems which arises most frequently in a Computer Laboratory is that of finding linear least squares solutions. These problems arise in a variety of contexts, e.g., statistical applications, numerical solution of integral equations of the first kind, etc. Linear least squares problems are particularly difficult to solve because they frequently involve large quantities of data, and they are ill-conditioned by their nature.

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

$$(1) \quad \|\mathbf{b} - A\hat{\mathbf{x}}\| = \min .$$

where  $\|\dots\|$  indicates the euclidean norm. It is well known (cf. [4]) that  $\hat{\mathbf{x}}$  satisfies the equation

$$(2) \quad A^T A \mathbf{x} = A^T \mathbf{b} .$$

Unfortunately, the matrix  $A^T A$  is frequently ill-conditioned [4] and influenced greatly by roundoff errors. The following example of LÄUCHLI [2] illustrates this well. Suppose

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

<sup>1)</sup> This report was supported in part by Office of Naval Research Contract Nonr — 225(37) (NR 044-11) at Stanford University.

then

$$(3) \quad A^T A = \begin{bmatrix} 1 + \varepsilon^2, & 1, & 1, & 1, & 1, \\ 1, & 1 + \varepsilon^2, & 1, & 1, & 1, \\ 1, & 1, & 1 + \varepsilon^2, & 1, & 1, \\ 1, & 1, & 1, & 1 + \varepsilon^2, & 1, \\ 1, & 1, & 1, & 1, & 1 + \varepsilon^2 \end{bmatrix}.$$

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  $\eta$  be the largest number on the computer such that  $\text{fl}(1.0 + \eta) \equiv 1.0$  where  $\text{fl}(\dots)$  indicates the floating point computation. Then if  $\varepsilon < \sqrt{\eta}/2$ , the rank of the computed representation of (3) will be one. Consequently, no matter how accurate the linear equation solver, it is impossible to solve the normal equations (2).

In [1], HOUSEHOLDER stressed the use of orthogonal transformations for solving linear least squares problems. In this paper, we shall exploit these transformations.

Since the euclidean norm of a vector is unitarily invariant,

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

where  $\mathbf{c} = Q\mathbf{b}$  and  $Q$  is an orthogonal matrix. We choose  $Q$  so that

$$(4) \quad Q A = R = \begin{pmatrix} \tilde{R} \\ \mathbf{0} \end{pmatrix}_{(m-n) \times n}$$

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

Clearly,

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

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

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

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 (4); 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 (4) is via Householder transformations [1]. Let  $A = A^{(1)}$ , and let  $A^{(2)}$ ,  $A^{(3)}$ , ...,  $A^{(n+1)}$  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 [5]. 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  $\sigma_k = \left(\sum_{i=k}^m (a_{i,k}^{(k)})^2\right)^{\frac{1}{2}}$ ,  $\beta_k = [\sigma_k(\sigma_k + |a_{k,k}^{(k)}|)]^{-1}$ , and

$$u_i^{(k)} = 0 \text{ for } i < k, \quad u_k^{(k)} = \text{sgn}(a_{k,k}^{(k)}) (\sigma_k + |a_{k,k}^{(k)}|), \quad u_i^{(k)} = a_{i,k}^{(k)} \text{ for } i > k.$$

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:

$$A^{(k+1)} = \begin{array}{c|c} \tilde{R}^{(k+1)} & \text{diagonal lines} \\ \hline & \text{diagonal lines} \end{array}$$

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.

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  $\mathbf{x}$  is sufficiently accurate. Let

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

Then if  $\|\mathbf{e}^{(1)}\|/\|\mathbf{x}^{(1)}\| > c$  and if  $c < \frac{1}{2}$ , i.e., "at least one bit of the initial solution is correct", one should not iterate since there is little likelihood that the iterative method will converge. Since convergence tends to be linear, one should terminate the procedure as soon as  $\|\mathbf{e}^{(k+1)}\| > c\|\mathbf{e}^{(k)}\|$ .

#### References

- [1] *A. S. Householder*: Unitary Triangularization of a Nonsymmetric Matrix, *J. Assoc. Comput. Mach.*, Vol. 5 (1958), pp. 339—342.
- [2] *P. Lauchli*: Jordan-Elimination und Ausgleichung nach kleinsten Quadraten, *Numer. Math.*, Vol. 3 (1961), pp. 226—240.
- [3] *Y. Linnik*: Method of Least Squares and Principles of the Theory of Observations, translated from Russian by R. C. Elandt, Pergamon Press, New York, 1961.
- [4] *E. E. Osborne*: On Least Squares Solutions of Linear Equations, *J. Assoc. Comput. Mach.*, Vol. 8 (1961), pp. 628—636.
- [5] *J. H. Wilkinson*: Householder's Method for the Solution of the Algebraic Eigenproblem, *Comput. J.*, Vol. 3 (1960), pp. 23—27.

*Gene H. Golub*, Computation Center, Stanford University, Stanford, California, U.S.A.