

Journal

Numerische Mathematik

in: Numerische Mathematik I Journal

473 page(s)

Terms and Conditions

The Göttingen State and University Library provides access to digitized documents strictly for noncommercial educational, research and private purposes and makes no warranty with regard to their use for other purposes. Some of our collections are protected by copyright. Publication and/or broadcast in any form (including electronic) requires prior written permission from the Goettingen State- and University Library. Each copy of any part of this document must contain there Terms and Conditions. With the usage of the library's online system to access or download a digitized document you accept there Terms and Conditions. Reproductions of material on the web site may not be made for or donated to other repositories, nor may be further reproduced without written permission from the Goettingen State- and University Library

For reproduction requests and permissions, please contact us. If citing materials, please give proper attribution of the source.

Contact:

Niedersächsische Staats- und Universitätsbibliothek
Digitalisierungszentrum
37070 Goettingen
Germany
Email: gdz@sub.uni-goettingen.de

Purchase a CD-ROM

The Goettingen State and University Library offers CD-ROMs containing whole volumes / monographs in PDF for Adobe Acrobat. The PDF-version contains the table of contents as bookmarks, which allows easy navigation in the document. For availability and pricing, please contact:

Niedersaechisische Staats- und Universitaetsbibliothek Goettingen - Digitalisierungszentrum
37070 Goettingen, Germany, Email: gdz@sub.uni-goettingen.de

Note on the Iterative Refinement of Least Squares Solution

G. H. GOLUB and J. H. WILKINSON

Received June 12, 1966

1. Introduction

In the standard least squares problem the vector x is required for which $\|b - Ax\|_2$ is a minimum. Here A is an $m \times n$ matrix ($m \geq n$) and b and x are vectors. In a recent paper GOLUB [4] described a method for solving this problem using orthogonal transformations of the type $P^{(r)} = I - 2w^{(r)}(w^{(r)})^T$ where

$$(w^{(r)})^T = (0, \dots, 0, w_{r+1}^{(r)}, \dots, w_n^{(r)}), \quad \|w^{(r)}\|_2 = 1. \quad (1)$$

Using pre-multiplications by $P^{(0)}, P^{(1)}, \dots, P^{(n-1)}$ (the last matrix is not required if $m = n$) the original problem is reduced to that of minimising $\|c - Bx\|_2$ where

$$B = P^{(n-1)} \dots P^{(1)} P^{(0)} A = \begin{bmatrix} U \\ \dots \\ 0 \end{bmatrix}_{m-n}^n \quad (2)$$

$$c = P^{(n-1)} \dots P^{(1)} P^{(0)} b = \begin{bmatrix} p \\ \dots \\ q \end{bmatrix}_{m-n}^n \quad (3)$$

Obviously the solution is given by

$$Ux = p \quad (4)$$

since writing $P = P^{(n-1)} \dots P^{(1)} P^{(0)}$ we have

$$\|b - Ax\|_2 = \|P(b - Ax)\|_2 = \left\| \begin{bmatrix} p - Ux \\ \dots \\ q \end{bmatrix} \right\|_2. \quad (5)$$

When $m = n$ the method reduces to HOUSEHOLDER'S triangularization [5] for solving linear equations; the vector q then has no components and provided A (and therefore U) is non-singular, x is the unique solution.

In practice, whether or not $m = n$, rounding errors contaminate the solution. If $x^{(1)}$ is the computed solution and $\delta^{(1)}$ is determined so that $\|(b - Ax^{(1)}) - A\delta^{(1)}\|_2$ is a minimum, then since this implies that $\|b - A(x^{(1)} + \delta^{(1)})\|_2$ is a minimum, $x^{(1)} + \delta^{(1)}$ is the least squares solution of the original problem. Hence if $r^{(1)}$ is the residual vector defined by

$$r^{(1)} = b - Ax^{(1)} \quad (6)$$

the required correction is the least squares solution of the system with matrix A and right-hand side $r^{(1)}$. The computed $P^{(r)}$ and U may be used to derive an approximation to this correction and the process may be continued to give suc-

cessive refinement of the solution. BUSINGER and GOLUB [2] described this process in detail and have published an ALGOL procedure embodying it.

Although the process is formally the same whether or not $m=n$, when $m>n$ the behaviour of the refinement procedure exhibits some features which are not present when $m=n$. These features are the main topic of this note.

2. The Basic Error Analysis

The error analysis depends to some extent on the details of the arithmetic operations used. Since it is *essential* that inner-products are accumulated to double-precision in the calculation of the residual vectors (otherwise no refinement is achieved) we shall assume that this is done throughout so that only one inner-product procedure is used.

WILKINSON [8] has given a fairly general error analysis of orthogonal transformations of this type based on the assumptions that floating-point arithmetic is used with a mantissa of t binary digits and that inner-products are accumulated wherever possible. This analysis applies directly to the above algorithm and although the bounds apply to a specific rounding procedure the differences arising from alternative procedures are of minor importance. The results may be summarised as follows:

The computed matrices \bar{B} and \bar{b} , which are of the forms

$$\bar{B} = \begin{bmatrix} U \\ \dots \\ 0 \end{bmatrix}, \quad \bar{b} = \begin{bmatrix} \bar{p} \\ \dots \\ \bar{q} \end{bmatrix} \quad (7)$$

are such that there exists an exactly orthogonal matrix Q (*not* the matrix corresponding to exact computation throughout) such that

$$\bar{B} = Q(A + E), \quad \bar{b} = Q(b + e). \quad (8)$$

Here E and e satisfy the bounds

$$\|E\|_F \leq 12.5 n 2^{-t} \|A\|_F, \quad \|e\|_2 \leq 12.5 n 2^{-t} \|b\|_2, \quad (9)$$

where the suffix F denotes the Frobenius norm (i.e. $\|A\|_F = (\sum \sum |a_{ij}|^2)^{1/2}$).

Further errors are made in the back-substitution $Ux = \bar{p}$ and it has been shown by WILKINSON [7, 8] that the computed solution $x^{(1)}$ satisfies the equation

$$(U + F^{(1)}) x^{(1)} = b \quad (10)$$

where

$$\|F^{(1)}\|_F \leq (2^{-t} + \frac{3}{2} n 2^{-2t}) \|U\|_F. \quad (11)$$

If we define $G^{(1)}$ by the relation $G^{(1)} = Q^T F^{(1)}$ then

$$\|G^{(1)}\|_F = \|F^{(1)}\|_F \quad (12)$$

and the computed $x^{(1)}$ is therefore the exact least squares solution of the system

$$(A + E + G^{(1)})b + e = (A + H^{(1)})b + e \quad (\text{say}) \quad (13)$$

where

$$\|H^{(1)}\|_F = \|E + G^{(1)}\|_F \leq \|E\|_F + \|F^{(1)}\|_F. \quad (14)$$

Having performed the orthogonal reduction of A we can treat any right-hand side. The matrix E is independent of the right-hand side, but this is not true of $G^{(1)}$ which is why we have included the superscript. However, the important point for the subsequent analysis is that we have the uniform bound (11) for $\|F^{(1)}\|_F$ and therefore for $\|H^{(1)}\|_F$.

3. The Linear Equation Case

When $m=n$ the analysis of the previous section implies that the computed solution satisfies the equation

$$(A + H^{(1)})x^{(1)} = b + e. \quad (15)$$

Obviously if A is too ill-conditioned $A + H^{(1)}$ could even be singular. We show that if

$$27n^{\frac{1}{2}}2^{-t}\|A\|_2\|A^{-1}\|_2 < 2^{-p} \quad (p \geq 0) \quad (16)$$

then

$$\|x - x^{(1)}\|_2 / \|x\|_2 < 2^{-p} / (1 - 2^{-p-1}) \quad (17)$$

where $x = A^{-1}b$ is the true solution.

From Eq. (15)

$$x^{(1)} = (A + H^{(1)})^{-1}(b + e) \quad (18)$$

providing $A + H^{(1)}$ is non-singular, giving

$$\|x^{(1)} - x\|_2 \leq (\|A^{-1}\|_2\|H^{(1)}\|_2\|x\|_2 + \|A^{-1}\|_2\|e\|_2) / (1 - \|A^{-1}\|_2\|H^{(1)}\|_2). \quad (19)$$

Writing $\alpha = 12.5n2^{-t}$ and remembering that $\|\bar{B}\|_F = \|\bar{U}\|_F$ we have from (14), (9), (11), (8) and (7)

$$\|H^{(1)}\|_F \leq \alpha\|A\|_F + (2^{-t} + \frac{3}{2}n2^{-2t})(1 + \alpha)\|A\|_F \quad (20)$$

and since $\|A\|_2\|A^{-1}\|_2 \geq 1$, it follows from (16) that $H^{(1)}$ certainly satisfies the bound

$$\|H^{(1)}\|_F < 13.5n2^{-t}\|A\|_F. \quad (21)$$

Using the relations $\|A\|_F \leq n^{\frac{1}{2}}\|A\|_2$ (a very weak inequality for most A), $\|H^{(1)}\|_2 \leq \|H^{(1)}\|_F$ and $\|e\|_2 \leq 12.5n2^{-t}\|b\|_2 \leq 12.5n2^{-t}\|A\|_2\|x\|_2$ we have finally

$$\frac{\|x^{(1)} - x\|_2}{\|x\|_2} \leq \frac{13.5n^{\frac{1}{2}}2^{-t}\|A^{-1}\|_2\|A\|_2 + 12.5n2^{-t}\|A^{-1}\|_2\|A\|_2}{1 - 13.5n^{\frac{1}{2}}2^{-t}\|A^{-1}\|_2\|A\|_2} \leq \frac{26n^{\frac{1}{2}}2^{-t}\kappa(A)}{1 - 13.5n^{\frac{1}{2}}2^{-t}\kappa(A)} \quad (22)$$

where $\kappa(A)$ is the usual spectral condition number. The condition (16) now gives the required result.

Since $H^{(1)}$ is uniformly bounded for all right-hand sides this shows that provided A is "not too ill-conditioned" one can be certain that the computed solution will have some correct figures. The bound (16) is extremely conservative. Experience suggest that (17) is true when p is defined by some such relation as

$$n^{\frac{1}{2}}2^{-t}\|A\|_2\|A^{-1}\|_2 = 2^{-p}. \quad (23)$$

Unless A is too ill-conditioned the iterative refinement procedure is certain to work with any right-hand side *provided the errors made in computing the residual are unimportant*. Indeed if the residual were computed exactly the analysis

guarantees that the s -th computed correction $\delta^{(s)}$ satisfies the relation

$$\|\delta^{(s)} - (x - x^{(s)})\|_2 / \|x - x^{(s)}\| \leq 2^{-p} / (1 - 2^{-p-1}). \quad (24)$$

If, further, no error were made in adding $\delta^{(s)}$ to $x^{(s)}$ we have

$$\|x^{(s+1)} - x\|_2 = \|x - x^{(s+1)}\|_2 \leq \|x - x^{(s)}\|_2 2^{-p} / (1 - 2^{-p-1}). \quad (25)$$

On the average $x^{(s)}$ gains roughly the same number of figures in each iteration until it is correct to working accuracy. Indeed by a slight modification of the procedure it is usually possible to obtain an $x^{(s)}$ of any accuracy using the original factorization of A . In practice it is adequate to compute $r^{(s)}$, the s -th residual, using accumulation of inner-products, since if $\bar{r}^{(s)}$ is the computed residual and $r^{(s)}$ the true residual then $\bar{r}^{(s)} = r^{(s)} + f^{(s)}$ where

$$\|f^{(s)}\|_2 \leq 2^{-t} \|r^{(s)}\|_2 + \frac{3}{2} n^{\frac{1}{2}} 2^{-2t} \|A\|_2 \|x^{(s)}\|_2. \quad (26)$$

(This result is an immediate consequence of the error analysis given by WILKINSON [7], pp. 82–85.)

The accumulation of inner-products is not vital at any other stage of the procedure. Failure to use it elsewhere merely leads to a strengthening of condition (16) thereby reducing the range of condition numbers for which the iterative refinement process will succeed and also slows down the rate of improvement to some extent. Provided $\kappa(A)$ satisfies the requisite bound an $x^{(s)}$ is ultimately attained for which

$$\|x - x^{(s)}\| / \|x\| \leq 2^{-t} \quad (27)$$

at which stage x is “correct to working accuracy”.

4. The Least Squares Case with $m > n$

A moment's reflection will make it obvious that we cannot obtain a result which is exactly analogous to (27) in the least squares case, because if $A^T b = 0$ the correct solution is $x = 0$. Nevertheless, we might expect that $\|x - x^{(s)}\|_2$ will show a progressive decrease in magnitude provided only that A satisfies some bound of the form shown in (16).

However, a detailed study of the linear equation case shows that even here the behaviour is not quite as simple as one might expect. As far as the computed $x^{(s)}$ are concerned the behaviour is analogous to that of $x^{(s)}$ defined exactly by the relations

$$r^{(s)} = b - A x^{(s)}, \quad \delta^{(s)} = (A + H^{(s)})^{-1} r^{(s)}, \quad x^{(s+1)} = x^{(s)} + \delta^{(s)} \quad (28)$$

where the $H^{(s)}$ are uniformly bounded. But the $\|r^{(s)}\|_2$ defined by Eqs. (28) diminish at much the same rate as the $\|x - x^{(s)}\|_2$ whereas in the practical iterative refinement procedure $\|r^{(s)}\|_2$ remains roughly constant. This behaviour is predicted by the detailed error analysis (see, for example, WILKINSON [7], pp. 121–126).

5. 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

$$\kappa(A) = \sigma_1 / \sigma_n, \quad \sigma_1 = \max_{x \neq 0} \|A x\|_2 / \|x\|_2, \quad \sigma_n = \min_{x \neq 0} \|A x\|_2 / \|x\|_2 \quad (29)$$

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 U is defined as in Eq. (2) then

$$\sigma_1(U) = \sigma_1(A), \quad \sigma_n(U) = \sigma_n(A), \quad \kappa(U) = \kappa(A), \quad (30)$$

while

$$\sigma_1(U) = \|U\|_2 \quad \text{and} \quad \sigma_n(U) = 1/\|U^{-1}\|_2. \quad (31)$$

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

$$A^T A x = A^T b. \quad (32)$$

The matrix $A^T A$ is square and we have

$$\kappa(A^T A) = \kappa^2(A). \quad (33)$$

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 the Eqs. (32). Now from Eq. (4) the method of orthogonal transformations replaces the least squares problem by the solution of the equations $Ux = p$ and $\kappa(U) = \kappa(A)$. 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

$$\sigma_1 = \|A\|_2 = \|b\|_2 = 1; \quad (34)$$

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

$$\kappa(A) = \kappa(U) = \|U^{-1}\|_2 = 1/\sigma_n. \quad (35)$$

Consider the perturbed system

$$(A + \varepsilon E)b + \varepsilon e, \quad \|E\|_2 = \|e\|_2 = 1, \quad (36)$$

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

$$(A + \varepsilon E)^T (A + \varepsilon E) \bar{x} = (A + \varepsilon E)^T (b + \varepsilon e). \quad (37)$$

If x is the exact solution of the original system and P is the exact orthogonal transformation corresponding to A we have

$$PA = \begin{bmatrix} U \\ 0 \end{bmatrix}, \quad P(A + \varepsilon E) = \begin{bmatrix} U + \varepsilon F \\ \varepsilon G \end{bmatrix}, \quad Pe = \begin{bmatrix} f \\ g \end{bmatrix} \quad (38)$$

and

$$r = b - Ax, \quad A^T r = 0. \quad (39)$$

Eq. (37) therefore becomes

$$(A + \varepsilon E)^T (A + \varepsilon E) (Ax + r + \varepsilon e) = (A + \varepsilon E)^T (b + \varepsilon e) \quad (40)$$

giving

$$\begin{bmatrix} U + \varepsilon F \\ \varepsilon G \end{bmatrix}^T \begin{bmatrix} U + \varepsilon F \\ \varepsilon G \end{bmatrix} \bar{x} = \begin{bmatrix} U + \varepsilon F \\ \varepsilon G \end{bmatrix}^T \left(\begin{bmatrix} U \\ 0 \end{bmatrix} x + \varepsilon \begin{bmatrix} f \\ g \end{bmatrix} \right) + \varepsilon E^T r. \quad (41)$$

Neglecting ε^2 where advantageous

$$\begin{aligned} (U + \varepsilon F)^T (U + \varepsilon F) \bar{x} &= (U + \varepsilon F)^T U x + \varepsilon (U + \varepsilon F)^T f + \varepsilon E^T r + O(\varepsilon^2) \\ \bar{x} &= (U + \varepsilon F)^{-1} U x + \varepsilon (U + \varepsilon F)^{-1} f + \\ &\quad + \varepsilon (U^T U)^{-1} E^T r + O(\varepsilon^2) \\ &= x - \varepsilon U^{-1} F x + \varepsilon U^{-1} f + \varepsilon (U^T U)^{-1} E^T r + O(\varepsilon^2) \end{aligned} \quad (42)$$

giving

$$\begin{aligned} \|\bar{x} - x\|_2 &\leq \varepsilon \|U^{-1}\|_2 \|F\|_2 \|x\|_2 + \varepsilon \|U^{-1}\|_2 \|f\|_2 + \varepsilon \|U^{-1}\|_2^2 \|E\|_2 \|r\|_2 + O(\varepsilon^2) \\ &\leq \varepsilon \kappa(A) \|x\|_2 + \varepsilon \kappa(A) + \varepsilon \kappa^2(A) \|r\|_2 + O(\varepsilon^2). \end{aligned} \quad (43)$$

We observe that the bounds include a term $\varepsilon \kappa^2(A) \|r\|_2$. It is easy to verify by means of a 3×2 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 of $\kappa^2(A)$ is still relevant to some extent.*

When the equations are compatible $\|r\| = 0$ and the term in $\kappa^2(A)$ disappears. In the non-singular linear equation case r is always null and hence it is always $\kappa(A)$ rather than $\kappa^2(A)$ which is relevant.

6. Correction Derived from the Correct Solution

An indication of the fundamental limitations of the algorithm may be obtained by considering the "correction" which corresponds to the correct solution. This correction will be the computed solution of the system $(A : r)$. We know that this computed solution y is the exact solution of some system $(A + H : r + e)$ where

$$\|H\|_F \leq 13.5 n 2^{-t} \|A\|_F, \quad \|e\|_2 \leq 12.5 n 2^{-t} \|r\|_2, \quad (44)$$

and we have as in section 5

$$\begin{aligned} \|y\| &\leq n 2^{-t} [(12.5)/\sigma_n + (13.5) \|A\|_F / \sigma_n^2] \|r\|_2 \\ &\leq 26 n^{\frac{1}{2}} 2^{-t} \kappa^2(A) \|r\|_2, \end{aligned} \quad (45)$$

the last inequality following trivially because $\kappa(A) \leq \kappa^2(A)$ and $1 < n^{\frac{1}{2}}$. Hence even when we start with the correct solution the procedure cannot recognize this. In the case when $x \neq 0$ the relative error in the corrected solution is bounded by $26 n^{\frac{1}{2}} 2^{-t} \kappa^2(A) \|r\|_2 / \|x\|_2$. We can guarantee that the correction will leave x unaffected to working accuracy only if

$$26 n^{\frac{1}{2}} \kappa^2(A) \|r\|_2 / \|x\|_2 < 1. \quad (46)$$

The factor $26 n^{\frac{1}{2}}$ need not be taken too seriously since it will, in general, be a severe overestimate, but even if it is replaced by unity we still require

$$\kappa^2(A) \|r\|_2 / \|x\|_2 < 1. \quad (47)$$

Hence whatever precision of computation is used there will be right-hand sides for which iterative refinement will never give solutions which are correct to working accuracy. This is in striking contrast with the linear equation case.

7. Modified Iterative Procedure

A further indication of the limitations of the iterative refinement procedure is provided by the following considerations. Suppose we were able to perform the iterative refinement procedure *exactly* using the *exact* matrices $P^{(0)}, \dots, P^{(n-1)}$ defined by the algorithm for the *computed* $A^{(n)}$ obtained in the reduction of A to B . Writing $P^{(n-1)} \dots P^{(0)} = P$ this hypothetical procedure would give

$$P r^{(s)} = \begin{bmatrix} \dot{p}^{(s)} \\ \dots \\ q^{(s)} \end{bmatrix}, \quad U \delta^{(s)} = \dot{p}^{(s)} \quad (48)$$

from which

$$\begin{bmatrix} \dot{p}^{(s)} \\ \dots \\ q^{(s)} \end{bmatrix} = P r^{(s)} = P(b - A x^{(s)}) \quad (49)$$

giving

$$\begin{bmatrix} \bar{U} \\ \dots \\ 0 \end{bmatrix}^T \begin{bmatrix} \dot{p}^{(s)} \\ \dots \\ q^{(s)} \end{bmatrix} = \begin{bmatrix} \bar{U} \\ \dots \\ 0 \end{bmatrix}^T P(b - A x^{(s)}). \quad (50)$$

Now $\delta^{(s)} = 0$ when $\dot{p}^{(s)} = 0$, that is from Eq. (50) when

$$\begin{bmatrix} \bar{U} \\ \dots \\ 0 \end{bmatrix}^T P(b - A x^{(s)}) = 0$$

or

$$(A + E)^T (b - A x^{(s)}) = 0. \quad (51)$$

This implies that even this hypothetical procedure would ultimately give an $x^{(s)}$ such that the residual is orthogonal to the columns of $A + E$ instead of A . In the linear equation case $(A + E)^T$ is a square non-singular matrix and Eq. (51) implies $b - A x^{(s)} = 0$. In other words in the linear equation case the hypothetical procedure gives a zero correction if and only if $x^{(s)}$ is the exact solution.

This has led KAHAN (oral remark) to suggest an alternative procedure which has some advantages. Ignoring rounding errors we have

$$A^T A = A^T P^T P A = \begin{bmatrix} U \\ \dots \\ 0 \end{bmatrix}^T \begin{bmatrix} U \\ \dots \\ 0 \end{bmatrix} = U^T U \quad (52)$$

and the usual normal equations are equivalent to

$$U^T U x = A^T b. \quad (53)$$

Hence using the computed \bar{U} we may carry out the iterative refinement procedure defined by

$$\begin{aligned} x^{(0)} &= 0, & r^{(s)} &= b - A x^{(s)} \\ \bar{U}^T \bar{U} \delta^{(s)} &= A^T r^{(s)}, & x^{(s+1)} &= x^{(s)} + \delta^{(s)}. \end{aligned} \quad (54)$$

At each stage we multiply $r^{(s)}$ by A^T and perform a forward and a backward substitution to obtain $\delta^{(s)}$. This certainly has the advantage that when b is exactly

orthogonal to the columns of A then $A^T r^{(0)} = A^T b = 0$, so that the computed correction is the null vector and the process "converges" to the correct answer. (This is not true of the procedure of BUSINGER and GOLUB.)

Unless we are prepared to go further than the mere accumulation of inner-products the limiting accuracy attainable with this algorithm is generally much the same as with that of BUSINGER and GOLUB. For if $x^{(s)}$ is, in fact, the correctly rounded solution then even if we accumulate inner-products when computing $r^{(s)}$, each of its components will have to be rounded, in general, before pre-multiplying with A^T . Instead of obtaining the true residual r , the best we can hope to obtain is $r + h$ where $\|h\| \leq 2^{-t} \|r\|_2$. (Even the correctly rounded r will have an error of this nature.) Hence Eq. (54) gives

$$\bar{U}^T \bar{U} \delta^{(s)} = A^T h \quad (55)$$

since $A^T r$ is null, and if $\|A\| = 1$ we have effectively

$$\|\delta^{(s)}\|_2 \leq 2^{-t} \kappa^2(A) \|r\|_2. \quad (56)$$

However, if we multiply the accumulated $r^{(s)}$ by A^T without previous rounding to form $A^T r^{(s)}$ (this involves the multiplication of a *double-precision* vector by a *single-precision* matrix, accumulating each element in double-precision and rounding only on completion) then, if $x^{(s)}$ is the correct solution, $\delta^{(s)}$ will be given by

$$\bar{U}^T \bar{U} \delta^{(s)} = A^T h \quad \text{where} \quad \|h\| \leq 2^{-2t} \|r\|_2. \quad (57)$$

With this modified arithmetic procedure we can expect $\delta^{(s)}$ to satisfy an inequality of the form

$$\|\delta^{(s)}\|_2 \leq 2^{-2t} \kappa^2(A) \|r\|_2. \quad (58)$$

It is therefore much more satisfactory but a substantial amount of extra work is involved.

8. Numerical Example

The predictions of the error analysis are well illustrated by the numerical examples given by BUSINGER and GOLUB [2]. The matrix A consists of the first five columns of the inverse of the 6×6 Hilbert matrix. A matrix of 39 binary digits was used and all inner-products were accumulated.

The right-hand side b_1 was chosen so that $Ax = b_1$ is compatible. The third iteration gave the correctly rounded solution to working accuracy as is to be expected.

The second right-hand side b_2 is such that $A^T b_2 = 0$ so that the correct solution is the null vector and $r = b_2$. The first computed solution had a component which was of the order of magnitude 10^{-3} and subsequent iterations produced random fluctuations in the sixth significant figure of the first iterate, the norm remaining roughly constant.

The successive residuals also remain roughly constant and differ only slightly from the original right-hand side which is, in any case, the exact residual corresponding to the correct solution. At first sight it may seem surprising that the original right-hand side gives a computed solution of the order of magnitude 10^{-3} , yet the subsequent corrections are of order 10^{-8} although they too are derived from a right-hand side which is much the same as the original. The application

of two stages of our error analysis shows that this is precisely what is to be expected. The first solution contains a term $(A^T A)^{-1}(H^{(1)})^T r$. This produces a contribution to the first residual which, though small, is just adequate to cancel out the term of order $(A^T A)^{-1}(H^{(2)})^T r$ which would otherwise arise in the first correction.

Example

A (exact)

3.60000 ₁₀ +1	-6.30000 ₁₀ +2	3.36000 ₁₀ +3	-7.56000 ₁₀ +3	7.56000 ₁₀ +3
-6.30000 ₁₀ +2	1.47000 ₁₀ +4	-8.82000 ₁₀ +4	2.11680 ₁₀ +5	-2.20500 ₁₀ +5
3.36000 ₁₀ +3	-8.82000 ₁₀ +4	5.64480 ₁₀ +5	-1.41120 ₁₀ +6	1.51200 ₁₀ +6
-7.56000 ₁₀ +3	2.11680 ₁₀ +5	-1.41120 ₁₀ +6	3.62880 ₁₀ +6	-3.96900 ₁₀ +6
7.56000 ₁₀ +3	-2.20500 ₁₀ +5	1.51200 ₁₀ +6	-3.96900 ₁₀ +6	4.41000 ₁₀ +6
-2.77200 ₁₀ +3	8.31600 ₁₀ +5	-5.8212 ₁₀ +5	1.55232 ₁₀ +6	-1.74636 ₁₀ +6

b_1 (exact)

$x^{(3)}$ (Third iteration)

4.63000 ₁₀ +2	1.00000 00000 ₁₀ +0
-1.38600 ₁₀ +4	5.00000 00000 ₁₀ -1
9.70200 ₁₀ +4	3.33333 33333 ₁₀ -1
-2.58720 ₁₀ +5	2.50000 00000 ₁₀ -1
2.91060 ₁₀ +5	2.00000 00000 ₁₀ -1
-1.16424 ₁₀ +5	

First iteration

b_2 (exact)	$x^{(1)}$	$r^{(1)}$	$\delta^{(1)}$
4.62000 ₁₀ +3	-1.34714 39142 ₁₀ -3	4.62000 19967 ₁₀ +3	1.11227 27028 ₁₀ -8
3.96000 ₁₀ +3	-4.48614 75412 ₁₀ -4	3.96000 02989 ₁₀ +3	3.43958 93111 ₁₀ -9
3.46500 ₁₀ +3	-1.91979 16562 ₁₀ -4	3.46499 95056 ₁₀ +3	1.37004 51340 ₁₀ -9
3.08000 ₁₀ +3	-8.39004 73880 ₁₀ -5	3.07999 91136 ₁₀ +3	5.62977 18071 ₁₀ -10
2.77200 ₁₀ +3	-2.98101 00501 ₁₀ -5	2.77199 89175 ₁₀ +3	1.89861 20958 ₁₀ -10
2.52000 ₁₀ +3		2.51999 88239 ₁₀ +3	

Second iteration

$x^{(2)}$	$r^{(2)}$	$\delta^{(2)}$
-1.34713 27915 ₁₀ -3	4.62000 19966 ₁₀ +3	-1.11227 27027 ₁₀ -8
-4.48611 31453 ₁₀ -4	3.96000 02988 ₁₀ +3	-5.05313 76285 ₁₀ -9
-1.91977 29557 ₁₀ -4	3.46499 95056 ₁₀ +3	-2.42938 00078 ₁₀ -9
-8.38999 10403 ₁₀ -5	3.07999 91136 ₁₀ +3	-1.12757 48695 ₁₀ -9
-2.98099 10639 ₁₀ -5	2.77199 89175 ₁₀ +3	-4.15470 31492 ₁₀ -10
	2.51999 88238 ₁₀ +3	

9. General Comments

For the accurate solution of linear equations iterative refinement procedures are extremely attractive whether used in connexion with the Householder triangularization or the Crout triangularization with pivoting [1, 3, 6]. For large systems a single-precision factorization plus iterative refinement has the advantage

over double-precision factorization without refinement that it requires less storage (only about half), is faster if only a few right-hand sides are involved and gives useful indication about the condition of A . It has the disadvantage that if A is too ill-conditioned iterative refinement will not work whereas double-precision computation might well give an answer of acceptable (but unknown) accuracy.

Turning to the least squares problem the choice is effectively between forming the normal equations $A^T A x = A^T b$ and using an LL^T decomposition of $A^T A$, all in double-precision or using the orthogonal reduction in single-precision with iterative refinement.

If, as is true in some fields of operation $m \gg n$, the matrix $A^T A$ is small compared with A . The storage requirement then ceases to be a valid argument against the use of double-precision. The matrix $A^T A$ has the condition number $\kappa^2(A)$ but we have shown that the performance of the iterative refinement process is dependent to some extent on $\kappa^2(A)$ (unless the system is compatible), in spite of the fact that PA and U have the same condition number $\kappa(A)$ as A itself. In some fields of operation overdetermined systems are obtained which are almost compatible so that $\|r\|_2$ is far smaller than κ . In such cases, particularly if m is not excessively large compared with n , the iterative refinement algorithm is very attractive.

In this paper we have discussed only the case when the matrix A has independent columns. Clearly the least squares problem merits further analysis, both theoretical and experimental.

Acknowledgements. The work of J. H. WILKINSON included here is part of the Research Programme of the National Physical Laboratory and is published by permission of the Director of the Laboratory.

The work of G. H. GOLUB was supported in part by the National Science Foundation.

References

- [1] BOWDLER, HILARY, J., R. S. MARTIN, G. PETERS, and J. H. WILKINSON: Solution of real and complex systems of linear equations. *Numerische Mathematik* **8**, 217—234 (1966).
- [2] BUSINGER, P., and G. H. GOLUB: Linear least squares solutions by Householder transformations. *Numerische Mathematik* **7**, 269—276 (1965).
- [3] FORSYTHE, G. E.: Crout with pivoting. *Commun. Ass. Comp. Mach.* **3**, 507—508 (1960).
- [4] GOLUB, G. H.: Numerical methods for solving linear least squares problems. *Numerische Mathematik* **7**, 206—216 (1965).
- [5] HOUSEHOLDER, A. S.: Unitary triangularization of a non-symmetric matrix. *J. Assoc. Comput. Mach.* **5**, 339—342 (1958).
- [6] McKEEMAN, W. M.: Crout with equilibration and iteration. *Commun. Ass. Comp. Mach.* **5**, 552—555 (1962).
- [7] WILKINSON, J. H.: Rounding errors in algebraic processes. Her Majesty's Stationery Office, London. New Jersey: Prentice-Hall 1963.
- [8] — The algebraic eigenvalue problem. London: Oxford University Press 1965.

Computation Center
Stanford University
Stanford, Calif. 94305 (USA)

National Physical Lab.
Teddington, Middlesex
Great Britain