

# SOLVING LINEAR LEAST SQUARES PROBLEMS BY GRAM-SCHMIDT ORTHOGONALIZATION

ÅKE BJÖRCK

## Abstract.

A general analysis of the condition of the linear least squares problem is given. The influence of rounding errors is studied in detail for a modified version of the Gram-Schmidt orthogonalization to obtain a factorization  $A = QR$  of a given  $m \times n$  matrix  $A$ , where  $R$  is upper triangular and  $Q^T Q = I$ . Let  $x$  be the vector which minimizes  $\|b - Ax\|_2$  and  $r = b - Ax$ . It is shown that if inner-products are accumulated in double precision then the errors in the computed  $x$  and  $r$  are less than the errors resulting from some simultaneous initial perturbation  $\delta A$ ,  $\delta b$  such that

$$\|\delta A\|_E / \|A\|_E \approx \|\delta b\|_2 / \|b\|_2 \approx 2 \cdot n^{3/2} \text{ machine units.}$$

No reorthogonalization is needed and the result is independent of the pivoting strategy used.

Key words: Least square, linear equations, orthogonalization.

## 1. Introduction.

Let  $A$  be a given  $m \times n$  real matrix of rank  $n$ ,  $m \geq n$ , and  $b$  a given  $m \times 1$  real vector. Then there exists a unique vector  $x$  which solves the least squares problem of minimizing

$$\|b - Ax\|_2.$$

It is well known that the solution  $x$  satisfies the condition

$$A^T(b - Ax) = 0, \quad (1.1)$$

i.e. the residual vector  $r = b - Ax$  is orthogonal to the columns of  $A$ . It follows that we can compute  $x$  from the normal equations

$$A^T A x = A^T b. \quad (1.2)$$

We define following Bauer [1] the (columnwise) condition of the rectangular matrix  $A$  to be

$$\text{cond}(A) = \max_{\|x\|=1} \|Ax\| / \min_{\|x\|=1} \|Ax\|.$$

The condition number corresponding to the  $L_2$ -norm is denoted by  $\kappa(A)$ . In section 7 we will show that under some restrictions,  $\kappa(A)$  can be con-

sidered as an approximate condition number for the problem of determining  $\mathbf{x}$ .<sup>1</sup> On the other hand

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

This shows that in general using  $t$ -digit binary arithmetic, we will not be able to obtain even an approximate solution to (1.2) unless

$$\kappa(\mathbf{A}) \leq 2^{t/2} .$$

Now let  $\mathbf{B} = \mathbf{A}\mathbf{S}$ , where  $\mathbf{S}$  is square and non-singular. Then from (1.1) follows  $\mathbf{B}^T(\mathbf{b} - \mathbf{A}\mathbf{x}) = 0$  and thus the equations

$$\mathbf{B}^T \mathbf{A} \mathbf{x} = \mathbf{B}^T \mathbf{b} \quad (1.3)$$

can be used instead of the normal equations. A natural question to ask is if we can choose  $\mathbf{B}$  in such a way that

$$\kappa(\mathbf{B}^T \mathbf{A}) = \kappa(\mathbf{A}) . \quad (1.4)$$

This is indeed possible. Since the columns of  $\mathbf{A}$  are linearly independent we have a factorization of  $\mathbf{A}$

$$\mathbf{A} = \mathbf{Q}\mathbf{R} \quad (1.5)$$

where  $\mathbf{Q}^T \mathbf{Q} = \mathbf{I}$  and  $\mathbf{R}$  is an upper triangle. Moreover, each row of  $\mathbf{R}$  and each column of  $\mathbf{Q}$  is uniquely determined up to a scalar factor of modulus one ([4] p. 7).

Choosing  $\mathbf{B} = \mathbf{A}\mathbf{R}^{-1} = \mathbf{Q}$  the matrix  $\mathbf{B}^T \mathbf{A} = \mathbf{Q}^T \mathbf{Q} \mathbf{R} = \mathbf{R}$  is triangular. The equations (1.3) are thus easily solved by back-substitution. Further the condition (1.4) is satisfied since

$$\kappa(\mathbf{R}) = \kappa(\mathbf{Q}\mathbf{R}) = \kappa(\mathbf{A}) . \quad (1.6)$$

This factorization can be obtained mainly in two different ways, and both have been proposed for solving least squares problems. One way is to annihilate the subdiagonal elements of  $\mathbf{A}$  by a sequence of elementary Hermitian matrices

$$\begin{aligned} \mathbf{H}^{(k)} &= (\mathbf{I} - 2\mathbf{w}_k \mathbf{w}_k^T), \quad k = 1, 2, \dots, n, \\ \mathbf{H}^{(n)} \dots \mathbf{H}^{(2)} \mathbf{H}^{(1)} \mathbf{A} &= \mathbf{Q}^T \mathbf{A} = \mathbf{R} . \end{aligned}$$

A detailed error estimate for this algorithm has been given by Wilkinson [7] pp. 152–161, and the application to least squares problems is treated by Golub in [2].

A second way to obtain the factorization is by applying the Gram–Schmidt orthogonalization procedure to the columns of  $\mathbf{A}$ . This is however known to have poor numerical stability.

---

<sup>1</sup> The reader may prefer to read this section at this stage.

In this paper we will use a slightly modified version of the Gram-Schmidt procedure, which is equivalent to the elimination with weighted row combinations in [1]. This will be applied to the solution of linear least squares problems. The propagation of rounding errors will be studied in detail and in particular an estimate of the deviation from orthogonality of the computed  $Q$  and an error bound for the computed solution will be derived.

## 2. Description of the algorithm.

In the actual computation it is preferable to use a different normalization of the factors in (1.5). We avoid computing square roots if we take instead

$$A = Q'R' \quad (2.1)$$

where  $R'$  is a *unit* upper triangle and  $(Q')^T Q'$  diagonal. For completeness we give below the algorithms both for the classical and the modified Gram-Schmidt procedure. In both cases we compute a sequence of matrices

$$A = A^{(1)}, A^{(2)}, \dots, A^{(n+1)} = Q'.$$

*Classical Gram-Schmidt:*

Here the elements of  $R'$  are computed one column at a time. We define

$$A^{(k)} = (q_1', \dots, q_{k-1}', a_k, \dots, a_n)$$

and assume

$$(q_p')^T q_r' = \delta_{pr} d_p, \quad 1 \leq p, r \leq k-1.$$

In step  $k$  we compute

$$r'_{jk} = (q_j')^T a_k / d_j, \quad 1 \leq j \leq k-1,$$

$$q_k' = a_k - \sum_{j=1}^{k-1} r'_{jk} q_j', \quad d_k = \|q_k'\|_2^2.$$

After step  $n$  we treat the right hand side  $b$  in a similar way and compute

$$y_j' = (q_j')^T b / d_j, \quad 1 \leq j \leq n, \quad r = b - \sum_{j=1}^n y_j' q_j'.$$

*Modified Gram-Schmidt:*

Here the elements of  $R'$  are computed one row at a time. We define

$$A^{(k)} = (q_1', \dots, q_{k-1}', a_k^{(k)}, \dots, a_n^{(k)})$$

and assume

$$(\mathbf{q}_p')^T \mathbf{q}_r' = \delta_{pr} \mathbf{d}_p, \quad (\mathbf{q}_p')^T \mathbf{a}_j^{(k)} = 0, \quad 1 \leq p, r \leq k-1, \quad k \leq j \leq n.$$

In step  $k$  we take  $\mathbf{q}_k' = \mathbf{a}_k^{(k)}$  and compute

$$\begin{aligned} d_k &= \|\mathbf{q}_k'\|_2^2, & r'_{kj} &= (\mathbf{q}_k')^T \mathbf{a}_j^{(k)} / d_k, \\ \mathbf{a}_j^{(k+1)} &= \mathbf{a}_j^{(k)} - r'_{kj} \mathbf{q}_k', & k+1 &\leq j \leq n. \end{aligned} \quad (2.2)$$

The right hand side is transformed in a similar way

$$\mathbf{b} = \mathbf{b}^{(1)}, \mathbf{b}^{(2)}, \dots, \mathbf{b}^{(n+1)} = \mathbf{r}$$

where in step  $k$  we compute

$$\mathbf{y}_k' = (\mathbf{q}_k')^T \mathbf{b}^{(k)} / d_k, \quad \mathbf{b}^{(k+1)} = \mathbf{b}^{(k)} - \mathbf{y}_k' \mathbf{q}_k'. \quad (2.3)$$

For both procedures we have taking  $r'_{kk} = 1$

$$\mathbf{a}_k = \sum_{j=1}^k r'_{jk} \mathbf{q}_j', \quad \mathbf{r} = \mathbf{b} - \sum_{j=1}^n \mathbf{y}_j' \mathbf{q}_j'.$$

Thus they both produce, if performed without rounding errors, the wanted unique factorizations

$$\mathbf{A} = \mathbf{Q}' \mathbf{R}', \quad \mathbf{b} = \mathbf{Q}' \mathbf{y}' + \mathbf{r},$$

where

$$(\mathbf{Q}')^T \mathbf{Q}' = \mathbf{D} = \text{diag}\{d_1, d_2, \dots, d_n\}.$$

The two procedures, however, have completely different numerical properties when  $n > 2$ . If  $\mathbf{A}$  is at all ill-conditioned, then using the classical procedure, the computed columns of  $\mathbf{Q}'$  will soon completely lose their orthogonality. Consequently this procedure should never be used without reorthogonalization, which greatly increases the amount of computation. As we shall show, reorthogonalization is never needed when using the modified procedure for solving least squares problems.

The following example due to Läuchli illustrates well the different behavior of the two procedures. Suppose

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

where  $\varepsilon$  is a small number such that due to rounding the result of the operation  $1 + \varepsilon^2$  is everywhere put equal to 1. It is easily verified, that if no other rounding errors are made, then the maximum deviation from orthogonality of the computed columns of  $\mathbf{Q}'$  is given by

$$\text{classical: } \frac{|(\mathbf{q}_3')^T \mathbf{q}_2'|}{\|\mathbf{q}_3'\|_2 \|\mathbf{q}_2'\|_2} = \frac{1}{2}; \quad \text{modified: } \frac{|(\mathbf{q}_3')^T \mathbf{q}_1'|}{\|\mathbf{q}_3'\|_2 \|\mathbf{q}_1'\|_2} = \sqrt{\frac{2}{3}} \varepsilon.$$

Note that  $\kappa(\mathbf{A}) \approx \sqrt{3} \varepsilon^{-1}$  and that  $\mathbf{A}^T \mathbf{A}$  (rounded) has rank 1, so that the normal equations are singular.

The modified procedure also has the advantage of easily allowing for instance one of the following pivoting strategies to be used. We can choose as  $\mathbf{q}_k'$  the column of  $\mathbf{A}^{(k)}$  for which  $\|\mathbf{a}_j^{(k)}\|_2$ ,  $k \leq j \leq n$ , is maximized. Then as is seen from (2.2) the elements of  $\mathbf{R}'$  will satisfy  $|r'_{kj}| \leq 1$ . If we wish to express  $\mathbf{b}$  in as few columns of  $\mathbf{A}$  as possible, we can choose instead as  $\mathbf{q}_k'$  the column for which  $|p_{kj}|$  is maximized, where

$$p_{kj} = (\mathbf{a}_j^{(k)})^T \mathbf{b}^{(k)} / \|\mathbf{a}_j^{(k)}\|_2, \quad k \leq j \leq n.$$

Both  $\|\mathbf{a}_j^{(k)}\|_2$  and  $p_{kj}$  are easily calculated by recursion. This is probably one reason why the modified procedure often has been preferred in practice.

We will show that the superiority of the modified procedure is due to other factors and independent of the pivoting strategy. This is in agreement with experimental results reported by Rice [5].

### 3. Basic definitions for the error analysis.

In the error analysis we shall assume that floating point arithmetic is used, and follow the technique and the notations introduced by Wilkinson [6]. Thus let 'op' denote any of the four arithmetic operators  $+$   $-$   $*$   $/$ . Then an equation of the form

$$z = \text{fl}(x' \text{op}' y)$$

will imply that  $x$ ,  $y$  and  $z$  are floating point numbers and  $z$  obtained from  $x$  and  $y$  using the appropriate floating point operation. We assume that the rounding errors in these operations are such that

$$\text{fl}(x' \text{op}' y) = (x' \text{op}' y)(1 + \varepsilon), \quad |\varepsilon| \leq 2^{-t}. \quad (3.1)$$

Let  $\mathbf{x}$  and  $\mathbf{y}$  be vectors of dimension  $n$ , where  $n \cdot 2^{-t} \leq 0.1$ . Then the following error bounds for the computed innerproduct of  $\mathbf{x}$  and  $\mathbf{y}$  are valid ([7] pp. 114, 117):

$$|\text{fl}(\mathbf{x}^T \mathbf{y}) - \mathbf{x}^T \mathbf{y}| \leq n \cdot 2^{-t_1} |\mathbf{x}^T| |\mathbf{y}|, \quad (3.2)$$

$$|\text{fl}_2(\mathbf{x}^T \mathbf{y}) - \mathbf{x}^T \mathbf{y}| \leq 2^{-t} |\mathbf{x}^T \mathbf{y}| + \frac{3}{2} \cdot n 2^{-2t_2} |\mathbf{x}^T| |\mathbf{y}| \quad (3.3)$$

where

$$t_1 = t - \log_2(1.06), \quad 2t_2 = 2t - \log_2(1.06). \quad (3.4)$$

Here  $|x|$  denotes a vector with components  $|x_i|$  and  $\text{fl}_2(\cdot)$  applies to the case when the inner-product is accumulated in double precision and then rounded.

Computed quantities will in the following be denoted by bars. Thus, we write

$$\bar{\mathbf{A}}^{(k)} = (\bar{\mathbf{q}}_1', \dots, \bar{\mathbf{q}}_{k-1}', \bar{\mathbf{a}}_k^{(k)}, \dots, \bar{\mathbf{a}}_n^{(k)})$$

and the formulas (2.2) and (2.3) become

$$\begin{aligned} \bar{d}_k &= \text{fl}(\|\bar{\mathbf{q}}_k'\|_2^2), & \bar{r}_{kj}' &= \text{fl}((\bar{\mathbf{q}}_k')^T \bar{\mathbf{a}}_j^{(k)} / \bar{d}_k), \\ \bar{\mathbf{a}}_j^{(k+1)} &= \text{fl}(\bar{\mathbf{a}}_j^{(k)} - \bar{r}_{kj}' \bar{\mathbf{q}}_k'), & k+1 \leq j \leq n, \end{aligned} \quad (3.5)$$

and

$$\bar{y}_k' = \text{fl}((\bar{\mathbf{q}}_k')^T \bar{\mathbf{b}}^{(k)} / \bar{d}_k), \quad \bar{\mathbf{b}}^{(k+1)} = \text{fl}(\bar{\mathbf{b}}^{(k)} - \bar{y}_k' \bar{\mathbf{q}}_k'). \quad (3.6)$$

This algorithm will not break down unless  $\bar{d}_k = 0$  for some  $k$ . This may happen due to rounding errors even when the rank of  $\mathbf{A}$  equals  $n$ . According to (3.2) and (3.3)  $\bar{d}_k = 0$  implies that  $\bar{\mathbf{q}}_k' = \mathbf{0}$ . To account for this case we add the (trivial) singular rule that when  $\bar{\mathbf{q}}_k' = \mathbf{0}$  we take

$$\bar{r}_{kj}' = 0, \quad k+1 \leq j \leq n, \quad \bar{y}_k' = 0. \quad (3.7)$$

Then  $\bar{\mathbf{R}}'$  and  $\bar{\mathbf{y}}'$  are always uniquely defined.

For the convenience of notation we also introduce the normalized quantities

$$\bar{\mathbf{q}}_k = d_k^{-1/2} \bar{\mathbf{q}}_k', \quad \bar{r}_{kj} = d_k^{1/2} \bar{r}_{kj}', \quad \bar{y}_k = d_k^{1/2} \bar{y}_k' \quad (3.8)$$

where

$$d_k^{1/2} = \begin{cases} \|\bar{\mathbf{q}}_k'\|_2, & \bar{\mathbf{q}}_k' \neq \mathbf{0} \\ 1, & \bar{\mathbf{q}}_k' = \mathbf{0} \end{cases}$$

Note that these quantities are never computed and thus (3.8) are exact relations.

#### 4. Errors in an elementary projection.

If  $\bar{\mathbf{q}}_k' \neq \mathbf{0}$ , then in the  $k$ :th step we compute the vectors  $\bar{\mathbf{a}}_j^{(k+1)}$ ,  $j = k+1, \dots, n$  as the projection of  $\bar{\mathbf{a}}_j^{(k)}$  on the subspace complementary to  $\bar{\mathbf{q}}_k'$ . If this is performed without rounding errors the result is

$$(\mathbf{I} - \bar{\mathbf{q}}_k \bar{\mathbf{q}}_k^T) \bar{\mathbf{a}}_j^{(k)} = \bar{\mathbf{a}}_j^{(k)} - r_{kj} \bar{\mathbf{q}}_k$$

where  $r_{kj}$  is the *exact* multiplier corresponding to the *computed*  $\bar{\mathbf{q}}_k$  and  $\bar{\mathbf{a}}_j^{(k)}$ . If using the *computed* multiplier  $\bar{r}_{kj}'$  the subtraction is performed exactly, then the result is

$$\bar{\mathbf{a}}_j^{(k)} - \bar{r}_{kj}' \bar{\mathbf{q}}_k' = \bar{\mathbf{a}}_j^{(k)} - \bar{r}_{kj} \bar{\mathbf{q}}_k.$$

We define the related errors  $\delta_j^{(k)}$  and  $\eta_j^{(k)}$  by

$$\bar{\mathbf{a}}_j^{(k+1)} = \bar{\mathbf{a}}_j^{(k)} - \bar{r}_{kj} \bar{\mathbf{q}}_k + \delta_j^{(k)} \quad (4.1)$$

$$\bar{\mathbf{a}}_j^{(k+1)} = (\mathbf{I} - \bar{\mathbf{q}}_k \bar{\mathbf{q}}_k^T) \bar{\mathbf{a}}_j^{(k)} + \eta_j^{(k)}. \quad (4.2)$$

In the singular case when  $\bar{\mathbf{q}}_k' = \mathbf{0}$  these relations are satisfied with

$$\bar{\mathbf{a}}_j^{(k+1)} = \bar{\mathbf{a}}_j^{(k)}, \quad \delta_j^{(k)} = \eta_j^{(k)} = \mathbf{0}. \quad (4.3)$$

In the non-singular case we shall prove the following estimates which are basic for our analysis

$$\|\delta_j^{(k)}\|_2 \leq 1.45 \cdot 2^{-t} \|\bar{\mathbf{a}}_j^{(k)}\|_2 \quad (4.4)$$

$$\|\eta_j^{(k)}\|_2 \leq \begin{cases} 3.23 \cdot 2^{-t} \|\bar{\mathbf{a}}_j^{(k)}\|_2 \text{ d.p. acc.} \\ (2m+3) \cdot 2^{-t} \|\bar{\mathbf{a}}_j^{(k)}\|_2 \text{ s.p.} \end{cases} \quad (4.5)$$

where (4.4) is valid also when single precision is used. For simplicity of notation we omit for a while the indices  $j$  and  $k$  and put

$$\bar{\mathbf{a}}_j^{(k+1)} = \mathbf{z}, \quad \bar{\mathbf{a}}_j^{(k)} = \mathbf{y}, \quad \bar{\mathbf{q}}_k = \mathbf{x}.$$

Then (3.5), (4.1) and (4.2) become

$$\bar{r}' = \text{fl}((\mathbf{x}')^T \mathbf{y} / (\mathbf{x}')^T \mathbf{x}'), \quad \mathbf{z} = \text{fl}(\mathbf{y} - \bar{r}' \mathbf{x}'), \quad (4.6)$$

$$\mathbf{r} = \mathbf{x}^T \mathbf{y}, \quad \mathbf{z} = \mathbf{y} - \bar{r} \mathbf{x} + \delta, \quad \mathbf{z} = \mathbf{y} - r \mathbf{x} + \eta. \quad (4.7)$$

From (3.1) and (3.8) follows

$$z_i = (y_i - \bar{r} x_i (1 + \varepsilon_1)) (1 + \varepsilon_2), \quad |\varepsilon_1| \leq 2^{-t}, \quad |\varepsilon_2| \leq 2^{-t}.$$

Using this to eliminate  $y_i$  from the definition of  $\delta$  we get

$$\delta_i = \frac{\varepsilon_2}{1 + \varepsilon_2} z_i - \varepsilon_1 \cdot \bar{r} x_i,$$

and hence since  $\|\mathbf{x}\|_2 = 1$

$$\|\delta\|_2 \leq \frac{2^{-t}}{1 - 2^{-t}} \|\mathbf{z}\|_2 + 2^{-t} |\bar{r}|. \quad (4.8)$$

Immediately from (4.7)

$$\|\eta\|_2 \leq \|\delta\|_2 + |r - \bar{r}|. \quad (4.9)$$

To estimate the error in the multiplier  $\bar{r}$  we use (3.1)–(3.3) and the identity

$$\frac{\text{fl}(\mathbf{x}^T \mathbf{y})}{\text{fl}(\mathbf{x}^T \mathbf{x})} - r = \frac{\mathbf{x}^T \mathbf{x}}{\text{fl}(\mathbf{x}^T \mathbf{x})} \left[ \frac{\text{fl}(\mathbf{x}^T \mathbf{y}) - \mathbf{x}^T \mathbf{y}}{\mathbf{x}^T \mathbf{x}} - r \frac{\text{fl}(\mathbf{x}^T \mathbf{x}) - \mathbf{x}^T \mathbf{x}}{\mathbf{x}^T \mathbf{x}} \right].$$

If we compute  $\bar{r}'$  as

$$\bar{d} = \text{fl}_2((x')^T x'), \quad \bar{r}' = \text{fl}_2((x')^T y / \bar{d})$$

accumulating inner-products in double precision and divide  $\bar{d}$  into the double precision mantissa before rounding, then

$$|\bar{r} - r| \leq \frac{2^{-t}(2 + \frac{3}{2} \cdot m \cdot 2^{-t_1})|r| + \frac{3}{2} \cdot m \cdot 2^{-t_2} \|y\|_2}{1 - 2^{-t}(1 + \frac{3}{2} \cdot m \cdot 2^{-t_1})}.$$

Using instead single precision we have

$$|\bar{r} - r| \leq \frac{(m+1)2^{-t_1}|r| + m(1+2^{-t})2^{-t_1}\|y\|_2}{1 - m \cdot 2^{-t_1}}.$$

In order to simplify these and later inequalities we make in the following the reasonable assumptions

$$m \geq 2, \quad 2(m+1)2^{-t_1} < 0.01. \quad (4.10)$$

Then we certainly have

$$|\bar{r} - r| < \begin{cases} (2.01 \cdot |r| + 0.01\|y\|_2)2^{-t} & \text{d.p. acc.} \\ ((m+1) \cdot |r| + m\|y\|_2)2^{-t_1} & \text{s.p.} \end{cases} \quad (4.11)$$

where we have taken into account that because of (4.10) the factor 1.06 in  $2^{-t_1}$  is now sufficiently generous to allow also for the factors  $1 - m \cdot 2^{-t_1}$  and  $1 + 2^{-t}$ .

Since  $(z - \eta)$  is orthogonal to  $x$  it follows

$$\|z\|_2 \leq (\|y\|_2^2 - r^2)^{\frac{1}{2}} + \|\eta\|_2. \quad (4.12)$$

Substituting this and (4.8) in (4.9) and solving for  $\|\eta\|_2$  gives after some rearranging

$$(1 - 2 \cdot 2^{-t})\|\eta\|_2 \leq 2^{-t}((\|y\|_2^2 - r^2)^{\frac{1}{2}} + |r|) + |r - \bar{r}|.$$

Hence using (4.11)

$$(1 - 2 \cdot 2^{-t})\|\eta\|_2 \leq \begin{cases} ((\|y\|_2^2 - r^2)^{\frac{1}{2}} + 3|r| + 0.02\|y\|_2)2^{-t}, & \text{d.p. acc.} \\ ((\|y\|_2^2 - r^2)^{\frac{1}{2}} + (m+2)|r| + m\|y\|_2)2^{-t_1}, & \text{s.p.} \end{cases} \quad (4.13)$$

Maximizing over  $|r|$ , where  $0 \leq |r| \leq \|y\|_2$  we have

$$(\|y\|_2^2 - r^2)^{\frac{1}{2}} + k|r| \leq (1 + k^2)^{\frac{1}{2}}\|y\|_2, \quad (4.14)$$

and finally (4.5) follows from (4.13) by invoking (4.10).

For the estimation of  $\|\delta\|_2$  (4.8) and (4.12) yields

$$\|\delta\|_2 \leq \frac{2^{-t}}{1 - 2^{-t}}((\|y\|_2^2 - r^2)^{\frac{1}{2}} + |r| + \|\eta\|_2 + |r - \bar{r}|).$$



Using (4.14) for the first two terms and single precision bounds for the last two we get

$$\|\delta\|_2 \leq \frac{2^{-t}}{1-2^{-t}} (2^{\frac{1}{2}} + 4(m+1) \cdot 2^{-t_1}) \|y\|_2.$$

and invoking (4.10) again this proves (4.4).

In exact computation we certainly have  $\|\bar{\mathbf{a}}_j^{(k+1)}\|_2 \leq \|\bar{\mathbf{a}}_j^{(k)}\|_2$ . This however need not be true when rounding errors are taken into account. Estimating  $\|\eta\|_2$  from (4.13) in (4.12) and again maximizing the resulting right hand side as a function of  $|r|$ , we get after some computations

$$\|\bar{\mathbf{a}}_j^{(k+1)}\|_2 < \begin{cases} (1 + 1.05 \cdot 2^{-t}) \|\bar{\mathbf{a}}_j^{(k)}\|_2, & \text{d.p. acc.} \\ (1 + 1.01(m+2) \cdot 2^{-t_1}) \|\bar{\mathbf{a}}_j^{(k)}\|_2, & \text{s.p.} \end{cases} \quad (4.15)$$

Note that these bounds can only be approached when  $\bar{\mathbf{a}}_j^{(k)}$  is almost orthogonal to  $\bar{\mathbf{q}}_k$ .

## 5. Errors in the factorization.

Using the error analysis of the elementary projection we can now derive bounds for the errors related to the factorization of  $\mathbf{A}$  and  $\mathbf{b}$ . *While we will be able to show that the error  $(\bar{\mathbf{Q}}\bar{\mathbf{R}} - \mathbf{A})$  is small independently of  $\kappa(\mathbf{A})$ , this will not be true of the error  $(\bar{\mathbf{R}} - \bar{\mathbf{Q}}^T \mathbf{A})$ .* This leads us to introduce the matrices

$$\tilde{\mathbf{Q}} = (\tilde{\mathbf{q}}_1, \tilde{\mathbf{q}}_2, \dots, \tilde{\mathbf{q}}_n), \quad \tilde{\tilde{\mathbf{Q}}} = (\tilde{\tilde{\mathbf{q}}}_1, \tilde{\tilde{\mathbf{q}}}_2, \dots, \tilde{\tilde{\mathbf{q}}}_n)$$

where

$$\begin{aligned} \tilde{\mathbf{q}}_k &= (\mathbf{I} - \bar{\mathbf{q}}_1 \bar{\mathbf{q}}_1^T)(\mathbf{I} - \bar{\mathbf{q}}_2 \bar{\mathbf{q}}_2^T) \dots (\mathbf{I} - \bar{\mathbf{q}}_{k-1} \bar{\mathbf{q}}_{k-1}^T) \bar{\mathbf{q}}_k, \\ \tilde{\tilde{\mathbf{q}}}_k &= (\mathbf{I} - \bar{\mathbf{q}}_n \bar{\mathbf{q}}_n^T)(\mathbf{I} - \bar{\mathbf{q}}_{n-1} \bar{\mathbf{q}}_{n-1}^T) \dots (\mathbf{I} - \bar{\mathbf{q}}_{k+1} \bar{\mathbf{q}}_{k+1}^T) \bar{\mathbf{q}}_k. \end{aligned} \quad (5.1)$$

When  $\bar{\mathbf{q}}_j \neq 0$  the matrix

$$\mathbf{P}^{(j)} = (\mathbf{I} - \bar{\mathbf{q}}_j \bar{\mathbf{q}}_j^T)$$

is the projector for the subspace complementary to  $\bar{\mathbf{q}}_j$ . Otherwise it is equal to the unit matrix. It follows that

$$\|\mathbf{P}^{(j)}\|_2 = 1, \quad \|\tilde{\mathbf{q}}_k\|_2 \leq 1, \quad \|\tilde{\tilde{\mathbf{q}}}_k\|_2 \leq 1. \quad (5.2)$$

We shall first prove a simple lemma concerning  $\tilde{\mathbf{Q}}$  and  $\tilde{\tilde{\mathbf{Q}}}$ .

LEMMA 5.1. *Let  $\mathbf{U}$  be the upper triangular matrix with elements*

$$u_{ij} = \bar{\mathbf{q}}_i^T \bar{\mathbf{q}}_j, \quad j > i, \quad u_{ij} = 0, \quad j \leq i.$$

*Then*

$$\bar{\mathbf{Q}} = \tilde{\mathbf{Q}}(\mathbf{I} + \mathbf{U}), \quad \bar{\mathbf{Q}} = \tilde{\tilde{\mathbf{Q}}}(\mathbf{I} + \mathbf{U}^T).$$

We first prove by induction in  $k$  that for  $k=1, 2, \dots, n$

$$(I - \bar{\mathbf{q}}_1 \bar{\mathbf{q}}_1^T) \dots (I - \bar{\mathbf{q}}_k \bar{\mathbf{q}}_k^T) = I - \tilde{\mathbf{q}}_k \bar{\mathbf{q}}_k^T - \dots - \tilde{\mathbf{q}}_1 \bar{\mathbf{q}}_1^T.$$

This is true for  $k=1$  and we have

$$(I - \bar{\mathbf{q}}_1 \bar{\mathbf{q}}_1^T) \dots (I - \bar{\mathbf{q}}_k \bar{\mathbf{q}}_k^T) (I - \bar{\mathbf{q}}_{k+1} \bar{\mathbf{q}}_{k+1}^T) = -\tilde{\mathbf{q}}_{k+1} \bar{\mathbf{q}}_{k+1}^T + \\ + (I - \bar{\mathbf{q}}_1 \bar{\mathbf{q}}_1^T) \dots (I - \bar{\mathbf{q}}_k \bar{\mathbf{q}}_k^T),$$

which proves the hypothesis.

Using this result we get from the definition of  $\tilde{\mathbf{q}}_k$

$$\tilde{\mathbf{q}}_k = (I - \tilde{\mathbf{q}}_{k-1} \bar{\mathbf{q}}_{k-1}^T - \dots - \tilde{\mathbf{q}}_1 \bar{\mathbf{q}}_1^T) \bar{\mathbf{q}}_k,$$

and it follows that for  $k=1, 2, \dots, n$

$$\bar{\mathbf{q}}_k = \tilde{\mathbf{q}}_k + (\bar{\mathbf{q}}_k^T \bar{\mathbf{q}}_{k-1}) \tilde{\mathbf{q}}_{k-1} + \dots + (\bar{\mathbf{q}}_k^T \bar{\mathbf{q}}_1) \tilde{\mathbf{q}}_1.$$

This is the  $k$ :th column of the equality  $\bar{\mathbf{Q}} = \tilde{\mathbf{Q}}(I + U)$  and the first part of the lemma is proved.

The second part is almost obvious from the symmetry of definition of  $\tilde{\mathbf{Q}}$  and  $\bar{\mathbf{Q}}$ , and is proved in the same way.

By taking the transpose of the induction hypothesis used in the proof of the lemma we obtain the

COROLLARY. *If  $\tilde{\mathbf{Q}}$  is defined by (5.1) we have the identity*

$$I - \bar{\mathbf{Q}} \tilde{\mathbf{Q}}^T = (I - \bar{\mathbf{q}}_n \bar{\mathbf{q}}_n^T) \dots (I - \bar{\mathbf{q}}_2 \bar{\mathbf{q}}_2^T) (I - \bar{\mathbf{q}}_1 \bar{\mathbf{q}}_1^T).$$

We now define

$$\mathbf{E}_1 = \bar{\mathbf{Q}} \bar{\mathbf{R}} - \mathbf{A}, \quad \mathbf{e}_1 = \bar{\mathbf{Q}} \bar{\mathbf{y}} + \bar{\mathbf{b}}^{(n+1)} - \mathbf{b}, \quad (5.3)$$

$$\mathbf{E}_2 = (I - \bar{\mathbf{Q}} \tilde{\mathbf{Q}}^T) \mathbf{A}, \quad \mathbf{e}_2 = (I - \bar{\mathbf{Q}} \tilde{\mathbf{Q}}^T) \mathbf{b} - \bar{\mathbf{b}}^{(n+1)}, \quad (5.4)$$

$$\mathbf{E}_3 = \bar{\mathbf{R}} - \tilde{\mathbf{Q}}^T \mathbf{A}, \quad \mathbf{e}_3 = \bar{\mathbf{y}} - \tilde{\mathbf{Q}}^T \mathbf{b}. \quad (5.5)$$

We shall prove the following estimates, which are valid if *inner-products are accumulated in double precision*<sup>1</sup>

$$\|\mathbf{E}_1\|_E \leq 1.5(n-1) \cdot 2^{-t} \|\mathbf{A}\|_E \\ \|\mathbf{e}_1\|_2 \leq 1.5n \cdot 2^{-t} \|\mathbf{b}\|_2 \quad (5.6)$$

$$\|\mathbf{E}_2\|_E \leq 3.25(n-1) \cdot 2^{-t} \|\mathbf{A}\|_E \\ \|\mathbf{e}_2\|_2 \leq 3.25n \cdot 2^{-t} \|\mathbf{b}\|_2 \quad (5.7)$$

$$\|\mathbf{E}_3\|_E \leq 1.9(n-1) \cdot 2^{-t} \|\mathbf{A}\|_E \\ \|\mathbf{e}_3\|_2 \leq 1.9n \cdot 2^{-t} \|\mathbf{b}\|_2 \quad (5.8)$$

---

<sup>1</sup> Here the suffix  $E$  denotes the Frobenius norm i.e.  $\|\mathbf{A}\|_E = (\sum \sum |a_{ij}|^2)^{\frac{1}{2}}$ .

Summing (4.1) for  $k=1, 2, \dots, j-1$  and using

$$\bar{\mathbf{a}}_j^{(j)} = \bar{r}_{jj} \bar{\mathbf{q}}_j, \quad \bar{\mathbf{a}}_j^{(1)} = \mathbf{a}_j$$

we get

$$\sum_{k=1}^j \bar{r}_{kj} \cdot \bar{\mathbf{q}}_k - \mathbf{a}_j = \sum_{k=1}^{j-1} \delta_j^{(k)} = \delta_j.$$

From (4.4) follows

$$\|\delta_j\|_2 \leq 1.45 \cdot 2^{-t} \sum_{k=1}^{j-1} \|\bar{\mathbf{a}}_j^{(k)}\|_2.$$

Using the  $\text{fl}_2(\cdot)$  bound in (4.15)

$$\|\bar{\mathbf{a}}_j^{(k)}\|_2 < (1 + 1.05 \cdot 2^{-t})^{k-1} \|\mathbf{a}_j\|_2$$

and thus invoking (4.10)

$$\|\mathbf{a}_j^{(k)}\|_2 < (1 + 1.05n \cdot 2^{-t}) \|\mathbf{a}_j\|_2 < 1.006 \|\mathbf{a}_j\|_2. \quad (5.9)$$

This inequality holds even when using single precision if (4.10) is replaced by the stronger condition

$$2n(m+2)2^{-t_1} < 0.01.$$

We then certainly have

$$\|\delta_j\|_2 < 1.5 \cdot (j-1)2^{-t} \|\mathbf{a}_j\|_2 \quad (5.10)$$

and

$$\|\mathbf{E}_1\|_E = \left( \sum_{j=1}^n \|\delta_j\|_2^2 \right)^{\frac{1}{2}} < 1.5 \cdot (n-1)2^{-t} \left( \sum_{j=1}^n \|\mathbf{a}_j\|_2^2 \right)^{\frac{1}{2}}$$

which proves the first part of (5.6). Since the right hand side  $\mathbf{b}$  is treated in the same way as the columns of  $\mathbf{A}$ , the second part follows immediately.

Next we solve the difference equation (4.2) for  $\bar{\mathbf{a}}_j^{(k)}$  and get for  $k \leq j$

$$\begin{aligned} \bar{\mathbf{a}}_j^{(k)} = & \mathbf{a}_j^{(k)} + \mathbf{P}^{(k-1)} \dots \mathbf{P}^{(3)} \mathbf{P}^{(2)} \eta_j^{(1)} + \dots \\ & + \mathbf{P}^{(k-1)} \eta_j^{(k-2)} + \eta_j^{(k-1)} \end{aligned} \quad (5.11)$$

where

$$\mathbf{a}_j^{(k)} = \mathbf{P}^{(k-1)} \dots \mathbf{P}^{(2)} \mathbf{P}^{(1)} \mathbf{a}_j.$$

Hence  $\mathbf{a}_j^{(k)}$  is the vector we should obtain if we performed the computation *exactly* with the *computed*  $\bar{\mathbf{q}}_k$ .

Since  $\bar{\mathbf{q}}_j' = \bar{\mathbf{a}}_j^{(j)}$  we have  $(\mathbf{I} - \bar{\mathbf{q}}_j \bar{\mathbf{q}}_j^T) \bar{\mathbf{a}}_j^{(j)} = \mathbf{0}$ . Consequently (5.11) holds also for  $j+1 \leq k \leq n+1$  with

$$\bar{\mathbf{a}}_j^{(k)} = \mathbf{0}, \quad \eta_j^{(j)} = \dots = \eta_j^{(k-1)} = \mathbf{0}.$$

Taking norms in (5.11) and using (5.2) we get

$$\|\bar{\mathbf{a}}_j^{(k)} - \mathbf{a}_j^{(k)}\|_2 \leq \sum_{i=1}^{s-1} \|\eta_j^{(i)}\|_2, \quad s = \min(j, k).$$

From this follows, using  $\text{fl}_2(\cdot)$  bounds from (4.5)

$$\|\bar{\mathbf{a}}_j^{(k)} - \mathbf{a}_j^{(k)}\|_2 \leq 3.23 \cdot 2^{-t} \sum_{i=1}^{s-1} \|\bar{\mathbf{a}}_j^{(i)}\|_2 \leq 3.25(s-1)2^{-t}\|\mathbf{a}_j\|_2. \quad (5.12)$$

Using the corollary of lemma 5.1 we can write this for  $k = n+1$  as

$$\|(I - \bar{\mathbf{Q}}\bar{\mathbf{Q}}^T)\mathbf{a}_j\|_2 \leq 3.25(j-1) \cdot 2^{-t}\|\mathbf{a}_j\|_2.$$

This and a similar relation for  $\mathbf{b}$  proves (5.7).

Using the identity

$$\tilde{\mathbf{q}}_k^T \mathbf{a}_j = \bar{\mathbf{q}}_k^T (I - \bar{\mathbf{q}}_{k-1} \bar{\mathbf{q}}_{k-1}^T) \dots (I - \bar{\mathbf{q}}_1 \bar{\mathbf{q}}_1^T) \mathbf{a}_j = \bar{\mathbf{q}}_k^T \mathbf{a}_j^{(k)}$$

we can write the component  $(k, j)$  of  $\bar{\mathbf{R}} - \tilde{\mathbf{Q}}^T \mathbf{A}$

$$\bar{r}_{kj} - \tilde{\mathbf{q}}_k^T \mathbf{a}_j = (\bar{r}_{kj} - r_{kj}) + \bar{\mathbf{q}}_k^T (\bar{\mathbf{a}}_j^{(k)} - \mathbf{a}_j^{(k)}), \quad k < j.$$

It follows that

$$|\bar{r}_{kj} - \tilde{\mathbf{q}}_k^T \mathbf{a}_j| \leq |\bar{r}_{kj} - r_{kj}| + \|\bar{\mathbf{a}}_j^{(k)} - \mathbf{a}_j^{(k)}\|_2.$$

Now (5.12) and (4.11) gives

$$|\bar{r}_{kj} - \tilde{\mathbf{q}}_k^T \mathbf{a}_j| = (2.02 + 3.25(k-1))2^{-t}\|\mathbf{a}_j\|_2 < 3.25k \cdot 2^{-t}\|\mathbf{a}_j\|_2.$$

The  $L_2$ -norm of the  $j$ :th column in  $\mathbf{E}_3$  is therefore bounded by

$$3.25 \cdot 2^{-t} \left( \sum_{k=1}^{j-1} k^2 \right)^{\frac{1}{2}} \|\mathbf{a}_j\|_2 < 1.9(j-1)^{\frac{1}{2}} \cdot 2^{-t}\|\mathbf{a}_j\|_2,$$

which finally proves (5.8).

Note that *when  $\mathbf{A}$  is ill-conditioned cancellation will occur, so that  $\|\bar{\mathbf{a}}_j^{(k)}\|_2 \ll \|\mathbf{a}_j\|_2$  often for a rather small value of  $k$* . Then the estimate (5.9) which we have used in (5.10) and (5.12) is very crude, and (5.6)–(5.8) will considerably overestimate the error.

If we use single precision (5.6) holds unchanged, but in (5.7) and (5.8) the bounds must be increased by a factor of  $(\frac{2}{3}m + 1)$ .

## 6. Orthogonality of the computed vectors.

All the error bounds given in section 5 apply also in the singular case when  $\bar{\mathbf{q}}_k' = 0$  for some  $k$ , i.e. when the rank of  $\bar{\mathbf{A}} = \bar{\mathbf{Q}}\bar{\mathbf{R}}$  is less than  $n$ . We now derive a sufficient condition for  $\bar{\mathbf{A}}$  to have rank  $n$ . Let the exact factorization of  $\mathbf{A}$  be  $\mathbf{A} = \mathbf{Q}\mathbf{R}$ . Then using (5.3) we can write

$$\bar{\mathbf{A}}^T \bar{\mathbf{A}} = (\mathbf{A} + \mathbf{E}_1)^T (\mathbf{A} + \mathbf{E}_1) = \mathbf{R}^T (\mathbf{I} + \mathbf{F}_1) \mathbf{R} \quad (6.1)$$

where

$$\mathbf{F}_1 = (\mathbf{Q}^T \mathbf{E}_1 \mathbf{R}^{-1})^T + \mathbf{Q}^T \mathbf{E}_1 \mathbf{R}^{-1} + (\mathbf{E}_1 \mathbf{R}^{-1})^T \mathbf{E}_1 \mathbf{R}^{-1}. \quad (6.2)$$

Taking norms in (6.2) we get

$$\|\mathbf{F}_1\|_2 \leq 2\|\mathbf{E}_1\|_E \|\mathbf{R}^{-1}\|_2 + \|\mathbf{E}_1\|_E^2 \|\mathbf{R}^{-1}\|_2^2 \quad (6.3)$$

and thus  $\|\mathbf{F}_1\|_2 < 1$  if

$$\|\mathbf{E}_1\|_E \|\mathbf{R}^{-1}\|_2 < \sqrt{2} - 1.$$

Using (5.6) it follows that  $\bar{\mathbf{A}}^T \bar{\mathbf{A}}$  is non-singular and  $\bar{\mathbf{A}}$  has rank  $n$  if

$$1.5(n-1)2^{-t} \|\mathbf{A}\|_E \|\mathbf{R}^{-1}\|_2 < \sqrt{2} - 1. \quad (6.4)$$

We assume in the following that (6.4) is satisfied.

The orthogonality of the computed vectors  $\bar{\mathbf{q}}_1, \bar{\mathbf{q}}_2, \dots, \bar{\mathbf{q}}_n$  can be measured by the norm of the matrix  $(\mathbf{I} - \bar{\mathbf{Q}}^T \bar{\mathbf{Q}})$ . Now

$$\mathbf{I} - \bar{\mathbf{Q}}^T \bar{\mathbf{Q}} = -(\mathbf{U} + \mathbf{U}^T) \quad (6.5)$$

where  $\mathbf{U}$  is the matrix introduced in lemma 5.1.

Summing (4.1) for  $k=i+1, i+2, \dots, j-1$  we get

$$\bar{\mathbf{a}}_j^{(i+1)} = \sum_{k=i+1}^j \bar{r}_{kj} \bar{\mathbf{q}}_k - \sum_{k=i+1}^{j-1} \delta_j^{(k)}. \quad (6.6)$$

From (4.2) follows

$$\bar{\mathbf{q}}_i^T \bar{\mathbf{a}}_j^{(i+1)} = \bar{\mathbf{q}}_i^T \eta_j^{(i)}$$

and hence multiplying (6.6) by  $\bar{\mathbf{q}}_i^T$  we get

$$s_{ij} = \sum_{k=i+1}^j \bar{r}_{kj} (\bar{\mathbf{q}}_i^T \bar{\mathbf{q}}_k) = \bar{\mathbf{q}}_i^T \left( \eta_j^{(i)} + \sum_{k=i+1}^{j-1} \delta_j^{(k)} \right), \quad (6.7)$$

where  $s_{ij}$  is the component  $(i, j)$  of the matrix  $\mathbf{S} = \mathbf{U} \bar{\mathbf{R}}$ .

Using the  $\text{fl}_2(\cdot)$  bound from (4.10), (4.11) and (5.9) we get

$$|s_{ij}| \leq (3.25 + 1.5(j-i-1)) \cdot 2^{-t} \|\mathbf{a}_j\|_2$$

and since

$$\sum_{i=1}^{j-1} (j-i+\frac{7}{6})^2 \leq \sum_{i=1}^{n-1} (i+\frac{7}{6})^2 = \frac{1}{3}(n-1)((n+\frac{3}{2})^2 + \frac{11}{6}) < \frac{1.0001}{3} n(n+1)^2$$

we can bound the  $L_2$ -norm of each column of  $\mathbf{S}$  by

$$0.87 \cdot n^{\frac{1}{2}} (n+1) 2^{-t} \|\mathbf{a}_j\|_2.$$

After a similar estimation using the bounds for single precision we obtain

$$\|\mathbf{U}\bar{\mathbf{R}}\|_E < \begin{cases} 0.87 \cdot n^{\frac{1}{2}}(n+1)2^{-t}\|\mathbf{A}\|_E & \text{d.p.acc.} \\ 0.87 \cdot n^{\frac{1}{2}}(n+1+2.5m)2^{-t}\|\mathbf{A}\|_E & \text{s.p.} \end{cases}$$

From (6.5) follows that when accumulating inner-products in double precision we have the a posteriori estimation

$$\|\mathbf{I} - \bar{\mathbf{Q}}^T \bar{\mathbf{Q}}\|_2 \leq 2\|\mathbf{U}\|_2 \leq 1.74n^{\frac{1}{2}}(n+1)2^{-t}\|\mathbf{A}\|_E\|\bar{\mathbf{R}}^{-1}\|_2. \quad (6.8)$$

To estimate  $\|\bar{\mathbf{R}}^{-1}\|_2$  in (6.8) we use the identity

$$\bar{\mathbf{R}}^T \bar{\mathbf{R}} = (\mathbf{A} + \mathbf{E}_1)^T (\mathbf{A} + \mathbf{E}_1) - \bar{\mathbf{R}}^T (\mathbf{U} + \mathbf{U}^T) \bar{\mathbf{R}}$$

which can be derived from (5.3) and (6.5). We write this

$$\bar{\mathbf{R}}^T \bar{\mathbf{R}} = \mathbf{R}^T (\mathbf{I} + \mathbf{F}_2) \mathbf{R} \quad (6.9)$$

where

$$\mathbf{F}_2 = \mathbf{F}_1 - \mathbf{R}^{-T} (\mathbf{U} \bar{\mathbf{R}})^T \bar{\mathbf{R}} \mathbf{R}^{-1} - \mathbf{R}^{-T} \bar{\mathbf{R}}^T (\mathbf{U} \bar{\mathbf{R}}) \mathbf{R}^{-1}.$$

By (6.9) follows

$$\|\bar{\mathbf{R}} \mathbf{R}^{-1}\|_2^2 \leq 1 + \|\mathbf{F}_2\|_2 \quad (6.10)$$

and hence using (6.3) we have the inequality

$$\|\mathbf{F}_2\|_2 \leq 2\|\mathbf{E}_1\|_E \|\mathbf{R}^{-1}\|_2 + \|\mathbf{E}_1\|_E^2 \|\mathbf{R}^{-1}\|_2^2 + 2\|\mathbf{U} \bar{\mathbf{R}}\|_E \|\mathbf{R}^{-1}\|_2 (1 + \|\mathbf{F}_2\|_2)^{\frac{1}{2}}.$$

Using the derived estimates for  $\|\mathbf{E}_1\|_E$  and  $\|\mathbf{U} \bar{\mathbf{R}}\|_E$  and the inequality  $(n-1) \leq 0.3n^{\frac{1}{2}}(n+1)$  we get, assuming d.p.acc.

$$\|\mathbf{F}_2\|_2 \leq 2 \cdot 0.3 \cdot 1.5c + (0.3 \cdot 1.5 \cdot c)^2 + 2 \cdot 0.87c(1 + \|\mathbf{F}_2\|_2)^{\frac{1}{2}}$$

where

$$c = n^{\frac{1}{2}}(n+1) \cdot 2^{-t} \|\mathbf{A}\|_E \|\mathbf{R}^{-1}\|_2.$$

It follows that  $\|\mathbf{F}_2\|_2 \leq b$  where  $b$  is the positive root to the equation

$$b = 0.9 \cdot c + (0.45 \cdot c)^2 + 1.74 \cdot c(1 + b)^{\frac{1}{2}}.$$

Since  $b=1$  corresponds to  $c=1/3.419 \dots$  it is readily verified that if we require

$$\beta = 3.42c = 3.42n^{\frac{1}{2}}(n+1)2^{-t}\|\mathbf{A}\|_E\|\mathbf{R}^{-1}\|_2 < 1, \quad (6.12)$$

which is a stronger condition than (6.4), then  $\|\mathbf{F}_2\|_2 < \beta$  and from (6.9)

$$\|\bar{\mathbf{R}}^{-1}\|_2^2 \leq \frac{1}{1-\beta} \|\mathbf{R}^{-1}\|_2^2. \quad (6.13)$$

This inequality holds for single precision if we define  $\beta$  by

$$\beta = 3.42n^{\frac{1}{2}}(n+1+2.5m)2^{-t}\|\mathbf{A}\|_E\|\mathbf{R}^{-1}\|_2.$$

Finally from (6.8) and (6.13) we get the a priori estimate

---

<sup>1</sup> The notation  $\mathbf{R}^{-T}$  is used for  $(\mathbf{R}^{-1})^T$ .

$$\|\mathbf{I} - \bar{\mathbf{Q}}^T \bar{\mathbf{Q}}\|_2 \leq \frac{1.74}{(1-\beta)^{\frac{1}{2}}} n^{\frac{1}{2}}(n+1) \cdot 2^{-t} \|\mathbf{A}\|_{\mathcal{E}} \|\mathbf{R}^{-1}\|_2. \quad (6.14)$$

In [5] Rice, after extensive computational experiments, reports that when using the modified Gram–Schmidt method “pivoting results in a perceptible, but small (even negligible) improvement” in the orthogonality of the computed vectors. Our analysis gives a theoretical explanation of this. *The bound in (6.14) is independent of pivoting, and since*

$$\|\mathbf{A}\|_{\mathcal{E}} \|\mathbf{R}^{-1}\|_2 \leq n^{\frac{1}{2}} \kappa(\mathbf{A})$$

*the deviation from orthogonality of  $\bar{\mathbf{Q}}$  can be expected to depend essentially only on  $\kappa(\mathbf{A})$ .*

## 7. Condition of the least squares problem.

Before estimating the errors in the computed solution, we shall study the sensitivity to perturbations in  $\mathbf{A}$  and  $\mathbf{b}$  of the least squares problem, which, following a suggestion by Golub, we write

$$\begin{pmatrix} \mathbf{r} + \mathbf{A}\mathbf{x} \\ \mathbf{A}^T \mathbf{r} \end{pmatrix} = \begin{pmatrix} \mathbf{I} & \mathbf{A} \\ \mathbf{A}^T & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{r} \\ \mathbf{x} \end{pmatrix} = \begin{pmatrix} \mathbf{b} \\ \mathbf{0} \end{pmatrix} \quad (7.1)$$

Let the perturbed system be

$$\begin{pmatrix} \mathbf{I} & \mathbf{A} + \delta\mathbf{A} \\ \mathbf{A}^T + \delta\mathbf{A}^T & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{r} + \delta\mathbf{r} \\ \mathbf{x} + \delta\mathbf{x} \end{pmatrix} = \begin{pmatrix} \mathbf{b} + \delta\mathbf{b} \\ \mathbf{0} \end{pmatrix}$$

and let

$$\tilde{\mathbf{A}} = \mathbf{A} + \delta\mathbf{A} = \tilde{\mathbf{V}}\tilde{\mathbf{R}}, \quad \tilde{\mathbf{V}}^T \tilde{\mathbf{V}} = \mathbf{I} \quad (7.2)$$

be the factorization of the perturbed matrix. Then

$$\begin{pmatrix} \delta\mathbf{r} \\ \delta\mathbf{x} \end{pmatrix} = \begin{pmatrix} \mathbf{I} & \tilde{\mathbf{A}} \\ \tilde{\mathbf{A}}^T & \mathbf{0} \end{pmatrix}^{-1} \left[ - \begin{pmatrix} \mathbf{0} & \delta\mathbf{A} \\ \delta\mathbf{A}^T & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{r} \\ \mathbf{x} \end{pmatrix} + \begin{pmatrix} \delta\mathbf{b} \\ \mathbf{0} \end{pmatrix} \right].$$

Using (7.2) the inverse can readily be calculated ([4] p. 17) and we get

$$\begin{pmatrix} \delta\mathbf{r} \\ \tilde{\mathbf{R}}\delta\mathbf{x} \end{pmatrix} = \begin{pmatrix} \mathbf{I} - \tilde{\mathbf{V}}\tilde{\mathbf{V}}^T & \tilde{\mathbf{V}} \\ \tilde{\mathbf{V}}^T & -\mathbf{I} \end{pmatrix} \begin{pmatrix} -\delta\mathbf{A}\mathbf{x} + \delta\mathbf{b} \\ -\tilde{\mathbf{R}}^{-T}\delta\mathbf{A}^T\mathbf{r} \end{pmatrix} \quad (7.3)$$

The eigenvalues of the symmetric matrix  $(\mathbf{I} - \tilde{\mathbf{V}}\tilde{\mathbf{V}}^T)$  are all equal to 0 or 1. Hence  $\|\mathbf{I} - \tilde{\mathbf{V}}\tilde{\mathbf{V}}^T\|_2 \leq 1$  and (7.3) yields

$$\begin{pmatrix} \|\delta\mathbf{r}\|_2 \\ \|\tilde{\mathbf{R}}\delta\mathbf{x}\|_2 \end{pmatrix} \leq \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} \|\tilde{\mathbf{R}}^{-1}\|_2 \|\delta\mathbf{A}\|_2 \|\mathbf{r}\|_2 \\ \|\delta\mathbf{A}\|_2 \|\mathbf{x}\|_2 + \|\delta\mathbf{b}\|_2 \end{pmatrix} \quad (7.4)$$

To estimate  $\|\tilde{\mathbf{R}}^{-1}\|_2$ , let the factorization of  $\mathbf{A}$  be  $\mathbf{A} = \mathbf{Q}\mathbf{R}$ . Then we can write

$$\tilde{\mathbf{R}}^T \tilde{\mathbf{R}} = \tilde{\mathbf{A}}^T \tilde{\mathbf{A}} = (\mathbf{A} + \delta \mathbf{A})^T (\mathbf{A} + \delta \mathbf{A}) = \mathbf{R}^T (\mathbf{I} + \mathbf{K}) \mathbf{R}$$

where

$$\mathbf{K} = (\mathbf{Q}^T \delta \mathbf{A} \mathbf{R}^{-1})^T + \mathbf{Q}^T \delta \mathbf{A} \mathbf{R}^{-1} + (\delta \mathbf{A} \mathbf{R}^{-1})^T \delta \mathbf{A} \mathbf{R}^{-1}.$$

Put

$$\alpha = (\sqrt{2} + 1) \|\mathbf{R}^{-1}\|_2 \|\delta \mathbf{A}\|_2, \quad (7.5)$$

and assume that  $\alpha < 1$ . Then, it is easily shown that

$$\|\mathbf{K}\|_2 \leq 2 \|\mathbf{R}^{-1}\|_2 \|\delta \mathbf{A}\|_2 + \|\mathbf{R}^{-1}\|_2^2 \|\delta \mathbf{A}\|_2^2 < \alpha < 1,$$

and hence

$$\|\tilde{\mathbf{R}}^{-1}\|_2^2 \leq \|\mathbf{R}^{-1}\|_2 \|(\mathbf{I} + \mathbf{K})^{-1}\|_2 \|\mathbf{R}^{-T}\|_2 < \frac{1}{1 - \alpha} \|\mathbf{R}^{-1}\|_2^2.$$

Hence by (1.6)

$$\|\mathbf{A}\|_2 \|\tilde{\mathbf{R}}^{-1}\|_2 \leq (1 - \alpha)^{-1} \|\mathbf{R}\|_2 \|\mathbf{R}^{-1}\|_2 = (1 - \alpha)^{-1} \kappa(\mathbf{A}).$$

Since  $\|\delta \mathbf{x}\|_2 \leq \|\tilde{\mathbf{R}}^{-1}\|_2 \|\tilde{\mathbf{R}} \delta \mathbf{x}\|_2$  it now follows from (7.4)

$$\left( \begin{array}{c} \|\delta \mathbf{r}\|_2 \\ \|\mathbf{A}\|_2 \|\delta \mathbf{x}\|_2 \end{array} \right) \leq \left( \begin{array}{cc} \frac{\kappa(\mathbf{A})}{\sqrt{1 - \alpha}} & 1 \\ \frac{\kappa^2(\mathbf{A})}{1 - \alpha} & \frac{\kappa(\mathbf{A})}{\sqrt{1 - \alpha}} \end{array} \right) \left( \begin{array}{c} \|\mathbf{r}\|_2 \frac{\|\delta \mathbf{A}\|_2}{\|\mathbf{A}\|_2} \\ \|\mathbf{A}\|_2 \|\mathbf{x}\|_2 \frac{\|\delta \mathbf{A}\|_2}{\|\mathbf{A}\|_2} + \|\delta \mathbf{b}\|_2 \end{array} \right) \quad (7.6)$$

It is possible to construct examples where these bounds are nearly obtained.

The result (7.6) has the following interesting interpretations. Notice that

$$\mathbf{x} = \mathbf{R}^{-1} \mathbf{R} \mathbf{x} = \mathbf{R}^{-1} \mathbf{Q}^T \mathbf{A} \mathbf{x}$$

and thus

$$\|\mathbf{A}\|_2 \|\mathbf{x}\|_2 + \kappa(\mathbf{A}) \|\mathbf{r}\|_2 \leq \kappa(\mathbf{A}) (\|\mathbf{A} \mathbf{x}\|_2 + \|\mathbf{r}\|_2) \leq \sqrt{2} \kappa(\mathbf{A}) \|\mathbf{b}\|_2.$$

Then from the first row of (7.6) follows

$$\frac{\|\delta \mathbf{r}\|_2}{\|\mathbf{b}\|_2} \leq \sqrt{2} \frac{\kappa(\mathbf{A})}{\sqrt{1 - \alpha}} \frac{\|\delta \mathbf{A}\|_2}{\|\mathbf{A}\|_2} + \frac{\|\delta \mathbf{b}\|_2}{\|\mathbf{b}\|_2}. \quad (7.7)$$

Hence as condition number for the decomposition of  $\mathbf{b}$  into orthogonal components  $\mathbf{A} \mathbf{x}$  and  $\mathbf{r}$  we can take  $\sqrt{2}(1 - \alpha)^{-1} \kappa(\mathbf{A})$ .

If we assume that  $\|\mathbf{x}\|_2 \neq 0$ , then we can also derive a condition number for the determination of  $\mathbf{x}$ . From the second row (7.6) follows



$$\frac{\|\delta \mathbf{x}\|_2}{\|\mathbf{x}\|_2} \leq \frac{\kappa(\mathbf{A})}{\sqrt{1-\alpha}} \left( 1 + \frac{\kappa(\mathbf{A})}{\sqrt{1-\alpha}} \frac{\|\mathbf{r}\|_2}{\|\mathbf{A}\|_2 \|\mathbf{x}\|_2} \right) \frac{\|\delta \mathbf{A}\|_2}{\|\mathbf{A}\|_2} + \frac{\kappa(\mathbf{A})}{\sqrt{1-\alpha}} \frac{\|\delta \mathbf{b}\|_2}{\|\mathbf{A}\|_2 \|\mathbf{x}\|_2}.$$

Notice in particular the importance of the ratio  $\|\mathbf{r}\|_2/\|\mathbf{A}\|_2\|\mathbf{x}\|_2$ . If the equations are nearly compatible in the sense that  $\|\mathbf{r}\|_2 \ll \|\mathbf{A}\|_2\|\mathbf{x}\|_2$ , then the situation is much the same when  $m > n$  as when  $m = n$ , and  $\kappa(\mathbf{A})$  is an approximate condition number for the determination of  $\mathbf{x}$ . However, if  $\|\mathbf{r}\|_2$  is of the same magnitude as  $\|\mathbf{A}\|_2\|\mathbf{x}\|_2$  then the condition number (in an ill-conditioned case) is more like  $\kappa^2(\mathbf{A})$ .

### 8. Errors in the computed solution.

As the computed solution we take

$$\bar{\mathbf{r}} = \bar{\mathbf{b}}^{(n+1)}, \quad \bar{\mathbf{x}} = fl((\bar{\mathbf{R}}')^{-1}\bar{\mathbf{y}}'). \quad (8.1)$$

If we define  $\tilde{\mathbf{R}}$  so that  $\bar{\mathbf{x}}$  is the exact solution of

$$\tilde{\mathbf{R}}\bar{\mathbf{x}} = (\bar{\mathbf{R}} + \delta\bar{\mathbf{R}})\bar{\mathbf{x}} = \bar{\mathbf{y}},$$

then using definitions (5.4) and (5.5) we can write (8.1)

$$\bar{\mathbf{r}} = (\mathbf{I} - \bar{\mathbf{Q}}\tilde{\mathbf{Q}}^T)\mathbf{b} - \mathbf{e}_2, \quad \bar{\mathbf{x}} = \tilde{\mathbf{R}}^{-1}(\tilde{\mathbf{Q}}^T\mathbf{b} + \mathbf{e}_3). \quad (8.2)$$

Substituting  $\mathbf{b} = \mathbf{r} + \mathbf{A}\mathbf{x}$  in (8.2) we get

$$\begin{pmatrix} \bar{\mathbf{r}} \\ \bar{\mathbf{x}} \end{pmatrix} = \begin{pmatrix} \mathbf{I} - \bar{\mathbf{Q}}\tilde{\mathbf{Q}}^T & (\mathbf{I} - \bar{\mathbf{Q}}\tilde{\mathbf{Q}}^T)\mathbf{A} \\ \tilde{\mathbf{R}}^{-1}\tilde{\mathbf{Q}}^T & \tilde{\mathbf{R}}^{-1}\tilde{\mathbf{Q}}^T\mathbf{A} \end{pmatrix} \begin{pmatrix} \mathbf{r} \\ \mathbf{x} \end{pmatrix} + \begin{pmatrix} -\mathbf{e}_2 \\ \tilde{\mathbf{R}}^{-1}\mathbf{e}_3 \end{pmatrix} \quad (8.3)$$

Now since  $\mathbf{A}^T\mathbf{r} = 0$  we have from (5.3)

$$\bar{\mathbf{Q}}^T\mathbf{r} = \bar{\mathbf{R}}^{-T}(\mathbf{A} + \mathbf{E}_1)^T\mathbf{r} = \bar{\mathbf{R}}^{-T}\mathbf{E}_1^T\mathbf{r}$$

and from lemma 5.1 it follows

$$\tilde{\mathbf{Q}}^T = (\mathbf{I} + \mathbf{U}^T)^{-1}\bar{\mathbf{Q}}^T, \quad \tilde{\mathbf{Q}} = \bar{\mathbf{Q}}(\mathbf{I} + \mathbf{U}^T)^{-1}.$$

Thus the errors in the computed solution can be written

$$\begin{pmatrix} \bar{\mathbf{r}} - \mathbf{r} \\ \bar{\mathbf{x}} - \mathbf{x} \end{pmatrix} = \begin{pmatrix} -\tilde{\mathbf{Q}}\bar{\mathbf{R}}^{-T}\mathbf{E}_1^T & (\mathbf{I} - \bar{\mathbf{Q}}\tilde{\mathbf{Q}}^T)\mathbf{A} \\ (\bar{\mathbf{R}}^T(\mathbf{I} + \mathbf{U}^T)\tilde{\mathbf{R}})^{-1}\mathbf{E}_1^T & \tilde{\mathbf{R}}^{-1}(\tilde{\mathbf{Q}}^T\mathbf{A} - \tilde{\mathbf{R}}) \end{pmatrix} \begin{pmatrix} \mathbf{r} \\ \mathbf{x} \end{pmatrix} + \begin{pmatrix} -\mathbf{e}_2 \\ \tilde{\mathbf{R}}^{-1}\mathbf{e}_3 \end{pmatrix} \quad (8.4)$$

We shall now estimate the  $L_2$ -norms of the submatrices appearing in (8.4). We assume for shortness that the  $fl_2(\cdot)$  mode of computation is used both in the decomposition and in the backsubstitution. The estimates relevant when using single precision are not very different and can be derived in a similar way.

Since according to (5.2)  $\|\tilde{\mathbf{Q}}\|_E \leq n^{\frac{1}{2}}$  we have immediately from (5.6) and (5.7)

$$\|-\tilde{\mathbf{Q}}\bar{\mathbf{R}}^{-T}\mathbf{E}_1^T\|_2 \leq 1.5 \cdot n^{\frac{1}{2}}(n-1) \cdot 2^{-t} \|\bar{\mathbf{R}}^{-1}\|_2 \|\mathbf{A}\|_E, \quad (8.5)$$

$$\|(\mathbf{I} - \bar{\mathbf{Q}}\tilde{\mathbf{Q}}^T)\mathbf{A}\|_2 \leq 3.25 \cdot (n-1) \cdot 2^{-t} \|\mathbf{A}\|_E. \quad (8.6)$$

We rewrite the remaining two submatrices in (8.4) as

$$(\bar{\mathbf{R}}^T(\mathbf{I} + \mathbf{U}^T)\tilde{\mathbf{R}})^{-1}\mathbf{E}_1^T = \tilde{\mathbf{R}}^{-1}\bar{\mathbf{R}}'(\bar{\mathbf{R}}^T(\mathbf{I} + \mathbf{U}^T)\bar{\mathbf{R}})^{-1}\mathbf{E}_1^T, \quad (8.7)$$

$$\tilde{\mathbf{R}}^{-1}(\tilde{\mathbf{Q}}^T\mathbf{A} - \tilde{\mathbf{R}}) = \tilde{\mathbf{R}}^{-1}\bar{\mathbf{R}}'(\bar{\mathbf{R}}^{-1}(\tilde{\mathbf{Q}}^T\mathbf{A} - \bar{\mathbf{R}}) - \bar{\mathbf{R}}^{-1}\delta\bar{\mathbf{R}}). \quad (8.8)$$

Now from (6.1) and (6.5) it follows

$$\bar{\mathbf{R}}^T(\mathbf{I} + \mathbf{U}^T)\bar{\mathbf{R}} = \bar{\mathbf{A}}^T\bar{\mathbf{A}} - \bar{\mathbf{R}}^T\mathbf{U}\bar{\mathbf{R}} = \mathbf{R}^T(\mathbf{I} + \mathbf{F}_3)\mathbf{R}$$

where

$$\mathbf{F}_3 = \mathbf{F}_1 - \mathbf{R}^{-T}\bar{\mathbf{R}}^T(\mathbf{U}\bar{\mathbf{R}})\mathbf{R}^{-1}.$$

Comparing this with (6.9) and (6.13) we observe that certainly  $\|\mathbf{F}_3\|_2 \leq \beta$  and

$$\|(\bar{\mathbf{R}}^T(\mathbf{I} + \mathbf{U}^T)\bar{\mathbf{R}})^{-1}\|_2 \leq \frac{1}{1-\beta} \|\mathbf{R}^{-1}\|_2^2$$

If we put

$$\|\tilde{\mathbf{R}}^{-1}\bar{\mathbf{R}}\|_2 = \|( \mathbf{I} + \bar{\mathbf{R}}^{-1}\delta\bar{\mathbf{R}} )^{-1}\|_2 = 1 + \eta, \quad (8.9)$$

then using (5.6) we get from (8.7)

$$\|(\bar{\mathbf{R}}^T(\mathbf{I} + \mathbf{U}^T)\tilde{\mathbf{R}})^{-1}\mathbf{E}_1^T\|_2 \leq \frac{1+\eta}{1-\beta} 1.5 \cdot (n-1) \cdot 2^{-t} \|\mathbf{R}^{-1}\|_2^2 \|\mathbf{A}\|_E, \quad (8.10)$$

and from (5.5) and (8.8)

$$\|\tilde{\mathbf{R}}^{-1}(\tilde{\mathbf{Q}}^T\mathbf{A} - \tilde{\mathbf{R}})\|_2 \leq (1+\eta)(\|\mathbf{E}_3\|_E \|\bar{\mathbf{R}}^{-1}\|_2 + \|\bar{\mathbf{R}}^{-1}\delta\bar{\mathbf{R}}\|_2). \quad (8.11)$$

Wilkinson shows in [6] pp. 99–104 that if double precision accumulation is used, then

$$\|\delta\bar{\mathbf{R}}'\|_E \leq (1.001 \cdot 2^{-t} + \frac{3}{2}(n+1)2^{-2t_2})\|\bar{\mathbf{R}}'\|_E \leq 2^{-t_1}\|\bar{\mathbf{R}}'\|_E. \quad (8.12)$$

Since  $\|\bar{\mathbf{R}}'\|_E \leq n^{\frac{1}{2}}\|\bar{\mathbf{R}}'\|_2$  we have

$$\|\bar{\mathbf{R}}^{-1}\delta\bar{\mathbf{R}}\|_2 = \|(\bar{\mathbf{R}}')^{-1}\delta\bar{\mathbf{R}}'\|_2 \leq n^{\frac{1}{2}} \cdot 2^{-t_1\kappa}(\bar{\mathbf{R}}'). \quad (8.13)$$

But (8.12) is also valid without primes and as

$$\|\bar{\mathbf{R}}\|_2 \leq (1+\beta)^{\frac{1}{2}}\|\mathbf{R}\|_2$$

when  $\beta < 1$ , it follows

$$\|\bar{\mathbf{R}}^{-1}\delta\bar{\mathbf{R}}\|_2 \leq n^{\frac{1}{2}} 2^{-t_1\kappa}(\bar{\mathbf{R}}) \leq (1+\beta)^{\frac{1}{2}} n^{\frac{1}{2}} 2^{-t_1} \|\bar{\mathbf{R}}^{-1}\|_2 \|\mathbf{A}\|_E. \quad (8.14)$$

Substituting in (8.11) from (5.8) and (8.14) we can show that

$$\|\tilde{\mathbf{R}}^{-1}(\tilde{\mathbf{Q}}^T \mathbf{A} - \tilde{\mathbf{R}})\|_2 \leq (1 + \eta) \cdot 1.9n^{\frac{1}{2}}(n+1) \cdot 2^{-t} \|\bar{\mathbf{R}}^{-1}\|_2 \|\mathbf{A}\|_E. \quad (8.15)$$

For the estimation of  $\eta$  (8.9) yields

$$\eta \leq \|\bar{\mathbf{R}}^{-1} \delta \bar{\mathbf{R}}\|_2 / (1 - \|\bar{\mathbf{R}}^{-1} \delta \bar{\mathbf{R}}\|_2). \quad (8.16)$$

From (8.14), (6.13) and (6.12) we get

$$\|\bar{\mathbf{R}}^{-1} \delta \mathbf{R}\|_2 \leq 1.06 \left( \frac{1+\beta}{1-\beta} \right)^{\frac{1}{2}} \cdot n^{\frac{1}{2}} 2^{-t} \|\mathbf{R}^{-1}\|_2 \|\mathbf{A}\|_E \leq \frac{1.06\beta}{3.42(n+1)} \left( \frac{1+\beta}{1-\beta} \right),$$

and if we make the additional assumption  $\beta \leq 0.9$  then

$$\|\bar{\mathbf{R}}^{-1} \delta \bar{\mathbf{R}}\|_2 \leq 1.22/(n+1) \leq 0.41.$$

Using this to estimate the denominator in (8.16) we can show that

$$\eta \leq 1.8n^{\frac{1}{2}} 2^{-t} \kappa(\bar{\mathbf{R}}'), \quad (8.17)$$

$$\eta \leq \frac{0.53}{n+1} \left( \frac{1+\beta}{1-\beta} \right)^{\frac{1}{2}} \beta. \quad (8.18)$$

Although it is not generally true that  $\kappa(\bar{\mathbf{R}}') \leq \kappa(\mathbf{A})$ , often  $\kappa(\bar{\mathbf{R}}') \ll \kappa(\bar{\mathbf{R}}) \approx \kappa(\mathbf{A})$ . In fact for many very ill-conditioned systems  $\kappa(\bar{\mathbf{R}}')$  will be of order unity. Then (8.14) and consequently (8.18) will considerably overestimate the error from the backsubstitution, which nearly always is of little importance.

We now take norms in (8.4) and use (8.5), (8.6), (8.10) and (8.15) together with the estimates for  $\|\mathbf{e}_2\|_2$  and  $\|\mathbf{e}_3\|_2$  in (5.7) and (5.8). Eliminating  $\|\bar{\mathbf{R}}^{-1}\|_2$  by means of the inequality

$$\|\bar{\mathbf{R}}^{-1}\|_2 \|\mathbf{A}\|_2 \leq (1-\beta)^{-1} \cdot \kappa(\mathbf{A})$$

and invoking the trivial inequalities

$$n-1 \leq 0.3n^{\frac{1}{2}}(n+1), \quad n \leq \frac{2}{3}n^{\frac{1}{2}}(n+1)$$

we arrive at the result

$$\left( \begin{array}{c} \|\bar{\mathbf{r}} - \mathbf{r}\|_2 \\ \|\mathbf{A}\|_2 \|\bar{\mathbf{x}} - \mathbf{x}\|_2 \end{array} \right) \leq \left( \begin{array}{cc} \frac{0.79\kappa(\mathbf{A})}{\sqrt{1-\beta}} & 0.81 \\ \frac{0.24\kappa^2(\mathbf{A})}{1-\beta} & \frac{\kappa(\mathbf{A})}{\sqrt{1-\beta}} \end{array} \right) \left( \begin{array}{c} \|\mathbf{r}\|_2 \frac{\|\mathbf{A}\|_E}{\|\mathbf{A}\|_2} \\ \|\mathbf{A}\|_2 \|\mathbf{x}\|_2 \frac{\|\mathbf{A}\|_E}{\|\mathbf{A}\|_2} + \|\mathbf{b}\|_2 \end{array} \right) \cdot f(n) 2^{-t} \quad (8.19)$$

where

$$f(n) = (1 + \eta) 1.9n^{\frac{1}{2}}(n+1).$$

We observe that (8.19) gives a smaller bound for the error than that obtained from (7.6) when

$$\|\delta \mathbf{A}\|_2 = f(n) \cdot 2^{-t} \|\mathbf{A}\|_E, \quad \|\delta \mathbf{b}\|_2 = f(n) 2^{-t} \|\mathbf{b}\|_2,$$

as this according to (6.12) and (7.5) implies

$$\alpha = (1 + \eta)(\sqrt{2} + 1) \cdot 1.9n^{\frac{1}{2}}(n + 1) \cdot 2^{-t} \|\mathbf{R}^{-1}\|_2 \|\mathbf{A}\|_E > \beta.$$

As the factor  $(1 + \eta)$  usually is of no importance, we can say that *the deviation of the computed solution from the true solution is roughly less than the deviation resulting from some perturbation  $\delta \mathbf{A}$ ,  $\delta \mathbf{b}$  such that*

$$\|\delta \mathbf{A}\|_E / \|\mathbf{A}\|_E \approx \|\delta \mathbf{b}\|_2 / \|\mathbf{b}\|_2 \approx 2 \cdot n^{3/2} \cdot 2^{-t}.$$

Moreover when  $\mathbf{A}$  is ill-conditioned we can, as remarked in the end of section 5, expect a much better result than this.

We finally observe that when the modified procedure is used for computing  $\mathbf{R}$ , then the classical procedure should on no account be used for calculating  $\mathbf{y}$ . If this was done we would have

$$\mathbf{y} \approx \bar{\mathbf{Q}}^T \mathbf{b} = (\bar{\mathbf{Q}} - \tilde{\mathbf{Q}})^T \mathbf{b} + \tilde{\mathbf{Q}}^T \mathbf{b}$$

and consequently

$$\|\mathbf{y} - \tilde{\mathbf{Q}}^T \mathbf{b}\|_2 \approx \|\bar{\mathbf{Q}} - \tilde{\mathbf{Q}}\|_2 \|\mathbf{b}\|_2 = \|\mathbf{U} \tilde{\mathbf{Q}}\|_2 \|\mathbf{b}\|_2 = \text{const.} \kappa(\mathbf{A}) 2^{-t} \|\mathbf{b}\|_2.$$

It can be verified that if  $\|\mathbf{r}\|_2 \ll \|\mathbf{b}\|_2$  this means that the errors in the computed solution usually would be multiplied by  $\kappa(\mathbf{A})!$

## 9. Conclusions.

The modified Gram-Schmidt procedure has been successfully used since 1960 as part of a data processing system for optical spectra [8]. It has about the same high numerical stability as other comparable algorithms based on nearly orthogonal transformations, [7] ch. 3, and generally should be preferred to the classical Gram-Schmidt procedure. Reorthogonalization is then never necessary, and the solution is obtained in approximately  $mn^2$  multiplications. Pivoting can easily be done, but is not required for stability. The storage requirement is approximately  $mn + n^2/2$ . This should be compared to  $2mn^2 + 4n^3/3$  s.p. multiplications and  $n^2$  storage locations needed for forming and solving the normal equations in double precision.

More accurate solutions to linear least squares problems can be obtained by using iterative refinement. This was first proposed by Golub [2] and used also in [1]. Different schemes for the refinement are discussed

in [3]. The results obtained so far, however, are only satisfactory when the equations are almost compatible, i.e. when

$$\kappa^2(\mathbf{A})\|\mathbf{r}\|_2 < \|\mathbf{A}\|_2\|\mathbf{x}\|_2.$$

In a forthcoming paper, where also numerical examples will be provided, the author will apply iterative refinement to linear least squares problems by calculating the residuals to the system (7.1) in double precision, and use the inverse

$$\begin{pmatrix} \mathbf{I} & \mathbf{A} \\ \mathbf{A}^T & \mathbf{0} \end{pmatrix}^{-1} \approx \begin{pmatrix} \mathbf{I} - \bar{\mathbf{Q}}\tilde{\mathbf{Q}}^T & \tilde{\mathbf{Q}}\bar{\mathbf{R}}^{-T} \\ \bar{\mathbf{R}}^{-1}\tilde{\mathbf{Q}}^T & -\bar{\mathbf{R}}^{-1}(\mathbf{I} + \mathbf{U})^{-T}\bar{\mathbf{R}}^{-T} \end{pmatrix}$$

to solve for the increments in  $\mathbf{r}$  and  $\mathbf{x}$ . This inverse results from a slight generalization of the modified Gram-Schmidt procedure. It will be shown that the iterates generally will converge to the exact solution correctly rounded even when  $\kappa^2(\mathbf{A})\|\mathbf{r}\|_2$  is of the same order of magnitude as  $2^t\|\mathbf{A}\|_2\|\mathbf{x}\|_2$ .

#### REFERENCES

1. Bauer, F. L., *Elimination with Weighted Row Combinations for Solving Linear Equations and Least Squares Problems*, Num. Math. 7 (1965), 338-352.
2. Golub, G. H., *Numerical Methods for Solving Linear Least Squares Problems*, Num. Math. 7 (1965), 206-216.
3. Golub, G. H. and Wilkinson, J. H., *Note on the Iterative Refinement of Least Squares Solution*, Num. Math. 9 (1966), 139-148.
4. Householder, A. S., *The Theory of Matrices in Numerical Analysis*, New York: Blaisdell 1964.
5. Rice, J. R., *Experiments on Gram-Schmidt Orthogonalization*, Math. Comp. 20 (1966), 325-328.
6. Wilkinson, J. H., *Rounding Errors in Algebraic Processes*, London: H.M.S.O. 1963.
7. Wilkinson, J. H.: *The Algebraic Eigenvalue Problem*, Oxford: Clarendon Press 1965.
8. Åslund, N., *A Data Processing System for Spectra of Diatomic Molecules*, Arkiv Fysik 30 (1965), 377-396.

INSTITUTE FOR INFORMATION PROCESSING  
ROYAL INSTITUTE OF TECHNOLOGY  
STOCKHOLM, SWEDEN