# ITERATIVE REFINEMENT
# OF LINEAR LEAST SQUARES SOLUTIONS I

ÅKE BJÖRCK

## Abstract.

An iterative procedure is developed for reducing the rounding errors in the computed least squares solution to an overdetermined system of equations $Ax = b$, where $A$ is an $m \times n$ matrix $(m \geqq n)$ of rank $n$. The method relies on computing accurate residuals to a certain augmented system of linear equations, by using double precision accumulation of inner products. To determine the corrections, two methods are given, based on a matrix decomposition of $A$ obtained either by orthogonal Householder transformations or by a modified Gram–Schmidt orthogonalization. It is shown that the rate of convergence in the iteration is independent of the right hand side, $b$, and depends linearly on the condition number, $\varkappa(A)$, of the rectangular matrix $A$. The limiting accuracy achieved will be approximately the same as that obtained by a double precision factorization.

In a second part of this paper the case when $x$ is subject to linear constraints and/or $A$ has rank less than $n$ is covered. Here also ALGOL-programs embodying the derived algorithms will be given.

## 1. Introduction.

Let $Ax = b$ be a given overdetermined system of linear equations where $A$ is an $m \times n$ matrix $(m \geqq n)$ and $b$ is a vector. A vector $x$ which minimizes $\|b - Ax\|_2$ is called a least squares solution to the system.

Least squares problems are often ill-conditioned. Rounding errors may then seriously contaminate the solution. For the linear equation case $(m = n)$ Wilkinson [8] has proposed the following process of iterative refinement for reducing the rounding errors: Compute the sequence of vectors $x^{(s)}$, $s = 0, 1, 2, \ldots$ defined by

$$\begin{aligned} x^{(0)} &= 0, \quad r^{(s)} = b - Ax^{(s)}, \\ \delta x^{(s)} &= A^{-1} r^{(s)}, \quad x^{(s+1)} = x^{(s)} + \delta x^{(s)}. \end{aligned} \qquad (1.1)$$

Here the residual vector $r^{(s)}$ is computed using double precision accumulation of inner products. Single precision is used in all other steps. In particular, the corrections $\delta x^{(s)}$ are computed using a suitable single precision factorization of $A$.

The performance of this process in floating point arithmetic has been analysed by Moler [7]. It has also been embodied in an ALGOL procedure by Martin, Peters and Wilkinson [6] for the special case when $A$ is positive definite. Now assume that $m > n$ and that $A$ has rank $n$. Then there exists a decomposition $A = QR$, where $R$ is upper triangular and $Q^T Q = I$, and it is well known that the least squares solution is given by

$$x = R^{-1} Q^T b \ .$$

Thus it is natural to use the refinement procedure (1.1) with

$$\delta x^{(s)} = R^{-1} Q^T r^{(s)} \tag{1.2}$$

for refining least squares solutions. This was first pointed out by Golub [4] and used also by Bauer in [1]. However it has been shown by Golub and Wilkinson in [5] that this process works satisfactorily only when the overdetermined system is nearly compatible.

In part I of this paper a procedure for the iterative refinement of least squares solutions without this restriction will be developed and analysed. In chapter 2 we describe the procedure as a special case of (1.1) and show, why the analysis in [7] is too general to be of any use here. In chapter 3 we formulate the assumptions on the arithmetic underlying our analysis. Any method for solving least squares problems can, after modification, be used in our procedure to solve for the corrections. In chapter 4 we analyse the errors which are independent of the particular method chosen. In recent papers it has been pointed out that methods related to an orthogonal triangularisation of the matrix $A$ either by Householder transformations [4] or by a modified Gram–Schmidt procedure [2], have several advantages over the classical method of solving the normal equations. In chapter 5 an algorithm is derived for the Householder method and a detailed error analysis for a single step is carried out. In chapter 6 these results are used to derive estimates for the rate of convergence and the limiting accuracy. Finally in chapter 7 the corresponding algorithm for the Gram–Schmidt method is analysed.

In part II the case when $x$ is subject to linear constraints and/or $A$ has rank less than $n$ is covered. Here ALGOL-programs embodying the derived algorithms will also be given.

## 2. The refinement procedure.

It is well known that a least squares solution is characterized by the property that the residual vector $r = b - Ax$ is orthogonal to the columns

of $A$. Thus for the unknowns $r$ and $x$ we have the system of $(m+n)$ equations

$$\begin{pmatrix} I & A \\ A^T & 0 \end{pmatrix} \begin{pmatrix} r \\ x \end{pmatrix} = \begin{pmatrix} b \\ 0 \end{pmatrix} \tag{2.1}$$

If we assume that the rank of $A$ equals $n$, then this system is non-singular and determines $r$ and $x$ uniquely.

We now propose to use the iterative procedure (1.1) for the refinement of the solution to the *augmented system* (2.1). As initial approximation we take

$$r^{(0)} = 0, \quad x^{(0)} = 0 .$$

The $s$th iteration consists of the three steps:

(i) Compute the residuals

$$\begin{pmatrix} f^{(s)} \\ g^{(s)} \end{pmatrix} = \begin{pmatrix} b \\ 0 \end{pmatrix} - \begin{pmatrix} I & A \\ A^T & 0 \end{pmatrix} \begin{pmatrix} r^{(s)} \\ x^{(s)} \end{pmatrix} \tag{2.2}$$

Here inner products are accumulated in double precision.

(ii) Solve for the corrections $\delta r^{(s)}$ and $\delta x^{(s)}$ from

$$\begin{pmatrix} I & A \\ A^T & 0 \end{pmatrix} \begin{pmatrix} \delta r^{(s)} \\ \delta x^{(s)} \end{pmatrix} = \begin{pmatrix} f^{(s)} \\ g^{(s)} \end{pmatrix} \tag{2.3}$$

Note that when $s \neq 0$ we generally have $g^{(s)} \neq 0$. Thus, modifications to the usual methods for solving linear least squares problems are necessary in order to solve (2.3).

(iii) Add the corrections

$$\begin{pmatrix} r^{(s+1)} \\ x^{(s+1)} \end{pmatrix} = \begin{pmatrix} r^{(s)} \\ x^{(s)} \end{pmatrix} + \begin{pmatrix} \delta r^{(s)} \\ \delta x^{(s)} \end{pmatrix} \tag{2.4}$$

If we put $r^{(s)} = 0$ for all $s$, then this procedure degenerates into

$$f^{(s)} = b - A x^{(s)}, \quad \delta x^{(s)} = R^{-1} Q^T f^{(s)} ,$$

which is precisely the scheme (1.2) proposed by Golub. This indicates that when the overdetermined system is compatible, the final performance of the two schemes should be the same.

Since the proposed procedure (2.2)–(2.4) is a special case of the general scheme (1.1), the analysis in [7] applies. Because of the special structure of the matrix

$$B = \begin{pmatrix} I & A \\ A^T & 0 \end{pmatrix} \tag{2.5}$$

this analysis, however, does not give a true assessment of the perfor-

mance. According to this general analysis the condition number $\varkappa(\boldsymbol{B})$ should play an essential role. We now make the following observation. If we scale the matrix $\boldsymbol{A}$ so that

$$\boldsymbol{A} := \alpha^{-1}\boldsymbol{A}, \quad \alpha = 2^{-q} \ (q \text{ integer})$$

then the fractional parts of the floating point numbers in our algorithm remain the same. The iterations are now, however, associated with the matrix

$$\boldsymbol{B}_\alpha = \begin{pmatrix} \alpha\boldsymbol{I} & \boldsymbol{A} \\ \boldsymbol{A}^T & 0 \end{pmatrix}.$$

*We will show that the condition number of $\boldsymbol{B}_\alpha$ varies considerably with $\alpha$.* Let $\lambda$ be an eigenvalue of $\boldsymbol{B}_\alpha$ and $(\boldsymbol{x},\boldsymbol{y})^T$ the corresponding eigenvector. Then

$$\alpha\boldsymbol{x} + \boldsymbol{A}\boldsymbol{y} = \lambda\boldsymbol{x}$$
$$\boldsymbol{A}^T\boldsymbol{x} \quad\ = \lambda\boldsymbol{y}$$

and it follows that

$$\alpha\lambda\boldsymbol{y} + \boldsymbol{A}^T\boldsymbol{A}\boldsymbol{y} = \lambda^2\boldsymbol{y} \ .$$

Now $\boldsymbol{y} \neq 0$ implies that $\boldsymbol{y}$ is an eigenvector and $(\lambda^2 - \alpha\lambda)$ is an eigenvalue of $\boldsymbol{A}^T\boldsymbol{A}$. On the other hand $\boldsymbol{y} = 0$ implies that

$$\boldsymbol{A}^T\boldsymbol{x} = 0, \quad \alpha\boldsymbol{x} = \lambda\boldsymbol{x}, \quad \boldsymbol{x} \neq 0 \ .$$

Thus the eigenvalues of $\boldsymbol{B}_\alpha$ are

$$\lambda(\boldsymbol{B}_\alpha) = \begin{cases} \dfrac{\alpha}{2} \pm \left( \dfrac{\alpha^2}{4} + \sigma_i{}^2 \right)^{\frac{1}{2}} \\[2mm] \alpha \end{cases}$$

where $\sigma_i{}^2$, $i = 1, 2, \ldots, n$ are the eigenvalues of $\boldsymbol{A}^T\boldsymbol{A}$ and the eigenvalue $\alpha$ has multiplicity $(m - n)$. From this it can be deduced that

$$\min_\alpha \varkappa(\boldsymbol{B}_\alpha) = \tfrac{1}{2} + \left( \tfrac{1}{4} + 2\frac{\sigma_{\max}^2}{\sigma_{\min}^2} \right)^{\frac{1}{2}} \leq 2\varkappa(\boldsymbol{A}) \ , \tag{2.6}$$

where the minimum is attained for $\alpha = 2^{-\frac{1}{2}}\sigma_{\min}$. Here $\varkappa(\boldsymbol{A})$ is the Euclidian condition number for the rectangular matrix $\boldsymbol{A}$,

$$\varkappa(\boldsymbol{A}) = \sigma_{\max}/\sigma_{\min} \geq 1 \ .$$

Furthermore, if we take $\alpha = 2^{-\frac{1}{2}}\sigma_{\max}$, then

$$\varkappa(\boldsymbol{B}_\alpha) > \varkappa^2(\boldsymbol{A}) \left( \frac{1}{2\sqrt{2}} + \frac{3}{2\sqrt{2}} \right) \left( \frac{1}{2\sqrt{2}} + \frac{1}{2\sqrt{2}} \right) = \varkappa^2(\boldsymbol{A}) \ . \tag{2.7}$$

This shows that a special analysis is needed.

## 3. Preliminaries for the error analysis.

We assume in the following that normalized floating point arithmetic with a single (double) precision mantissa of $t_1$ ($t_2$) digits of base $\beta$ is used. By accumulating an inner product in double precision we mean that multiplications which produce a $t_2$-digit product from $t_1$-digit numbers and double precision addition is used, the result finally being rounded to single precision.

More precisely, we assume that if $x$ and $y$ are single precision numbers and $z$ a double precision number then

$$\begin{aligned}
\mathrm{fl}(x'op'y) &= (x'op'y)(1+\delta_1), \quad 'op' = +,-,\times,/ , \\
\mathrm{fl}_2(x \times y + z) &= (x \times y + z)(1+\delta_2)
\end{aligned} \tag{3.1}$$

where

$$|\delta_i| \leqq \varepsilon_i, \quad \varepsilon_i = \beta^{1-t_i}, \quad i = 1,2 . \tag{3.2}$$

If the single precision machine operations are rounded rather than chopped, then $\varepsilon_1$ can usually be halved. For a detailed discussion, see Wilkinson [8]. We furthermore assume that all quantities remain within the permitted range of the computer.

It has been pointed out that the accumulation of inner products is essential only in the calculation of the residuals. For convenience we will, however, assume that this is done also when computing the decomposition of $A$, and when solving the systems (2.3). This is not an essential restriction. If single precision inner products are used in these steps, most derived error bounds will only increase by a factor less than $m$.

## 4. Rounding errors in the residuals.

We consider here the rounding errors introduced in the calculation of the residuals (2.2). Here and in the following we distinguish computed quantities by using a bar. To make for easier accumulation of errors we assume that $m\varepsilon_2 \leqq 0.1$ and define

$$\varepsilon_2' = 1.06(1+\varepsilon_1)\varepsilon_2 . \tag{4.1}$$

Following Moler [7] (p. 318) we have for the computed $i$th component of $f^{(s)} = b - \bar{r}^{(s)} - A\bar{x}^{(s)}$

$$\bar{f}_i^{(s)} = (1+\delta)[-\bar{r}_i^{(s)}(1+\eta_1) + (b_i(1+\eta_2) - t(1+\eta_3))(1+\eta_4)]$$

$$t = \sum_{j=1}^{n} a_{ij}\bar{x}_j^{(s)}(1+\gamma_j) ,$$

where

$$|\delta| \leqq \varepsilon_1, \ |\eta_i| \leqq \varepsilon_2, \ i = 1,2,3,4, \ |\gamma_j| \leqq 1.06(n-j+2)\varepsilon_2 .$$

If we denote the errors in the computed residuals by

$$\delta f^{(s)} = \bar{f}^{(s)} - f^{(s)} ,$$

then combining these results

$$\delta f_i^{(s)} = (b_i - \bar{r}_i^{(s)} - \sum_{j=1}^{n} a_{ij}\bar{x}_j^{(s)})\delta + (1+\delta)$$

$$[-\bar{r}_i^{(s)}\eta_1 + b_i((1+\eta_2)(1+\eta_4)-1) - \sum_{j=1}^{n} a_{ij}\bar{x}_j^{(s)}((1+\gamma_j)(1+\eta_3)(1+\eta_4)-1)] ,$$

and thus

$$\|\delta f^{(s)}\|_2 \leq \varepsilon_1 \|b - \bar{r}^{(s)} - A\bar{x}^{(s)}\|_2 +$$
$$(1+\varepsilon_1)\varepsilon_2(\|\bar{r}^{(s)}\|_2 + 2 \cdot 1.06\|b\|_2 + 1.06(n+3)\|A\|_2\|\bar{x}^{(s)}\|_2) .$$

Using $b = r + Ax$ and the estimate

$$\|b\|_2 \leq \|r\|_2 + \|A\|_2\|x\|_2$$

we obtain

$$\|\delta f^{(s)}\|_2 \leq 3\varepsilon_2'\|r\|_2 + (n+5)\varepsilon_2'\|A\|_2\|x\|_2 + \tag{4.2}$$
$$(\varepsilon_1+\varepsilon_2')\|r - \bar{r}^{(s)}\|_2 + (\varepsilon_1 + (n+3)\varepsilon_2')\|A\|_2\|x - \bar{x}^{(s)}\|_2 .$$

In the same way we can bound the rounding error in the computed residual vector $g^{(s)} = -A^T\bar{r}^{(s)}$. We obtain

$$\delta g_i^{(s)} = -\delta \sum_{i=1}^{m} a_{ij}\bar{r}_i^{(s)} - (1+\delta) \sum_{i=1}^{m} a_{ij}\bar{r}_i^{(s)}\gamma_i ,$$

where

$$|\gamma_i| \leq 1.06(m-i+2)\varepsilon_2, \quad |\delta| \leq \varepsilon_1 .$$

From this it follows that

$$\|\delta g^{(s)}\|_2 \leq \|A\|_2[(m+1)\varepsilon_2'\|r\|_2 + (\varepsilon_1 + (m+1)\varepsilon_2')\|r - \bar{r}^{(s)}\|_2] . \tag{4.3}$$

## 5. Error analysis of Householder's method.

In [4] Golub has described a method for solving linear least squares problems using a sequence of elementary orthogonal transformations of Householder type:

$$P^{(r)} = I - 2w^{(r)}w^{(r)T} .$$

Here $w^{(r)} = (0, \ldots, 0, w_{r+1}^{(r)}, \ldots, w_m^{(r)})^T$, $r = 0, 1, \ldots, n-1$ is chosen so that

$$P^{(n-1)} \ldots P^{(1)}P^{(0)}A = \left(\frac{U}{0}\right){\begin{matrix}\}n \\ \}m-n\end{matrix}} \tag{5.1}$$

where $U$ is an upper triangular matrix.

We now show how to adopt this method to solve the system (2.3) which we write, for the moment,

$$r + Ax = f$$
$$A^T r = g \tag{5.2}$$

Multiplying the first set of these equations from the left with $Q = P^{(n-1)} \dots P^{(1)} P^{(0)}$ and using (5.1) we get

$$Qr + \left(\frac{U}{0}\right) x = Qf \tag{5.3}$$
$$(U^T | 0) Qr = g$$

Thus, it is easily seen that $r$ and $x$ can be computed by the following algorithm:

$$h = U^{-T} g, \quad d = Qf = \left(\frac{d_1}{d_2}\right){}^{\}n}_{\}m-n}$$
$$r = Q^T \left(\frac{h}{d_2}\right), \quad x = U^{-1}(d_1 - h) \tag{5.4}$$

Wilkinson [9] has, under the assumption that inner products are accumulated with $t_2 = 2t_1$ digits, given an error analysis of orthogonal transformations of the type used here. We state below in (5.5)–(5.8) those of his results needed here.

Consider the computed sequence of transformed matrices

$$A = \bar{A}^{(0)}, \quad \bar{A}^{(r+1)} = \mathrm{fl}_2(\bar{P}^{(r)}\bar{A}^{(r)}), \quad r = 0, 1, \dots, n-1 .$$

For a certain prescribed method of computation there exists a sequence of elementary orthogonal transformations $P^{(0)}, P^{(1)}, \dots, P^{(n-1)}$ (not the matrices corresponding to exact computation throughout) such that

$$P^{(n-1)} \dots P^{(1)} P^{(0)}(A + E) = \bar{A}^{(n)} = \left(\frac{\bar{U}}{0}\right) \tag{5.5}$$

where[1]

$$\|E\|_E \leq n\beta(1+\beta)^{n-1}\|A\|_E, \quad \beta = 12.36\, \varepsilon_1 \tag{5.6}$$

and $\bar{U}$ is the computed upper triangular matrix.

If the computed transformation $\bar{Q} = \bar{P}^{(n-1)} \dots \bar{P}^{(1)} \bar{P}^{(0)}$ is applied to a vector $b$,

$$b = \bar{b}^{(0)}, \quad \bar{b}^{(r+1)} = \mathrm{fl}_2(\bar{P}^{(r)}\bar{b}^{(r)}), \quad r = 0, 1, \dots, n-1 .$$

then

$$\bar{b}^{(r+1)} = P^{(r)}\bar{b}^{(r)} + f^{(r)}, \quad \|f^{(r)}\|_2 \leq \beta\|\bar{b}^{(r)}\|_2 . \tag{5.7}$$

---

[1] Here the suffix $E$ denotes the Frobenius norm i.e. $\|A\|_E = (\sum\sum a_{ij}^2)^{\frac{1}{2}}$.

From (5.7) it follows by induction that

$$\|\bar{b}^{(r)}\|_2 \leq (1+\beta)\|\bar{b}^{(r-1)}\|_2 \leq (1+\beta)^r\|b\|_2$$

and

$$d^{(r+1)} = P^{(r)}d^{(r)}+f^{(r)}, \ d^{(r)} = \bar{b}^{(r)} - P^{(r-1)} \dots P^{(0)}b \ .$$

Since $d^{(0)}=0$ and $d^{(n)}=\bar{b}^{(n)}-Qb$, it follows that

$$\|\bar{b}^{(n)}-Qb\|_2 \leq \sum_{r=0}^{n-1} \beta(1+\beta)^r\|b\|_2 < n\beta(1+\beta)^{n-1}\|b\|_2 \ . \tag{5.8}$$

We now derive an only slightly different result for the error when $\bar{Q}^T=\bar{P}^{(0)}\bar{P}^{(1)} \dots \bar{P}^{(n-1)}$ is applied to a vector $c$;

$$c = \bar{c}^{(0)}, \ \bar{c}^{(r+1)} = \text{fl}_2(\bar{P}^{(n-r-1)}\bar{c}^{(r)}), \quad r = 0,1,\dots,n-1 \ .$$

In analogy to (5.7) we obviously have

$$\bar{c}^{(r+1)} = P^{(n-r-1)}\bar{c}^{(r)}+g^{(r)}, \ \|g^{(r)}\|_2 \leq \beta\|\bar{c}^{(r)}\|_2 \ . \tag{5.9}$$

As $P^{(n-r-1)}$ is orthogonal and symmetric we get

$$\bar{c}^{(r)} = P^{(n-r-1)}(\bar{c}^{(r+1)}+g^{(r)})$$

and thus

$$\|\bar{c}^{(r)}\|_2 \leq (1-\beta)^{-1}\|\bar{c}^{(r+1)}\|_2 \leq (1-\beta)^{-(n-r)}\|\bar{c}^{(n)}\|_2 \ .$$

By induction from (5.9)

$$e^{(r+1)} = P^{(n-r-1)}e^{(r)}+g^{(r)}, \ e^{(r)} = \bar{c}^{(r)} - P^{(n-r)} \dots P^{(n-1)}c \ ,$$

where $e^{(0)}=0$ and $e^{(n)}=\bar{c}^{(n)}-Q^Tc$. Hence

$$\|\bar{c}^{(n)}-Q^Tc\|_2 \leq \sum_{r=0}^{n-1} \beta(1-\beta)^{-(n-r)}\|\bar{c}^{(n)}\|_2 < n\beta(1-\beta)^{-n}\|\bar{c}^{(n)}\|_2 \tag{5.10}$$

For convenience we assume in the following that

$$12.5 \cdot n \cdot \varepsilon_1 \leq 0.01 \ , \tag{5.11}$$

which, as is easily shown, implies that

$$n\beta(1+\beta)^{n-1} < n\beta(1-\beta)^{-n} \leq 12.485 \cdot n \cdot \varepsilon_1 \ . \tag{5.12}$$

From (5.8), (5.10) and (5.12) it now follows that the computed quantities $\bar{d}$ and $\bar{r}$, in the algorithm (5.4), satisfy

$$\bar{d} = Q(f+e_1), \ \bar{r} = Q^T\begin{pmatrix}\bar{h} \\ \bar{d}_2\end{pmatrix}+e_2 \tag{5.13}$$

where

$$\|e_1\|_2 \leq 12.5n\varepsilon_1\|f\|_2, \ \|e_2\|_2 \leq 12.5n\varepsilon_1\|\bar{r}\|_2 \ . \tag{5.14}$$

Further errors are made in the two back-substitutions

$$\bar{U}^T h = g, \quad \bar{U} x = \bar{d}_1 - \bar{h} .$$

Wilkinson [8], pp. 99–104, has shown that the computed quantity $\bar{h}$ satisfies exactly

$$(\bar{U} + F_2)^T \bar{h} = g \tag{5.15}$$

where

$$\|F_2\|_E \leqq (\varepsilon_1 + n\varepsilon_2)\|\bar{U}\|_E . \tag{5.16}$$

A similar result holds for $\bar{x}$, which satisfies

$$(\bar{U} + F_1)\bar{x} = \bar{d}_1 - \bar{h} , \tag{5.17}$$

where

$$\|F_1\|_E \leqq (1 - \varepsilon_1)^{-1}(2\varepsilon_1 + n\varepsilon_2)\|\bar{U}\|_E . \tag{5.18}$$

The slightly greater bound for $F_1$ accounts for the rounding of the difference $(\bar{d}_1 - \bar{h})$. From (5.5), (5.6) and (5.12) we have

$$\|\bar{U}\|_E = \|A + E\|_E \leqq \|A\|_E + \|E\|_E \leqq 1.01\|A\|_E .$$

When $t_2 \geqq 2t_1$ we certainly have $\varepsilon_2 \leqq 2\varepsilon_1{}^2$. Hence if we assume $n \geqq 2$, from (5.16) and (5.18) we have

$$\|F_i\|_E \leqq (1 - \varepsilon_1)^{-1}(1 + n\varepsilon_1)1.01n\varepsilon_1\|A\|_E, \quad i = 1, 2 ,$$

or, after using (5.12) twice,

$$\|F_i\|_E \leqq 1.012n\varepsilon_1\|A\|_E, \quad i = 1, 2 . \tag{5.19}$$

Now define

$$H_i = E + Q^T F_i, \quad i = 1, 2 \tag{5.20}$$

where $E$ is the matrix in (5.5); then the relation

$$Q(A + H_i) = \left(\frac{\bar{U} + F_i}{0}\right), \quad i = 1, 2 \tag{5.21}$$

holds exactly. From (5.6), (5.12), (5.19) and (5.20) it follows that

$$\|H_i\|_E \leqq 13.5n\varepsilon_1\|A\|_E, \quad i = 1, 2 . \tag{5.22}$$

We summarize the results obtained in this section:

*Assume that the solutions $r$ and $x$ of the system of equations (5.2) are computed by the algorithm (5.4), using a certain method of computation described in [9]. Then, provided $12.5n\varepsilon_1 \leqq 0.01$ and $n \geqq 2$, the computed solution $\bar{r}$ and $\bar{x}$ is the exact solution to the perturbed system*

$$\left(\begin{array}{c|c} I & A+H_1 \\ \hline (A+H_2)^T & 0 \end{array}\right)\left(\begin{array}{c} \bar{r}-e_2 \\ \bar{x} \end{array}\right) = \left(\begin{array}{c} f-e_1 \\ g \end{array}\right), \qquad (5.23)$$

where $e_1$, $e_2$, $H_1$ and $H_2$ satisfies the bounds given in (5.14) and (5.22).

Perturbed systems of type (5.23) with symmetric perturbations $H_1 = H_2$ have been studied by Björck in [2]. To obtain from (5.23), an estimate of the errors in $\bar{r}$ and $\bar{x}$ we need a slightly more general result, which we state in the following theorem:

THEOREM 1. *Let $\bar{r}$ and $\bar{x}$ satisfy a perturbed system of equations*

$$\left(\begin{array}{c|c} I & A+H_1 \\ \hline (A+H_2)^T & 0 \end{array}\right)\left(\begin{array}{c} \bar{r} \\ \bar{\bar{x}} \end{array}\right) = \left(\begin{array}{c} f+d_1 \\ g+d_2 \end{array}\right)$$

*where $A$ is a given $m \times n$ matrix of rank $n$, and let $r$ and $x$ denote the solutions to the corresponding unperturbed system.*

*Assume that there exists an orthogonal matrix $Q$ such that the matrices $Q(A+H_i)$ $i = 1, 2$ are upper triangular, and let the perturbations satisfy the bounds*

$$\|H_i\|_2 \leqq \eta \|A\|_2, \quad i = 1, 2,$$
$$\|d_1\|_2 \leqq \tau_2, \quad \|d_2\|_2 \leqq \tau_1 \|A\|_2.$$

*Then, provided*

$$\alpha = (\sqrt{2}+1)\varkappa(A)\cdot\eta < 1,$$

*the following estimate holds*

$$\left(\begin{array}{c} \|\bar{r}-r\|_2 \\ \|A\|_2\|\bar{x}-x\|_2 \end{array}\right) \leqq \left(\begin{array}{cc} \varkappa' & 1 \\ (\varkappa')^2 & \varkappa' \end{array}\right)\left\{\eta\left(\begin{array}{c} \|r\|_2 \\ \|A\|_2\|x\|_2 \end{array}\right)+\left(\begin{array}{c} \tau_1 \\ \tau_2 \end{array}\right)\right\} \qquad (5.24)$$

*where*

$$\varkappa' = (1-\alpha)^{-1}\varkappa(A). \qquad (5.25)$$

PROOF. The theorem is proved by the same technique as used in [2] pp. 15–16, only trivial changes being necessary.

## 6. Convergence and limiting accuracy.

Denote by $\mathscr{M}(\varkappa')$ the set of $2 \times 2$ non-negative matrices of the type

$$M = \left(\begin{array}{c} 1 \\ \varkappa' \end{array}\right)(a,b). \qquad (6.1)$$

Then, obviously, $M_1, M_2 \in \mathscr{M}(\varkappa')$ implies that $M_1 + M_2 \in \mathscr{M}(\varkappa')$ and $M_1 M_2 \in \mathscr{M}(\varkappa')$. More generally, if $M \in \mathscr{M}(\varkappa')$ and $B$ is an arbitrary non-negative $2 \times 2$ matrix then $MB \in \mathscr{M}(\varkappa')$.

We note that the matrix

$$M_0 = \begin{pmatrix} \varkappa' & 1 \\ (\varkappa')^2 & \varkappa' \end{pmatrix} = \begin{pmatrix} 1 \\ \varkappa' \end{pmatrix} (\varkappa', 1)$$

in theorem 1 (5.24), belongs to $\mathscr{M}(\varkappa')$. This set of matrices will play an essential part in the analysis of the refinement procedure. We now (in lemmas 1 and 2) prove some simple properties of these matrices.

LEMMA 1. *Let* $M \in \mathscr{M}(\varkappa')$ *be defined by* (6.1). *Then the spectral radius of* $M$ *is given by*

$$\varrho(M) = a + \varkappa' b .$$

PROOF. $M$ has by definition rank one and thus only one eigenvalue different from zero. From the non-negativity it follows that

$$\varrho(M) = \mathrm{trace}(M) = a + \varkappa' b .$$

LEMMA 2. *Let* $M_1, M_2 \in \mathscr{M}(\varkappa')$. *Then the following multiplication rule holds*

$$M_1 M_2 = \varrho(M_1) M_2 . \tag{6.2}$$

PROOF. We have

$$M_1 M_2 = \begin{pmatrix} 1 \\ \varkappa' \end{pmatrix} (a_1, b_1) \begin{pmatrix} 1 \\ \varkappa' \end{pmatrix} (a_2, b_2) = (a_1 + \varkappa' b_1) \begin{pmatrix} 1 \\ \varkappa' \end{pmatrix} (a_2, b_2) = \varrho(M_1) M_2 .$$

COROLLARY.     $$M^n = \varrho(M)^{n-1} M, \quad n \geq 1 . \tag{6.3}$$

The estimate (5.24), moreover, suggests that we define a pseudo-norm in the space $R^{m+n}$ by

$$\|z\|_A = \begin{pmatrix} \|r\|_2 \\ \|A\|_2 \|x\|_2 \end{pmatrix}, \quad z = \begin{pmatrix} r \\ x \end{pmatrix} \begin{matrix} \}m \\ \}n \end{matrix} \in R^{m+n} . \tag{6.4}$$

Note that with the pseudo-distance in $R^{m+n}$

$$d_A(z_1, z_2) = \begin{pmatrix} \|r_1 - r_2\|_2 \\ \|A\|_2 \|x_1 - x_2\|_2 \end{pmatrix}, \quad z_1, z_2 \in R^{m+n} ,$$

$R^{m+n}$ becomes a pseudo-metric space cf. Collatz [3] (p. 40).

With the notation introduced in (6.1) and (6.4) we can write (5.24) in the simple form

$$\|\bar{z} - z\|_A \leq M_0(\eta \|z\|_A + \tau) . \tag{6.5}$$

We will now analyse the iterative refinement procedure defined by (2.2)–(2.4) assuming that the Householder method is used for the calculation of the corrections. Let $f^{(s)}$, $g^{(s)}$ and $\delta r^{(s)}$, $\delta x^{(s)}$ denote the exact

residuals and corrections corresponding to the computed approximations $\bar{r}^{(s)}$, $\bar{x}^{(s)}$. Then we have

$$r = \bar{r}^{(s)} + \delta r^{(s)}, \quad x = \bar{x}^{(s)} + \delta x^{(s)} .$$

From (5.23) it follows that the computed corrections satisfy

$$\left( \begin{array}{c|c} I & A + H_1^{(s)} \\ \hline (A + H_2^{(s)})^T & 0 \end{array} \right) \left( \begin{array}{c} \delta \bar{r}^{(s)} - e_2^{(s)} \\ \delta \bar{x}^{(s)} \end{array} \right) = \left( \begin{array}{c} f^{(s)} + e_1^{(s)} \\ g^{(s)} \end{array} \right) + \left( \begin{array}{c} \delta f^{(s)} \\ \delta g^{(s)} \end{array} \right) \quad (6.6)$$

where

$$\|e_1^{(s)}\|_2 \leqq 12.5 n \varepsilon_1 \|\bar{f}^{(s)}\|_2, \quad \|e_2^{(s)}\|_2 \leqq 12.5 n \varepsilon_1 \|\delta \bar{r}^{(s)}\|_2 .$$

Let $\delta \tilde{r}^{(s+1)}$ and $\delta \tilde{x}^{(s+1)}$ denote the errors in the exact sums

$$\tilde{r}^{(s+1)} = \bar{r}^{(s)} + \delta \bar{r}^{(s)}, \quad \tilde{x}^{(s+1)} = \bar{x}^{(s)} + \delta \bar{x}^{(s)} .$$

Then it follows that

$$\delta \tilde{r}^{(s+1)} = \delta \bar{r}^{(s)} - \delta r^{(s)}, \quad \delta \tilde{x}^{(s+1)} = \delta \bar{x}^{(s)} - \delta x^{(s)} . \quad (6.7)$$

Using (6.7) and the identity $f^{(s)} = \delta r^{(s)} + A \delta x^{(s)}$ we have

$$\|e_1^{(s)}\|_2 \leqq 12.5 n \varepsilon_1 (\|\delta r^{(s)}\|_2 + \|A\|_2 \|\delta x^{(s)}\|_2 + \|\delta f^{(s)}\|_2) ,$$
$$\|e_2^{(s)}\|_2 \leqq 12.5 n \varepsilon_1 (\|\delta r^{(s)}\|_2 + \|\delta \tilde{r}^{(s+1)}\|_2) .$$

From (5.22) and (the often very weak) inequality $\|A\|_E \leqq n^{\frac{1}{2}} \|A\|_2$ it follows that we can apply theorem 1 with

$$\alpha = (\sqrt{2} + 1) 13.5 n^{3/2} \varepsilon_1 \varkappa(A) . \quad (6.8)$$

Assuming $n \geqq 2$ it certainly follows that

$$12.5 n < \tfrac{2}{3} 13.5 n^{3/2} .$$

Taking (5.11) into account we obtain

$$\left( \begin{array}{c} 0.99 \|\delta \tilde{r}^{(s+1)}\|_2 \\ \|A\|_2 \|\delta \tilde{x}^{(s+1)}\|_2 \end{array} \right) \leqq 13.5 n^{3/2} \varepsilon_1 \left\{ \left( \begin{array}{c} 1 \\ \varkappa' \end{array} \right) [(\varkappa' \; 1) + \tfrac{2}{3}(1 \; 1)] + \right.$$

$$\left. \tfrac{2}{3} \left( \begin{array}{cc} 1 & 0 \\ 0 & 0 \end{array} \right) \right\} \left( \begin{array}{c} \|\delta r^{(s)}\|_2 \\ \|A\|_2 \|\delta x^{(s)}\|_2 \end{array} \right) + \left( \begin{array}{c} 1 \\ \varkappa' \end{array} \right) (\varkappa', 1) \left( \begin{array}{c} 1.01 \|\delta f^{(s)}\|_2 \\ \|\delta g^{(s)}\|_2 / \|A\|_2 \end{array} \right) .$$

From this follows, using a more compact notation,

$$0.99 \|\delta \tilde{z}^{(s+1)}\|_A \leqq 13.5 n^{3/2} \varepsilon_1 M' \|\delta z^{(s)}\|_A + 1.01 M_0 \tau \quad (6.9)$$

where

$$M' = \left( \begin{array}{cc} \varkappa' + \tfrac{4}{3} & \tfrac{5}{3} \\ \varkappa'(\varkappa' + \tfrac{2}{3}) & \varkappa' \tfrac{5}{3} \end{array} \right) < \left( \begin{array}{c} 1 \\ \varkappa' \end{array} \right) ((\varkappa' + \tfrac{4}{3}), \tfrac{5}{3}) \quad (6.10)$$

and from (4.2) and (4.3) we have

$$\tau \leq \left\{ \varepsilon_1 \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix} + \varepsilon_2' \begin{pmatrix} m+1 & 0 \\ 1 & n+3 \end{pmatrix} \right\} \|\delta z^{(s)}\|_A + \varepsilon_2' \begin{pmatrix} m+1 & 0 \\ 3 & n+5 \end{pmatrix} \|z\|_A \, . \quad (6.11)$$

Assuming $n \geq 2$ we have

$$M_0 \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix} = \begin{pmatrix} 1 \\ \varkappa' \end{pmatrix} ((\varkappa'+1),1) < M' \, ,$$

$$M_0 \begin{pmatrix} m+1 & 0 \\ 1 & n+3 \end{pmatrix} = \begin{pmatrix} 1 \\ \varkappa' \end{pmatrix} \left( \left( \varkappa' + \frac{1}{m+1} \right), \frac{n+3}{m+1} \right) < M' \quad (6.12)$$

Now consider the errors made when adding the corrections. For the error in the $i$th component of $\bar{r}^{(s+1)}$

$$\bar{r}_i^{(s+1)} - r_i = (1+\delta)\tilde{r}_i^{(s+1)} - r_i = (1+\delta)(\tilde{r}_i^{(s+1)} - r_i) + r_i \delta \, ,$$

where $|\delta| \leq \varepsilon_1$. Hence

$$\|\bar{r}^{(s+1)} - r\|_2 \leq (1+\varepsilon_1)\|\tilde{r}^{(s+1)} - r\|_2 + \varepsilon_1\|r\|_2 \, .$$

Using a similar result for $x^{(s+1)}$ we get

$$\|\delta z^{(s+1)}\|_A \leq (1+\varepsilon_1)\|\delta\tilde{z}^{(s+1)}\|_A + \varepsilon_1\|z\|_A \, . \quad (6.13)$$

We now summarize the results obtained in (6.9)–(6.13):
Define the matrices $C$ and $D$ by

$$C = cM_1, \quad D = dM_2 \quad (6.14)$$

where

$$c = [(13.64n^{3/2} + 1.021)\varepsilon_1 + 1.021(m+1)\varepsilon_2'](1+\varepsilon_1) \, ,$$
$$d = 1.021(m+1)\varepsilon_2'(1+\varepsilon_1) \, , \quad (6.15)$$

and

$$M_1 = \begin{pmatrix} 1 \\ \varkappa' \end{pmatrix} ((\varkappa' + \tfrac{4}{3}), \tfrac{5}{3}) \, ,$$

$$M_2 = \begin{pmatrix} 1 \\ \varkappa' \end{pmatrix} \left( \left( \varkappa' + \frac{3}{m+1} \right), \frac{n+5}{m+1} \right). \quad (6.16)$$

Then the errors in $\tilde{z}^{(s+1)}$ and $\bar{z}^{(s+1)}$ satisfy the recurrence relations

$$\|\delta\tilde{z}^{(s+1)}\|_A \leq C\|\delta z^{(s)}\|_A + D\|z\|_A$$
$$\|\delta z^{(s+1)}\|_A \leq C\|\delta z^{(s)}\|_A + (D + \varepsilon_1 I)\|z\|_A \, . \quad (6.17)$$

Thus, by induction,

$$\|\delta z^{(s)}\|_A \leq C^s\|z^{(0)}\|_A + (I + C + \ldots + C^{s-1})(D + \varepsilon_1 I)\|z\|_A$$
$$\|\delta\tilde{z}^{(s)}\|_A \leq C^s\|z^{(0)}\|_A + [(C + C^2 + \ldots + C^{s-1})(D + \varepsilon_1 I) + D]\|z\|_A \, . \quad (6.18)$$

Since $C$, $D \in \mathcal{M}(\varkappa')$, we can use the multiplication rule (6.2) from lemma 2. Hence, if we put

$$\varrho = \varrho(C) = \tfrac{8}{3}c(\varkappa' + \tfrac{1}{2}) \tag{6.19}$$

and assume $\varrho < 1$, then, for $s \geq 1$, $C^s = \varrho^{s-1}C$ and

$$(I + C + \ldots + C^{s-1})D = (1 + \varrho + \ldots + \varrho^{s-1})D < (1 - \varrho)^{-1}D$$
$$(C + C^2 + \ldots + C^{s-1}) = (1 + \varrho + \ldots + \varrho^{s-2})C < (1 - \varrho)^{-1}C .$$

Substituting this into (6.18) and noting that $z^{(0)} = z$ we finally get

$$\|\delta z^{(s)}\|_A < \big(\varrho^{s-1}C + (1-\varrho)^{-1}(D + \varepsilon_1 C) + \varepsilon_1 I\big)\|z\|_A$$
$$\|\delta \tilde{z}^{(s)}\|_A < \big(\varrho^{s-1}C + (1-\varrho)^{-1}(D + \varepsilon_1 C)\big)\|z\|_A . \tag{6.20}$$

If we make the reasonable assumption

$$(m+1)\varepsilon_2' < \varepsilon_1$$

then, as a consequence of (5.11) and the assumption $n \geq 2$

$$c < 1.0004(13.64 + 2 \cdot 2^{-3/2}1.021)n^{3/2}\varepsilon_1 ,$$

or

$$c < 14.4n^{3/2}\varepsilon_1 . \tag{6.21}$$

In the initial stages the term $\varrho^{s-1}C$ in (6.20) dominates and we have approximately

$$\|\delta z^{(s)}\|_A \lesssim \varrho^{s-1}C\|z\|_A .$$

*This justifies calling $\varrho$ the initial rate of convergence. From (6.19) and (6.21) we have the estimate*

$$\varrho < 38.4n^{3/2}(\varkappa' + \tfrac{1}{2})\varepsilon_1 . \tag{6.22}$$

*We note that this bound for $\varrho$ is independent of the right hand side and roughly proportional to $\varkappa(A)$. Thus, $\varkappa^2(A)$ does not enter, which, remembering (2.7), might have been conjectured from the general analysis.*

*When $\varrho < 1$, the term $\varrho^{s-1}C\|z\|_A$ in (6.20) approaches zero as $s \to \infty$ and the limiting accuracy in $\tilde{z}^{(s)}$ is*

$$\lim_{s \to \infty} \|\delta \tilde{z}^{(s)}\|_A \lesssim (1-\varrho)^{-1}(D + \varepsilon_1 C)\|z\|_A = \binom{1}{\varkappa'}(1-\varrho)^{-1}K \tag{6.23}$$

*where*

$$K = 1.022\varepsilon_2'(\varkappa'(m+4)\|r\|_2 + (n+5)\|A\|_2\|x\|_2) +$$
$$14.4n^{3/2}\varepsilon_1^2((\varkappa' + \tfrac{4}{3})\|r\|_2 + \tfrac{5}{3}\|A\|_2\|x\|_2) . \tag{6.24}$$

The first term in $K$ is proportional to $\varepsilon_2$ and comes from the errors made when computing the residuals. The second term, which is propor-

tional to $\varepsilon_1{}^2$, comes from the errors introduced when solving the system of equations defining the corrections. *To get full benefit from the refinement it is obviously necessary to have* $t_2 \geqq 2t_1$. The first term in $K$ will then usually be negligible, and the limiting accuracy will approximately be $(1-\varrho)^{-1}\varepsilon_1 C\|z\|_A$. This can be compared to the bound for the error $\|\bar{z}^{(1)} - z\|_A$ derived from (6.18), which is $C\|z\|_A$. *Hence, we can expect to gain* $t_1$ *figures during the refinement, and to achieve almost the same accuracy as if, without any refinement,* $2t_1$ *digit precision had been used throughout the computation.*

Since either $r$ or $x$ can be equal to zero, we obviously can not always expect to achieve a small relative error in either $\bar{r}^{(s)}$ or $\bar{x}^{(s)}$. We now derive simple sufficient conditions for the relations

$$\lim_{s\to\infty} \|\bar{r}^{(s)} - r\|_2 \leqq 2\varepsilon_1 \cdot \|r\|_2 , \tag{6.25}$$

$$\lim_{s\to\infty} \|\bar{x}^{(s)} - x\|_2 \leqq 2\varepsilon_1 \cdot \|x\|_2 \tag{6.26}$$

to be satisfied. We assume that $\varrho \leqq \frac{1}{4}$ and that the second term in $K$ can be neglected, since the most important case in practice is $t_2 \geqq 2t_1$. From (6.14), (6.20) and (6.23) it follows that (6.25) holds provided

$$\tfrac{4}{3}c(\varkappa' + \tfrac{4}{3})(1 + \tfrac{5}{3}\gamma^{-1}) \leqq 1$$

where we have put

$$\gamma = \frac{(\varkappa' + \tfrac{4}{3})\|r\|_2}{\|A\|_2\|x\|_2} . \tag{6.27}$$

Substituting for $c$ from (6.19) we obtain

$$1 + \tfrac{5}{3}\gamma^{-1} \leqq \frac{2}{\varrho}\frac{\varkappa' + \tfrac{1}{2}}{\varkappa' + \tfrac{4}{3}} .$$

Since $\varkappa' \geqq 1$ this relation is satisfied if

$$\gamma^{-1} \leqq \tfrac{3}{5}\left(\frac{2\cdot 9}{\varrho\cdot 14} - 1\right) .$$

Remembering that $\varrho \leqq \frac{1}{4}$ it follows that (6.25) certainly holds if

$$\gamma^{-1} \leqq \tfrac{87}{140}\varrho^{-1} . \tag{6.28}$$

Similarly (6.26) holds if

$$\gamma + \tfrac{5}{3} \leqq \frac{2}{\varrho}\frac{\varkappa' + \tfrac{1}{2}}{\varkappa'}$$

or if

$$\gamma \leqq \frac{1}{\varrho}(2 - \tfrac{5}{12}) = \tfrac{19}{12}\varrho^{-1} . \tag{6.29}$$

*We conclude that if $\gamma$ satisfies*

$$1.61\varrho \leqq \gamma \leqq 1.58\frac{1}{\varrho}$$

*then, for sufficiently large s, both $\bar{r}^{(s)}$ and $\bar{x}^{(s)}$ will have a relative accuracy better than $2\varepsilon_1$.* It is obvious that, for a given matrix $A$, there will always exist right hand sides for which this is not true either for $\bar{r}^{(s)}$ or for $\bar{x}^{(s)}$. However, since $\varrho \leqq \frac{1}{4}$ implies $1.61\varrho < 1.58\varrho^{-1}$, *at least one of the relations (6.25) and (6.26) is always satisfied.*

## 7. Error analysis of the modified Gram-Schmidt method.

This method is based on a decomposition of $A$ obtained in the following way: Let $A^{(1)} = A$ and $A^{(k)}$, $k = 2, 3, \ldots, n+1$ be defined by

$$r_k{}^T = d_k^{-1} q_k{}^T A^{(k)}, \quad d_k = \|q_k\|_2^2 ,$$
$$A^{(k+1)} = A^{(k)} - q_k r_k{}^T \tag{7.1}$$

where $q_1, q_2, \ldots, q_n$ is a suitable sequence of linearly independent vectors. By induction it follows that

$$A^{(n+1)} = A - QR \tag{7.2}$$

where $R^T = (r_1, r_2, \ldots, r_n)$, and for $k = 1, 2, \ldots, n$, that

$$q_k{}^T(q_1 r_1{}^T + \ldots + q_{k-1} r_{k-1}^T + q_k r_k{}^T) = q_k{}^T A .$$

Thus, if we define the lower triangular matrix $L$ by

$$L = \{l_{kj}\}, \qquad l_{kj} = \begin{cases} q_k{}^T q_j, & k \geqq j \\ 0 & , & k < j \end{cases}, \tag{7.3}$$

then

$$LR = Q^T A . \tag{7.4}$$

We now choose $q_k$ so that in step $k$ the $k$th column of $A^{(k)}$ is annihilated, i.e. so that

$$A^{(k)} = (0, \ldots, 0, a_k^{(k)}, \ldots, a_n^{(k)}) .$$

This is obviously achieved if we take $q_k = a_k^{(k)}$. Then, by (7.2),

$$A = QR \tag{7.5}$$

and $R$ becomes unit upper triangular. If the calculations are performed *exactly* the columns of $Q$ will be mutually orthogonal and thus $L$ will be a diagonal matrix.

The method described can be interpreted as an elimination with

weighted row combinations (Bauer [1]) or as a modification of the classical Gram–Schmidt orthogonalization method (Björck [2]).

We now adopt this method to solve the system (5.2). This must be done carefully, since, as is well known, [2], the computed columns of $Q$ may deviate considerably from orthogonality. In the derivation we, therefore, do not assume $L$ to be diagonal. If we define

$$y = (y_1, y_2, \ldots, y_n)^T, \quad h = (h_1, h_2, \ldots, h_n)^T$$

by

$$y = Rx, \quad h = R^{-T}g \tag{7.6}$$

and use (7.5), then we can write (5.2) as

$$r + Qy = f$$
$$Q^T r = h .$$

Multiplying the first set of equations in (5.2) from the left by $Q^T$ and using (7.4) we get

$$Q^T r + Ly = Q^T f .$$

Thus $r$ and $y$ are determined by

$$Ly = Q^T f - h, \quad r = f - Qy . \tag{7.7}$$

Now let $f^{(1)} = f$ and define $f^{(k)}$, $k = 2, 3, \ldots, n+1$, and $y$ by

$$y_k = (q_k^T f^{(k)} - h_k)/d_k, \quad f^{(k+1)} = f^{(k)} - q_k y_k . \tag{7.8}$$

By induction it follows that

$$f^{(k)} = f - (q_1 y_1 + \ldots + q_{k-1} y_{k-1}) ,$$
$$q_k^T(q_1 y_1 + \ldots + q_{k-1} y_{k-1} + q_k y_k) = q_k^T f - h_k ,$$

and thus $y$ and $r = f^{(n+1)}$ satisfy (7.7). Hence, (7.6) and (7.8) is the desired algorithm for solving (5.2).

In the error analysis below, we again assume that (5.11) holds and that $t_2 \geq 2t_1$. Let the computed factors in the decomposition (7.5) of $A$ be $\bar{R}$ and $\bar{Q} = (\bar{q}_1, \bar{q}_2, \ldots, \bar{q}_n)$. In the back-substitutions (7.6) the computed quantities $h$ and $x$ satisfy, cf. (5.16),

$$\bar{R}_1 \bar{x} = \bar{y}, \quad \bar{R}_2^T \bar{h} = g ,$$

where $\bar{R}_i = \bar{R} + G_i$, and

$$\|G_i\|_E \leq (\varepsilon_1 + n\varepsilon_2)\|\bar{R}\|_E \leq 1.002\varepsilon_1\|\bar{R}\|_E, \quad i = 1, 2 . \tag{7.9}$$

Let $\tilde{y}$ and $\tilde{f}^{(1)}, \tilde{f}^{(2)}, \ldots, \tilde{f}^{(n+1)} = \tilde{r}$ denote the *exact* results when performing (7.8) using the *computed* quantities $\bar{q}_k$ and $\bar{h}_k$. If we put

$$e_1 = \tilde{r} - \overline{r}, \quad e_2 = \tilde{y} - \overline{y},$$

then the following expression for the errors in $\overline{r}$ and $\overline{x}$ holds:

$$\begin{pmatrix} r - \overline{r} \\ x - \overline{x} \end{pmatrix} = (I - \overline{B}^{-1}B) \begin{pmatrix} r \\ x \end{pmatrix} + \begin{pmatrix} e_1 \\ \overline{R}_1^{-1}e_2 \end{pmatrix} \tag{7.10}$$

where

$$B = \begin{pmatrix} I & A \\ \hline A^T & 0 \end{pmatrix}, \quad \overline{B}^{-1} = \begin{pmatrix} I & -\overline{Q} \\ \hline 0 & \overline{R}_1^{-1} \end{pmatrix} \begin{pmatrix} I & 0 \\ \hline 0 & \overline{L}^{-1} \end{pmatrix} \begin{pmatrix} I & 0 \\ \hline \overline{Q}^T & -\overline{R}_2^{-T} \end{pmatrix}.$$

Note that, in order to make the errors small, $\overline{B}^{-1}$, which is not symmetric, has to be a good *left-hand* inverse. It is readily verified that

$$I - \overline{B}^{-1}B = \begin{pmatrix} \overline{Q}F_1 & -E_2 \\ \hline -\overline{R}_1^{-1}F_1 & \overline{R}_1^{-1}(E_3 + G_1) \end{pmatrix} \tag{7.11}$$

where

$$\begin{aligned} F_1 &= \overline{L}^{-1}\overline{R}_2^{-T}(E_1 + \overline{Q}G_2)^T, & E_1 &= \overline{Q}\,\overline{R} - A, \\ E_2 &= (I - \overline{Q}\,\overline{L}^{-1}\overline{Q}^T)A, & E_3 &= \overline{R} - \overline{L}^{-1}\overline{Q}^T A, \end{aligned} \tag{7.12}$$

and from [2] p. 10 we have the following results:

$$\begin{aligned} \|E_1\|_E &\leq 1.5(n-1)\varepsilon_1\|A\|_E, \quad \|E_2\|_E \leq 3.25(n-1)\varepsilon_1\|A\|_E \\ \|E_3\|_E &\leq 1.9(n-1)^{\frac{1}{2}}n\varepsilon_1\|A\|_E. \end{aligned} \tag{7.13}$$

We now estimate the rounding errors made in (7.8). To simplify the analysis, we assume that

$$\|\overline{q}_k\|_2 = 1, \quad k = 1, 2, \ldots, n.$$

Since this can be achieved by a proper scaling of the columns of $A$, this assumption will not influence the derived floating point error bounds. Then we have

$$\tilde{f}^{(k+1)} = (I - \overline{q}_k\overline{q}_k{}^T)\tilde{f}^{(k)} + \overline{q}_k\overline{h}_k.$$

If we define the error vector $\eta^{(k)}$ by

$$\overline{f}^{(k+1)} = (I - \overline{q}_k\overline{q}_k{}^T)\overline{f}^{(k)} + \overline{q}_k\overline{h}_k + \eta^{(k)}, \tag{7.14}$$

then, subtracting,

$$\overline{f}^{(k+1)} - \tilde{f}^{(k+1)} = (I - \overline{q}_k\overline{q}_k{}^T)(\overline{f}^{(k)} - \tilde{f}^{(k)}) + \eta^{(k)}.$$

Since $\overline{f}^{(1)} = \tilde{f}^{(1)} = f$ it follows that

$$\|\overline{f}^{(k+1)} - \tilde{f}^{(k+1)}\|_2 \leq \sum_{\nu=1}^{k} \|\eta^{(\nu)}\|_2, \quad k = 1, 2, \ldots, n. \tag{7.15}$$

For the error in $\overline{y}_k$ we write

$$\bar{y}_k - \tilde{y}_k = \bar{y}_k - y_k' + \bar{q}_k^T(\bar{f}^{(k)} - \tilde{f}^{(k)})$$

where $y_k'$ is the exact multiplier corresponding to the computed vector $\bar{f}^{(k)}$. Then

$$|\bar{y}_k - \tilde{y}_k| \leqq |\bar{y}_k - y_k'| + \|\bar{f}^{(k)} - \tilde{f}^{(k)}\|_2 \qquad (7.16)$$

and

$$|y_k'| \leqq \|\bar{f}^{(k)}\|_2 + |\bar{h}_k| .$$

Now, if we make the reasonable assumption that $2(m+1)\varepsilon_1 \leqq 0.01$, then we can use results from [2] p. 7–8, duly modified, to obtain

$$|\bar{y}_k - y_k'| < 2.02\varepsilon_1(\|\bar{f}^{(k)}\|_2 + |\bar{h}_k|) \qquad (7.17)$$

and

$$\|\eta^{(k)}\|_2 < (1 - \varepsilon_1)^{-1}\varepsilon_1\|\bar{f}^{(k+1)}\|_2 + \varepsilon_1|\bar{y}_k| + |\bar{y}_k - y_k'| .$$

Consequently

$$(1 - \varepsilon_1)\|\eta^{(k)}\|_2 \leqq \varepsilon_1\|\bar{f}^{(k+1)}\|_2 + 3.02\varepsilon_1(\|\bar{f}^{(k)}\|_2 + |\bar{h}_k|) \qquad (7.18)$$

From (7.14) it follows that

$$\|\bar{f}^{(k+1)}\|_2 \leqq \|\bar{f}^{(k)}\|_2 + |\bar{h}_k| + \|\eta^{(k)}\|_2$$

and using (7.18) we have

$$(1 - 2\varepsilon_1)\|\bar{f}^{(k+1)}\|_2 \leqq (1 + 2.02\varepsilon_1)(\|\bar{f}^{(k)}\|_2 + |\bar{h}_k|) .$$

The assumption (5.11) certainly implies that

$$(1 + 2.02\varepsilon_1)/(1 - 2\varepsilon_1) < 1.002^{1/n}$$

and thus, for $k = 1, 2, \ldots, n$,

$$\|\bar{f}^{(k+1)}\|_2 < 1.002(\|f\|_2 + s_k), \quad s_k = \sum_{j=1}^{k} |\bar{h}_j| .$$

Using this bound for the growth of the computed vectors $\bar{f}^{(k)}$ in (7.17) and (7.18) we obtain

$$\begin{aligned} |\bar{y}_k - y_k'| &< 4.05\tfrac{1}{2}\varepsilon_1(\|f\|_2 + s_k) , \\ \|\eta^{(k)}\|_2 &< 4.05\varepsilon_1(\|f\|_2 + s_k) . \end{aligned} \qquad (7.19)$$

Using Schwarz's inequality we get

$$\sum_{\nu=1}^{k} s_\nu \leqq \left(\sum_{\nu=1}^{k} \nu^2\right)^{\tfrac{1}{2}} \|\bar{h}\|_2 < \frac{1}{\sqrt{3}} k^{\tfrac{1}{2}}(k+1)\|\bar{h}\|_2$$

and

$$\tfrac{1}{2}s_k + \sum_{\nu=1}^{k-1} s_\nu < \frac{1}{\sqrt{3}} k^{3/2}\|\bar{h}\|_2 .$$

Hence, from (7.15), (7.16) and (7.19) we have

$$\|e_1\|_2 = \|\bar{\tilde{f}}^{(n+1)} - \tilde{f}^{(n+1)}\|_2 < 4.05(n+1)\varepsilon_1\varphi , \qquad (7.20)$$

$$|\bar{y}_k - \tilde{y}_k| < 4.05k\varepsilon_1\varphi ,$$

where

$$\varphi = \|f\|_2 + \frac{1}{\sqrt{3}} n^{\frac{1}{2}} \|\bar{h}\|_2 . \qquad (7.21)$$

Using Schwarz's inequality again we obtain

$$\|e_2\|_2 = \|\bar{y} - \tilde{y}\|_2 < \frac{1}{\sqrt{3}} 4.05 n^{\frac{1}{2}}(n+1)\varepsilon_1\varphi . \qquad (7.22)$$

From (5.2) we have $h = R^{-T}A^T r = Q^T r$, and thus

$$\|f\|_2 \leqq \|r\|_2 + \|A\|_2\|x\|_2, \quad \|h\|_2 \leqq \|r\|_2 . \qquad (7.33)$$

Since the computed matrix $\overline{Q}$ is not exactly orthogonal, a further strict analysis becomes very cumbersome. In [2] pp. 12–15 we have shown, however, that

$$\|I - \overline{Q}^T \overline{Q}\|_2 \leqq \text{const.}\delta + 0(\delta^2) , \qquad (7.24)$$

where

$$\delta = n^2 \cdot \varepsilon_1 \varkappa(A) .$$

*In the following we therefore assume that $\delta$ is a small quantity.* From (7.24) it follows

$$\|\overline{Q}\|_2 = 1 + 0(\delta), \quad \|\overline{L}^{-1}\|_2 = 1 + 0(\delta) .$$

Let $A = QR$ where $Q$ is normalized so that $Q^T Q = I$. Then we also have ([2] p. 14)

$$\|\overline{R}\|_E = \|R\|_E(1 + 0(\delta)) = \|A\|_E(1 + 0(\delta)) .$$

Using these relations, we get from (7.9) and (7.12)

$$\begin{align}
\|E_1 + \overline{Q}\, G_2\|_E &\leqq 1.5n\varepsilon_1\|A\|_E(1 + 0(\delta)) , \\
\|E_3 + G_1\|_E &\leqq 1.9 \cdot n^{3/2}\varepsilon_1\|A\|_E(1 + 0(\delta))
\end{align} \qquad (7.25)$$

Neglecting powers of $\delta$ we derive from the inequality

$$\|A\|_E\|\overline{R}_i^{-1}\|_2 \leqq n^{\frac{1}{2}}\|A\|_2\|\overline{R}_i^{-1}\|_2 = n^{\frac{1}{2}}\varkappa(A)(1 + 0(\delta))$$

$i = 1, 2,$ (7.10)–(7.12) and (7.25) the estimate

$$\|z - \bar{z}\|_A \leqq 1.5n^{\frac{1}{2}}(n+1)\varepsilon_1 \begin{pmatrix} \varkappa & 1.27\sqrt{3} \\ \varkappa^2 & \varkappa \cdot 1.27 \cdot n^{\frac{1}{2}} \end{pmatrix} \|z\|_A + \begin{pmatrix} \|e_1\|_2 \\ \varkappa\|e_2\|_2 \end{pmatrix}.$$

Further, from $\|\bar{h}\|_2 = \|h\|_2(1 + 0(\delta))$, (7.21) and (7.23) it follows that

$$\varphi \leq \|r\|_2 \left(1 + \frac{n^{\frac{1}{2}}}{\sqrt{3}}\right) + \|A\|_2\|x\|_2 + 0(\delta) .$$

Hence, we obtain from (7.20) and (7.22), neglecting powers of $\delta$

$$\begin{pmatrix} \|e_1\|_2 \\ \varkappa\|e_2\|_2 \end{pmatrix} \leq 4.05(n+1)\varepsilon_1 \begin{pmatrix} 1 \\ \varkappa\dfrac{n^{\frac{1}{2}}}{\sqrt{3}} \end{pmatrix} \left(1 + \frac{n^{\frac{1}{2}}}{\sqrt{3}}, 1\right) \|z\|_A .$$

After some manipulation we have finally

$$\|z - \bar{z}\|_A \leq 1.5n^{\frac{1}{2}}(n+1)\varepsilon_1 \begin{pmatrix} \varkappa + 0.9(\sqrt{3} + 3n^{-\frac{1}{2}}) & 1.27\sqrt{3}(1 + 1.25n^{-\frac{1}{2}}) \\ \varkappa(\varkappa + 0.9(n^{\frac{1}{2}} + \sqrt{3})) & \varkappa \cdot 1.27(n^{\frac{1}{2}} + 1.25) \end{pmatrix} \|z\|_A .$$

*If $n \geq 3$, which we therefore assume, we can obviously write this relation in the form*

$$\|z - \bar{z}\|_A \leq 1.5n^{\frac{1}{2}}(n+1)\varepsilon_1 M_{GS}\|z\|_A + 0(\delta^2) \tag{7.26}$$

*where*

$$M_{GS} = \begin{pmatrix} 1 \\ \varkappa \end{pmatrix} (\varkappa + 0.9(n^{\frac{1}{2}} + \sqrt{3}), \ 1.27(n^{\frac{1}{2}} + 1.25)) \in \mathcal{M}(\varkappa) .$$

Using this as a starting point, the iterative refinement with the Gram–Schmidt method can be analysed in a way similar to that used for the Householder method. Obviously, the effects of the errors in the residuals will, in the first approximation, be the same.

It will again be true, that if $t_2 \geq 2t_1$, we can expect to gain $t_1$ figures during the refinement. To compare the limiting accuracy of the methods it is therefore sufficient to look at the error bounds for the solution of the system of equations (5.2).

A result corresponding to (7.26) for the Householder method follows from (6.9), if we put $s = 0$ and $\tau = 0$. Then, we get

$$\|z - z\|_A \leq 13.65n^{3/2}\varepsilon_1 M_H\|z\|_A + 0(\delta^2) \tag{7.27}$$

where

$$M_H = \begin{pmatrix} 1 \\ \varkappa \end{pmatrix} (\varkappa + \tfrac{4}{3}, \tfrac{5}{3}) .$$

Comparing (7.26) and (7.27) it is seen that when $n \leq 110$ the bound for the Gram–Schmidt method is always smaller. In particular for large values of $\varkappa$ and $\gamma$ (defined in (6.26)), the ratio is approximately $1:9$.

By lemma 1 and 2 the rate of convergence for the Gram–Schmidt method is approximately equal to

$$\varrho_{GS} = 1.5n^{\frac{1}{2}}(n+1)\varepsilon_1\varrho(M_{GS}) \tag{7.28}$$

where

$$\varrho(M_{GS}) = \varkappa + 0.9(n^{\frac{1}{2}}+\sqrt{3}) + 1.27(n^{\frac{1}{2}}+1.25)$$
$$< 1.27(n^{\frac{1}{2}}+2)(\varkappa+0.71) .$$

This can be compared to the corresponding expression for the House-holder method which is

$$\varrho_H = 13.65n^{3/2}\varepsilon_1\tfrac{8}{3}(\varkappa+\tfrac{1}{2}) . \tag{7.29}$$

For reasonable values of $n$, the Gram–Schmidt method is again seen to be more favourable. The advantage when $\varkappa \gg 1$ and $\gamma \gg 1$ is however less marked.

## 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. Björck, Å., *Solving Linear Least Squares Problems by Gram–Schmidt Orthogonalization*, BIT 7 (1967), 1–21.
3. Collatz, L., *Funktionalanalysis und Numerische Mathematik*, Berlin: Springer-Verlag (1964).
4. Golub, G. H., *Numerical Methods for Solving Linear Least Squares Problems*, Num. Math. 7 (1965), 206–216.
5. Golub, G. H. and Wilkinson, J. H., *Note on the Iterative Refinement of Least Squares Solution*, Num. Math. 9 (1966), 139–148.
6. Martin, R. S., Peters, G. and Wilkinson, J. H., *Iterative Refinement of the Solution of a Positive Definite System of Equations*, Num. Math. 8, (1966), 203–216.
7. Moler, C. B., *Iterative Refinement in Floating Point*, J. Assoc. Comput. Mach. 14, (1967), 316–321.
8. Wilkinson, J. H., *Rounding Errors in Algebraic Processes*, London: Her Majesty's Stationary Office; Englewood Cliffs, N.J.: Prentice–Hall 1963.
9. Wilkinson, J. H., *Error Analysis of Transformations Based on the Use of Matrices of the Form $I-2ww^T$*, Error in Digital Computations. Volume II. Rall, L. B. ed. New York: John Wiley (1965), 77–101.