

AN IMPLEMENTATION OF THE QR FACTORIZATION FOR SOLVING OVERDETERMINED SYSTEMS OF LINEAR EQUATIONS

L. F. ESCUDERO

Different procedures for solving an overdetermined system of linear equations are revisited and new technical details in its computer implementations are described. The new procedures produce more stable solutions than the implementation of the direct theoretical approaches.

1. INTRODUCTION AND MOTIVATION

In this paper we describe an implementation of a procedure to factorize a $n \times t$ full rank matrix A so that we may solve with the maximum accuracy an overdetermined system of linear equations. In this paper maximum accuracy means maximum accuracy in today available computers. Consider the problem

$$A\mu = g \quad (1.1)$$

where g is an n -column vector, μ is the unknown t -column vector and $n > t$. If the system is not compatible, vector μ has to minimize the norm $\|b\|^2$

$$\|b\|^2 = \|g - A\mu\|^2 \quad (1.2)$$

If the system is compatible, vector μ that satisfies (1.1), say μ_L is termed minimal least square solution.

The algorithms described in this paper calculate least square solutions to overdetermined systems of linear equations through the use of modified Gram-Schmidt orthogonalization with iterative refinement. The least square problem is a very general one which arises in a variety of contexts in many scientific disciplines. Frequently, least squares problems are ill-conditioned, and not all algorithms for computing solutions are numerically stable. See Escudero /1/. Among the variety of cases in which this problem arises, there are two very important cases; one is the multiple linear regression in which the data are the n observed values

(vector g) of a dependent variable and the observed values (matrix A) in the same n experiments corresponding to t independent variables; the problem consists in obtaining the vector μ of the regression coefficients such that the residual $\|b\|^2$ in eq. (1.2) is to be minimized. (In this case the partial regression coefficient μ_i takes the influence of independent variable, say X_i on the dependent variable, say Y).

The other important case with direct applicability of this work is the general non-linear constrained problem: $\min f(x)$ subject to $c(X) \geq 0$, where $f(X)$ and $c(X)$ are differentiable non-linear functions. The point X is stationary (see Escudero /1/) if $\|b\|^2 = 0$ in eq. (1.2), where A is the Jacobian matrix of the active constraints ($c(X)=0$), g is the gradient vector of objective function $f(X)$ - evaluated at point X and μ is the corresponding vector of Lagrange multipliers. At an intermediate iteration of the algorithm for solving the problem, point X needs that, at least, $\|b\|^2 = 0$ for satisfying the optimality conditions. If $\|b\|^2 \neq 0$, that means that the vector μ_L for which $\|b\|^2$ is minimized is the so-called first-order estimate of the Lagrange multipliers. See /1/.

The necessary and sufficient condition /3 p. 309/ for the minimal least square solution μ_L is that

$$A^t(g - A\mu) = 0 \quad (1.3)$$

where A^t is the transpose of matrix A .

- L.F. Escudero. IBM Scientific Center. Palo Alto, California. On leave from IBM España, Pº de la Castellana, 4, Madrid-1.
- Article rebut 1'Abril del 1980, nova versió.

Sufficient

Suppose $A^t(g - A\mu) = 0$. The residual $r(\epsilon)$ corresponding to any other vector $\mu + \epsilon$, where ϵ is a vector, is

$$\begin{aligned} \|r(\epsilon)\|^2 &= \|g - A\mu - A\epsilon\|^2 \\ &= \|g - A\mu\|^2 + \|A\epsilon\|^2 - 2(g - A\mu)^t A\epsilon \\ &= \|g - A\mu\|^2 + \|A\epsilon\|^2 \end{aligned} \quad (1.4)$$

Since $\|A\epsilon\|^2$ is non-negative, the residual $\|g - A\mu\|^2$ is a minimum. Unless $A\epsilon = 0$, the residual corresponding to $\mu + \epsilon$ is greater than the residual corresponding to μ ; the only other vectors with minimum residual are those vectors $\mu + \epsilon$ for which $A\epsilon = 0$. But since A is a full rank matrix, $A\epsilon \neq 0$.

Necessary

Suppose μ has a residual (1.2) that is minimum, but

$$A^t(g - A\mu) = s \neq 0 \quad (1.5)$$

The residual corresponding to $\mu + \epsilon s$ where ϵ is a scalar, is

$$\begin{aligned} \|g - A\mu - A\epsilon s\|^2 &= \|g - A\mu\|^2 - \\ &- 2\epsilon s^t A^t(g - A\mu) + \epsilon^2 \|As\|^2 \end{aligned} \quad (1.6)$$

For sufficiently small positive ϵ , the right hand-side of (1.6) is less than $\|g - A\mu\|^2$; - but it contradicts the hypothesis that the residual of μ is a minimum.

Then the vector μ_L that satisfies eq. (1.3) is

$$\mu_L = A^+ g \quad (1.7)$$

where A^+ is the pseudoinverse of matrix A . There are several computational methods to calculate matrix A^+ . A well known method - computes A^+ such that

$$\mu_L = (A^t A)^{-1} A^t g \quad (1.8)$$

assuming that $(A^t A)$ is non-singular; it will be non-singular if A is a full rank matrix. But from a practical point of view even in this case (since the computer has finite precision) the rounding errors produce a solution μ_L that does not satisfy eq. (1.3).

In order to avoid the computation of matrix A^+ , we may factorize matrix A such that

$$Q^t A = \begin{pmatrix} R \\ 0 \end{pmatrix} \quad (1.9)$$

where Q^t is a $n \times n$ non-symmetric orthogonal ($Q^t Q = I$) matrix and R is a $t \times t$ upper triangular matrix (not necessarily with identity diagonal). Matrix Q may be partitioned so that

$$Q = (Q_1 Q_2) \quad (1.10)$$

where Q_1 is a $n \times t$ matrix and Q_2 is a $n \times (n - t)$ matrix, so that

$$A = Q_1 R \quad (1.11)$$

Since A is a full rank matrix, it results

(a) R is a non-singular matrix

$$(b) I = Q Q^t = (Q_1 Q_2) \begin{pmatrix} Q_1^t \\ Q_2^t \end{pmatrix} = Q_1 Q_1^t + Q_2 Q_2^t \quad (1.12)$$

$$I = Q^t Q; Q_1^t Q_1 = I; Q_1^t Q_2 = 0; Q_2^t Q_1 = 0; Q_2^t Q_2 = I$$

$$(c) \begin{pmatrix} R \\ 0 \end{pmatrix} = \begin{pmatrix} Q_1^t \\ Q_2^t \end{pmatrix} \cdot A \equiv \begin{pmatrix} Q_1^t A = R \\ Q_2^t A = 0 \end{pmatrix} \quad (1.13)$$

(d) The minimal least square solution μ_L is found by solving

$$R \mu_L = Q_1^t g \quad (1.14)$$

In fact, considering eqs. (1.2) and (1.9) - (1.4), and in a similar way to eq. (1.4), it results

$$\begin{aligned} b &= g - A \mu_L = g - Q_1 R R^{-1} Q_1^t g = g - Q_1 Q_1^t g = \\ &= (I - Q_1 Q_1^t) g = Q_2 Q_2^t g \end{aligned} \quad (1.15)$$

being

$$A \mu_L = Q_1 Q_1^t g \quad (1.16)$$

where μ is calculated by formula (1.14). - For any other vector, say $\mu_L + \epsilon$ (where ϵ is a vector) its norm $\|b_1\|^2$ is greater than $\|b\|^2$. In effect,

$$b_1 = g - A(\mu_L + \epsilon) = b - A\epsilon \quad (1.17)$$

$$\|b_1\|^2 = \|b - A\epsilon\|^2 = \|b\|^2 + \|A\epsilon\|^2 - 2b^t A\epsilon$$

$$= \|b\|^2 + \|A\epsilon\|^2 - 2g^t Q_2 Q_2^t A\epsilon \\ = \|b\|^2 + \|A\epsilon\|^2 > \|b\|^2 \quad (1.18)$$

since $Q_2^t A = 0$, and being A a full rank matrix it results that $A\epsilon \neq 0$ for any ϵ .

For the rest of the paper we will not use - matrix Q_2 . In order to simplify notation - let substitute Q_1 by Q . Then the new notation is $A = QR$ and it is termed QR-factorization of matrix A . Formula (1.14) is more - stable than formula (1.8), since from (1.14) it results

$$\mu_{L_t} = Q_t^t g / R_{tt}$$

$$\mu_{L_i} = (Q_i^t g - \sum_{r=i+1}^t R_{ir} \mu_{L_r}) / R_{ii} \quad (1.19)$$

for $i = t-1, \dots, 1$

where Q_i is the i -th column of matrix Q . - Then formula (1.19) is not losing precision if the computation of $Q^t g$ is correctly made. It is interesting to point out that in this computation it is implicit the orthogonality of matrix Q (in this case it means that $Q^t Q = I$, but $Q Q^t \neq I$); if while obtaining matrix Q - there are rounding errors, formula (1.19) is also unstable although it is better than formula (1.8).

Then the goal is to calculate the orthogonal matrix Q and the upper triangular matrix R (Sec.2), the minimum residual vector b (Sec. 3) and the unknown vector μ_L (Sec.4) with - the maximum accuracy that it is possible in the computer. We will use the Modified Gram-Schmidt QR factorization. See e.g. /2 pp. - 381-385/, /3 p. 313/, /4/ and /5/, although the technical details make a strong difference in the performance of any procedure.

2. CALCULATION OF MATRICES Q AND R

Matrix Q (dimension $n \times t$) is calculated in t iterations in a recurrence way, so that at iteration 0 we assign $Q^{(0)} = A$ and at the end of iteration k , matrix $Q^{(k)}$ has the following expression

$$Q^{(k)} = (q_1^{(k)}, \dots, q_k^{(k)}, \dots, q_t^{(k)}) \quad (2.1)$$

where q_i is the i -th column vector of matrix

Q , so that vectors $q_1^{(k)}$ to $q_k^{(k)}$ take the definitive vectors $q_1^{(t)}$ to $q_k^{(t)}$ of - matrix Q ($\equiv Q^{(t)}$).

In successive iterations (from $k+1$ to t) vectors $q_{k+1}^{(k)}$ to $q_t^{(t)}$ will be transformed in - the definitive vectors $q_{k+1}^{(t)}$ to -

$q_t^{(t)}$ of matrix Q . Then at iteration k , vectors $q_1^{(k)}$ to $q_k^{(k)}$ are orthonormals, that is

$$q_i^{(k)t} q_j^{(k)} = \begin{cases} 1 & i=j \\ 0 & i \neq j \end{cases} \quad (2.2)$$

for $i, j = 1, 2, \dots, k$

Matrix R is calculated in a similar way, so that $R_{ij} = 0$ for $i > j$ and $i, j = 1, 2, \dots, t$. The row vector R_k of matrix R is obtained - at the same iteration k where vector q_k is obtained. At iteration k the procedure is - as follows.

Normalized Q algorithm (NOA):

(1) Obtain norms

$$N_i^{(k)} = \left(\sum_{h=1}^n q_{hi}^{(k)2} \right)^{1/2} \quad (2.3)$$

for $i = k, k+1, \dots, t$

Note that the square root is not calculated in an exact way because of the limited representation of floating-point numbers in the computer. See below a slightly modification that avoids this possible instability.

(2) Let $N_{k'}^{(k)}$ be the maximum of norms (2.3). If there are several norms with the same maximum value, we arbitrarily take the norm - with the first subscript. Since we assume that A is a full rank matrix, $N_{k'}^{(k)} \neq 0$.

(3) If $k' > k$, the k' -th columns of matrix $Q^{(k)}$ are interchanged. Also we interchange the elements $R_{ik'}$ and R_{ik} of row vectors R_i (for $i = 1, 2, \dots, k-1$). After the interchange we will continue calling $q_i^{(k)}$ to the i -th column of matrix $Q^{(k)}$ and R_i to the i -th row of matrix R .

(4) Obtain vector q_k for the iteration $k+1$ (it will be the definitive vector q_k of matrix Q).

$$q_k^{(k+1)} = q_k^{(k)} / N_k^{(k)} \quad (2.4)$$

(5) Obtain elements R_{ki} ($i=k, k+1, \dots, t$) of row vector R_k

$$R_{kk} = q_k^{(k+1)T} q_k^{(k)} = (q_k^{(k)T} q_k^{(k)}) / N_k^{(k)} = N_k^{(k)} \quad (2.5a)$$

$$R_{ki} = q_k^{(k+1)T} q_i^{(k)} = (q_k^{(k)T} q_i^{(k)}) / N_k^{(k)} \quad (2.5b)$$

for $i=k+1, \dots, t$

(6) Calculate vectors q for iteration $k+1$

$$q_i^{(k+1)} = q_i^{(k)} - R_{ki} q_k^{(k+1)} \text{ for } i=k+1, \dots, t \quad (2.6)$$

The orthogonality (2.2) implies that column vectors $q_i^{(k)}$ ($i=1, 2, \dots, k$) are theoretically normalized at the end of iteration k - (i.e. $\sum_{h=1}^n q_{hi}^{(k)2}$ should be equal to 1), but

eq. (2.4) is not very stable since the calculations of $N_i^{(k)}$ (2.3) are not exact. And - this instability is augmented while calculating $q_i^{(k+1)}$ (2.6). In order to avoid it we modify the above algorithm mainly by not executing step (4) and storing in scalar d_k the norm $N_k^{(k)2}$ that was calculated in steps (2)-(3); at the end of the new algorithm, the t -column vector d stores the norms of vectors q_1, q_2, \dots, q_t that are no longer normalized. Then the expression of the orthogonality of matrix Q will be

$$(q_i^T q_j) / (N_i N_j) = \begin{cases} 1, & i=j \\ 0, & i \neq j \end{cases} \quad (2.7a)$$

for $i, j = 1, 2, \dots, t$; that is

$$Q^T Q \delta^T \delta = I \quad (2.7b)$$

where the i -th element of the t -column vector δ is $1/d_i^{1/2}$.

Non-Normalized Q algorithm (NNQA):

At iteration k of the new algorithm:

(1) Obtain elements d_i ($i=k, k+1, \dots, t$) of vector d

$$d_i = N_i^{(k)2} \equiv \sum_{h=1}^n q_{hi}^{(k)2} \quad (2.8)$$

(2) -(3) as in NQA

(4) Assign directly

$$q_k^{(k+1)} = q_k^{(k)} \quad (2.9)$$

(5) As in NQA, but normalizing the row vector R_k with $N_k^{(k)}$. See eq. (2.12).

Then,

$$R_{kk} = ((q_k^{(k)T} q_k^{(k)}) / N_k^{(k)}) / N_k^{(k)} = 1 \quad (2.10a)$$

$$R_{ki} = ((q_k^{(k)T} q_i^{(k)}) / N_k^{(k)}) / N_k^{(k)} = (q_k^{(k)T} q_i^{(k)}) / d_k \equiv (q_k^{(k+1)T} q_i^{(k)}) / d_k \quad (2.10b)$$

(6) Obtain vectors q for iteration $k+1$. We use formula (2.6), but considering that vector $q_k^{(k+1)}$ is not normalized. Note that now element R_{ki} is normalized in formula (2.10b). Then, we have the same result

$$q_i^{(k+1)} = q_i^{(k)} - R_{ki} N_k^{(k)} (q_k^{(k+1)} / N_k^{(k)}) = q_i^{(k)} - R_{ki} q_k^{(k+1)} \quad (2.11)$$

for $i=k+1, \dots, t$. Note that since NNQA is not using $N_k^{(k)}$, but d_k (and it is used only at the last calculations), there are fewer rounding errors. Although in NNQA we calculate Q and R without normalizing the columns of matrix Q and normalizing the rows of matrix R , the QR factorization is the same. In fact, we have reduced the rounding errors by affecting the norm of each column only once to the whole column vector (see eq. (2.10)). If we consider that Q and R are the matrices calculated by NNQA, and Q' and R' are the matrices calculated by NQA, theoretically - $Q'R' = QR$. In effect, element A_{ij} of matrix A can be written

$$A_{ij} = \sum_{h=1}^j Q'_{ih} R'_{hj} = \sum_{h=1}^j (Q_{ih} / N_h) (R_{hj} N_h) = \sum_{h=1}^j Q_{ih} R_{hj} \quad (2.12)$$

where $R_{jj}=1$ (eq. (2.10a)). Then $A=QR=Q'R'$. We may see that when Q is a non-normalized matrix and R is a normalized matrix with the identity diagonal, the procedure is more stable than when Q is normalized and R is non-normalized.

3. CALCULATION OF RESIDUAL VECTOR $B=G-AU$

Since $A=QR$ (being Q and R calculated by NNQA)

it is possible to obtain μ with formula -- (1.19). But it is based on the orthogonality of matrix Q ($A\mu=g$; $A\mu=QQ^t g$; $Q^t A\mu=Q^t Q Q^t g$; $R\mu=Q^t g$); then if while calculating matrix Q there are some rounding errors, it is possible that Q is not orthogonal and μ may be unstable. Alternatively, it is better the algorithm described below. The main idea is to avoid any possibility of rounding errors by substituting any intermediate operation -- $q_i^t q_j$ for $i \neq j$ by zero (where zero is its theoretical value). Firstly we obtain the residual b and second the unknown vector μ is obtained. The n -column vector b is obtained in t iterations (recall that Q is a $n \times t$ matrix) in recurrence way, so that at iteration 0 we assign $b^{(0)}=g$, and at the end of iteration t vector $b^{(t)}$ takes the residual vector b .

The justification of the procedure is as follows. Consider eqs. (1.15) and (1.16); note that Q' and R' is the notation of the matrices calculated by NQA and Q and R is the notation of matrices calculated by NNQA. Eqs. (1.15) and (1.16) assume the orthogonality $Q'^t Q' = I$ (or $Q'^t Q \delta^t \delta = I$). Then

$$\begin{aligned} b &= g - Q' Q'^t g = g - \sum_{i=1}^t (q_i^t q_i) q_i' \\ &= g - \sum_{i=1}^{t-1} ((q_i^t g)/d_i) q_i - \\ &\quad - ((q_t^t g)/d_t) q_t \end{aligned} \quad (3.1)$$

Defining

$$b^{(t-1)} = g - \sum_{i=1}^{t-1} ((q_i^t g)/d_i) q_i \quad (3.2)$$

results

$$b = b^{(t-1)} - ((q_t^t g)/d_t) q_t \quad (3.3)$$

Similarly we may substitute eq. (3.2) by

$$b^{(k)} = b^{(k-1)} - ((q_k^t g)/d_k) q_k \quad (3.4a)$$

for $k=1, 2, \dots, t$

where by definition

$$b^{(k-1)} = g - \sum_{i=1}^{k-1} ((q_i^t g)/d_i) q_i \quad (3.4b)$$

where $b^{(k)}$ for $k=t$ is the residual vector b (3.3).

The stability of procedure (3.4) depends on the accuracy of matrix Q calculated by NNQA.

In effect, if we premultiply q_k^t/N_k by eq. (3.4b), it results

$$\begin{aligned} (q_k^t b^{(k-1)})/N_k &= (q_k^t g)/N_k - \\ &= (q_k^t g)/N_k - \sum_{i=1}^{k-1} ((q_i^t g)/N_i) (q_i^t/N_i) \\ &= (q_k^t g)/N_k - \sum_{i=1}^{k-1} ((q_i^t g)/N_i) ((q_k^t q_i)/(N_k N_i)) \end{aligned} \quad (3.5a)$$

The calculation $(q_k^t g)/d_k$ (where $d_k \equiv N_k^2$) of formula (3.4a) will be as precise as it is eq. (3.5a); if in this eq., $q_k^t q_i \neq 0$ ($i=1, 2, \dots, k-1$, then $i \neq k$), $(q_k^t g)$ will not be exact and $b^{(k)}$ calculated by formula (3.4a) will be unstable. Since we know that without rounding errors, $q_k^t q_i = 0$, it is substituted by zero and eq. (3.5a) gives

$$\begin{aligned} (q_k^t b^{(k-1)})/N_k &= (q_k^t g)/N_k \therefore \\ q_k^t b^{(k-1)} &= q_k^t g \end{aligned} \quad (3.5b)$$

Then the final formula to calculate $b^{(k)}$ is

$$b^{(k)} = b^{(k-1)} - ((q_k^t b^{(k-1)})/d_k) q_k \quad (3.6)$$

for $k=1, 2, \dots, t$; and $b^{(0)}=g$, $b \equiv b^{(t)}$

4. CALCULATION OF THE UNKNOWN VECTOR μ

With the same notation Q' , R' , Q and R used in the previous sections, the vector μ that minimizes norm $\|b\|^2$ (1.2) is found by solving eqs. (1.14): $R'\mu = Q'^t g$, that for element μ_k gives

$$R'_{kk} \mu_k + \sum_{i=k+1}^t R'_{ki} \mu_i = q_k^t g \quad (4.1a)$$

for $k = 1, 2, \dots, t-1$

$$R'_{tt} \mu_t = q_t^t g \quad (4.1b)$$

By using matrices Q and R calculated by NNQA (see Sec. 2), eqs. (4.1) give

$$N_k \mu_k = (q_k^t g)/N_k - \sum_{i=k+1}^t N_k R_{ki} \mu_i \quad (4.2a)$$

for $k=1, 2, \dots, t-1$

$$N_t \mu_t = (q_t^t g)/N_t \quad (4.2b)$$

Since $q_k^t g = q_k^t b^{(k-1)}$ (eq. 3.5b), eqs. (4.2) give

$$\mu_t = (q_t^t b^{(t-1)}) / d_t \quad (4.3a)$$

Transactions on Mathematical Software 5, --
1979, pp. 457-465.

$$\mu_k = (q_k^t b^{(k-1)}) / d_k - \sum_{i=k+1}^t R_{ki} \mu_i \quad (4.3b)$$

for $k=1, 2, \dots, t-1$; where $d_k \equiv N_k^2$.

We may note that formulae (3.6) and (4.3) - use the same intermediate operations; then we may avoid to repeat these calculations. Also note that square root calculations are avoided (this operation produces the most of the rounding errors in alternate procedures).

5. CONCLUSION

A procedure for solving an overdetermined - system of linear equations is revisited and some new technical details are introduced. - It calculates the orthogonal matrix Q (NNQA in Sec. 2), the residual vector b (3.6) and the unknown vector μ (4.6) with the maximum accuracy. It is much more stable than procedures that use the pseudoinverse of the system matrix.

6. ACKNOWLEDGEMENT

I thank Dr. Augustin Drubulle and one of the referees for their helpful suggestions.

7. REFERENCES

- /1/ ESCUDERO, L.F., "Lagrange multipliers estimates for constrained optimization", IBM Scientific Center Report, 1980, Palo Alto, California.
- /2/ GOLUB, G.H., "Matrix decomposition and - statistical calculations", Academic Press London, 1969, pp. 381-385.
- /3/ PETERS, G. and WILKINSON, J.H., "The -- least-square problem and pseudoinverses", Computer Journal 13, 1970, pp. 309-316.
- /4/ RICE, J.R., "Experiments on Gram-Schmidt orthogonalization", Mathematical Computations 20, 1966, pp. 325-328.
- /5/ WAMPLER, R.H., "Solutions to weighted least squares problems by Modified Gram-Schmidt with iterative refinement, ACM -