

A Direct Method for the Solution of Sparse Linear Least Squares Problems

Å. Björck*

*Department of Mathematics
Linköping University
S-581 83 Linköping, Sweden*

and

I. S. Duff

*Computer Science and Systems Division
Building 8.9
AERE Harwell
Didcot, Oxon, U.K., OX11 0RA.*

Submitted by Robert J. Plemmons

ABSTRACT

We examine a direct method based on an LU decomposition of the rectangular coefficient matrix for the solution of sparse linear least squares problems. We wish to develop a method which is also stable for weighted systems and at the same time preserves much of the original sparsity. We also describe a general updating scheme for modifying the solution when extra equations are added. This is particularly useful in the case when these added equations are nonsparse. We illustrate our description and analyses by runs on test examples.

1. INTRODUCTION

We consider an algorithm for the direct solution of the linear least squares problem

$$\min_x \|b - Ax\|_2 \quad (1.1)$$

*Research for the first author was supported in part by National Science Foundation Grant No. MCS78-17697.

where A is an $m \times n$ sparse matrix. Duff and Reid [7] compared variants of four methods for solving (1.1) from the point of view of preserving sparsity. They concluded from numerical examples that when stability is important, the method of Peters and Wilkinson [14] was the best in this respect. In particular, they found that methods based on orthogonal transformations which hold these transformations explicitly are particularly bad at preserving sparsity. However, for problems which do not have very heavily weighted rows (see Sec. 4), recent work of George and Heath [10] avoids this problem by discarding these orthogonal factors and working with the normal equations where the Cholesky factors are obtained using orthogonal transformations.

In this paper, we develop and extend the method of Peters and Wilkinson. We discuss the basic method and some minor modifications to it in Sec. 2. We allow the rows of $b - Ax$ to have widely differing weights and also allow linear constraints for which we do not have to introduce an artificial weighting. In Sec. 3 we discuss the stability of our techniques and in Sec. 4 we consider weighted problems. Pivoting strategies for preserving sparsity are discussed in Sec. 5.

In some problems, the matrix A can have a few nonsparse rows, which, if included, would cause severe fill-in in the sparse factorizations. For such problems, we suggest in Sec. 6 the use of an updating algorithm which efficiently updates the solution when a few equations and/or constraints are added. A new feature of this algorithm (compare [9]) is that the initial problem is allowed to be rank deficient.

Throughout this paper, we illustrate our discussion with results of runs on test examples. The examples we use are similar to those in [7], to which we have added two problems discussed in [17] and some small illustrative examples described in this text. All our programs have been written in FORTRAN and compiled using the IBM FORTRAN H extended (enhanced) compiler with OPT=2. All arithmetic is in IBM double precision (approximately 16 decimals), and timings are in milliseconds on an IBM 3033.

2. THE METHOD OF PETERS AND WILKINSON

We now describe a slightly modified version of the method of Peters and Wilkinson [14]. The modified method has advantages for solving problems where the system $Ax=b$ is consistent or nearly consistent. We also discuss how the method can be easily extended to handle linear constraints.

The first step is to compute an LU factorization of the rectangular matrix A , using Gaussian elimination with both row and column interchanges. This

is equivalent to multiplying a permutation of A from the left by the product, M , of a sequence of elementary elimination matrices.

If we carry out the same transformations on the right hand side, we obtain

$$MP_1AP_2 = \begin{pmatrix} U \\ 0 \end{pmatrix} \}_r, \quad (2.1)$$

$$MP_1b = \begin{pmatrix} c \\ d \end{pmatrix} \}_r, \quad (2.2)$$

where P_1, P_2 are permutation matrices, and U is an $r \times n$ upper trapezoidal matrix with $r = \text{rank}(A)$.

If we look at this in terms of an LU decomposition of A ,

$$P_1AP_2 = LU, \quad (2.3)$$

with L a unit lower trapezoidal $m \times r$ matrix, then we have

$$P_1b = Lc + \begin{pmatrix} 0 \\ d \end{pmatrix}. \quad (2.4)$$

Now if x_c is any solution of the system

$$UP_2^T x = c, \quad (2.5)$$

the residual norm corresponding to it is given by (cf. [1])

$$\begin{aligned} \|b - Ax_c\|_2 &= \|P_1(b - Ax_c)\|_2 \\ &= \left\| Lc + \begin{pmatrix} 0 \\ d \end{pmatrix} - Lc \right\|_2 \quad [\text{from (2.3), (2.4), and (2.5)}] \\ &= \|d\|_2. \end{aligned} \quad (2.6)$$

Thus, if $\|d\|_2 < \epsilon$ (ϵ some suitable tolerance), then x_c is a solution to (1.1) with a slightly perturbed right hand side b , and we can immediately accept x_c as the solution to our problem at the cost of a simple forward elimination (2.2) and back substitution (2.5). We show in Table 1 the savings which can be achieved in such a case.

TABLE 1
RESULTS COMPARING THE MODIFIED PETERS-WILKINSON ALGORITHM AGAINST THE
ORIGINAL FOR CONSISTENT TABLE FLATNESS EXAMPLES^a

M , number of rows	84	144	220
N , number of columns	64	100	144
NZ, number of nonzeros	252	432	660
Time for initial decomposition (2.3)	40	78	126
Time for forward elimination (2.2)	1	2	3
Total time for solution (full P&W)	114	226	381
Total time for solution (using consistency)	41	80	131
Final residual (P&W)	9×10^{-15}	5×10^{-14}	7×10^{-14}
Final residual (consistent method)	5×10^{-14}	1×10^{-13}	5×10^{-14}

^aThreshold parameter in the initial decomposition was 0.01.

However, if $\|d\|_2$ is larger, we would like to be able to solve the least squares problem for our initial decomposition (2.1) and (2.2). For an arbitrary x we have that

$$\begin{aligned}
 P_1(b - Ax) &= Lc + \begin{pmatrix} 0 \\ d \end{pmatrix} - LUP_2^T x \\
 &= \begin{pmatrix} 0 \\ d \end{pmatrix} - Lz,
 \end{aligned} \tag{2.7}$$

where $UP_2^T x = c + z$. (2.8)

Hence, x is a least squares solution of (1.1) if it satisfies (2.8), where z is the solution of

$$\min \left\| \begin{pmatrix} 0 \\ d \end{pmatrix} - Lz \right\|_2 \quad [\text{from (2.7)}] \tag{2.9}$$

Since we formed L by pivoting on A , the least squares problem (2.9) can be solved by the method of normal equations:

$$L^T Lz = L^T \begin{pmatrix} 0 \\ d \end{pmatrix}, \tag{2.10}$$

the solution to which can be obtained using the normal equations

$$L_2^T L_2 z_2 = L_2^T \begin{pmatrix} 0 \\ d \end{pmatrix}. \quad (2.14)$$

Then the solution $x = P_2 \tilde{x}$ to the constrained, rank deficient least squares problem is obtained by the back substitution

$$\begin{pmatrix} U_1 & U_2 \end{pmatrix} \begin{pmatrix} \tilde{x}_1 \\ \tilde{x}_2 \end{pmatrix} = \begin{pmatrix} c_1 \\ c_2 + z_2 \end{pmatrix}, \quad (2.15)$$

where the undetermined components \tilde{x}_2 are normally set to zero.

We remark that if linearly dependent constraints occur, this will be detected during the factorization of A , and it is simple to add a check for consistency of such constraints.

3. STABILITY

One of the principal reasons for studying the method of Peters and Wilkinson is the known stability problems in using the normal equations approach. This is highlighted in an example in the following section on weighted least squares problems. In this section we examine the stability of our modification of Peters and Wilkinson.

Clearly, we would expect to require as numerically satisfactory a pivoting strategy for the initial LU decomposition (2.3) as in the solution of linear equations by Gaussian elimination. However, in the least squares case, there are additional features which make life a little more difficult. Clearly, we cannot scale the rows of the coefficient matrix without changing the least squares objective function. Also, any ill-conditioning in L will affect us to a worse degree than for linear equations, since here we are forming $L^T L$ and thus squaring the condition number. Furthermore, it is not true we can immediately take nonzeros in singleton rows (that is rows with only one nonzero in them) as pivot, since the situation

$$\begin{pmatrix} \varepsilon & 0 \\ 1 & 1 \\ 1 & 2 \end{pmatrix}, \quad \varepsilon \ll 1$$

is quite possible when solving a least squares problem (even of full rank), while such an occurrence in the solution of linear equations implies the

matrix is singular. Obviously, pivoting on ε would be disastrous, and in general partial pivoting either by rows or columns is not sufficient.

For these reasons, Peters and Wilkinson [14] use complete pivoting in their LU factorization. We cannot afford to do this in the sparse case, because then we would not be able to preserve sparsity. We therefore relax the stability criterion by allowing a nonzero $a_{kk}^{(k)}$ to be the pivot at the k th step if

$$|a_{kk}^{(k)}| \geq u \max \left\{ \max_{k < i \leq m} |a_{ik}^{(k)}|, \max_{k < j \leq n} |a_{kj}^{(k)}| \right\}, \quad (3.1)$$

where u is a preset parameter in the range 0 to 1. As in the solution of linear equations, we call u the threshold (parameter) and the pivoting criterion (3.1) threshold pivoting, and will normally recommend a compromise value for u of about 0.1. We are, of course, assuming that the whole of the active matrix can be held in core in order that we can implement this pivot selection technique efficiently.

This threshold criterion (3.1) will cause the elements in L to satisfy the bound

$$|l_{ik}| \leq u^{-1}, \quad i > k, \quad (3.2)$$

and will also ensure that for $k = 1, 2, \dots, n$

$$\max_{k < j \leq n} |a_{ij}^{(k+1)}| \leq (1 + u^{-1}) \max_{k < j \leq n} |a_{ij}^{(k)}|, \quad k < i \leq m, \quad (3.3)$$

where $A_k = \{a_{ij}^{(k)}\}$ is the reduced matrix at the start of the k th elimination step. By an obvious extension of the result of Reid [16], we can see that the factorization we obtain is exact for the perturbed matrix

$$A + E$$

with $E = \{e_{ij}\}$, where

$$\max_{1 \leq j \leq n} |e_{ij}| \leq 3.01 \varepsilon \xi_i (1 + u^{-1})^{\xi_i} \max_{1 \leq j \leq n} |a_{ij}|, \quad (3.4)$$

where ξ_i is no greater than the number of nonzeros in row i of L and corresponds to the maximum number of operations performed on any columns up to stage $\min(i-1, j)$. This gives a bound for the relative

perturbation in each row of A and shows that the factorization (2.3) is stable even when the row norms of A differ widely.

It is well known that even when $u = 1$, the inequality (3.2) does not mean that L (and therefore $L^T L$) will be well conditioned. However, in practice, as a result of our pivotal strategy, any ill-conditioning of A will usually be fully reflected in U , and L will be well-conditioned (see [14]). In the next section we show that the situation is essentially the same for weighted least squares problems, even when the weights differ greatly in magnitude.

4. WEIGHTED LEAST SQUARES PROBLEMS

We now consider the weighted least squares problem

$$\min_x \|W(\hat{b} - \hat{A}x)\|_2 \quad (4.1)$$

where W is a nonsingular (often diagonal) weighting matrix. In the following, we will assume \hat{A} to be row equilibrated, so that all rows of \hat{A} have equal norm. It follows from [18] that such a row scaling will approximately minimize the condition number of \hat{A} . It is then $\kappa(\hat{A})$ rather than $\kappa(W\hat{A})$ that determines the conditioning of the problem (4.1).

Note that for the statistical linear model

$$\hat{b} = \hat{A}x + r,$$

where r is an unknown noise vector of zero mean and covariance S , the correct choice of weighting is $W = S^{-1/2}$. Here, we will consider only the case

$$W = D = \text{diag}\{d_1, d_2, \dots, d_m\}.$$

We remark that in practice the weights d_i may vary widely in size, e.g. when certain linear combinations of the unknowns are known to high accuracy.

To illustrate the stability of our methods we will study the particular weighting

$$D = \text{diag} \left\{ \underbrace{w, w, w, \dots, w}_p, 1, 1, \dots, 1 \right\}$$

with $w \gg 1$ and $p < n$. We denote the corresponding scaled quantities by

$$A_w = \begin{pmatrix} w\hat{A}_1 \\ \hat{A}_2 \end{pmatrix}, \quad b_w = \begin{pmatrix} w\hat{b}_1 \\ \hat{b}_2 \end{pmatrix} \quad (4.2)$$

where each row of \hat{A}_1 and \hat{b}_1 is multiplied by w .

For the largest singular value of A we immediately get the estimate $\sigma_{\max}(A_w) \geq w\sigma_{\max}(\hat{A}_1)$. If we assume that the singular vector corresponding to the smallest singular value of \hat{A}_2 lies in the nullspace of \hat{A}_1 , we further have $\sigma_{\min}(A_w) = \sigma_{\min}(\hat{A}_2)$ and the inequality

$$\kappa(A_w) \geq w \frac{\sigma_{\max}(\hat{A}_1)}{\sigma_{\min}(\hat{A}_2)}.$$

We conclude that even when the matrix \hat{A} is well conditioned, the method of normal equations can break down if $w > (\text{macheps})^{-1/2}$, where macheps is the machine precision. Indeed, formation of the normal equations gives

$$(w^2\hat{A}_1^T\hat{A}_1 + \hat{A}_2^T\hat{A}_2)x = w^2\hat{A}_1^T\hat{b}_1 + \hat{A}_2^T\hat{b}_2,$$

thus causing loss of all information in \hat{A}_2 and \hat{b}_2 unless \hat{A}_1 only involves $\text{rank}(\hat{A}_1)$ of the variables. We illustrate this by running our codes on the consistent system

$$\begin{bmatrix} w & w & w \\ 1 & & \\ & 1 & \\ & & 1 \end{bmatrix} x = \begin{bmatrix} 3w \\ 1 \\ 1 \\ 1 \end{bmatrix} \quad (4.3)$$

with varying values of w . The results are shown in Table 2.

TABLE 2
FAILURE OF NORMAL EQUATIONS ON WEIGHTED PROBLEMS

w	10^6	10^9	10^{12}
Normal equations residual	2(-4)	Fails	Fails
$P + W$ residual	5(-10)	1(-7)	3(-4)

Powell and Reid [15] have shown how pivoting can be used with orthogonal methods to produce satisfactory algorithms even in the presence of widely differing weights. We illustrate the similar good behavior of the Peters-Wilkinson technique on the example (4.2). In the following, we always indicate explicitly the dependence on w . The LU factorization of A_w will have the form

$$A_w = \begin{pmatrix} L_{11} & \\ w^{-1}L_{21} & L_{22} \end{pmatrix} \begin{pmatrix} wU_{11}wU_{12} \\ U_{22} \end{pmatrix},$$

where L_{11} is p by p , and so the matrix of the reduced normal equations, L^TL , is of the form

$$L^TL = \begin{pmatrix} L_{11}^TL_{11} + w^{-2}L_{21}^TL_{21} & w^{-1}L_{21}^TL_{22} \\ w^{-1}L_{22}^TL_{21} & L_{22}^TL_{22} \end{pmatrix},$$

which is almost block diagonal and with condition number almost independent of w . The Cholesky factorization of L^TL will yield

$$L^TL = \tilde{L}\tilde{L}^T$$

where

$$\tilde{L}^T = \begin{pmatrix} \tilde{L}_1^T + O(w^{-2}) & w^{-1}\tilde{L}_1^{-1}L_{21}^TL_{22} + O(w^{-3}) \\ 0 & \tilde{L}_2^T + O(w^{-2}) \end{pmatrix},$$

where \tilde{L}_1 and \tilde{L}_2 are the Cholesky factors of $L_{11}^TL_{11}$ and $L_{22}^TL_{22}$, respectively.

In the original method of Peters and Wilkinson, we must solve

$$\tilde{L}\tilde{L}^Ty = L^T \begin{pmatrix} w\hat{b}_1 \\ \hat{b}_2 \end{pmatrix},$$

but the heavy weighting on the rows of \hat{b}_1 will not cause us any instability because the elements in the $(2, 1)$ block of \tilde{L} are of order w^{-1} . There is even less of a problem with our modified method, since the right hand side of our reduced equations will have order unity.

One way of handling linear constraints (for example, [12, Chapter 22]) is to assign large weights to these equations, where "large" is defined to be a

weighting at least $(\text{macheps})^{-1/2}$ greater than for any nonconstraint equation. We do not pursue this approach, for two reasons. Clearly, as we illustrated in Sec. 2, our algorithm simplifies in the presence of constraints (giving us a smaller least squares problem to solve), and secondly, in order to choose a sensible weighting for constraints, we need to know the other weights *a priori* and also be assured that \hat{A} is well scaled. Thus, we effectively assign infinite weights to constraints so that they are always pivoted on first, or at least before any other equations whose pivot variable also appears in a constraint equation. The use of these "infinite" weights does not cause us any trouble in our implementation, because there we do not work with explicitly scaled equations. Thus, our threshold criterion (3.1) becomes

$$d_k |\hat{a}_{kk}^{(k)}| \geq u \max \left\{ \max_{k < i < m} |d_i \hat{a}_{ik}^{(k)}|, \max_{k < j < n} |d_k \hat{a}_{kj}^{(k)}| \right\}, \quad (4.4)$$

and we obtain a factorization

$$\hat{A} = \hat{L} \hat{U}$$

of the unscaled matrix \hat{A} .

If $D = \text{diag}\{d_1, d_2, \dots, d_n\}$, and $A = D\hat{A}$ etc., then this factorization is related to our earlier one (2.3)–(2.4) through

$$\hat{L} = D^{-1}LD, \quad \hat{U} = D^{-1}U,$$

and

$$\begin{pmatrix} \hat{c} \\ \hat{d} \end{pmatrix} = D^{-1} \begin{pmatrix} c \\ d \end{pmatrix},$$

and if $\hat{z} \equiv D^{-1}z$, Eqs. (2.14) and (2.15) can be written

$$\hat{L}^T D^2 \hat{L} \hat{z} = \hat{L}^T D^2 \begin{pmatrix} 0 \\ \hat{d} \end{pmatrix}$$

and

$$\hat{U} P_2^T x = \hat{c} + \hat{z}.$$

Provided that the scaling itself does not introduce any roundoff, we will get the same numerical results whether or not the equations are explicitly scaled.

This property of scaling invariance was shown to hold for Gaussian elimination by Bauer [2], and our result above is an easy generalization of this.

We end this section by establishing an *a posteriori* error analysis for the weighted problem. This is similar to the result in [19] for unweighted problems. We state our result as a theorem.

THEOREM 4.1. *Let \bar{x} and \bar{r} be the computed solution and residual to the least squares problem*

$$\min_x \|W(\hat{b} - \hat{A}x)\|_2.$$

Then, if $b = W\hat{b}$ and $A = W\hat{A}$, \bar{x} and \bar{r} are the exact solution and residual for the perturbed problem

$$\min_x \|(b + e) - (A + E)x\|_2 \quad (4.5)$$

with

$$E = - \frac{WW^T \bar{r} \bar{r}^T A}{\|W^T \bar{r}\|_2^2} \quad (4.6)$$

and

$$e = E\bar{x}. \quad (4.7)$$

Proof. Because of the relationship (4.7) between e and E , the residual corresponding to taking \bar{x} as the solution of (4.5) is

$$(b + e) - (A + E)\bar{x} = b - A\bar{x} = \bar{r},$$

and so to prove the theorem we have only to show that

$$(A + E)^T \bar{r} = 0.$$

But

$$(A + E)^T \bar{r} = A^T \bar{r} - \frac{A^T \bar{r} \bar{r}^T W W^T \bar{r}}{\|W^T \bar{r}\|_2^2} \quad [\text{from (4.6)}]$$

$$= A^T 0 = 0. \quad \blacksquare$$

One consequence of this result is that we can provide a criterion for estimating when our computed solution is satisfactory. This will be when the norm of \hat{E} ($=W^{-1}E$) is small relative to the norm of \hat{A} . If we consider a diagonal scaling $W=D$, then

$$\|D^{-1}E\|_2 = \frac{\|A^T \hat{r}\|_2}{\|D \hat{r}\|_2},$$

and so the relevant test is in terms of the unscaled quantities

$$\|\hat{A}^T D^2 \hat{r}\|_2 \leq \text{tol} \|\hat{A}\|_2 \|D^2 \hat{r}\|_2,$$

where we have put

$$\hat{r} = \hat{b} - \hat{A}\bar{x} = D^{-1}\tilde{r}.$$

5. PRESERVATION OF SPARSITY

There are three stages in Peters and Wilkinson's method at which we can lose sparsity present in the original system. These are

- (i) the original LU decomposition,
- (ii) the formation of the reduced normal equations $L^T L$, and
- (iii) the Cholesky factorization of $L^T L$.

We preserve sparsity when factoring $L^T L$, (iii), by choosing as pivot at each stage that nonzero from the diagonal with the least number of nonzeros in its row. This minimum degree criterion has proved itself effective over a wide range of matrices, and we know of no other technique competitive for general systems. Fill-in (when a zero element becomes nonzero) in steps (i) and (ii) is controlled by sparsity pivoting in the initial decomposition, and it is this which we now consider.

We compare three methods of pivoting for sparsity in the LU decomposition (i). The first is the well tried and tested criterion due to Markowitz [13], where the nonzero chosen as pivot minimizes the product of the number of other nonzeros in its row and the number of other nonzeros in its column, subject to the threshold criterion (3.1) [or (4.4)] being satisfied. However, since we are more concerned with keeping the number of nonzeros in L low, there are strong arguments for adopting a simpler strategy where the pivot at each stage is chosen from the column with the least number of nonzero elements, and within that column we select the nonzero with the least

number of nonzeros in its row, subject to our threshold criterion again being satisfied. We call this criterion c_j .

A third criterion aimed at avoiding creating nonzeros in $L^T L$ was suggested by Björck [4]. Let s_i , $k \leq i \leq m$, be the number of nonzeros in the i th row of the computed part of L . An upper bound on the number of nonzeros added to the lower triangular part of $L^T L$ when the pivot is chosen from the j th column is

$$k_j = \sum_{i \in B_j} (s_i + 1),$$

where B_j is the set of row indices of the nonzero elements in the j th column.

TABLE 3
A COMPARISON OF DIFFERENT ORDERINGS FOR INITIAL LU DECOMPOSITION^a

Problem	P1	P2	P3	P4	P5	P6	P7
Number of equations	197	220	100	1850	958	313	313
Number of variables	150	144	50	713	292	176	176
Nonzeros in A	546	660	300	10608	1916	1557	1557
<i>Nonzeros in L + U</i>							
Markowitz	546	1504	300	18270	1936	1818	1877
c_j	868	1522	386	19578	3234	2803	3142
k_j	825	1540	369	19355	3226	3005	3083
<i>Nonzeros in L</i>							
Markowitz	394	998	250	14935	1634	1374	1428
c_j	515	1032	307	15984	2701	2047	2329
k_j	484	1050	294	15732	2698	2204	2251
<i>Nonzeros in $L^T L$</i>							
Markowitz	628	1712	512	10687	1253	1572	1604
c_j	1499	1685	725	10985	1922	2447	2701
k_j	1239	1722	651	10554	1930	2555	2577
<i>Time for LU decomp.</i>							
Markowitz	29	195	15	3230	227	166	164
c_j	61	114	25	3210	234	228	268
k_j	71	126	27	3149	263	284	295
<i>Total time for soln.</i>							
Markowitz	145	446	81	5266	542	399	399
c_j	272	375	100	5856	688	660	745
k_j	248	391	100	5801	710	756	751

^aThe threshold parameter was set to 0.01 except for the last column, where it was 0.25. The test problems are: P1: random columns from 199; P2: table flatness; P3: random matrix; P4: [17]; P5–P7: surveying.

The column which minimizes k_j is chosen as pivotal column, and then the nonzero in this column with least nonzeros in its row is taken as pivot, subject to satisfying the numerical criterion.

It is hard to show the superiority of any one of these strategies, particularly since there is, in general, no monotonicity between the number of nonzeros in L , the number in $L^T L$, and the number in \tilde{L} , the Cholesky factor of $L^T L$. We have therefore run these three strategies on all of our test examples with several different values for u , the threshold parameter. We show results from a representative sample of these runs in Table 3. Problem P4 is from [17]. The other problems are similar to those described in [7]; however, their order is different, because different random number generators were used in this case, the surveying problems are from slightly different sources, and the table flatness example is larger than those used in that paper.

As is quite clear from these results, the Markowitz criterion is, in general, by far the best, and we have no hesitation in recommending its use as sparsity criterion for the initial LU decomposition of A . All of our other runs for this paper use this criterion. The results of Table 3 also confirm the findings of Duff and Reid [7], who compared the first two criteria described in this section, viz. Markowitz and c_j .

One practical bonus of this result is that the well-tested code of Duff [6] can be used as the basis for producing this initial decomposition. We do need, however, to make some minor changes, which we now discuss.

The MA28 code [6] first permutes the matrix to block lower triangular form; if this were done on the example (4.3), the following form could result:

$$\begin{pmatrix} 1 & & & \\ & 1 & & \\ & & 1 & \\ w & w & w & 0 \end{pmatrix},$$

so that the heavily weighted equation is not included in the solution process. Although this is fine for the consistent system given, it would wreak havoc for an inconsistent set. Although there is a switch to ignore this option, we choose to modify the code, since the resulting program is considerably simpler: not only are calls to the block triangularization routines removed, but also one level of loop structure is removed from the main subroutine, MA30A.

Our pivoting strategy (3.1) or (4.4) requires that we test for stability by columns as well as by rows. In order to avoid excessive searching during the pivot selection stage, we include an additional real array holding the current

maximum element in each column. This makes it very easy to test for stability by columns, but we must, of course, keep updating this maximum element array. Such care is not required in the solution of square nonsingular equations and so in MA28 we pivoted by rows only. As in MA28, we can obtain an *a posteriori* bound on the growth to give us some indication of the accuracy of our decomposition.

Finally, so that we can handle weighted problems with constraints (see Sec. 4), we introduce a real vector of weights and use the stability criterion (4.4). We set $W(I) = -d_i$ for constraint equations and $W(I) = 0$ for equations on which we do not wish to pivot (see Sec. 6 on updating solution); otherwise $W(I)$ is set to d_i , the weight for equation I .

6. UPDATING THE SOLUTION

It is quite common in the formulation of a least squares problem for the matrix A to have a few very dense rows. If these were used as pivot rows during the initial decomposition, severe fill-in would result. However, even if pivots were not chosen from these rows, there would still be a grave fill-in problem when forming and factorizing $L^T L$. In this section, we propose a method which avoids this problem.

An example of such a problem is given by the following set of equations:

$$x_i - x_j = b_{ij}, \quad i = 1, 2, \dots, n, \quad j > i, \quad \sum_{i=1}^n x_i = 1. \quad (6.1)$$

Obviously the first $n(n-1)/2$ equations, which are very sparse, only determine the x_i up to a constant, and a constraint is added to normalize the solution. This corresponds to a dense row in A and will cause $L^T L$ to fill in completely.

Our solution is to treat problems of this type in two steps, where the dense equations are added in the second "updating" step. It is important to realize that we update the solution rather than the factorization, since we wish to preserve the sparsity of the original decomposition and to preserve the ability to update our solution more than once. We note that a similar idea is used in the capacitance matrix method for solving the discrete Poisson equation on an irregular region described by Buzbee et al. [5]. There, the solution to be updated corresponds to a problem on a rectangular region for which fast direct methods exist.

We assume, in the following, that the data are partitioned as

$$A = \begin{pmatrix} A_1 \\ A_2 \end{pmatrix} \begin{matrix} \} M_1 \\ \} M_2 \end{matrix}, \quad b = \begin{pmatrix} b_1 \\ b_2 \end{pmatrix} \begin{matrix} \} M_1 \\ \} M_2 \end{matrix},$$

so that in the second step we will add m_2 equations and/or constraints. We refer to the least squares problem

$$\min \|b_1 - A_1 x\|_2 \quad (6.1)$$

as the unaugmented problem, and (1.1) as the full least squares problem.

Although updating schemes for the least squares problem have been with us since the days of Gauss (see the survey by Farebrother [9]), the novel features of the method we now describe (in addition to being sparsity oriented) are that we allow A_1 to be rank deficient and A_2 to contain a mixture of equations and constraints.

We assume that p_1 and p_2 rows of A_1 and A_2 respectively correspond to linear constraints, and denote by r the rank of A_1 . In the initial LU factorization A_2 and b_2 are adjoined and operated on, but we never pivot on rows in A_2 . After r elimination steps the transformed matrix and vector are structured as in Fig. 2. For ease of notation we assume, in the following, that the rows and columns in A are in pivotal order.

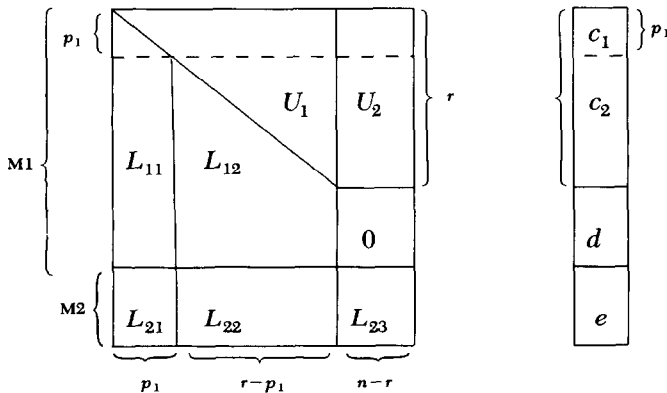


FIG. 2. Initial LU factorization performed on A and b .

The solution

$$x^{(0)} = \left(\begin{array}{c} x_1^{(0)} \\ 0 \end{array} \right) \begin{array}{l} \} r \\ \} n-r \end{array}$$

to the unaugmented least squares problem can be computed as in Sec. 2, where the last $n-r$ components of the solution are taken to be zero. Thus, with the notation of Fig. 2, we have the following equations corresponding to Eqs. (2.15) and (2.14) respectively:

$$U_1 x_1^{(0)} = \left(\begin{array}{c} c_1 \\ c_2 + z^{(0)} \end{array} \right) \begin{array}{l} \} p_1 \\ \} r \end{array} \quad (6.2)$$

with $z^{(0)}$ given by the solution of the $(r-p_1) \times (r-p_1)$ set of equations

$$L_{12}^T L_{12} z^{(0)} = L_{12}^T \begin{pmatrix} 0 \\ d \end{pmatrix}, \quad (6.3)$$

which we solve by using the Cholesky decomposition

$$L_{12}^T L_{12} = \tilde{L} \tilde{L}^T. \quad (6.4)$$

The complete residual vector corresponding to the above solution $x^{(0)}$ is given by [cf. Eq. (2.13)]

$$\begin{aligned} r_1^{(0)} &= b_1 - A_1 x^{(0)} = \begin{pmatrix} 0 \\ r_{12}^{(0)} \end{pmatrix}, & r_{12}^{(0)} &= \begin{pmatrix} 0 \\ d \end{pmatrix} - L_{12} z^{(0)}, \\ r_2^{(0)} &= b_2 - A_2 x^{(0)} = e - L_{22} z^{(0)} \end{aligned} \quad (6.5)$$

Now let

$$x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$

be any vector which satisfies the p_1 constraints in block A_1 , and define $z = z^{(0)} + \Delta z$ by

$$(U_1 \quad U_2) \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \begin{array}{l} \} r \\ \} n-r \end{array} = \begin{pmatrix} c_1 \\ c_2 + z \end{pmatrix} \begin{array}{l} \} p_1 \\ \} r-p_1 \end{array}. \quad (6.6)$$

It follows that we can write the corresponding residual vector as

$$r_1 = b_1 - A_1 x = \begin{pmatrix} 0 \\ r_{12} \end{pmatrix}, \quad r_{12} = r_{12}^{(0)} - L_{12} \Delta z \quad (6.7a)$$

and

$$r_2 = b_2 - A_2 x = r_2^{(0)} - \begin{pmatrix} L_{22} & L_{23} \end{pmatrix} \begin{pmatrix} \Delta z \\ x_2 \end{pmatrix}. \quad (6.7b)$$

From (6.3) and (6.5) it follows that $r_{12}^{(0)}$ is orthogonal to the columns of L_{12} and therefore

$$\|r_1\|_2^2 = \|r_1^{(0)}\|_2^2 + \|L_{12} \Delta z\|_2^2 \quad (6.8)$$

directly from Eq. (6.7a). If we change variables to

$$u = \tilde{L}^T \Delta z, \quad (6.9)$$

where \tilde{L} is the Cholesky factor defined in Eq. (6.4), then because of the orthogonality of $L_{12} \tilde{L}^{-T}$ we have from (6.8) that

$$\|r_1\|_2^2 = \|r_1^{(0)}\|_2^2 + \|u\|_2^2.$$

Therefore the problem of choosing x to solve the full least squares problem, including the p_2 constraints in block A_2 , is equivalent to

$$\min \left\| \begin{pmatrix} u \\ r_2 \end{pmatrix} \right\|_2^2, \quad (6.10)$$

subject to the constraints (6.7b), which we now write

$$\begin{pmatrix} C & J_{M2} & B \end{pmatrix} \begin{pmatrix} u \\ r_2 \\ x_2 \end{pmatrix} \begin{matrix} \}^{r-p_1} \\ \}^{M2} \\ \}^{n-r} \end{matrix} = r_2^{(0)}. \quad (6.11)$$

Here

$$C = L_{22} \tilde{L}^{-T} \quad \text{of dimension } M2 \text{ by } r - p_1,$$

$$B = L_{23} \quad \text{of dimension } M2 \text{ by } n - r,$$

and J_{M2} is the identity matrix of order $M2$ with p_2 elements equal to zero corresponding to constraints among equations in block A_2 . We now perform $r' = \text{rank}(C J_{M2})$ elimination steps on the equations in (6.11) with the vector $r_2^{(0)}$ appended. Since this system is small $[M2 \times (n - p_1 + M2)]$, we use complete pivoting but do not pivot on B . We show the factors diagrammatically in Fig. 3, where U_c is an $r' \times (r - p_1 + M2)$ unit trapezoidal matrix.

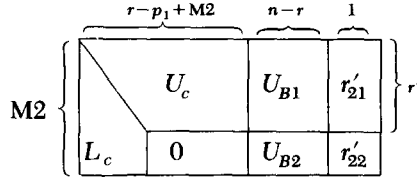


FIG. 3. LU factorization performed on (6.11).

The constraints (6.11) then become (using notation in Fig. 3)

$$U_c \begin{pmatrix} u \\ r_2 \end{pmatrix} + U_{B1} x_2 = r'_{21}, \quad (6.12)$$

$$U_{B2} x_2 = r'_{22},$$

For any fixed x_2 satisfying the last $M2 - r'$ constraint equations above, the problem (6.10)–(6.11) becomes

$$\min \left\| \begin{pmatrix} u \\ r_2 \end{pmatrix} \right\|_2^2 \quad \text{subject to} \quad U_c \begin{pmatrix} u \\ r_2 \end{pmatrix} = s_{21},$$

where

$$s_{21} = r'_{21} - U_{B1} x_2. \quad (6.13)$$

This is a minimum norm problem, with the solution

$$\begin{pmatrix} u \\ r_2 \end{pmatrix} = U_c^T (U_c U_c^T)^{-1} s_{21} = U_c^T \tilde{L}_c^{-T} \tilde{L}_c^{-1} s_{21}, \quad (6.14)$$

where we have used the Cholesky decomposition

$$U_c U_c^T = \tilde{L}_c \tilde{L}_c^T. \quad (6.15)$$

Then from (6.14), using the orthogonality of $U_c^T \tilde{L}_c^{-T}$, we get

$$\left\| \begin{pmatrix} u \\ r_2 \end{pmatrix} \right\|_2^2 = \|\tilde{L}_c^{-1} s_{21}\|_2^2.$$

Thus we can now find x_2 by solving the reduced $r' \times (n-r)$ constrained least squares problem

$$\min_{x_2} \|\tilde{L}_c^{-1}(r'_{21} - U_{B1}x_2)\|_2 \quad \text{subject to} \quad U_{B2}x_2 = r'_{22}. \quad (6.16)$$

We then compute u using (6.14) as

$$u = \begin{pmatrix} I_{r-p_1} & 0 \end{pmatrix} U_c^T \tilde{L}_c^{-T} \tilde{L}_c^{-1} s_{21}, \quad (6.17)$$

Δz from (6.9), and finally x_1 from (6.6).

Note that the dimensions of the matrices in (6.16) and (6.17) are small: \tilde{L}_c is an $r' \times r'$ matrix, and the least squares problem (6.16) has dimension $r' \times (n-r)$, where $n-r$ is the rank deficiency of A_1 .

We thus obtain:

OUTLINE FOR UPDATING ALGORITHM.

(1) Decompose A_1 to obtain L_{12} , and perform forward elimination on the right hand side and on the equations of block A_2 to obtain L_{22} , L_{23} , c , d , and e (see Fig. 2). Factor the reduced normal equations to get the Cholesky factor \tilde{L} and $z^{(0)}$, the solution to that problem, according to (6.2)–(6.4).

(2) Compute

$$r_2^{(0)} = e - L_{22}z^{(0)} \quad [\text{Eq. (6.5)}].$$

(3) Solve $\tilde{L}C^T = L_{22}^T$ to obtain

$$C = L_{22}\tilde{L}^{-T}.$$

(4) Perform elimination steps, using full pivoting on columns of C and J_{M2} , on

$$(C J_{M2} B r_2^{(0)}),$$

where $B = L_{23}$ to obtain

$$U_c, U_{B1}, U_{B2}, \text{ and } r'_{21}/r'_{22} \quad \text{as in Fig. 3.}$$

(5) Form $U_c U_c^T$ and perform Cholesky decomposition to obtain \tilde{L}_c [Eq. (6.15)].

(6) Form $\tilde{L}_c^{-1}(U_{B1} \ r'_{21})$ and compute x_2 as the solution to the constrained least squares problem (6.16).

(7) Compute first $s_{21} = r'_{21} - U_{B1}x_2$ and then u from Eq. (6.17).

(8) Solve

$$\tilde{L}^T \Delta z = u \quad [\text{Eq. (6.9)}]$$

for Δz , and put $z = z^{(0)} + \Delta z$.

(9) Hence the solution to the full least squares problem is

$$x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix},$$

where [Eq. (6.6)]

$$U_1 x_1 = c + \begin{pmatrix} 0 \\ z \end{pmatrix} \begin{matrix} p_1 \\ r-p_1 \end{matrix} - U_2 x_2$$

with

U_2 an $r \times (n-r)$ matrix,
 U_1 an $r \times r$ nonsingular matrix.

$$M \left\{ \begin{array}{cccccc} & & & & & N \\ \hline 1 & -1 & & & & \\ 2 & -1 & & & & \\ & 2 & -1 & & & \\ & & & -1 & & \\ & & & 2 & -1 & \\ \hline 1 & 1 & 1 & \dots & 1 & 1 \end{array} \right.$$

FIG. 4. Form of matrix used for results in Table 4.

We have written a computer program to implement this updating scheme, and we illustrate the effectiveness of this strategy by some results in Table 4. We have chosen a sequence of very simple matrices of the form shown in Fig. 4, where the matrix of the unaugmented system corresponds to the first $M-1$ rows of the matrix in that figure. We have shown runs with N set to 10, 20, and 30 where in each case $M=N+1$. In all cases, the right

TABLE 4
ILLUSTRATION OF EFFECTIVENESS OF UPDATING SCHEME

Number of variables, $N=$	10	20	30
Number of nonzeros in reduced normal equations matrix ($L_{12}^T L_{12}$):			
Using $P+W$	55	210	465
Using updating scheme	11	31	51
Time for solution of complete least squares problem:			
Using $P+W$	7	20	44
Using updating scheme	5	10	15
Residual measure (see Sec. 4)			
Using $P+W$	6×10^{-15}	4×10^{-14}	7×10^{-13}
Using updating scheme	7×10^{-15}	9×10^{-15}	6×10^{-14}

hand sides were generated by a pseudo-random-number generator. It is appropriate to give just the number of nonzeros in the factors of $L_{12}^T L_{12}$, because this is proportional to the total storage required for the solution of the reduced equations. The storage for the full matrix code update is minimal and can be overlaid with storage used in the sparse matrix code parts of the algorithm.

7. CONCLUSIONS

We have developed a stable direct method for solving sparse linear least squares problems. One useful feature of the method is that it detects consistency and treats consistent problems by solving a set of linear equations by a standard method. Also, a general updating algorithm is given, which allows adding a small number of nonsparse equations after the sparse solution phase. We remark that it is possible to let the algorithm select these nonsparse equations automatically, and will develop this idea in a later paper. In many practical applications the user may wish to obtain explicit values of entries in the variance-covariance matrix $[(A^T A)^{-1}]$ in our terminology. Since this matrix is generally full, it would be impractical to generate it in cases with many variables. Additionally, our partitioning and updating scheme discussed in Sec. 6 causes further complications. We plan at a later date to devise efficient methods for obtaining specific variances or covariances based on generalizations of work by Erisman and Tinney [8].

It should be pointed out that for well-conditioned sparse problems, especially using double precision, the method of normal equations usually gives a solution of sufficient accuracy. Although nothing can be said in general, the method of normal equations is then often faster and uses less storage. However, note that our algorithm is also stable for weighted systems, even when the weights vary widely in size. For such problems the method of normal equations should be avoided if possible.

The authors would like to thank Mike Powell for stimulating discussions during this work. In particular the decision to include the updating step is due to one of his comments.

REFERENCES

- 1 I. Barrodale and G. Stuart, A new variant of Gaussian elimination, *J. Inst. Math. Appl.* 19: 39–47 (1977).
- 2 F. L. Bauer, Optimal scaling of matrices and the importance of the minimal condition, in *Proceedings of the IFIP Conference* (C. M. Popplewell, Ed.), North-Holland, 1962, pp. 198–201.
- 3 Å. Björck, Methods for sparse linear least squares problems, in *Sparse Matrix Computations* (J. Bunch and D. Rose, Eds.), Academic, 1976, pp. 177–199.
- 4 Å. Björck, Numerical algorithms for linear least squares problems, Report No. 2/78, Dept. of Mathematics, Univ. of Trondheim, 1978.
- 5 B. L. Buzbee, F. W. Dorr, J. A. George, and G. H. Golub, The direct solution of the discrete Poisson equation on irregular regions, *SIAM J. Numer. Anal.* 8:722–736 (1971).
- 6 I. S. Duff, MA28—a set of FORTRAN subroutines for sparse unsymmetric linear equations, Harwell Report AERE-R 8730, July 1977.
- 7 I. S. Duff and J. K. Reid, A comparison of some methods for the solution of sparse overdetermined systems of linear equations, *J. Inst. Math. Appl.* 17:267–280 (1976).
- 8 A. M. Erisman and W. F. Tinney, On computing certain elements of the inverse of a sparse matrix, *Comm. ACM* 18:177–179 (1975).
- 9 R. W. Farebrother, An historical note on the least squares updating formulas, Report, Univ. of Manchester, 1978.
- 10 A. George and M. T. Heath, Solution of sparse linear least squares problems using Givens rotations, *Linear Algebra and Appl.*, 34:69–83 (1980).
- 11 K. H. Haskell and R. J. Hanson, An algorithm for linear least squares problems with equality and nonnegativity constraints, Report SAND 77-0552, Sandia Laboratories, Albuquerque, N.Mex., 1978.
- 12 C. L. Lawson and R. J. Hanson, *Solving Least Squares Problems*, Prentice-Hall, 1974.
- 13 H. M. Markowitz, The elimination form of the inverse and its application to linear programming, *Management Sci.* 3:255–269 (1957).

- 14 G. Peters and J. H. Wilkinson, The least squares problem and pseudoinverses, *Comput. J.* 13:309–316 (1970).
- 15 M. J. D. Powell and J. K. Reid, On applying Householder transformations to linear least squares problems, in *Information Processing 68* (A. J. H. Morell, Ed.), North-Holland, 1969, pp. 122–126.
- 16 J. K. Reid, A note on the stability of Gaussian elimination, *J. Inst. Math. Appl.* 8:374–375 (1971).
- 17 M. Saunders, Sparse least squares by conjugate gradients: a comparison of preconditioning methods, in *Proceedings of Computer Science and Statistics, 12th Annual Symposium on the Interface*, 1979.
- 18 A. van der Sluis, Condition numbers and equilibration of matrices, *Numer. Math.* 14:14–23 (1969).
- 19 G. W. Stewart, Research, development and LINPACK, in *Mathematical Software III* (J. R. Rice, Ed.), Academic, New York, 1977, pp. 1–14.
- 20 J. H. Wilkinson, *The Algebraic Eigenvalue Problem*, Oxford U.P., 1965.

Received January 1980; revised 12 April 1980