# Downdating the Cholesky factorization[1]

A. Bojanczyk, R. Brent

*Centre for Mathematical Analysis*
*The Australian National University*
*GPO Box 4, Canberra ACT 2601, Australia*


P. van Dooren

*Philips Research Laboratory*
*Avenue Van Becelaere 2, Box 8*
*B-1170 Brussels, Belgium*


F. de Hoog

*Division of Mathematics and Statistics*
*CSIRO, GPO Box 1965*
*Canberra ACT 2601, Australia*

We analyse and compare three algorithms for "downdating" the Cholesky factorization of a positive definite matrix. Although the algorithms are closely related, their numerical properties differ. Two algorithms are stable in a certain "mixed" sense while the other is unstable. As well as comparing the numerical properties of the algorithms, we compare their computational complexity and their suitability for implementation on parallel or vector computers.

## 1. INTRODUCTION

The Cholesky downdating problem considered in this note is: given an upper triangular matrix $R \in \mathfrak{R}^{n \times n}$ and a vector $x \in \mathfrak{R}^n$ such that $R^T R - xx^T$ is positive definite, find an upper triangular matrix $U$ such that

$$U^T U = R^T R - xx^T.$$

By our assumption of positive definiteness, $U$ exists and is unique up to the signs of its diagonal elements (GOLUB and VAN LOAN [10]).

The Cholesky downdating problem is closely related to that of downdating a $QR$-factorization. To show this, let

$$A = \begin{bmatrix} x^T \\ \tilde{A} \end{bmatrix} = Q \begin{bmatrix} R \\ 0 \end{bmatrix},$$

where $A \in \mathfrak{R}^{m \times n}$, $m > n$ and $Q \in \mathfrak{R}^{m \times m}$ is an orthogonal matrix. The problem is

---

1. This work was done while the third author was visiting the Centre for Mathematical Analysis at The Australian National University

recursive method, "Algorithm B", which has some computational advantages over algorithm A but does not have comparable stability properties. The recursive algorithm is then modified in Section 2c so that its stability properties are improved. The resulting algorithm will be called "Algorithm C".

The main results of this note are given in Section 3, where we present an error analysis of Algorithms B and C. In particular, we show that Algorithm C is not backward stable in the classical sense but does satisfy a "mixed" error bound which shows that it gives forward errors of the same order as an algorithm which is backward stable in the classical sense (see, e.g., DE JONG [4]). A similar result was obtained by STEWART [13] for Algorithm A. The best error bounds that could be obtained for Algorithm B depend on the data and can be arbitrarily large, from which we conclude that it is numerically unstable [4].

In Section 4 we show that Algorithms B and C are preferable to Algorithm A from the point of view of computational complexity, as they require about 20 percent fewer floating point multiplications ($2n^2 + O(n)$ versus $5n^2/2 + O(n)$). We also show that Algorithms B and C are more readily implemented on parallel or vector computers than is Algorithm A. Finally, in Section 5 some numerical results verifying our stability analysis are presented.

To summarize, Algorithm C appears the best overall, it is faster than Algorithm A and more stable than Algorithm B.

## 2. THREE DIFFERENT APPROACHES

In this section we compare and relate three different approaches that can be used for solving the downdating problem.

### 2.1. *The LINPACK method (Algorithm A)*

This method is the one implemented in the LINPACK package and is described by STEWART in [13]. A sequence of $(n+1) \times (n+1)$ Givens rotations of the form

$$J_k = \begin{bmatrix} \cos\theta_k & 0 & \cdots & \sin\theta_k & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & & & \vdots & & & \vdots \\ -\sin\theta_k & 0 & \cdots & \cos\theta_k & 0 & \cdots & 0 \\ 0 & 0 & \cdots & 0 & 1 & \cdots & 0 \\ \vdots & & & \vdots & & & \vdots \\ 0 & 0 & \cdots & 0 & 0 & \cdots & 1 \end{bmatrix} \leftarrow (k+1)\text{st row}$$

$(k+1)\text{st col}$

with $k = 1, \cdots, n$, are used to compute $U$ via the relationship

$$J_1 \cdots J_n \begin{bmatrix} 0^T \\ R \end{bmatrix} = \begin{bmatrix} x^T \\ U \end{bmatrix}. \tag{1}$$

---

to find an orthogonal matrix $\tilde{Q} \in \mathbb{R}^{(m-1)\times(m-1)}$ such that

$$\tilde{A} = \tilde{Q} \begin{bmatrix} U \\ 0 \end{bmatrix}.$$

Thus we have

$$R^T R = A^T A = xx^T + \tilde{A}^T \tilde{A} = xx^T + U^T U$$

which is precisely the Cholesky downdating problem described previously.

The problem of downdating a $QR$-factorization occurs for example when an observation (such as an outlier) is deleted from a regression. An algorithm for computing $\tilde{Q}$ and $U$ has been described by GOLUB [9] and GILL, GOLUB, MURRAY and SAUNDERS [8] and it is not difficult to show that the algorithm is backward stable in the classical sense (WILKINSON [9]). This rather satisfactory state of affairs is due to the fact that we know precisely which row (in our case, the first row $x^T$) of $A$ is to be deleted. However, this is not the case when downdating the Cholesky factorization. Here $x^T$ may bear no relation to rows of $R$ and the only requirement is that the matrix $R^T R - xx^T$ is positive definite. Thus, the best we can expect when computations are performed with finite precision is that the computed $U$ is the exact Cholesky factor of $\tilde{R}^T \tilde{R} - \tilde{x}\tilde{x}^T$ where $\tilde{R}$ and $\tilde{x}$ are close to $R$ and $x$ respectively. However, STEWART [13] has shown that under certain circumstances, small perturbations in $R$ and $x$ may lead to large perturbations in $U$. Thus, we do not recommend that algorithms based on downdating Cholesky factors (as described in the present note) be used to downdate the triangular factor in the $QR$ decomposition.

Nevertheless, the Cholesky downdating problem is important in its own right. Let $B$ be a positive definite matrix for which a Cholesky factorization has been obtained, that is

$$B = R^T R.$$

Now suppose we wish to compute the Cholesky factors of the positive definite matrix

$$\tilde{B} := B - xx^T =: U^T U.$$

Then,

$$U^T U = R^T R - xx^T.$$

This occurs quite often as a subproblem when a positive definite matrix is perturbed by a matrix of low rank. Often the problem can be solved by a sequence of rank-1 updates followed by a sequence of rank-1 downdates. Such a situation occurs for example in structural problems when elements are added and deleted from a design (see for example ARGYRIS and ROY [1]). For another application of Cholesky downdating see BOJANCZYK, BRENT and DE HOOG [3].

In Section 2a, we describe a method implemented in the LINPACK package, "Algorithm A", analysed by STEWART [13]. In Section 2b we describe a

and since $J_1, \ldots, J_n$ are orthogonal matrices, it is easy to verify that

$$R^T R = xx^T + U^T U. \qquad (2)$$

Although the right hand side of (1) will always be an upper Hessenberg matrix, we clearly require some restrictions on $J_1, \ldots, J_n$ to ensure that the first row is indeed the required $x^T$. Let

$$J_n^T \cdots J_1^T e_1 = q =: \begin{bmatrix} \alpha \\ u \end{bmatrix}$$

where $e_1 = (1.0, \ldots, 0)^T$. Then, from (1),

$$a^T R = e_1^T J_1 \cdots J_n \begin{bmatrix} 0^T \\ R \end{bmatrix} = x^T.$$

Thus we must have

$$R^T a = x \qquad (3)$$

and

$$\alpha^2 = 1 - a^T a. \qquad (4)$$

We note that a necessary and sufficient condition for $R^T R - xx^T$ to be positive definite is $\alpha^2 > 0$.

The basic steps in the algorithm are therefore as follows. First the $n$-vector $a$ is obtained by forward substitution from (3) and $\alpha$ is calculated from (4) (the sign of $\alpha$ is immaterial and the positive root of (4) is usually taken). Givens rotations $J_1, \ldots, J_n$ satisfying

$$J_1 \cdots J_n q = e_1$$

are constructed and $U$ is then calculated from (1).

The above algorithm has been analysed by STEWART [13] and his error analysis is summarized in Section 3. Roughly speaking, the algorithm performs as well as can be expected given the potentially ill-conditioned nature of the downdating problem. However, as discussed in Section 4, the algorithm is not particularly well suited to parallel implementation, the main problem being that the forward substitution (3) needs to be completed before the calculation of (1) can commence.

### 2.2. A recursive method (Algorithm B)

We now describe another standard algorithm for downdating which goes back to GOLUB [9]. Although we derive the method from the previous algorithm to show the close connection between them, the usual derivation is based on completely different ideas.

We begin rewriting (1) as

$$\begin{bmatrix} 0^T \\ R \end{bmatrix} = J_n^T \cdots J_1^T \begin{bmatrix} x^T \\ U \end{bmatrix} \qquad (5)$$

and define

$$\begin{bmatrix} x^{(k)T} \\ R^{(k)} \end{bmatrix} := J_k^T \cdots J_1^T \begin{bmatrix} x^T \\ U \end{bmatrix}, \quad k = 0.1, \ldots, n \qquad (6)$$

where we have used the convention that $J_1^T \cdots J_1^T$ is the identity when $k=0$. It is easy to verify that $x^{(0)} = x$, $R^{(0)} = U$, $x^{(n)} = 0$ and $R^{(n)} = R$. Furthermore, the first $k$ components of $x^{(k)}$ are zeros and the first $k$ rows of $R^{(k)}$ and $R$ are the same as are the last $n - k$ rows of $R^{(k)}$ and $U$. From (5) and (6).

$$\begin{bmatrix} x^{(k)T} \\ R^{(k)} \end{bmatrix} = J_k^T \begin{bmatrix} x^{(k-1)T} \\ R^{(k-1)} \end{bmatrix}, \quad k = 1, \ldots, n \qquad (7)$$

which can be shown to be equivalent to ([2]):

$$\begin{bmatrix} x^{(k)T} \\ R^{(k-1)} \end{bmatrix} = S_k \begin{bmatrix} x^{(k-1)T} \\ R^{(k)} \end{bmatrix}, \quad k = 1, \ldots, n \qquad (8)$$

where

$$S_k = \begin{bmatrix} \sec\theta_k & 0 & \cdots & -\tan\theta_k & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & & & & & & \vdots \\ -\tan\theta_k & 0 & \cdots & \sec\theta_k & 0 & \cdots & 0 \\ 0 & 0 & \cdots & 0 & 1 & \cdots & 0 \\ \vdots & & & & & & \vdots \\ 0 & 0 & \cdots & 0 & 0 & \cdots & 1 \end{bmatrix} \quad \leftarrow (k+1)\text{st row}$$

$$(k+1)\text{st col}$$

From (2). $\cos^2\theta_k > \alpha^2$, $k = 1, \ldots, n$ and hence $S_k$, $k = 1, \ldots, n$ are well defined.

To see that (8) forms the basis of a recursive algorithm to compute $U$, let us focus attention on the first and $(k+1)$-st row of (8). As noted previously, the $k$-th row of $R^{(k)}$ is the same as the $k$-th row of $R$ and is therefore known. Let us suppose that in the $k$-th step we also know $x^{(k-1)}$. The role played by $\theta_k$ in $S_k$ (or equivalently $J_k$) is to ensure that the $k$-th component of $x^{(k)}$ is zero. This can be achieved by setting

$$\cos\theta_k = \sqrt{r_{kk}^2 - (x_k^{(k-1)})^2}/r_{kk} \qquad (9)$$

and

$$\sin\theta_k = x_k^{(k-1)}/r_{kk}. \qquad (10)$$

We now calculate $x^{(k)}$ and the $k$-th row of $U$ (or equivalently the $k$-th row of $R^{(k-1)}$) from (8) as

$$x_j^{(4)} = (x_j^{(k-1)} - r_{kj}\sin\theta_k)/\cos\theta_k, \quad j = k+1, \ldots, n \tag{11}$$

and

$$u_{kk} = \sqrt{\tilde{r}_k - (x_k^{(k-1)})^2} \tag{12}$$

$$u_{kj} = (-x_j^{(k-1)}\sin\theta_k + r_{kj})/\cos\theta_k, \quad j = k+1, \ldots, n. \tag{13}$$

Since $x^{(0)} = x$ is known, the relations (9-13) yield a recursive method for calculating $U$.

Using (8), it can be shown that

$$\begin{bmatrix} 0^T \\ U \end{bmatrix} = S_n \cdots S_1 \begin{bmatrix} x^T \\ R \end{bmatrix}$$

and it easy to verify that $S_k, k = 1, \ldots, n$ are $\Sigma$-unitary with

$$\Sigma = \begin{bmatrix} -1 & 0^T \\ 0 & I_n \end{bmatrix},$$

where a matrix $A$ is $\Sigma$-unitary if

$$A^T \Sigma A = \Sigma.$$

The product

$$S_c = S_n \cdots S_1$$

is also $\Sigma$-unitary and hence

$$U^T U = [0|U^T]\Sigma\begin{bmatrix} 0^T \\ U \end{bmatrix} = [x|R^T]S^T\Sigma S\begin{bmatrix} x^T \\ R \end{bmatrix}$$

which is an independent verification that $U$ is the required Cholesky factor. Thus, the method can be regarded as that of finding a product of planar $\Sigma$-unitary transformations to reduce $[x|R^T]^T$ to $[0|U^T]^T$ in the same way that Givens rotations are used to reduce the matrix in the updating problem (see for example GOLUB [9] and GILL, GOLUB, MURRAY and SAUNDERS [8]).

The recursive method has a number of advantages over Algorithm A. Since there is no forward substitution phase, the operation count is somewhat less. In addition the algorithm is quite well suited to parallel implementation. However, the error analysis in Section 3 indicates that the stability properties of the method are substantially inferior to those of Algorithm A.

2.3. A modified recursive algorithm (Algorithm C)
It turns out that Algorithm B can be modified so that its stability properties are substantially improved.

At the $k$-th step, consider the $(k+1)$-st row of (7). In component form, this can be written as

$$u_{kk} = \sqrt{\tilde{r}_{kk} - (x_k^{(k-1)})^2} \tag{14}$$

$$u_{kj} = (-x_j^{(k-1)}\sin\theta_k + r_{kj})/\cos\theta_k, \quad j = k+1, \ldots, n \tag{15}$$

where $\cos\theta_k$ and $\sin\theta_k$ have been calculated using (9) and (10) respectively. Now, consider the first row of (7) and rewrite it as

$$x_j^{(k)} = \cos\theta_k x_j^{(k-1)} - \sin\theta_k u_{kj}, \quad j = k+1, \ldots, n. \tag{16}$$

This is clearly a well defined calculation as $u_{kj}, j = k+1, \ldots, n$ have already been calculated via (15).

Thus, (9) and (15-16) form the basis of an algorithm for calculating $U$ which differs only slightly from Algorithm B. However, as we shall see in Section 3, this small modification enables us to establish stability estimates that are comparable to those for Algorithm A. In addition, we still retain all of the desirable features of Algorithm B. Note that this method was also suggested in LAWSON and HANSON [12] but without any comments on its stability.

3. ERROR ANALYSIS
We denote quantities stored in the computer with a tilde. In addition, $\epsilon$ is the relative precision of the machine considered and terms of order $\epsilon^2$ are neglected.

3.2. Algorithm A
An error analysis of this method is to be found in [13]. It is shown there that there exists an exactly orthogonal matrix Q (which is not computed) such that

$$Q\begin{bmatrix} 0^T \\ R \end{bmatrix} = \begin{bmatrix} x^T + \Delta x^T \\ \tilde{U} + \Delta U \end{bmatrix} \tag{17}$$

where

$$\|[\Delta\tilde{U}]_i\|_2 \leq 6n\epsilon\|[R]_i\|_2 \tag{18}$$

$$|\Delta x_i| \leq [(13n+5)/2 + (i+2)\sqrt{i}]\epsilon\|[R]_i\|_2.$$

Here we use $[\cdot]_i$ to denote the $i$-th column of a matrix. From this, one easily obtains a bound for the total error matrix.

$$\left\| \begin{bmatrix} \Delta x^T \\ \Delta\tilde{U} \end{bmatrix} \right\|_F \leq [n^2/2 + 9n\sqrt{n} + O(n)]\epsilon\|R\|_F. \tag{20}$$

Notice that errors are not only superimposed on the data — $R$ and $x$ in this case — but also on the result $\tilde{U}$. Hence, the bound (20) does not guarantee forward or backward stability of the algorithm. but rather proves what could be called "mixed" stability. As argued in [13], the forward part of the error. $\Delta\tilde{U}$, is in fact unimportant since the "true" forward error $(U - \tilde{U})$ is usually much larger and mainly depends on $\Delta x$. the backward part of the error in (17). The result of mixed stability is thus as satisfactory in practice as backward stability since error bounds in both cases are comparable.

A. Bojanczyk, R. Brent, P. van Dooren, F. de Hoog

### 3.2. Algorithm B

The error analysis of this method depends on how it is implemented. We therefore write the order of computations for one step, which, without loss of generality, can be the first step as given in (9-13).

Let us denote the elements of the rows involved in this step as follows.

$$1/c_1 \begin{bmatrix} 1 & -s_1 \\ -s_1 & 1 \end{bmatrix} \begin{bmatrix} x_1^{(0)} x_2^{(0)} \cdots x_n^{(0)} \\ r_{11} \ r_{12} \cdots r_{1n} \end{bmatrix} = \begin{bmatrix} 0 \ x_2^{(1)} \cdots x_n^{(1)} \\ u_{11} \ u_{12} \cdots u_{1n} \end{bmatrix} \tag{21}$$

Now the operations are performed in the following order (where we use the $fl(\cdot)$ notation of WILKINSON [14]).

$$\begin{aligned}
\bar{u}_{11} &= fl(\sqrt{(r_{11}-x_1^{(0)})(r_{11}+x_1^{(0)})}) \\
\bar{c}_1 &= fl(\bar{u}_{11}/r_{11}) \\
\bar{s}_1 &= fl(x_1^{(0)}/r_{11}) \\
\text{for} \quad i &= 2 \text{ until } n \\
\bar{u}_{1i} &= fl((r_{1i}-\bar{s}_1 x_i^{(0)})/\bar{c}_1) \\
\bar{x}_i^{(1)} &= fl((x_i^{(0)}-\bar{s}_1 r_{1i})/\bar{c}_1).
\end{aligned} \tag{22}$$

In the Appendix we show that for this implementation of (21), the following equality holds.

$$1/c_1 \begin{bmatrix} 1 & -s_1 \\ -s_1 & 1 \end{bmatrix} \begin{bmatrix} x_i^{(0)}+\Delta x_i^{(0)} \\ r_{1i}^{(0)}+\Delta r_{1i} \end{bmatrix} = \begin{bmatrix} \bar{x}_i^{(1)} \\ \bar{u}_{1i} \end{bmatrix} \tag{23}$$

where

$$\left\| \begin{bmatrix} \Delta x_i^{(0)} \\ \Delta r_{1i} \end{bmatrix} \right\|_2 \le (8\epsilon/|c_1|) \left\| \begin{bmatrix} x_i^{(0)} \\ \bar{u}_{1i} \end{bmatrix} \right\|_2. \tag{24}$$

If we rewrite (23) as

$$1/c_1 \begin{bmatrix} 1 & -s_1 \\ -s_1 & 1 \end{bmatrix} \begin{bmatrix} x_i^{(0)}+\Delta x_i^{(0)} \\ r_{1i} \end{bmatrix} = \begin{bmatrix} \bar{x}_i^{(1)} \\ \bar{u}_{1i}+\Delta\bar{u}_{1i} \end{bmatrix} \tag{25}$$

then for the mixed error one obtains (see the Appendix).

$$\left\| \begin{bmatrix} \Delta x_i^{(0)} \\ \Delta\bar{u}_{1i} \end{bmatrix} \right\|_2 \le (8(1+|s_1|)\epsilon/|c_1|) \left\| \begin{bmatrix} x_i^{(0)} \\ \bar{u}_{1i} \end{bmatrix} \right\|^2. \tag{26}$$

Since this 2-vector is now a subcolumn of the $i$-th column of the matrix $\begin{bmatrix} x^T \\ U \end{bmatrix}$.

we have

$$J_1^T \begin{bmatrix} x^T+\Delta^{(1)}x^T \\ U+\Delta^{(1)}U \end{bmatrix} = \begin{bmatrix} \bar{x}^{(1)T} \\ R^{(1)} \end{bmatrix} \tag{27}$$

with $J_1$ exactly orthogonal and

$$= \left\| \begin{bmatrix} \Delta^{(1)}x^T \\ \Delta^{(1)}U \end{bmatrix} \right\|_2 \le (8(1+|s_1|)\varkappa/c_1) \| [R]_i \|_2,$$

since the $i$-the column of $R$ and $[x.U^T]^T$ have the same norms.

Errors in the subsequent steps are similarly bounded in terms of the deflated problem $[\bar{x}^{(1)}, R^{(1)T}]^T$. Because corresponding columns in (27) are related by an orthogonal transformation, the errors of step 2 can be mapped backwards onto $[x.U^T]^T$ without altering their norms. (This would not have been possible if we had used backward errors as in (23), (24)). This now gives

$$J_2^T J_1^T \begin{bmatrix} x^T+\Delta^{(1)}x^T+\Delta^{(2)}x^T \\ U+\Delta^{(1)}U+\Delta^{(2)}U \end{bmatrix} = \begin{bmatrix} \bar{x}^{(2)T} \\ R^{(2)} \end{bmatrix} \tag{29}$$

where

$$\left\| \begin{bmatrix} \Delta^{(1)}x^T+\Delta^{(2)}x^T \\ \Delta^{(1)}U+\Delta^{(2)}U \end{bmatrix} \right\|_2$$
$$\le 8((1+|s_1|)/|c_1|+(1+|s_2|)/|c_2|)\epsilon \|[R]_i\|_2 \tag{30}$$

for all columns $i$ except the first one (since that one is not affected anymore). By induction we finally have

$$J_n^T \cdots J_1^T \begin{bmatrix} x^T+\Delta x^T \\ U+\Delta U \end{bmatrix} = \begin{bmatrix} 0^T \\ R \end{bmatrix} \tag{31}$$

where $\Delta U = \sum_{i=1}^n \Delta^{(i)}U$ and $\Delta x = \sum_{i=1}^n \Delta^{(i)}x$ are bounded columnwise by

$$\left\| \begin{bmatrix} \Delta x^T \\ \Delta U \end{bmatrix} \right\|_2 \le 8i \max_{1\le i} \{(1+|s_i|)/|c_i|\} \epsilon \|[R]_i\|_2 \tag{32}$$

which is significantly worse than (18-19) and problem dependent. For the total error we obtain

$$\left\| \begin{bmatrix} \Delta x^T \\ \Delta U \end{bmatrix} \right\|_F \le 5n\sqrt{n} \max_{1\le j\le n} \{(1+|s_j|)/|c_j|\} \epsilon \|R\|_2 \tag{33}$$

### 3.3. Algorithm C

Since the method is essentially a rearrangement of the previous method, one might expect similar bounds for the numerical errors, but this is not the case because operations are performed in a different order.

Let us again denote the elements of the rows involved in step 1 as

$$\begin{bmatrix} c_1 & -s_1 \\ s_1 & c_1 \end{bmatrix} \begin{bmatrix} x_1^{(0)} x_2^{(0)} \cdots x_n^{(0)} \\ u_{11} u_{12} \cdots u_{1n} \end{bmatrix} = \begin{bmatrix} 0 \ x_2^{(1)} \cdots x_n^{(1)} \\ r_{11} \ r_{12} \cdots r_{1n} \end{bmatrix} \tag{34}$$

then $c_1, s_1$ and $u_{1j}, x_j^{(i)}$ are constructed as follows

...eded for the construction of the $J_i$. Both steps could thus be merged on an array of the type of Figure 1 and the number of time steps would then be ...lved.

On vector machines all these operations are also easy to exploit. The triangular system (3) can be solved via typical $A \cdot X(\cdot) + Y(\cdot)$ vector operations and the application of the rotations $J_i$ ammounts to so-called "pipelined ...ivens" (see e.g [5]). Both of these operations are known to perform well on ...ctor machines [5].

*4.2. The recursive Algorithms B and C*

Here in each step one essentially constructs an elementary rotation $J_k$ (or equivalently $S_k$) and applies it to determine $x^{(k)}$ and the $k$-th row of $U$ from ...(4.1) and the $k$-th row of $R$. This is easily seen to require $4(n - i + 1)$ multiplications for each transformation $J_i$ or $S_i$, resulting in an operation count of

$$2n^2 + O(n) \quad (40)$$

multiplications for the whole process.

Parallel implementation of these methods would require $n$ processors and $2n$ ...teps as shown above for the second stage of the LINPACK method. Since ...sentially, the first stage is skipped here, these algorithms compare favourably ...th the LINPACK method for both systolic arrays and vector machines. ...otice also that by using modified elementary rotations [6] or modified $\Sigma$-...rthogonal transformations, the work could be reduced and square roots could ...e avoided, just as in the second stage of the previous method.

*4.3. Consecutive up- and downdates*

...s explained in the introduction, it is often required to perform a series of ...onsecutive up- and/or downdates of Cholesky factorizations. In this case a ...ifferent scheme than the one suggested in Figure 1 is recommended. Indeed, ...he $\Sigma$-orthogonal method of Algorithm C is completely analogous to the ...rthogonal method used in recursive updating of Cholesky factors (see [7]). ...he array of processors shown in Figure 2 used for updating in GENTLEMAN ...nd KUNG [7] can thus be used also for downdating.

...Here the factor $R$ resides now in the $n(n+1)/2$ processors shown in Figure ...and $x^T$ is "passed through" it leaving the updated factor behind him in the ...riangular array (see [7] for details). When connecting a "flag" (say $\pm 1$) to ...ach row $x^T$ in order to indicate if an updating or downdating is requested ...hen the programs of the processors are easily adapted to take this into ...ccount and to "switch" from one type of operations to another. In this ...ashion consecutive up- and/or downdates can easily be pipelined through the ...riangular array depicted in Figure 2. This in fact constitutes a more efficient ...se of the processors than in Figure 1, since there the left most processors gra- ...ually become inactive as the time steps proceed. In the second scheme, shown ...n Figure 2, all processors are active at any time step, for as long as new rows ... are fed into the array.
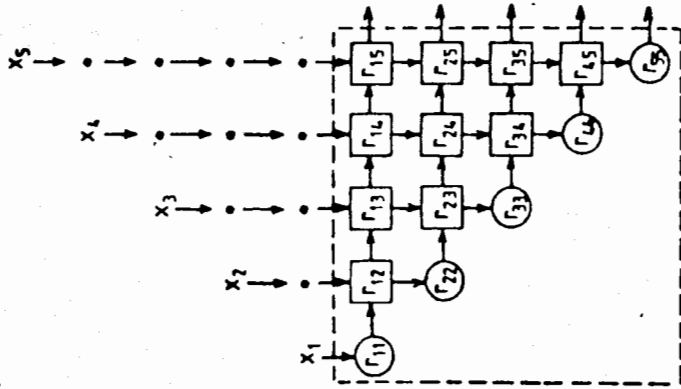


FIGURE 2. A triangular systolic array for up- and downdating a Cholesky factor

*5. NUMERICAL RESULTS*

The results of STEWART [13] on the conditioning of the downdating problem show that we cannot expect that $\tilde{U}$, the calculated factor, will be close to $U$. However a mixed stability result of the form

$$R^T R - (x + \Delta x)(x + \Delta x)^T = (\tilde{U} + \Delta \tilde{U})(\tilde{U} + \Delta \tilde{U})^T$$

where $\Delta x$ and $\Delta \tilde{U}$ are small does ensure that

$$R^T R - xx^T = \tilde{U}\tilde{U}^T + E$$

where $E$ is small. To demonstrate the superior stability properties of Algorithms A and C over Algorithm B, we have computed the quantity

$$e = \frac{\|R^T R - xx^T - \tilde{U}\tilde{U}^T\|_F}{\|U^T U\|_F}$$

for the problem

$$R = \begin{bmatrix} 1 & \sin(\theta/2) \\ 0 & \sqrt{2}\cos(\theta/2) \end{bmatrix}, \quad x = \begin{bmatrix} \sin\theta \\ \cos(\theta/2) \end{bmatrix},$$

which has the solution

$$U = \begin{bmatrix} \cos\theta & -\sin(\theta/2) \\ 0 & \cos(\theta/2) \end{bmatrix}.$$

results obtained on a MacIntosh which has working accuracy of between 7 and 8 significant digits are tabulated in Table 1 for $\cos\theta = 2^{-k}$, $k = 3,6,9,12$.

| $\cos\theta$ | Alg A | Alg B | Alg C |
|---|---|---|---|
| $2^{-3}$ | $1.337E-8$ | $2.367E-7$ | $1.183E-7$ |
| $2^{-6}$ | $1.210E-7$ | $1.073E-6$ | $6.939E-8$ |
| $2^{-9}$ | $1.788E-7$ | $3.218E-5$ | $2.946E-8$ |
| $2^{-12}$ | $1.264E-7$ | $1.011E-4$ | $2.467E-8$ |

TABLE 1.

The numerical results are consistent with the stability analysis of Section 3 and clearly demonstrate the superior properties of Algorithms A and C. Moreover, these results also show that Algorithm B is unstable (in de Jong's sense [4]) when $\cos\theta$ becomes very small.

APPENDIX

Here we derive the bounds (24), (26) and (36) for the numerical errors incurred in (22) and (35).

Let us denote by $\delta_i$ and $\epsilon_i$ quantities that are smaller in absolute value than the relative precision of the machine used. Then we have, according to WILKINSON [14].

$$\tilde{u}_{11} = \sqrt{\{[(r_{11} - x_1^{(0)})(1+\delta_1)(r_{11} + x_1^{(0)})(1+\delta_2)](1+\delta_3)\}(1+\delta_4)} =$$    (41)

$$= u_{11}(1 + 2.5\epsilon_1)$$

$$\tilde{c}_1 = (\tilde{u}_{11}/r_{11})(1+\delta_5) = (u_{11}/r_{11})(1+3.5\epsilon_2) = c_1(1+3.5\epsilon_2)$$    (42)

$$\tilde{s}_1 = (x_1^{(0)}/r_{11})(1+\delta_6) = (x_1^{(0)}/r_{11})(1+\epsilon_3) = s_1(1+\epsilon_3)$$    (43)

and for each $i$

$$\tilde{x}_i^{(1)} = [x_i^{(0)} - \tilde{s}_1 r_{1i}(1+\delta_{10})](1+\delta_{11})/[\tilde{c}_1(1+\delta_{12})]$$    (44)

$$= [x_i^{(0)} - \tilde{s}_1 r_{1i}(1+\epsilon_6)]/[\tilde{c}_1(1+2\epsilon_7)]$$

$$\tilde{u}_{1i} = [r_{1i} - \tilde{s}_1 x_i^{(0)}(1+\delta_7)](1+\delta_8)/[\tilde{c}_1(1+\delta_9)]$$    (45)

$$= [r_{1i} - \tilde{s}_1 x_i^{(0)}(1+\epsilon_4)]/[\tilde{c}_1(1+2\epsilon_5)].$$

Multiplying these with their respective denominators yields

$$x_i^{(0)} = \tilde{s}_1 r_{1i}(1+\epsilon_6) + \tilde{x}_i^{(1)}\tilde{c}_1(1+2\epsilon_7)$$    (46)

$$r_{1i} = \tilde{s}_1 x_i^{(0)}(1+\epsilon_4) + \tilde{u}_{1i}\tilde{c}_1(1+2\epsilon_5).$$

Using (42) and (43) we finally obtain

$$x_i^{(0)} - s_1 r_{1i} = c_1 \tilde{x}_i^{(1)} + 2\epsilon_8 s_1 r_{1i} + 5.5\epsilon_9 \tilde{x}_i^{(1)} c_1$$

$$r_{1i} - s_1 x_i^{(0)} = c_1 \tilde{u}_{1i} + 2\epsilon_{10} s_1 x_i^{(0)} + 5.5\epsilon_{11} c_1 \tilde{u}_{1i}$$    (47)

or

$$\frac{1}{c_1}\begin{bmatrix} 1 & -s_1 \\ -s_1 & 1 \end{bmatrix}\begin{bmatrix} x_i^{(0)} \\ r_{1i} \end{bmatrix} = \begin{bmatrix} \tilde{x}_i + \Delta\tilde{x}_i^{(1)} \\ \tilde{u}_{1i} + \Delta\tilde{u}_{1i} \end{bmatrix}$$

$$= \begin{bmatrix} \tilde{x}_i^{(1)} + 2\epsilon_8(s_1/c_1)r_{1i} + 5.5\epsilon_9\tilde{x}_i^{(1)} \\ \tilde{u}_{1i} + 2\epsilon_{10}(s_1/c_1)x_i^{(0)} + 5.5\epsilon_{11}\tilde{u}_{1i} \end{bmatrix}$$    (48)

Identifying the error terms in the right hand side with the vector $[\Delta\tilde{x}_i^{(1)} , \Delta\tilde{u}_{1i}]^T$ we then obtain from the Cauchy-Schwartz inequality the bounds

$$|\Delta\tilde{x}_i^{(1)}| \leq (1/|c_1|)\left\| \begin{bmatrix} 5.5\epsilon_9\cdot\tilde{x}_i^{(1)} \\ 2\epsilon_8 r_{1i} \end{bmatrix} \right\|_2 \leq (5.5\epsilon/|c_1|)\left\| \begin{bmatrix} \tilde{x}_i^{(1)} \\ r_{1i} \end{bmatrix} \right\|_2$$

$$= (5.5\epsilon/|c_1|)\left\| \begin{bmatrix} x_i^{(0)} \\ \tilde{u}_{1i} \end{bmatrix} \right\|_2 + O(\epsilon^2)$$    (49)

$$|\Delta\tilde{u}_{1i}| \leq (1/|c_1|)\left\| \begin{bmatrix} 2\epsilon_{10} x_i^{(0)} \\ 5.5\epsilon_{11}\tilde{u}_{1i} \end{bmatrix} \right\|_2 \leq (5.5\epsilon/|c_1|)\left\| \begin{bmatrix} x_i^{(0)} \\ \tilde{u}_{1i} \end{bmatrix} \right\|_2$$

and subsequently

$$= \left\| \begin{bmatrix} \Delta\tilde{x}_i^{(1)} \\ \Delta\tilde{u}_{1i} \end{bmatrix} \right\|_2 = (5.5\sqrt{2}\epsilon/|c_1|)\left\| \begin{bmatrix} x_i^{(0)} \\ \tilde{u}_{1i} \end{bmatrix} \right\|_2$$

$$\leq (8\epsilon/|c_1|)\left\| \begin{bmatrix} x_i^{(0)} \\ \tilde{u}_{1i} \end{bmatrix} \right\|_2.$$    (50)

For the backward and mixed errors one obtains via similar techniques the following bounds

$$= \left\| \begin{bmatrix} \Delta\tilde{x}_i^{(1)} \\ \Delta r_{1i} \end{bmatrix} \right\|_2 \leq (8\epsilon/|c_1|)\left\| \begin{bmatrix} x_i^{(0)} \\ \tilde{u}_{1i} \end{bmatrix} \right\|_2$$    (51)

$$= \left\| \begin{bmatrix} \Delta\tilde{x}_i^{(1)} \\ \Delta\tilde{u}_{1i} \end{bmatrix} \right\|_2 \leq (8\epsilon(1+|s_1|)/|c_1|)\left\| \begin{bmatrix} x_i^{(0)} \\ \tilde{u}_{1i} \end{bmatrix} \right\|_2.$$    (52)

Notice that all three bounds depend on $1/|c_1|$ which is large when the condition number $\kappa$ of $S_1$,

$$\kappa(S_1) = (1+|s_1|)/(1-|s_1|) = (1+|s_1|)^2/|c_1|^2$$    (53)

is large.

Now we consider Algorithm C. A similar analysis to t? used for

Statistical Computation. Eds. R.C. Milton. J.A. Nelder. pp. 365-397. Academic Press. New York.

10. G.H. GOLUB, C. VAN LOAN. (1983). Matrix Computations. The John Hopkins Press. Baltimore.

11. H.T. KUNG, C.E. LEISERSON. (1980). Algorithms for VLSI processor arrays, in Introduction to VLSI systems (C. Mead and L. Conway. eds.), Addison-Wesley, pp. 271-292.

12. C.L. LAWSON, R.J. HANSON, (1974). Solving least squares problems. Prentice Hall, Englewood Cliffs, N.J.

13. G.W. STEWART. (1979). The effect of rounding errors on an algorithm for downdating a Cholesky factorization. JIMA. vol. 23. pp. 203-213.

14. J.H. WILKINSON. (1965). The Algebraic Eigenvalue Problem. Clarendon Press. Oxford.

Algorithm B yields

$$\tilde{x}_i^{(1)} = [\bar{c}_1 x_i^{(0)}(1+\delta_{10}) - \bar{s}_1 \bar{u}_{1_i}(1+\delta_{11})](1+\delta_{12})$$
$$= \bar{c}_1 \tilde{x}_i^{(0)}(1+2\epsilon_{12}) - \bar{s}_1 \bar{u}_{1_i}(1+2\epsilon_{13}). \quad (54)$$

Together with the second equation of (46) and using (42). (43). this now leads

$$\begin{bmatrix} \tilde{x}_i^{(1)} \\ r_{1_i} \end{bmatrix} = \begin{bmatrix} c_1 & -s_1 \\ s_1 & c_1 \end{bmatrix} \begin{bmatrix} x_i^{(0)} \\ \bar{u}_{1_i} \end{bmatrix} + \begin{bmatrix} -3\epsilon_{14}s_1\bar{u}_{1_i} + 5.5\epsilon_{15}c_1 x_i^{(0)} \\ 5.5\epsilon_{11}c_1 \bar{u}_{1_i} + 2\epsilon_{10}s_1 x_i^{(0)} \end{bmatrix}. \quad (55)$$

Identifying the error terms on the right hand side with the vector $[\Delta \tilde{x}_i, -\Delta r_{1_i}]^T$ and using a similar argument as in the proof of (49) and (50), one now finds

$$= \left\| \begin{bmatrix} \Delta \tilde{x}_i^{(1)} \\ \Delta r_{1_i} \end{bmatrix} \right\|_2 \le 6.5\epsilon \left\| \begin{bmatrix} x_i^{(0)} \\ \bar{u}_{1_i} \end{bmatrix} \right\|_2. \quad (56)$$

From the relation

$$\begin{bmatrix} \Delta \tilde{x}_i^{(0)} \\ \Delta \bar{u}_{1_i} \end{bmatrix} = \begin{bmatrix} c_1 & s_1 \\ -s_1 & c_1 \end{bmatrix} \begin{bmatrix} \Delta \tilde{x}_i^{(1)} \\ \Delta r_{1_i} \end{bmatrix} \quad (57)$$

one then obtains a similar bound for the vector $[\Delta x_i^{(0)}, \Delta \bar{u}_{1_i}]^T$. which completes the proof.

REFERENCES

J.H. ARGYRIS, J.R. ROY. (1972). General treatments of structural modifications. J. Structural Division. A.S.C.E. 98. ST2. Proc. Paper 8732. pp. 465-492.

V. BELEVITCH. (1968). Classical network theory. Holden Day. San Francisco.

A.W. BOJANCZYK, R.P. BRENT, F.R. DE HOOG. (May 1985). QR factorization of Toeplitz matrices. Report CMA-R07-1985. Centre for Math. Anal. The Australian National University.

L. DE JONG. (1977). Towards a formal definition of numerical stability. Numerische Mathematik. Vol. 28, pp. 211-220.

J.J. DONGARRA, D.C. SORENSEN. (1986). Linear algebra on high performance computers, Appl. Math. & Comp.. Vol. 20. pp. 57-88.

W.M. GENTLEMAN. (1973). Least squares computations by Givens transformations without square roots. JIMA. vol. 12. pp. 329-336.

W.M. GENTLEMAN, H.T. KUNG. (1981). Matrix triangularization by systolic arrays, Proceedings SPIE Symp. (1981), Vol. 298. Real Time Signal Processing IV, pp. 19-26.

P.E. GILL, G.H. GOLUB, W. MURRAY, M.A. SAUNDERS, (1974). Methods for modifying matrix factorizations. Mathematics of Computation, vol. 28, pp. 505-535.

G.H. GOLUB, (1969). Matrix decompositions and statistical calculations, in: