

## ROUNDING ERRORS IN SOLVING BLOCK HESSENBERG SYSTEMS

URS VON MATT AND G. W. STEWART

ABSTRACT. A rounding error analysis is presented for a divide-and-conquer algorithm to solve linear systems with block Hessenberg matrices. Conditions are derived under which the algorithm computes a stable solution. The algorithm is shown to be stable for block diagonally dominant matrices and for M-matrices.

### 1. INTRODUCTION

In [9] a recursive algorithm was proposed for the solution of the linear system

$$(1) \quad AX = B,$$

where  $A$  is a block Hessenberg matrix. Its development was motivated by the attempt to find the steady-state of certain Markov chains. In this paper we will present an error analysis to explain the accurate results obtained by the algorithm.

Our analysis is a rounding error analysis in the style of Wilkinson [13, 14]. We will see that the computed matrix  $X$  can be regarded as the exact solution of a nearby linear system. In particular we will show that the computed  $X$  satisfies

$$AX = B + \Delta B.$$

We call the matrix  $X$  a stable solution if

$$\|\Delta B\| \leq \eta \|A\| \|X\|,$$

where  $\eta$  denotes a small multiple of the unit roundoff  $\varepsilon$ . This is an example of residual stability. Note that residual stability is the same as backward stability if the right hand side  $B$  is a vector (cf. [5]).

A stable solution is not to be confused with an accurate solution. The accuracy of  $X$  is usually limited by the condition number  $\kappa(A) := \|A\| \|A^{-1}\|$ . The relative error of  $X$  can be bounded by

$$\frac{\|X - A^{-1}B\|}{\|A^{-1}B\|} \leq \frac{\eta\kappa(A)}{1 - \eta\kappa(A)},$$

provided that  $\eta\kappa(A) < 1$ . Thus, we can only compute an accurate solution  $X$  if we use a stable algorithm to solve a well-conditioned problem.

---

Received by the editor August 22, 1994 and, in revised form, January 10, 1995.

1991 *Mathematics Subject Classification*. Primary 65G05; Secondary 65F05.

*Key words and phrases*. Rounding error analysis, linear systems, block Hessenberg matrices, block diagonally dominant matrices, M-matrices.

This work was supported in part by the National Science Foundation under grant CCR 9115568.

Our paper is organized as follows. In §2 we give a concise description of the algorithm to be analyzed. This algorithm consists of a few basic building blocks for which we will cite error bounds in §3. Since our algorithm calls itself recursively we have to make an assumption about the structure of the errors after each invocation. This is the purpose of §4, where we also analyze the local errors in each stage. We combine these local errors to give a global error bound in §5. The structure of this global error bound reveals a potential instability of our algorithm. This is discussed in §6. In §§7 and 8 we identify two classes of matrices for which our algorithm computes a stable solution. We conclude our presentation with some numerical examples in §9.

Throughout our analysis we will use the 2-norm, except where otherwise noted. Its main advantage is that the norm of an orthogonal matrix is one.

## 2. ALGORITHM

We assume that the matrix  $A$  in (1) has the following block Hessenberg structure:

$$A = \begin{bmatrix} A_{11} & A_{12} & \cdots & \cdots & A_{1n} \\ A_{21} & \ddots & & & \vdots \\ & \ddots & \ddots & & \vdots \\ & & & \ddots & A_{n-1,n} \\ & & & A_{n,n-1} & A_{nn} \end{bmatrix}.$$

The diagonal blocks  $A_{ii}$  are assumed to be square nonsingular matrices of order  $p_i$ . The total size of  $A$  is given by

$$N := \sum_{i=1}^n p_i.$$

If  $n > 1$  we can select a tear index  $k$  with  $1 \leq k < n$  and partition the matrix  $A$  as follows:

$$A = \begin{bmatrix} A_{nw} & A_{ne} \\ A_{sw} & A_{se} \end{bmatrix}.$$

The submatrix  $A_{nw}$  contains the first  $k$  diagonal blocks of  $A$ , and  $A_{se}$  contains the last  $n - k$  diagonal blocks of  $A$ . Note that  $A_{k+1,k}$  is the only nonzero block in  $A_{sw}$ . This partitioning is also shown as Figure 1 (next page). The dimensions  $n_{nw}$  and  $n_{se}$  are given by

$$n_{nw} = \sum_{i=1}^k p_i,$$

$$n_{se} = \sum_{i=k+1}^n p_i.$$

Let  $E$  be the last  $n_{se}$  columns of the  $N$ -by- $N$  identity matrix, and let  $F$  consist of the first  $n_{nw}$  columns of the  $N$ -by- $N$  identity matrix. Then we can also define

$$\hat{A} := \begin{bmatrix} A_{nw} & A_{ne} \\ 0 & A_{se} \end{bmatrix} = A - EA_{sw}F^T.$$

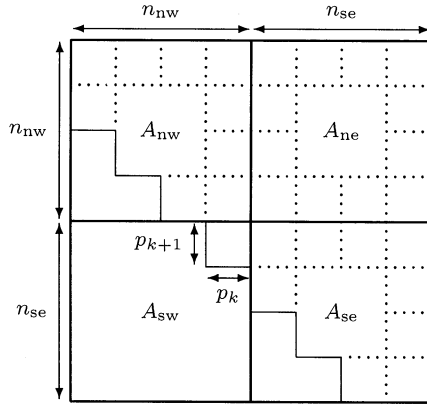


FIGURE 1. Structure of  $A$

In order to solve the linear system (1) for  $X$  we first compute the solution  $Y$  of

$$\hat{A}Y = B.$$

This step requires the solution of linear systems with the matrices  $A_{se}$  and  $A_{nw}$ , which can be solved recursively by the same divide-and-conquer algorithm. Afterwards we obtain  $X$  from  $Y$  by means of an updating formula. The well-known Sherman-Morrison-Woodbury formula (cf. [4, p. 51]) would give us

$$\begin{aligned} A^{-1} &= (\hat{A} + EA_{sw}F^T)^{-1} = \hat{A}^{-1} - \hat{A}^{-1}E(I + A_{sw}F^T\hat{A}^{-1}E)^{-1}A_{sw}F^T\hat{A}^{-1} \\ &= \hat{A}^{-1} - \hat{A}^{-1}EA_{sw}(I + F^T\hat{A}^{-1}EA_{sw})^{-1}F^T\hat{A}^{-1}. \end{aligned}$$

Unfortunately, this formula does not take advantage of the many zeros in  $A_{sw}$ , requiring the solution of a large intermediate linear system. We can reduce the size of this linear system with the help of the URV-decomposition (cf. [10])

$$A_{sw} = URV^T.$$

Let  $r$  denote the rank of  $A_{sw}$  as it is determined by the URV-decomposition. Then  $U$  will be an orthogonal  $n_{se}$ -by- $r$  matrix with  $p_{k+1}$  nonzero rows. Also  $R$  is a square  $r$ -by- $r$  matrix, and  $V$  is an orthogonal  $n_{nw}$ -by- $r$  matrix with  $p_k$  nonzero rows. Now we can express the inverse of  $A$  by

$$A^{-1} = \hat{A}^{-1} - \hat{A}^{-1}EU(I + RV^T F^T \hat{A}^{-1} EU)^{-1}RV^T F^T \hat{A}^{-1}.$$

In order to avoid the multiple evaluation of the same expressions, we introduce the following intermediate quantities:

$$\begin{aligned} G &:= \hat{A}^{-1}EU, \\ S &:= RV^T F^T G, \\ T &:= I + S, \\ \hat{R} &:= T^{-1}R, \\ P &:= G\hat{R}. \end{aligned}$$

## ALGORITHM 1. Solution of block Hessenberg systems

```

function  $X = solve(A, B)$ 
if not at the bottom then
  Compute the orthogonal URV-decomposition  $A_{sw} = URV^T$ .
   $G_s := solve(A_{se}, U)$ 
   $G_n := solve(A_{nw}, -A_{ne}G_s)$ 
   $S := RV^T G_n$ 
   $T := I + S$ 
  Solve  $T\hat{R} = R$  for  $\hat{R}$  by Gaussian elimination.
   $P := G\hat{R}$ 
   $Y_s := solve(A_{se}, B_s)$ 
   $Y_n := solve(A_{nw}, B_n - A_{ne}Y_s)$ 
   $X := Y - PV^TY_n$ 
else
  Solve  $AX = B$  for  $X$  by Gaussian elimination.
end

```

Note that these matrices are independent of the right-hand side  $B$ . The overall recursive procedure to solve the linear system (1) is also presented as Algorithm 1.

In [9] this algorithm is refined further by introducing the auxiliary procedures “patchgen” and “topsolve”. These refinements are critical for the efficiency of the algorithm, but they are not necessary for the purpose of this error analysis. Further implementation details may be found in [11].

The solution of the linear system (1) can also be described by the tear tree of Figure 2. Each node represents a linear system to be solved. The node on the top level ( $k = 3$ ) stands for the system (1), whereas the leaf nodes are the linear systems that are not divided any further but solved by Gaussian elimination. The number  $n$  of diagonal blocks in the matrix  $A$ , which is equal to the number of leaf nodes, and the height  $h$  of the tear tree are connected by the inequalities

$$n \leq 2^h,$$

$$h \geq \log_2 n.$$

These inequalities become equalities if the tear tree of Figure 2 is a complete binary tree.

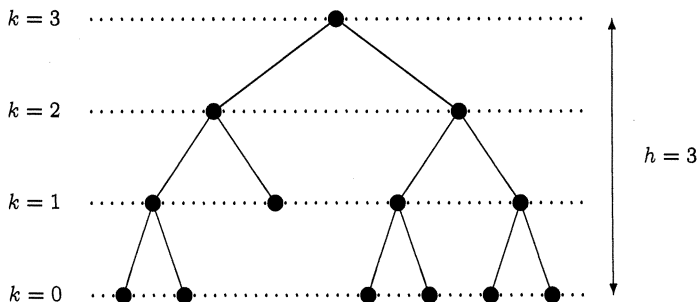


FIGURE 2. Tear tree

## 3. BASIC OPERATIONS

Algorithm 1 is composed of a few basic building blocks. These are the addition and multiplication of matrices, the calculation of a URV-decomposition, and the solution of linear systems by Gaussian elimination. We will now state bounds for the errors associated with these operations.

In the case of the addition of two matrices, we have

$$\text{fl}(A + B) = A + B + E,$$

where

$$\|E\| \leq \eta_1 \|A + B\|.$$

The quantity  $\eta_1$  is on the order of the unit roundoff  $\varepsilon$  and slowly increases with the size of the matrices  $A$  and  $B$ . See also [14, pp. 114–115 and p. 57].

If we multiply two matrices in floating point, we have

$$\text{fl}(AB) = AB + E,$$

where

$$\|E\| \leq \eta_2 \|A\| \|B\|.$$

Again,  $\eta_2$  is a small multiple of the unit roundoff and slowly grows with the size of the matrices  $A$  and  $B$  (cf. [14, pp. 115–116]).

The size of the error in computing an orthogonal URV-decomposition depends on the specifics of the decomposition. One may choose a QR-decomposition [4, Chapter 5], a rank-revealing decomposition [2, 10], or the singular value decomposition [4, §8.3]. All of these factorizations have in common that they can be expressed as a sequence of orthogonal transformations applied from the left and the right to the initial matrix. The orthogonal transformations are also accumulated to give the matrices  $U$  and  $V$ . Wilkinson shows in [14, pp. 160–161] that there are orthogonal matrices  $U_0$  and  $V_0$  and an  $\eta_3$  such that

$$(2) \quad \|R - U_0^T A V_0\| \leq 2\eta_3 \|A\|,$$

$$(3) \quad \|U - U_0\| \leq \eta_3,$$

$$(4) \quad \|V - V_0\| \leq \eta_3.$$

The quantity  $\eta_3$  is on the order of the unit roundoff and slowly grows with the size of the matrix  $A$  and the number of the orthogonal transformations applied to  $A$ . From (3) and (4) it immediately follows that

$$\|U\| \leq 1 + \eta_3,$$

$$\|V\| \leq 1 + \eta_3.$$

We can also show, by increasing  $\eta_3$  slightly as necessary, that

$$\|URV^T - A\| \leq 4\eta_3 \|A\|,$$

$$\|(U^T U)^{-1} U^T\| \leq 1 + 3\eta_3,$$

$$\|V(V^T V)^{-1}\| \leq 1 + 3\eta_3.$$

We assume that small linear systems are solved by Gaussian elimination. In [13, p. 108] and [14, p. 252] Wilkinson shows that this process can be described by the equation

$$(5) \quad A\mathbf{x} = \mathbf{b} + \Delta\mathbf{b},$$

where

$$(6) \quad \|\Delta \mathbf{b}\| \leq \eta'_4 \|A\| \|\mathbf{x}\|.$$

The value of  $\eta'_4$  is on the order of the unit roundoff and slowly increases with the size of  $A$ . It also depends on the pivoting strategy used. See [6] for a more recent survey.

Note that the bound (6) is only applicable if the right-hand side of (5) is a vector. If we solve several linear systems with the same matrix  $A$  we get

$$(7) \quad AX = B + \Delta B,$$

where

$$(8) \quad \|\Delta B\| \leq \eta'_4 \sqrt{r} \|A\| \|X\|,$$

and  $r$  denotes the number of columns in the matrix  $B$ .

Let  $r_{\max}$  be the maximum number of right-hand sides in a linear system that is solved by Gaussian elimination in Algorithm 1. If we define

$$\eta_4 := \eta'_4 \sqrt{r_{\max}},$$

then we can always bound the residual  $\Delta B$  in (7) by

$$\|\Delta B\| \leq \eta_4 \|A\| \|X\|.$$

Thanks to this convention, our error bounds will become somewhat simpler.

#### 4. ANALYSIS OF ONE STAGE

In the following we will give expressions for the rounding errors incurred at one stage of Algorithm 1. We assume that we are not at the bottom of the tear tree, and we use the assumptions of §3 to bound the size of the rounding errors.

In what follows the matrix  $A$  denotes the system matrix of an arbitrary interior node of the tear tree. In an attempt to keep the notation simple, we do not introduce an index to indicate the corresponding node. We also assume that the four submatrices  $A_{\text{nw}}$ ,  $A_{\text{ne}}$ ,  $A_{\text{sw}}$ , and  $A_{\text{se}}$  are predefined by the tearing strategy.

We make the inductive assumption that the solution  $X$  computed at level  $k$  satisfies

$$(9) \quad AX = B + \Delta B,$$

where the residual  $\Delta B$  can be expressed by

$$(10) \quad \Delta B = \Delta L_X + \Delta M_X X.$$

We use the index  $X$  for the matrices  $\Delta L_X$  and  $\Delta M_X$  to indicate that they depend on the solution  $X$ .

We assume that at level  $k$  we always have

$$(11) \quad \|\Delta L_X\| \leq \xi_k \|D^{-1} X\|,$$

$$(12) \quad \|\Delta M_X D\| \leq \zeta_k,$$

for all matrices  $X$ . The quantity  $D$  denotes a nonsingular block diagonal matrix, which is partitioned commensurably with  $A$ . In particular, we always have  $D_{\text{ne}} = 0$  and  $D_{\text{sw}} = 0$  for all the nodes in the tear tree. The matrix  $D$  will give us additional flexibility in bounding the norm of the residual  $\Delta B$ . We will discuss this issue in more detail in §§7–9.

Since we solve the systems at the bottom level by Gaussian elimination, we define

$$\begin{aligned}\xi_0 &:= \eta_4 \|A\| \|D\|, \\ \zeta_0 &:= 0,\end{aligned}$$

The purpose of the next sections will be to compute  $\xi_k$  and  $\zeta_k$  if  $\xi_{k-1}$  and  $\zeta_{k-1}$  are known.

Throughout our analysis we will assume that the rounding errors remain small compared to the norm of the computed quantities. This means that the computed and the exact quantities will agree to at least a few digits. We will use a factor of 1.01 in (20,23,25,26,29,30,36,40,49,50,52,61,62) to simplify our bounds.

**4.1. Calculation of the URV-decomposition.** The result of the initial URV-decomposition of Algorithm 1 can be described by

$$(13) \quad A_{\text{sw}} = URV^T + \Delta A_{\text{sw}},$$

where

$$(14) \quad \|\Delta A_{\text{sw}}\| \leq 4\eta_3 \|A_{\text{sw}}\|.$$

The matrices  $U$  and  $V$  are nearly orthogonal, and they satisfy

$$(15) \quad \|U\| \leq 1 + \eta_3, \quad \|(U^T U)^{-1} U^T\| \leq 1 + 3\eta_3,$$

$$(16) \quad \|V\| \leq 1 + \eta_3, \quad \|V(V^T V)^{-1}\| \leq 1 + 3\eta_3.$$

The expression  $RV^T$ , which we will also use later on, can be written as

$$(17) \quad RV^T = (U^T U)^{-1} U^T (A_{\text{sw}} - \Delta A_{\text{sw}}).$$

Therefore, we have the bound

$$(18) \quad \|RV^T\| \leq (1 + 3\eta_3)(1 + 4\eta_3) \|A_{\text{sw}}\|.$$

Because of

$$R = (U^T U)^{-1} U^T (A_{\text{sw}} - \Delta A_{\text{sw}}) V (V^T V)^{-1},$$

we can also bound the norm of  $R$  by

$$(19) \quad \|R\| \leq (1 + 3\eta_3)^2 (1 + 4\eta_3) \|A_{\text{sw}}\|.$$

**4.2. Calculation of  $G$ .** The calculation of the matrix  $G$  proceeds in three steps that can be described by the equations

$$\begin{aligned}A_{\text{se}} G_{\text{s}} &= U + \Delta G_{\text{s}}, \\ U_{\text{n}} &= -A_{\text{ne}} G_{\text{s}} + \Delta U_{\text{n}}, \\ A_{\text{nw}} G_{\text{n}} &= U_{\text{n}} + \Delta G_{\text{n}}.\end{aligned}$$

The error matrix  $\Delta U_{\text{n}}$  is bounded by

$$(20) \quad \|\Delta U_{\text{n}}\| \leq \eta_2 \|A_{\text{ne}}\| \|G_{\text{s}}\| \leq 1.01\eta_2 \|A_{\text{ne}}\| \|A_{\text{se}}^{-1}\|.$$

The residuals  $\Delta G_{\text{s}}$  and  $\Delta G_{\text{n}}$  have the expansion

$$\begin{aligned}\Delta G_{\text{s}} &= \Delta L_{G_{\text{s}}} + \Delta M_{G_{\text{s}}} G_{\text{s}}, \\ \Delta G_{\text{n}} &= \Delta L_{G_{\text{n}}} + \Delta M_{G_{\text{n}}} G_{\text{n}}.\end{aligned}$$

We can also write these equations in matrix terms as

$$(21) \quad \hat{A}G = EU + \Delta G,$$

where

$$(22) \quad \Delta G := \begin{bmatrix} \Delta G_n + \Delta U_n \\ \Delta G_s \end{bmatrix} = \begin{bmatrix} \Delta L_{G_n} + \Delta U_n \\ \Delta L_{G_s} \end{bmatrix} + \begin{bmatrix} \Delta M_{G_n} & \\ & \Delta M_{G_s} \end{bmatrix} G.$$

Since the matrix  $G$  can be written as

$$G = \hat{A}^{-1}(EU + \Delta G),$$

we also have

$$(23) \quad \|G\| \leq 1.01 \|\hat{A}^{-1}\|.$$

**4.3. Calculation of  $S$ .** We can express the matrix  $S$  by

$$(24) \quad S = RV^T G_n + \Delta S,$$

where

$$\|\Delta S\| \leq (2\eta_2 + \eta_2^2) \|\hat{R}\| \|V\| \|G_n\|.$$

This bound applies regardless of the sequence in which the two multiplications are performed. In view of (16,19,23) we can also bound  $\|\Delta S\|$  by

$$(25) \quad \begin{aligned} \|\Delta S\| &\leq 1.01(2\eta_2 + \eta_2^2)(1 + \eta_3)(1 + 3\eta_3)^2(1 + 4\eta_3) \|A_{sw}\| \|\hat{A}^{-1}\| \\ &\leq 2 \cdot 1.01^2 \eta_2 \|A_{sw}\| \|\hat{A}^{-1}\|. \end{aligned}$$

By means of the equations (13,21,24) we can derive the following more explicit expression for  $S$ :

$$S = (U^T U)^{-1} U^T (A_{sw} F^T \hat{A}^{-1} EU - \Delta A_{sw} F^T G + A_{sw} F^T \hat{A}^{-1} \Delta G) + \Delta S.$$

Obviously, the norm of  $S$  can be bounded by

$$(26) \quad \|S\| \leq 1.01 \|A_{sw} F^T \hat{A}^{-1}\|.$$

**4.4. Calculation of  $T$ .** The matrix  $T$  satisfies the equation

$$(27) \quad T = I + S + \Delta T,$$

where

$$\|\Delta T\| \leq \eta_1 \|I + S\|.$$

By means of some straightforward manipulations, using (13,21,24), we can see that

$$(28) \quad EU(I + S) = A\hat{A}^{-1}EU - E\Delta A_{sw}F^TG + EA_{sw}F^T\hat{A}^{-1}\Delta G + EU\Delta S.$$

Consequently, we can express  $I + S$  as

$$I + S = (U^T U)^{-1} U^T E^T (A\hat{A}^{-1}EU - E\Delta A_{sw}F^TG + EA_{sw}F^T\hat{A}^{-1}\Delta G) + \Delta S,$$

and  $T$  is given by

$$T = (U^T U)^{-1} U^T E^T (A\hat{A}^{-1}EU - E\Delta A_{sw}F^TG + EA_{sw}F^T\hat{A}^{-1}\Delta G) + \Delta S + \Delta T.$$

If  $\Delta A_{sw}$ ,  $\Delta G$ , and  $\Delta S$  are sufficiently small, we have

$$(29) \quad \|\Delta T\| \leq 1.01\eta_1 \|A\hat{A}^{-1}\|.$$

Similarly, we also have

$$(30) \quad \|T\| \leq 1.01 \|A\hat{A}^{-1}\|.$$



4.5. **Calculation of  $\hat{R}$ .** If we solve the linear system for  $\hat{R}$  by Gaussian elimination, we get

$$(31) \quad T\hat{R} = R + \Delta\hat{R},$$

where

$$(32) \quad \|\Delta\hat{R}\| \leq \eta_4 \|T\| \|\hat{R}\|.$$

In order to bound  $\|\Delta\hat{R}\|$  differently, we need an alternative expression for  $\hat{R}$ . It is useful to consider the quantity  $I - G\hat{R}V^T F^T$  first. By using (13,21,24,27,31) we have

$$(33) \quad \begin{aligned} A(I - G\hat{R}V^T F^T) &= \hat{A} + E\Delta A_{\text{sw}} F^T (I - G\hat{R}V^T F^T) \\ &\quad - \Delta G\hat{R}V^T F^T + EU((\Delta S + \Delta T)\hat{R} - \Delta\hat{R})V^T F^T, \end{aligned}$$

which is equivalent to

$$(34) \quad \begin{aligned} I - G\hat{R}V^T F^T &= A^{-1}\hat{A} + A^{-1}E\Delta A_{\text{sw}} F^T (I - G\hat{R}V^T F^T) \\ &\quad - A^{-1}\Delta G\hat{R}V^T F^T + A^{-1}EU((\Delta S + \Delta T)\hat{R} - \Delta\hat{R})V^T F^T. \end{aligned}$$

Consequently, we can represent  $\hat{R}$  as

$$\begin{aligned} \hat{R} &= (U^T U)^{-1} U^T E^T (\hat{A} - \hat{A}A^{-1}\hat{A}) FV (V^T V)^{-1} \\ &\quad - (U^T U)^{-1} U^T E^T \left( \hat{A}A^{-1} E\Delta A_{\text{sw}} F^T (I - G\hat{R}V^T F^T) FV (V^T V)^{-1} \right. \\ &\quad \left. + EA_{\text{sw}} F^T A^{-1} \Delta G\hat{R} + \hat{A}A^{-1} EU((\Delta S + \Delta T)\hat{R} - \Delta\hat{R}) \right). \end{aligned}$$

Note that we can also write  $\hat{A} - \hat{A}A^{-1}\hat{A}$  as

$$(35) \quad \hat{A} - \hat{A}A^{-1}\hat{A} = EA_{\text{sw}} F^T A^{-1} \hat{A} = \hat{A}A^{-1} EA_{\text{sw}} F^T.$$

If we assume the rounding errors to be bounded we can show that

$$(36) \quad \|\hat{R}\| \leq 1.01 \|\hat{A} - \hat{A}A^{-1}\hat{A}\|.$$

By combining this result with (30,32) we can bound  $\|\Delta\hat{R}\|$  by

$$(37) \quad \|\Delta\hat{R}\| \leq 1.01^2 \eta_4 \|A\hat{A}^{-1}\| \|\hat{A} - \hat{A}A^{-1}\hat{A}\|.$$

4.6. **Calculation of  $P$ .** The calculation of  $P$  can be described by the equation

$$(38) \quad P = G\hat{R} + \Delta P,$$

where

$$\|\Delta P\| \leq \eta_2 \|G\| \|\hat{R}\|.$$

By using (23,36) we can bound  $\|\Delta P\|$  by

$$(39) \quad \|\Delta P\| \leq 1.01^2 \eta_2 \|\hat{A}^{-1}\| \|\hat{A} - \hat{A}A^{-1}\hat{A}\|.$$

Because of (34) the following alternative expression for  $P$  applies:

$$\begin{aligned} P &= A^{-1} EA_{\text{sw}} V (V^T V)^{-1} - A^{-1} E\Delta A_{\text{sw}} F^T (I - G\hat{R}V^T F^T) FV (V^T V)^{-1} \\ &\quad + A^{-1} \Delta G\hat{R} - A^{-1} EU((\Delta S + \Delta T)\hat{R} - \Delta\hat{R}) + \Delta P. \end{aligned}$$

Provided that the rounding errors remain bounded we certainly have

$$(40) \quad \|P\| \leq 1.01\|A^{-1}EA_{sw}\|.$$

**4.7. Calculation of  $Y$ .** The matrix  $Y$  is computed in three steps as follows:

$$(41) \quad A_{se}Y_s = B_s + \Delta Y_s,$$

$$(42) \quad \widehat{B}_n = B_n - A_{ne}Y_s + \Delta \widehat{B}_n,$$

$$(43) \quad A_{nw}Y_n = \widehat{B}_n + \Delta Y_n.$$

The error in computing  $\widehat{B}_n$  can be bounded by

$$(44) \quad \begin{aligned} \|\Delta \widehat{B}_n\| &\leq \eta_1\|B_n - A_{ne}Y_s\| + (1 + \eta_1)\eta_2\|A_{ne}\| \|Y_s\| \\ &\leq \eta_1\|B_n\| + (\eta_1 + \eta_2 + \eta_1\eta_2)\|A_{ne}\| \|Y_s\|. \end{aligned}$$

On the other hand, we can assume the following expressions for the residuals  $\Delta Y_s$  and  $\Delta Y_n$ :

$$(45) \quad \Delta Y_s = \Delta L_{Y_s} + \Delta M_{Y_s}Y_s,$$

$$(46) \quad \Delta Y_n = \Delta L_{Y_n} + \Delta M_{Y_n}Y_n.$$

The equations (41,42,43,45,46) can also be written in matrix terms as

$$(47) \quad \widehat{A}Y = B + \Delta Y,$$

where

$$(48) \quad \Delta Y := \begin{bmatrix} \Delta Y_n + \Delta \widehat{B}_n \\ \Delta Y_s \end{bmatrix} = \begin{bmatrix} \Delta L_{Y_n} + \Delta \widehat{B}_n \\ \Delta L_{Y_s} \end{bmatrix} + \begin{bmatrix} \Delta M_{Y_n} & \\ & \Delta M_{Y_s} \end{bmatrix} Y.$$

It will turn out to be useful to eliminate  $B_n$  from the right-hand side of (44). Observe that

$$(49) \quad \|B_n\| \leq \|B\| \leq \|\widehat{A}\| \|Y\| + \|\Delta Y\| \leq 1.01\|\widehat{A}\| \|Y\|.$$

Consequently, we can also bound  $\|\Delta \widehat{B}_n\|$  by

$$(50) \quad \|\Delta \widehat{B}_n\| \leq (2.01\eta_1 + \eta_2 + \eta_1\eta_2)\|\widehat{A}\| \|Y\| \leq 1.01(2\eta_1 + \eta_2)\|\widehat{A}\| \|Y\|.$$

**4.8. Calculation of  $X$ .** The computed value of  $X$  satisfies

$$(51) \quad X = Y - PV^T Y_n + \Delta X,$$

where

$$\|\Delta X\| \leq \eta_1\|Y - PV^T Y_n\| + (1 + \eta_1)(2 + \eta_2)\eta_2\|P\| \|V\| \|Y_n\|.$$

This bound can be processed further by using (16,40), and we get

$$\|\Delta X\| \leq \left( \eta_1 + 1.01(1 + \eta_3)(\eta_1 + (1 + \eta_1)(2 + \eta_2)\eta_2)\|A^{-1}EA_{sw}\| \right) \|Y\|.$$

A simpler bound is given by

$$(52) \quad \|\Delta X\| \leq (\eta_1 + 1.01^2(\eta_1 + 2\eta_2)\|A^{-1}EA_{sw}\|) \|Y\|.$$

## 5. RESIDUAL ERROR

We are now ready to give an expression for the residual  $\Delta B$  in (9). By using the equations (33,38,47,51), we get

$$(53) \quad \begin{aligned} \Delta B = & \left( E\Delta A_{\text{sw}}F^T(I - G\hat{R}V^TF^T) - \Delta G\hat{R}V^TF^T \right. \\ & \left. + EU((\Delta S + \Delta T)\hat{R} - \Delta\hat{R})V^TF^T - A\Delta PV^TF^T \right) Y \\ & + \Delta Y + A\Delta X. \end{aligned}$$

The value of  $\Delta B$  can be processed further. By combining (38) and (51) we get

$$(54) \quad (I - G\hat{R}V^TF^T)Y = X + \Delta PV^TF^TY - \Delta X.$$

If we use (9,13,21,34,35,47,53,54) we can also derive the following alternative expression for  $\hat{R}V^TF^TY$ :

$$(55) \quad \begin{aligned} \hat{R}V^TF^TY = & RV^TF^TX - ((\Delta S + \Delta T)\hat{R} - \Delta\hat{R})V^TF^TY \\ & + RV^TF^T\Delta PV^TF^TY - RV^TF^T\Delta X. \end{aligned}$$

Because of (22,48,54,55) we are led to the following expression for  $\Delta B$ :

$$(56) \quad \begin{aligned} \Delta B = & E\Delta A_{\text{sw}}F^TX \\ & - \begin{bmatrix} \Delta L_{G_n} + \Delta U_n \\ \Delta L_{G_s} \end{bmatrix} RV^TF^TX - \begin{bmatrix} \Delta M_{G_n} & \Delta M_{G_s} \end{bmatrix} GRV^TF^TX \\ & + EU(\Delta S + \Delta T)RV^TF^TX - EU\Delta\hat{R}V^TF^TY - A\Delta PV^TF^TY \\ & + \begin{bmatrix} \Delta L_{Y_n} + \Delta\hat{B}_n \\ \Delta L_{Y_s} \end{bmatrix} + \begin{bmatrix} \Delta M_{Y_n} & \Delta M_{Y_s} \end{bmatrix} Y + A\Delta X + O(\varepsilon^2). \end{aligned}$$

A key expression in (56) is  $GRV^TF^T$ . With the help of (13,21) we can transform this matrix into

$$GRV^TF^T = \hat{A}^{-1}EA_{\text{sw}}F^T - \hat{A}^{-1}E\Delta A_{\text{sw}}F^T + \hat{A}^{-1}\Delta GRV^TF^T.$$

Similarly, we use (9,47) to write  $Y$  as

$$(57) \quad \begin{aligned} Y = & \hat{A}^{-1}AX - \hat{A}^{-1}(\Delta B - \Delta Y) \\ = & D(D^{-1}\hat{A}^{-1}AD)D^{-1}X - \hat{A}^{-1}(\Delta B - \Delta Y). \end{aligned}$$

Now it is straightforward to write the residual  $\Delta B$  as

$$\Delta B = \Delta L_X + \Delta M_X X,$$

where

$$(58) \quad \Delta L_X = \begin{bmatrix} \Delta L_{Y_n} + \Delta \hat{B}_n \\ \Delta L_{Y_s} \end{bmatrix} + A\Delta X + O(\varepsilon^2),$$

$$(59) \quad \begin{aligned} \Delta M_X &= E\Delta A_{sw}F^T - \begin{bmatrix} \Delta L_{G_n} + \Delta U_n \\ \Delta L_{G_s} \end{bmatrix} RV^T F^T \\ &\quad - \begin{bmatrix} \Delta M_{G_n} \\ \Delta M_{G_s} \end{bmatrix} \hat{A}^{-1} EA_{sw}F^T \\ &\quad + EU(\Delta S + \Delta T)RV^T F^T - EU\Delta \hat{R}V^T F^T \hat{A}^{-1}A \\ &\quad - A\Delta PV^T F^T \hat{A}^{-1}A + \begin{bmatrix} \Delta M_{Y_n} \\ \Delta M_{Y_s} \end{bmatrix} \hat{A}^{-1}A. \end{aligned}$$

Let us give bounds for the norms of  $\Delta L_X$  and  $\Delta M_X D$ . From (50,52,57) we have

$$(60) \quad \begin{aligned} \|\Delta L_X\| &\leq 1.01 \left( \sqrt{2}\xi_{k-1} + 1.01(2\eta_1 + \eta_2)\|\hat{A}\| \|D\| \right. \\ &\quad \left. + (\eta_1 + 1.01^2(\eta_1 + 2\eta_2)\|A^{-1}EA_{sw}\|)\|A\| \|D\| \right) \\ &\quad \cdot \|D^{-1}\hat{A}^{-1}AD\| \|D^{-1}X\|. \end{aligned}$$

In order to bound  $\|\Delta M_X D\|$ , we need to multiply (59) by  $D$  from the right:

$$\begin{aligned} \Delta M_X D &= E\Delta A_{sw}F^T D - \begin{bmatrix} \Delta L_{G_n} + \Delta U_n \\ \Delta L_{G_s} \end{bmatrix} RV^T F^T D \\ &\quad - \begin{bmatrix} \Delta M_{G_n} D_{nw} \\ \Delta M_{G_s} D_{se} \end{bmatrix} D^{-1}\hat{A}^{-1}EA_{sw}F^T D \\ &\quad + EU(\Delta S + \Delta T)RV^T F^T D - EU\Delta \hat{R}V^T F^T DD^{-1}\hat{A}^{-1}AD \\ &\quad - A\Delta PV^T F^T DD^{-1}\hat{A}^{-1}AD \\ &\quad + \begin{bmatrix} \Delta M_{Y_n} D_{nw} \\ \Delta M_{Y_s} D_{se} \end{bmatrix} D^{-1}\hat{A}^{-1}AD. \end{aligned}$$

Because of (17) we can write the expression  $RV^T F^T D$  also as

$$RV^T F^T D = (U^T U)^{-1}U^T(A_{sw} - \Delta A_{sw})F^T D.$$

Consequently, its norm is certainly bounded by

$$(61) \quad \|RV^T F^T D\| \leq 1.01\|A_{sw}D_{nw}\|.$$

If we use (11,12,14,15,16,18,20,25,29,37,39,61) we can bound  $\|\Delta M_X D\|$  as follows:

$$(62) \quad \begin{aligned} \|\Delta M_X D\| &\leq 4\eta_3\|A_{sw}\| \|D_{nw}\| + 1.01^4(\eta_1 + 3\eta_2)\|A_{sw}D_{nw}\| \|A\| \|\hat{A}^{-1}\| \\ &\quad + 1.01^3(\eta_2 + \eta_4)\|A\| \|\hat{A}^{-1}\| \|\hat{A} - \hat{A}A^{-1}\hat{A}\| \|D_{nw}\| \|D^{-1}\hat{A}^{-1}AD\| \\ &\quad + 1.01^2\sqrt{2}\|A_{sw}D_{nw}\| \|\hat{A}^{-1}\|\xi_{k-1} \\ &\quad + (\|D^{-1}\hat{A}^{-1}EA_{sw}D_{nw}\| + \|D^{-1}\hat{A}^{-1}AD\|)\zeta_{k-1}. \end{aligned}$$

The inequalities (60) and (62) contain the quantities  $\xi_{k-1}$  and  $\zeta_{k-1}$  which are used to bound the errors at the previous level of the tear tree. Consequently, we can

use (60) and (62) to give recurrence relationships for  $\xi_k$  and  $\zeta_k$ . Let us define the quantities

$$(63) \quad f_\xi := \max 1.01\sqrt{2}\|D^{-1}\hat{A}^{-1}AD\|,$$

$$(64) \quad f_\zeta := \max \|D^{-1}\hat{A}^{-1}EA_{\text{sw}}D_{\text{nw}}\| + \|D^{-1}\hat{A}^{-1}AD\|,$$

$$(65) \quad g := \max 1.01^2\sqrt{2}\|A_{\text{sw}}D_{\text{nw}}\| \|\hat{A}^{-1}\|,$$

$$(66) \quad c_\xi := \max 1.01 \left( 1.01(2\eta_1 + \eta_2)\|\hat{A}\| \|D\| \right. \\ \left. + (\eta_1 + 1.01^2(\eta_1 + 2\eta_2)\|A^{-1}EA_{\text{sw}}\|)\|A\| \|D\| \right) \\ \cdot \|D^{-1}\hat{A}^{-1}AD\|,$$

$$(67) \quad c_\zeta := \max 4\eta_3\|A_{\text{sw}}\| \|D_{\text{nw}}\| \\ + 1.01^4(\eta_1 + 3\eta_2)\|A_{\text{sw}}D_{\text{nw}}\| \|A\| \|\hat{A}^{-1}\| \\ + 1.01^3(\eta_2 + \eta_4)\|A\| \|\hat{A}^{-1}\| \|\hat{A} - \hat{A}A^{-1}\hat{A}\| \\ \cdot \|D_{\text{nw}}\| \|D^{-1}\hat{A}^{-1}AD\|,$$

where the maximum is to be taken over all the nodes in the tear tree. The sequences  $\{\xi_k\}$  and  $\{\zeta_k\}$  thus satisfy the following recurrence relationships:

$$\begin{aligned} \xi_0 &= \eta_4\|A\| \|D\|, \\ \xi_k &= f_\xi\xi_{k-1} + c_\xi, & k \geq 1, \\ \zeta_0 &= 0, \\ \zeta_k &= f_\zeta\zeta_{k-1} + g\xi_{k-1} + c_\zeta, & k \geq 1. \end{aligned}$$

Their explicit solutions are given by<sup>1</sup>

$$(68)$$

$$\xi_k = \xi_0 f_\xi^k + c_\xi \frac{f_\xi^k - 1}{f_\xi - 1},$$

$$(69)$$

$$\begin{aligned} \zeta_k &= \xi_0 g \frac{f_\xi^k - f_\zeta^k}{f_\xi - f_\zeta} + c_\zeta \frac{f_\zeta^k - 1}{f_\zeta - 1} \\ &\quad + c_\xi g \left( \frac{1}{(f_\xi - 1)(f_\zeta - 1)} + \frac{f_\xi^k}{(f_\xi - 1)(f_\xi - f_\zeta)} - \frac{f_\zeta^k}{(f_\zeta - 1)(f_\xi - f_\zeta)} \right). \end{aligned}$$

These explicit expressions for  $\xi_k$  and  $\zeta_k$  are only valid if  $f_\xi \neq 1$ ,  $f_\zeta \neq 1$ , and  $f_\xi \neq f_\zeta$ . It would be possible to give similar expressions for these special cases, too. However these formulas would not give us more insight than (68) and (69).

At the top level of the tear tree the residual  $\Delta B$  is given by

$$\Delta B = \Delta L_X + \Delta M_X X,$$

where

$$\begin{aligned} \|\Delta L_X\| &\leq \xi_h \|D^{-1}X\|, \\ \|\Delta M_X D\| &\leq \zeta_h. \end{aligned}$$

Consequently,  $\|\Delta B\|$  can be bounded by

<sup>1</sup>We used Maple [3] to derive this result.

$$(70) \quad \|\Delta B\| \leq (\xi_h + \zeta_h) \|D^{-1}X\|.$$

The expression  $\|D^{-1}X\|$  is nothing else than the size of the solution  $X$  expressed in another norm.

## 6. STABILITY CRITERION

Let us now analyze the conditions under which Algorithm 1 computes a stable solution. We assume that all the matrices  $A$ ,  $\hat{A}$ , and  $D$  in the tear tree are nonsingular and only moderately ill-conditioned. Without this assumption the quantities  $c_\xi$  and  $c_\zeta$  could become arbitrarily large, like in the case of the matrices

$$A := \begin{bmatrix} \epsilon & 1 \\ 1 & \epsilon \end{bmatrix}, \quad \hat{A} := \begin{bmatrix} \epsilon & 1 \\ 0 & \epsilon \end{bmatrix},$$

with  $\epsilon \downarrow 0$ .

If  $c_\xi$  and  $c_\zeta$  are only small multiples of  $\epsilon\|A\|$ , the norm of the residual  $\Delta B$  in (70) will be on the order of  $\epsilon\|A\| \|D^{-1}X\|$  for  $f_\xi \approx 1$  and  $f_\zeta \approx 1$ . This condition is equivalent to the requirement that there exists a nonsingular block diagonal matrix  $D$ , partitioned commensurably with  $A$ , such that

$$(71) \quad \|D^{-1}\hat{A}^{-1}AD\| \approx 1$$

for all the matrices in the tear tree. If this stability criterion is met, Algorithm 1 is guaranteed to compute a stable solution provided that all the matrices  $A$ ,  $\hat{A}$ , and  $D$  in the tear tree are only moderately ill-conditioned.

It should be noted that (71) is a sufficient but not a necessary stability condition. Since we use the quantities  $\xi_k$  and  $\zeta_k$  to bound  $\|\Delta L_X\|$  and  $\|\Delta M_X D\|$  at each level in the tear tree, these bounds may grow even if  $\Delta L_X$  and  $\Delta M_X D$  remain bounded. On the positive side, we get a manageable error analysis and a simple stability criterion.

In the next two sections we will identify two classes of matrices for which the criterion (71) is always satisfied.

## 7. BLOCK DIAGONALLY DOMINANT MATRICES

An important class of matrices, for which the condition (71) is always satisfied, is given by the set of nonsingular block diagonally dominant matrices. In order to see this, we need the following theorem.

**Theorem 7.1.** *Let  $A = (A_{ij})$  be an  $m$ -by- $n$  block matrix with  $m < n$ . Furthermore let  $A$  be nonsingular and block diagonally dominant, i.e.,  $A$  has square nonsingular diagonal blocks and*

$$\|A_{ii}^{-1}\|_\infty \sum_{\substack{j=1 \\ j \neq i}}^n \|A_{ij}\|_\infty \leq 1, \quad i = 1, \dots, m.$$

We partition  $A$  into

$$A = \begin{bmatrix} A_1 & A_2 \end{bmatrix},$$

where  $A_1$  is a square  $m$ -by- $m$  block matrix, and  $A_2$  is an  $m$ -by- $(n-m)$  block matrix. We also assume that  $A_1$  is nonsingular. Under these assumptions the inequality

$$\|A_1^{-1}A_2\|_\infty \leq 1$$

applies.

*Proof.* Let  $\mathbf{y} = A_1^{-1} A_2 \mathbf{x}$ , where  $\mathbf{x} = (\mathbf{x}_i)$  and  $\mathbf{y} = (\mathbf{y}_i)$  are partitioned commensurably with  $A_2$  and  $A_1$ , respectively. Let  $\|\mathbf{y}_i\|_\infty = \|\mathbf{y}\|_\infty$ . From  $A_1 \mathbf{y} = A_2 \mathbf{x}$  we deduce that

$$\sum_{j=1}^m A_{ij} \mathbf{y}_j = \sum_{j=m+1}^n A_{ij} \mathbf{x}_{j-m}.$$

Consequently we can write  $\mathbf{y}_i$  as

$$\mathbf{y}_i = A_{ii}^{-1} \left( \sum_{j=m+1}^n A_{ij} \mathbf{x}_{j-m} - \sum_{\substack{j=1 \\ j \neq i}}^m A_{ij} \mathbf{y}_j \right).$$

Hence,

$$\|\mathbf{y}\|_\infty = \|\mathbf{y}_i\|_\infty \leq \|A_{ii}^{-1}\|_\infty \sum_{j=m+1}^n \|A_{ij}\|_\infty \|\mathbf{x}\|_\infty + \|A_{ii}^{-1}\|_\infty \sum_{\substack{j=1 \\ j \neq i}}^m \|A_{ij}\|_\infty \|\mathbf{y}\|_\infty,$$

or

$$\|\mathbf{y}\|_\infty \leq \frac{\|A_{ii}^{-1}\|_\infty \sum_{j=m+1}^n \|A_{ij}\|_\infty}{1 - \|A_{ii}^{-1}\|_\infty \sum_{\substack{j=1 \\ j \neq i}}^m \|A_{ij}\|_\infty} \|\mathbf{x}\|_\infty,$$

which yields  $\|\mathbf{y}\|_\infty \leq \|\mathbf{x}\|_\infty$  in view of the block diagonal dominance. □

In [7] the proof of this theorem for a point diagonally dominant matrix will appear as an exercise.

As a straightforward application of Theorem 7.1 we consider the norm of  $\hat{A}^{-1}A$  when  $A$  is block diagonally dominant. It is easily verified that

$$\hat{A}^{-1}A = \begin{bmatrix} I - A_{nw}^{-1}A_{ne}A_{se}^{-1}A_{sw} & 0 \\ A_{se}^{-1}A_{sw} & I \end{bmatrix}.$$

The two matrices  $[A_{nw} \ A_{ne}]$  and  $[A_{se} \ A_{sw}]$  satisfy the assumptions of Theorem 7.1, and we have

$$\begin{aligned} \|A_{nw}^{-1}A_{ne}\|_\infty &\leq 1, \\ \|A_{se}^{-1}A_{sw}\|_\infty &\leq 1. \end{aligned}$$

It is now easy to see that

$$\|\hat{A}^{-1}A\|_\infty \leq 2.$$

This bound is tight. It is attained, for instance, by the matrices

$$A = \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}, \quad \hat{A} = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}.$$

If a matrix  $A$  is block diagonally dominant, then so are all its submatrices in the tear tree. Consequently, if we set  $D = I$ , the condition (71) is satisfied, and Algorithm 1 will compute a stable solution of the linear system (1), provided that all the matrices  $A$  and  $\hat{A}$  in the tear tree are only moderately ill-conditioned.

## 8. M-MATRICES

A nonsingular  $n$ -by- $n$  matrix  $A$  is called an M-matrix if  $a_{ij} \leq 0$  for  $i \neq j$  and all the entries in  $A^{-1}$  are nonnegative. Many alternative characterizations of an M-matrix may be found in [1, Chapter 6] and [8, §6.4]. For the purpose of this error analysis the following definition is the most useful: A nonsingular  $n$ -by- $n$  matrix  $A$  is called an M-matrix if  $a_{ij} \leq 0$  for  $i \neq j$ ,  $a_{ii} > 0$ , and there exists a positive diagonal matrix  $D$  such that  $AD$  is strictly diagonally dominant, i.e.,

$$a_{ii}d_{ii} > \sum_{j \neq i} |a_{ij}|d_{jj}, \quad i = 1, \dots, n.$$

Without loss of generality we may assume that  $\|D\| = 1$ . Note that this definition is equivalent to the condition (M<sub>35</sub>) in [1, p. 137].

In view of the preceding discussion in §7 the stability criterion (71) is satisfied for this particular choice of the matrix  $D$ . Therefore, Algorithm 1 will compute a stable solution for linear systems with M-matrices.

## 9. NUMERICAL RESULTS

In this section we will present numerical results for three classes of test matrices to illustrate our error bounds. As our first example we choose the matrix

$$A := \begin{bmatrix} \epsilon & 1 & & & \\ & 1 & 1 & & \\ & 1 & \epsilon & 1 & \\ & & & 1 & 1 \\ & & & & 1 & \epsilon \\ & & & & & & 1 \end{bmatrix}.$$

We will see in a moment that the condition (71) is not satisfied for this matrix. The basic building block of  $A$  is given by

$$M := \begin{bmatrix} 1 & 1 \\ 1 & \epsilon \end{bmatrix},$$

which is a well-conditioned matrix for  $0 < \epsilon < 1/2$ . Specifically we have

$$\kappa_{\infty}(M) = \frac{4}{1 - \epsilon}.$$

On the other hand the matrix

$$\hat{M} := \begin{bmatrix} 1 & 1 \\ 0 & \epsilon \end{bmatrix}$$

becomes increasingly ill-conditioned as  $\epsilon$  tends to zero:

$$\kappa_{\infty}(\hat{M}) = 2 \frac{1 + \epsilon}{\epsilon}.$$

Consequently, if we partition the matrix  $A$  as indicated we can expect large values of  $\|D^{-1}\hat{A}^{-1}AD\|$  for all nodes in the tear tree.



TABLE 1. Properties of the first class of test matrices

$N$	$\ A\ $	$\ A^{-1}\hat{A}\ $	$\ \hat{A}^{-1}\ $	$\ D^{-1}\hat{A}^{-1}AD\ $	$\ \hat{A} - \hat{A}A^{-1}\hat{A}\ $
4	1.735	$1.010 \cdot 10^2$	$1.000 \cdot 10^4$	$9.901 \cdot 10^1$	$1.010 \cdot 10^{-2}$
6	1.906	$1.021 \cdot 10^2$	$1.010 \cdot 10^4$	$1.414 \cdot 10^2$	$1.010 \cdot 10^{-2}$
8	1.954	$1.031 \cdot 10^2$	$1.021 \cdot 10^4$	$1.744 \cdot 10^2$	$1.010 \cdot 10^{-2}$
10	1.974	$1.041 \cdot 10^2$	$1.031 \cdot 10^4$	$2.026 \cdot 10^2$	$1.010 \cdot 10^{-2}$
12	1.984	$1.052 \cdot 10^2$	$1.042 \cdot 10^4$	$2.278 \cdot 10^2$	$1.010 \cdot 10^{-2}$
14	1.990	$1.063 \cdot 10^2$	$1.053 \cdot 10^4$	$2.509 \cdot 10^2$	$1.010 \cdot 10^{-2}$
16	1.994	$1.074 \cdot 10^2$	$1.064 \cdot 10^4$	$2.724 \cdot 10^2$	$1.010 \cdot 10^{-2}$
18	1.996	$1.085 \cdot 10^2$	$1.074 \cdot 10^4$	$2.928 \cdot 10^2$	$1.010 \cdot 10^{-2}$
20	1.998	$1.096 \cdot 10^2$	$1.085 \cdot 10^4$	$3.122 \cdot 10^2$	$1.010 \cdot 10^{-2}$
22	1.999	$1.107 \cdot 10^2$	$1.097 \cdot 10^4$	$3.308 \cdot 10^2$	$1.010 \cdot 10^{-2}$

In order to avoid large matrices we choose an unusual tearing strategy: If the size  $N$  of the matrix  $A$  is two we solve the linear system by Gaussian elimination, otherwise we set  $n_{\text{nw}} = N - 2$  and  $n_{\text{se}} = 2$  (cf. Figure 1). This strategy leads to a degenerate tear tree with a height of  $h = N/2 - 1$ . The point of our example, however, does not depend on a particular tearing strategy. For any given strategy we can construct a matrix  $A$  that exhibits the same problems.

For our numerical example we choose  $\epsilon = 0.01$ . If we set  $D = I$  the value of the quantity  $\|D^{-1}\hat{A}^{-1}AD\|$  is on the order of  $\epsilon^{-2} = 10^4$ . However, if we choose

$$D := \begin{bmatrix} 1 & & & & \\ & c & & & \\ & & 1 & & \\ & & & c & \\ & & & & 1 \\ & & & & & c \end{bmatrix},$$

where  $c := \epsilon^2 = 10^{-4}$ , this norm can be reduced to the order of  $\epsilon^{-1} = 10^2$ .

In Table 1 we give the values of some key quantities from (63,64,65,66,67). For any  $A$  with  $N \geq 4$  we obtain the corresponding  $\hat{A}$  by setting  $a_{N-1,N-2}$  to zero. This is consistent with our special tearing strategy. We can see from the values of Table 1 that all the subproblems in the tear tree are only moderately ill-conditioned.

The entries of the right-hand side  $\mathbf{b}$  are given by  $b_i = i$ . We computed our results on a DECstation 3100 using a MATLAB [12] implementation of Algorithm 1. The unit roundoff is given by  $\epsilon = 2^{-52} \approx 2.2204 \cdot 10^{-16}$ . We used the singular value decomposition as our URV-decomposition.

In Table 2 (next page) we present the absolute and relative residuals for increasing matrix sizes  $N$ . Owing to our special tearing strategy, the height  $h$  of the tear tree increases linearly with  $N$ . We also compare the residuals obtained from Algorithm 1 with those from Gaussian elimination (backslash operator in MATLAB). The reader should observe the exponential error growth of the residual  $\|Ax - \mathbf{b}\|$ , which is due to the large value of  $\|D^{-1}\hat{A}^{-1}AD\|$ .

TABLE 2. Residuals for the first class of test matrices

N	h	$\ Ax - b\ $		$\frac{\ Ax - b\ }{\ A\  \ x\ }$	
		Algorithm 1	Gaussian elimination	Algorithm 1	Gaussian elimination
4	1	$3.634 \cdot 10^{-14}$	$4.965 \cdot 10^{-16}$	$1.042 \cdot 10^{-16}$	$1.423 \cdot 10^{-18}$
6	2	$5.298 \cdot 10^{-12}$	$4.441 \cdot 10^{-16}$	$9.169 \cdot 10^{-15}$	$7.687 \cdot 10^{-19}$
8	3	$1.523 \cdot 10^{-10}$	$8.882 \cdot 10^{-16}$	$1.918 \cdot 10^{-13}$	$1.119 \cdot 10^{-18}$
10	4	$3.020 \cdot 10^{-8}$	$6.280 \cdot 10^{-16}$	$2.996 \cdot 10^{-11}$	$6.231 \cdot 10^{-19}$
12	5	$7.708 \cdot 10^{-6}$	$8.882 \cdot 10^{-16}$	$6.307 \cdot 10^{-9}$	$7.268 \cdot 10^{-19}$
14	6	$2.044 \cdot 10^{-3}$	$1.538 \cdot 10^{-15}$	$1.422 \cdot 10^{-6}$	$1.070 \cdot 10^{-18}$
16	7	$9.639 \cdot 10^{-2}$	$2.176 \cdot 10^{-15}$	$5.754 \cdot 10^{-5}$	$1.315 \cdot 10^{-18}$
18	8	$6.382 \cdot 10^0$	$1.776 \cdot 10^{-15}$	$2.538 \cdot 10^{-3}$	$9.483 \cdot 10^{-19}$
20	9	$1.462 \cdot 10^3$	$1.776 \cdot 10^{-15}$	$4.403 \cdot 10^{-3}$	$8.483 \cdot 10^{-19}$
22	10	$2.708 \cdot 10^5$	$1.986 \cdot 10^{-15}$	$1.958 \cdot 10^{-2}$	$8.572 \cdot 10^{-19}$

TABLE 3. Properties of the class of diagonally dominant matrices

N	$\ A\ $	$\ A^{-1}\hat{A}\ $	$\ \hat{A}^{-1}\ $	$\ D^{-1}\hat{A}^{-1}AD\ $	$\ \hat{A} - \hat{A}A^{-1}\hat{A}\ $
4	1.990	$7.143 \cdot 10^1$	$1.611 \cdot 10^0$	1.407	$5.024 \cdot 10^1$
6	1.995	$1.060 \cdot 10^2$	$1.010 \cdot 10^2$	1.572	$6.677 \cdot 10^1$
8	1.997	$1.328 \cdot 10^2$	$2.025 \cdot 10^2$	1.752	$7.505 \cdot 10^1$
10	1.998	$1.554 \cdot 10^2$	$3.458 \cdot 10^2$	1.922	$8.003 \cdot 10^1$
12	1.999	$1.752 \cdot 10^2$	$5.303 \cdot 10^2$	2.081	$8.335 \cdot 10^1$
14	1.999	$1.931 \cdot 10^2$	$7.558 \cdot 10^2$	2.230	$8.573 \cdot 10^1$
16	1.999	$2.095 \cdot 10^2$	$1.022 \cdot 10^3$	2.371	$8.751 \cdot 10^1$
18	1.999	$2.248 \cdot 10^2$	$1.330 \cdot 10^3$	2.503	$8.889 \cdot 10^1$
20	2.000	$2.391 \cdot 10^2$	$1.678 \cdot 10^3$	2.630	$9.000 \cdot 10^1$
22	2.000	$2.526 \cdot 10^2$	$2.067 \cdot 10^3$	2.751	$9.091 \cdot 10^1$

As our second test matrix we choose the diagonally dominant matrix

$$A := \begin{bmatrix} 1 & a & & & & \\ a & 1 & b & & & \\ & b & 1 & a & & \\ & & a & 1 & b & \\ & & & b & 1 & a \\ & & & & a & 1 \end{bmatrix},$$

where we set  $a := 0.01$  and  $b := 0.99$ . All the subproblems in the tear tree are relatively well-conditioned. To illustrate this, we present in Table 3 the values of some key quantities in the tear tree.

TABLE 4. Residuals for the class of diagonally dominant matrices

N	h	Ax - b		Ax - b   / (  A     x  )	
		Algorithm 1	Gaussian elimination	Algorithm 1	Gaussian elimination
4	1	8.379 · 10 <sup>-15</sup>	2.665 · 10 <sup>-15</sup>	6.066 · 10 <sup>-17</sup>	1.929 · 10 <sup>-17</sup>
6	2	1.801 · 10 <sup>-14</sup>	1.779 · 10 <sup>-14</sup>	4.593 · 10 <sup>-17</sup>	4.535 · 10 <sup>-17</sup>
8	3	7.121 · 10 <sup>-14</sup>	3.815 · 10 <sup>-14</sup>	8.786 · 10 <sup>-17</sup>	4.708 · 10 <sup>-17</sup>
10	4	8.175 · 10 <sup>-14</sup>	1.113 · 10 <sup>-13</sup>	5.759 · 10 <sup>-17</sup>	7.842 · 10 <sup>-17</sup>
12	5	2.283 · 10 <sup>-13</sup>	1.088 · 10 <sup>-13</sup>	1.018 · 10 <sup>-16</sup>	4.851 · 10 <sup>-17</sup>
14	6	2.863 · 10 <sup>-13</sup>	1.632 · 10 <sup>-13</sup>	8.677 · 10 <sup>-17</sup>	4.947 · 10 <sup>-17</sup>
16	7	4.768 · 10 <sup>-13</sup>	1.835 · 10 <sup>-13</sup>	1.034 · 10 <sup>-16</sup>	3.981 · 10 <sup>-17</sup>
18	8	5.265 · 10 <sup>-13</sup>	3.521 · 10 <sup>-13</sup>	8.504 · 10 <sup>-17</sup>	5.687 · 10 <sup>-17</sup>
20	9	1.019 · 10 <sup>-12</sup>	5.198 · 10 <sup>-13</sup>	1.265 · 10 <sup>-16</sup>	6.450 · 10 <sup>-17</sup>
22	10	8.965 · 10 <sup>-13</sup>	4.299 · 10 <sup>-13</sup>	8.763 · 10 <sup>-17</sup>	4.201 · 10 <sup>-17</sup>

We use the same tearing strategy as in the first example. Again, the vector **b** is given by  $b_i = i$ . The norms of the residuals for different matrix sizes are presented as Table 4.

As our last example we choose the M-matrix

$$A := \begin{bmatrix} 1 & b & & & & \\ a & 1 & a & & & \\ & b & 1 & b & & \\ & & a & 1 & a & \\ & & & b & 1 & b \\ & & & & a & 1 \\ & & & & & b & 1 & b \\ & & & & & & a & 1 \end{bmatrix},$$

where we set  $a := -0.00049$  and  $b := -490$ . The matrix  $AD$  is diagonally dominant for

$$D := \begin{bmatrix} c & & & & & \\ & 1 & & & & \\ & & c & & & \\ & & & 1 & & \\ & & & & c & \\ & & & & & 1 \\ & & & & & & c \\ & & & & & & & 1 \end{bmatrix},$$

where  $c := 1000$ . The properties of  $A$  are listed as Table 5 (next page). We choose the same tearing strategy and the same right-hand side **b** as before, and the resulting residuals are presented as Table 6 (next page).

For the second and third class of test problems the condition (71) always holds. For these cases, the results in Tables 4 and 6 confirm that Algorithm 1 can compute a solution with a residual whose norm is on the same order as the norm of the residual from Gaussian elimination.

TABLE 5. Properties of the class of M-matrices

$N$	$\ A\ $	$\ A^{-1}\hat{A}\ $	$\ \hat{A}^{-1}\ $	$\ D^{-1}\hat{A}^{-1}AD\ $	$\ \hat{A} - \hat{A}A^{-1}\hat{A}\ $
4	$7.928 \cdot 10^2$	$1.158 \cdot 10^3$	$7.547 \cdot 10^2$	1.303	$8.387 \cdot 10^2$
6	$8.830 \cdot 10^2$	$1.445 \cdot 10^3$	$2.167 \cdot 10^3$	1.349	$9.549 \cdot 10^2$
8	$9.209 \cdot 10^2$	$1.614 \cdot 10^3$	$4.028 \cdot 10^3$	1.396	$1.001 \cdot 10^3$
10	$9.403 \cdot 10^2$	$1.712 \cdot 10^3$	$6.077 \cdot 10^3$	1.431	$1.021 \cdot 10^3$
12	$9.515 \cdot 10^2$	$1.768 \cdot 10^3$	$8.130 \cdot 10^3$	1.453	$1.029 \cdot 10^3$
14	$9.586 \cdot 10^2$	$1.799 \cdot 10^3$	$1.007 \cdot 10^4$	1.466	$1.033 \cdot 10^3$
16	$9.633 \cdot 10^2$	$1.815 \cdot 10^3$	$1.182 \cdot 10^4$	1.474	$1.035 \cdot 10^3$
18	$9.666 \cdot 10^2$	$1.823 \cdot 10^3$	$1.338 \cdot 10^4$	1.478	$1.035 \cdot 10^3$
20	$9.691 \cdot 10^2$	$1.828 \cdot 10^3$	$1.474 \cdot 10^4$	1.480	$1.036 \cdot 10^3$
22	$9.709 \cdot 10^2$	$1.830 \cdot 10^3$	$1.591 \cdot 10^4$	1.481	$1.036 \cdot 10^3$

TABLE 6. Residuals for the class of M-matrices

$N$	$h$	$\ Ax - b\ $		$\frac{\ Ax - b\ }{\ A\  \ x\ }$	
		Algorithm 1	Gaussian elimination	Algorithm 1	Gaussian elimination
4	1	$1.017 \cdot 10^{-12}$	$4.441 \cdot 10^{-16}$	$1.570 \cdot 10^{-19}$	$6.859 \cdot 10^{-23}$
6	2	$4.547 \cdot 10^{-12}$	$9.095 \cdot 10^{-13}$	$2.087 \cdot 10^{-19}$	$4.173 \cdot 10^{-20}$
8	3	$5.457 \cdot 10^{-12}$	$1.819 \cdot 10^{-12}$	$1.100 \cdot 10^{-19}$	$3.667 \cdot 10^{-20}$
10	4	$7.500 \cdot 10^{-12}$	$1.819 \cdot 10^{-12}$	$8.214 \cdot 10^{-20}$	$1.992 \cdot 10^{-20}$
12	5	$1.373 \cdot 10^{-11}$	$8.702 \cdot 10^{-15}$	$9.347 \cdot 10^{-20}$	$5.923 \cdot 10^{-23}$
14	6	$1.925 \cdot 10^{-11}$	$3.638 \cdot 10^{-12}$	$8.923 \cdot 10^{-20}$	$1.686 \cdot 10^{-20}$
16	7	$4.002 \cdot 10^{-11}$	$3.638 \cdot 10^{-12}$	$1.349 \cdot 10^{-19}$	$1.226 \cdot 10^{-20}$
18	8	$4.426 \cdot 10^{-11}$	$5.087 \cdot 10^{-14}$	$1.139 \cdot 10^{-19}$	$1.309 \cdot 10^{-22}$
20	9	$7.570 \cdot 10^{-11}$	$7.276 \cdot 10^{-12}$	$1.544 \cdot 10^{-19}$	$1.484 \cdot 10^{-20}$
22	10	$8.678 \cdot 10^{-11}$	$2.183 \cdot 10^{-11}$	$1.444 \cdot 10^{-19}$	$3.632 \cdot 10^{-20}$

## 10. CONCLUSIONS

We have presented an error analysis for a divide-and-conquer algorithm to solve linear systems with block Hessenberg matrices. Our error analysis corresponds closely to the recursive nature of this algorithm. The key to our analysis is equation (10) which gives a representation for the residuals  $\Delta B$  in the tear tree. Another important equation is (53) which shows the structure of the residuals in terms of local errors and errors from previous nodes in the tear tree. By combining (10) and (53) we can derive the linear recurrence relationships for  $\xi_k$  and  $\zeta_k$  which lead to the final error bound (70).

The precise value of the bound (70) is of minor importance. Rather we can show that Algorithm 1 computes a stable solution if the condition (71) is satisfied for all the matrices in the tear tree. This condition ensures that all the quantities in Algorithm 1 remain bounded. In particular equation (57) shows that the matrix  $D^{-1}Y$  can never be much larger than  $D^{-1}X$ . On the other hand, if  $\|D^{-1}\hat{A}^{-1}AD\|$  is significantly larger than one for all the nodes in the tear tree, we may encounter very large matrices  $D^{-1}Y$  during the execution of Algorithm 1, even if the final result  $D^{-1}X$  is small. This is exactly what happens in our first test case in §9, leading to a large residual  $\Delta B$ .

We have also shown that the condition (71) is always satisfied in the case of block diagonally dominant matrices and M-matrices. This explains the accurate results of the algorithm for this type of problems.

#### ACKNOWLEDGMENT

We would like to thank Nick Higham for reading and commenting on the manuscript.

#### REFERENCES

1. A. Berman and R. J. Plemmons, *Nonnegative matrices in the mathematical sciences*, Academic Press, New York, 1979. MR **82b**:15013
2. T. F. Chan, *Rank revealing QR factorizations*, *Linear Algebra Appl.* **88/89** (1987), 67–82. MR **88c**:15011
3. B. Char, K. Geddes, G. Gonnet, B. Leong, M. Monagan, and S. Watt, *Maple V language reference manual*, Springer, New York, 1991.
4. G. H. Golub and C. F. Van Loan, *Matrix computations*, 2nd ed., The Johns Hopkins University Press, Baltimore, MD, 1989. MR **90d**:65055
5. D. J. Higham and N. J. Higham, *Componentwise perturbation theory for linear systems with multiple right-hand sides*, *Linear Algebra Appl.* **174** (1992), 111–129. MR **93e**:65041
6. N. J. Higham, *How accurate is Gaussian elimination?*, *Numerical Analysis 1989*, Proceedings of the 13th Dundee Conference (D. F. Griffiths and G. A. Watson, eds.), Longman Scientific and Technical, 1990, pp. 137–154. CMP 91:17
7. ———, *Stability and accuracy of numerical algorithms (provisional title)*, 1994, in preparation.
8. H. Minc, *Nonnegative matrices*, Wiley, New York, 1988. MR **89i**:15001
9. G. W. Stewart, *On the solution of block Hessenberg systems*, *Numerical Linear Algebra with Applications* **2** (1995), 287–296.
10. ———, *An updating algorithm for subspace tracking*, *IEEE Trans. Signal Processing* **40** (1992), 1535–1541.
11. ———, *Implementing an algorithm for solving block Hessenberg systems*, Tech. Report CS-TR-3295, Department of Computer Science, University of Maryland, June 1994.
12. The MathWorks Inc., *MATLAB, high-performance numeric computation and visualization software*, Natick, Massachusetts, 1992.
13. J. H. Wilkinson, *Rounding errors in algebraic processes*, Prentice-Hall, Englewood Cliffs, NJ, 1963. MR **28**:4661
14. ———, *The algebraic eigenvalue problem*, Clarendon Press, Oxford, 1965. MR **32**:1894

INSTITUTE FOR ADVANCED COMPUTER STUDIES, UNIVERSITY OF MARYLAND, COLLEGE PARK, MARYLAND 20742

*E-mail address:* vonmatt@na-net.ornl.gov

DEPARTMENT OF COMPUTER SCIENCE AND INSTITUTE FOR ADVANCED COMPUTER STUDIES, UNIVERSITY OF MARYLAND, COLLEGE PARK, MARYLAND 20742

*E-mail address:* stewart@cs.umd.edu