

Validated Solution of Large Linear Systems¹

Siegfried M. Rump

Dedicated to U. Kulisch on the occasion of his 60th birthday

Abstract

Some new methods will be presented for computing verified inclusions of the solution of large linear systems. The matrix of the linear system is typically of band or sparse structure. There are no prerequisites to the matrix such as being M-matrix, symmetric, positive definite or diagonally dominant. For general band matrices of lower, upper bandwidth p, q of dimension n the computing time is less than $n \cdot (pq + p^2 + q^2)$. Examples with up to 1.000.000 unknowns will be presented.

Zusammenfassung

Es werden neuartige Methoden vorgestellt zur Berechnung sicherer Schranken der Lösung großer linearer Gleichungssysteme. Die Matrix des Gleichungssystems hat typischerweise Bandstruktur oder ist spärlich besetzt. Es werden keinerlei Voraussetzungen an die Matrix gestellt wie etwa M-Matrix, symmetrisch, positiv definit oder diagonal dominant. Für Bandmatrizen von oberer bzw. unterer Bandbreite p bzw. q der Dimension n ist die Rechenzeit kleiner als $n \cdot (pq + p^2 + q^2)$. Es werden Beispiele bis Dimension 1.000.000 diskutiert.

0 Notation

Let \mathbb{R} denote the set of real numbers, \mathbb{R}^n vectors and $\mathbb{R}^{n \times n}$ matrices over those. The letter n is only used for the dimension of vectors and matrices, others then n -vectors and $n \times n$ -matrices do not occur in this paper.

IPT denotes the power set over T , IT the interval extension for $T \in \{\mathbb{R}, \mathbb{R}^n, \mathbb{R}^{n \times n}\}$. Usually hyperrectangulars are used but others are not excluded. It should be stressed that interval operations producing validated bounds are rigorously and very efficiently implementable on digital computers, see [25], [1], [5], [28] for details.

¹published in R. Albrecht et al. (eds.): Validation numerics: theory and applications, vol. 9 of Computing Supplementum , pp. 191–212, Springer 1993

Intervals are written in brackets $a \pm b$ denotes the interval $[a - b, a + b]$, for some interval $[X]$ is $\| [X] \| := \max \{ |x| \mid x \in [X] \}$, $\text{mid}([X])$ denotes the midpoint, $\text{rad}([X])$ the radius of an interval $[X]$. Those terms apply to vectors and matrices componentwise.

The interior of a set is denoted by int , ρ denotes the spectral radius of a matrix and $\rho([A]) := \max \{ \rho(A) \mid A \in [A] \}$ for $[A] \in \mathbb{IR}^{n \times n}$. An interval linear system is sometimes written in short notation $[A] \cdot x = [b]$, solving it means to compute bounds for

$$\Sigma([A], [b]) := \{ x \in \mathbb{R}^n \mid \exists A \in [A], b \in [b] \quad \text{with} \quad Ax = b \}.$$

$\sigma_1, \dots, \sigma_n$ denote the singular values of a matrix in nonincreasing order such that $\sigma_1 = \|A\|_2$.

If not stated otherwise all operations are real or floating-point operations. We use operations \triangle with upwardly directed rounding, $* \in \{+, -, \cdot, /\}$ having the property

$$a \triangle b \geq a * b$$

where the latter operation $*$ is the real operation. In case a, b are vectors or matrices the \leq -sign applies componentwise.

1 Introduction

Few papers are known dealing with the problem of finding validated inclusions for the solution of sparse linear systems without calculating an approximate inverse of the system matrix. All of the papers known to the author not using an approximate inverse require special properties of the system matrix which is essentially being an M-matrix. The approximate inverse of a sparse matrix is in general full thus limiting the size of the tractable problems significantly. This is because of limitations in memory and because for banded systems the computing time depends quadratically on n . Our goal is to go large sizes, that is 100.000 unknowns and beyond and to keep the computing time for banded systems linearly dependant on n .

There are very interesting papers for condition estimation of sparse matrices (cf. [4], [13]). However, these are yielding estimations rather than verified bounds.

In this paper we describe our method for banded linear systems. The numerical examples are for banded systems, too. For sparse systems the techniques for reducing bandwidth, symbolic pivoting and others (cf. [13]) can be applied. The resulting linear system of reduced bandwidth can be treated by our methods.

There is one yet unpublished method without using an approximate inverse and without prerequisites on the matrix by Jansson [20]. Despite that there are essentially two different approaches known in the literature. The first is the direct extension of some numerical decomposition algorithm by means of replacing every real operation by the corresponding

interval operation. It has been shown that, for example, the interval version of Gaussian elimination is executable in this way for diagonally dominant matrices or M-matrices. In the general case intervals tend to grow in diameter rapidly due to data dependencies such that soon a pivot column only consists of intervals containing zero and the algorithm stops prematurely. This effect depends mainly on the dimension, not on the condition number. As a rule of thumb for general matrices with floating-point input data, for example for random matrices, the range of application of this approach is limited to dimension 50 when calculating in double precision which is roughly 17 decimals. The dimension is even more limited for interval input data.

The other approach uses fixed point methods. We shortly describe this ansatz because it gives insight in the problems we have to deal with.

Let a linear system $Ax = b$ with matrix $A \in \mathbb{R}^{n \times n}$ and right hand side $b \in \mathbb{R}^n$ be given together with some $\tilde{x} \in \mathbb{R}^n, R \in \mathbb{R}^{n \times n}$. \tilde{x} is considered to be an approximate solution to the linear system, R an approximate inverse of A . Krawczyk [21], [22] defines for $X \in \mathbb{IR}^n$ the following operator

$$K(X) := \tilde{x} - R \cdot (A\tilde{x} - b) + (I - RA) \cdot (X - \tilde{x}). \quad (1)$$

He shows that

$$\|I - RA\| < 1 \quad \text{and} \quad K(X) \subseteq X \quad \text{implies} \quad \exists \hat{x} \in X : A\hat{x} = b$$

(see also [26], [27]). In [29] it has been shown that the assumption $\|I - RA\| < 1$ can be replaced by $K(X) \subseteq \text{int}(X)$. Algorithms were designed to compute validated inclusions for the solution of general nonlinear systems [30]. There are a number of specializations to specific problems such as polynomial zeros [7], algebraic eigenproblems [29], evaluation of arithmetic expressions [8] and others taking advantage of the special situation. A basic theorem for linear systems is as follows.

Theorem 1.1. Let $\mathcal{A} \in \mathbb{IPR}^{n \times n}, \mathcal{B} \in \mathbb{IPR}^n$ be given and let $\tilde{x} \in \mathbb{R}^n, R \in \mathbb{R}^{n \times n}, \emptyset \neq X \in \mathbb{IPR}^n, X$ being compact. Define

$$\mathcal{Z} := R \cdot (\mathcal{B} - \mathcal{A}\tilde{x}) \quad \text{and} \quad \mathcal{C} := I - R \cdot \mathcal{A}, \quad (2)$$

$$L(X) := \mathcal{Z} + \mathcal{C} \cdot X, \quad (3)$$

all operations being power set operations. If

$$L(X) \subseteq \text{int}(X) \quad (4)$$

then R and every $A \in \mathbb{R}^{n \times n}$, $A \in \mathcal{A}$ is nonsingular and for every $b \in \mathcal{B}$ the unique solution $\hat{x} := A^{-1}b$ satisfies

$$\hat{x} \in \tilde{x} + L(X). \quad (5)$$

The proof consists of three basic steps. First take fixed but arbitrary $A \in \mathcal{A}, b \in \mathcal{B}$ thus reducing the problem to a point problem. Second, show that $C := I - RA \in \mathcal{C}$ is convergent ($\rho(C) < 1$) and therefore A and R are nonsingular. Moreover, the iteration $x^{k+1} := R(b - Ax^k) + C \cdot x^k$ has a unique fixed point $\hat{x} \in X$. Third show that this fixed point is the (unique) solution of $Ax = b$.

Thus theorem 1 already verifies the solvability of the linear system and gives a sufficient criterion for some $X \in \mathbb{I}\mathbb{R}^n$ for including the solution. To devise an algorithm for *finding* a validated inclusion $[X]$ we have to solve two problems. First the operations have to become executable on the computer and second we need a constructive way to obtain a suitable $[X]$. The first problem is solved by using interval operations rather than power set operations. On the computer floating-point bounds for the intervals are used. Then systems with $[A] \in \mathbb{I}\mathbb{R}^{n \times n}, [b] \in \mathbb{I}\mathbb{R}^n$ can be attacked. This includes for example point matrices the entries of which not being exactly representable on the computer by replacing those by the smallest enclosing machine interval (see [1], [27], [5], [28]).

For the second problem we use an iteration with a so-called ϵ -inflation (see [29], [31]). In this technique for a starting interval $[X] := [Z] := R \cdot ([b] - [A] \cdot \tilde{x})$ the iterated interval is made “fatter” in every step. This is used in combination with an Einzelschrittverfahren. It can be shown [31] that a validated inclusion will be found

- for a point system $Ax = b$ and *power set operations*
iff $\rho(I - R \cdot A) < 1$
- for an interval linear system $[A]x = [b]$ and *interval operations*
iff $\rho(|I - R \cdot [A]|) < 1$.

All of the fixed point methods known in the literature basically use theorem 1, especially (1.2) - (1.4), in one or the other way. Thus in our discussions for sparse linear systems we may concentrate on how to satisfy those conditions.

For simplicity let a point linear system $Ax = b, A \in \mathbb{R}^{n \times n}, b \in \mathbb{R}^n$ be given. We do not impose restrictions on A or b . For large banded or sparse linear systems the original approach cannot be used because it needs an approximate inverse R of A which is in general full. We

may omit this by using some decomposition of A . For $A = LU$ and $R = U^{-1}L^{-1}$ we obtain for $x \in \mathbb{R}^n$

$$R \cdot (b - A\tilde{x}) + (I - RA)x = U^{-1}L^{-1} \cdot (b - A\tilde{x} + (LU - A) \cdot x). \quad (6)$$

L and U preserve a banded structure of A . In a practical application we would think of replacing U^{-1} and L^{-1} by an efficient algorithm for solving triangular systems. From a mathematical point of view L and U are arbitrary. If for some $L, U \in \mathbb{R}^{n \times n}$ and $[X] \in \mathbb{IR}^n$ we can show that

$$M([X]) := U^{-1}L^{-1}(b - A\tilde{x} + (LU - A) \cdot [X]) \subseteq \text{int}([X]) \quad (7)$$

then theorem 1 implies that A is nonsingular and the unique solution $\hat{x} = A^{-1} \cdot b$ satisfies $\hat{x} \in \tilde{x} + M([X])$. $LU - A$ can be estimated during the decomposition of A , most simple and without additional cost for example using Crout's variant. Thus we have reduced our problem to computing a validated inclusion of the solution of a linear system with triangular point matrix and interval right hand side.

In (1.7) $b - A\tilde{x}$ is of order $\epsilon \cdot \|A\| \cdot \|\tilde{x}\|$ if \tilde{x} is a reasonable approximate solution, for example the one computed by floating-point Gaussian elimination. Also, numerical error analysis tells us that $LU - A$ will be of the order $\epsilon \cdot \|A\|$. $[X]$ shall contain the error of the approximate solution \tilde{x} which means that $(LU - A) \cdot [X]$ will be an interval vector of small magnitude. Thus we would not loose too much accuracy going to intervals being symmetric to the origin. This saves us half of the storage per interval vector. Clearly, for $0 < x \in \mathbb{R}^n$

$$U^{-1}L^{-1} \cdot (|b - A\tilde{x}| + |LU - A| \cdot x) < x \quad (8)$$

implies A being nonsingular and $A^{-1}b \in \tilde{x} \pm x$.

This reduces our problem to solving a triangular system with right hand side $[b]$ symmetric to the origin and we may further simplify it to $[b] := [-1, 1]$. In other words find

$$\text{validated bounds for } S := \{L^{-1} \cdot b \mid -1 \leq b \leq 1\}, L \in \mathbb{R}^{n \times n} \text{ lower triangular.} \quad (9)$$

All of the papers [11], [12], [23] using the fixed point approach solve (1.9) using interval backward substitution:

$$\text{for } i = 1 : n \text{ do } [x]_i = ([-1, +1] - \sum_{j=1}^{i-1} L_{ij} \cdot [x]_j) / L_{ii} \quad (10)$$

all operations in (1.10) being interval operations. Thus the intervals $[x]_j$ are symmetric to the origin and (1.10) can be written using absolute values

$$\text{for } i = 1 : n \text{ do } x_i = (1 + \sum_{j=1}^{i-1} |L_{ij}| \cdot x_j) / |L_{ii}| \quad (11)$$

yielding a true inclusion $S \subseteq [-x, +x]$. The overestimation can be estimated observing $x = \langle L \rangle^{-1} \cdot e$ where $e \in \mathbb{R}^n$, $e_i = 1$ for $1 \leq i \leq n$ and $\langle L \rangle$ is Ostrowski's comparison matrix (see [28]):

$$\langle L \rangle_{ij} := \begin{cases} |L_{ii}| & \text{for } i = j \\ -|L_{ij}| & \text{otherwise.} \end{cases}$$

For our special right hand side the maximal overestimation is the ratio

$$\|\langle L \rangle^{-1}\|_\infty / \|L^{-1}\|_\infty. \quad (12)$$

If we could estimate $\|L^{-1}\|_\infty$ then our problem (1.9) would be solved. In practical applications the ratio (1.12) is exponentially increasing with n unless L has special properties. Such properties are A and therefore L and U being M-matrices in which case $L = \langle L \rangle, U = \langle U \rangle$. This is the reason why M-matrices can be solved using interval Gaussian elimination without overestimation. To further illustrate the effect consider the following example due to Neumaier:

$$L = \begin{pmatrix} 1 & & & & & & & & & & \\ 1 & 1 & & & & & & & & & \\ 1 & 1 & 1 & & & & & & & & \\ & 1 & 1 & 1 & & & & & & & \\ & & & & \ddots & & & & & & \\ & & & & & & 1 & 1 & 1 & & \\ & & & & & & & & & & \end{pmatrix}, [b]_i = [-1, +1]. \quad (13)$$

Using interval backward substitution we obtain with $E := [-1, 1]$

$$\begin{aligned} [x]_1 &= E \\ [x]_2 &= E - [x]_1 = 2 \cdot E \\ [x]_3 &= E - [x]_1 - [x]_2 = 4 \cdot E \\ [x]_4 &= E - [x]_2 - [x]_3 = 7 \cdot E \end{aligned}$$

with exponentially growing diameter of $[x]_i$. This can also be seen from $\langle L \rangle^{-1}$ which we show for $n = 7$:

$$\langle L \rangle^{-1} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 2 & 1 & 1 & 0 & 0 & 0 & 0 \\ 3 & 2 & 1 & 1 & 0 & 0 & 0 \\ 5 & 3 & 2 & 1 & 1 & 0 & 0 \\ 8 & 5 & 3 & 2 & 1 & 1 & 0 \\ 13 & 8 & 5 & 3 & 2 & 1 & 1 \end{pmatrix}.$$

Thus $[x]$ computed by (1.10) is a huge overestimation of the true solution set which computes to

$$(L^{-1} \cdot [b])_i = \pm |L^{-1}| \cdot E = (i - [i/3]) \cdot E \subseteq n \cdot E.$$

This can be seen from

$$L^{-1} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ -1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & -1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & -1 & 1 & 0 & 0 & 0 \\ -1 & 1 & 0 & -1 & 1 & 0 & 0 \\ 0 & -1 & 1 & 0 & -1 & 1 & 0 \\ 1 & 0 & -1 & 1 & 0 & -1 & 1 \end{pmatrix}.$$

Unfortunately, this behaviour is typical for practical examples with matrices without special properties.

Methods based on the first approach (replacing floating-point operations by their corresponding interval operations in some numerical decomposition algorithm) are by their nature essentially restricted to diagonally dominant or inverse positive matrices (see for example [1], [28]). See also [33] for an interval version of Bunemann's algorithm for Poisson equation.

As we have just seen the fixed point approach as described in the literature is restricted to a similar class of matrices. This approach is used in [3], [11], [12], [23].

Using a coded version [3] of this algorithm the effect can be demonstrated. We used algorithm DSSSB with $IWK = 5$ which means the maximum possible amount of work is invested. We used $A = 0.1 \cdot LL^T$ with L from (1.13) and right hand side $(1, 0, \dots, 0)^T$. The factor 0.1 is used to make the factors of A not exactly representable on the computer. Then using double precision floating-point format which is approximately 17 decimal digits the algorithm fails for $n \geq 41$. For $n = 41$ we have $cond(A) = 2.3e3$. Taking the matrix (4.20) from [16] with $a = 1$ and the same right hand side $(1, 0, \dots, 0)^T$ the algorithm fails for $n \geq 48$. For $n = 48$ we have $cond(A) = 42$.

The amount of overestimation (1.12) is displayed in the following table.

n	10	20	30	40	50
$\ \langle L \rangle^{-1}\ _\infty / \ L^{-1}\ _\infty$	20.4	1265	1.0e5	9.9e6	9.7e8

Table 1.1 Overestimation of interval Gaussian elimination for L from (1.13)

The figures demonstrate the exponential behaviour of the overestimation.

2 The method

In order to bound (1.9) we may look for the singular values of L . Let U_r be the unit disk of radius r . Then $\|L^{-1} \cdot u\|_2, u \in U_r$ is bounded by $\sigma_n(L)^{-1} \cdot \|u\|_2 = \sigma_n(L)^{-1} \cdot r$. Thus a validated lower bound on the smallest singular value of a triangular matrix would solve the problem. This, in turn, would also yield a validated condition estimator. The problem of finding fast and reliable (although not validated) condition estimators has been attacked by many authors ([9], [10], [15], [17], [18], [2], [6]).

Given an approximation $\tilde{\lambda}$ of $\sigma_n(L)$, $\tilde{\lambda}^2$ is an approximate eigenvalue of LL^T . If for some $\kappa \in \mathbb{R}$ being slightly less than one we could prove that $LL^T - \kappa\tilde{\lambda}^2 \cdot I$ is positive definite then $\kappa^{1/2} \cdot \tilde{\lambda}$ proved to be a lower bound of $\sigma_n(L)$.

L is a Cholesky factor of LL^T . The change of the Cholesky factor L into G with $GG^T = LL^T - \tilde{\lambda}^2 I$ is given by the following formulas:

$$\begin{aligned} \sum_{\nu=1}^i G_{i\nu}^2 &= \sum_{\nu=1}^i L_{i\nu}^2 - \tilde{\lambda}^2 & \text{for } i = j \\ \sum_{\nu=1}^j G_{i\nu}G_{j\nu} &= \sum_{\nu=1}^j L_{i\nu}L_{j\nu} & \text{for } i > j. \end{aligned} \tag{14}$$

We need, however, a validation for the fact that $LL^T - \tilde{\lambda}^2 I$ is positive semidefinite. When performing an exact Cholesky factorization of $LL^T - \tilde{\lambda}^2 I$ this is true if the algorithm is executable, i.e. if the diagonal elements stay nonnegative. Using floating-point operations we have to estimate the rounding errors during the computation. Rather than estimating them a priori by replacing the floating-point operations by the corresponding interval operations we estimate them a posteriori by estimating the difference of GG^T and $LL^T - \tilde{\lambda}^2 I$ for the computed Cholesky factor G and by using perturbation theory.

For the diagonal elements this means

$$\text{computing } G_{ii} := \left(\sum_{\nu=1}^i L_{i\nu}^2 - \sum_{\nu=1}^i G_{i\nu}^2 - \tilde{\lambda}^2 \right)^{1/2} \quad \text{approximately and}$$

estimating $|(LL^T - \tilde{\lambda}^2 I - GG^T)_{ii}| = |\sum_{\nu=1}^i L_{i\nu}^2 - \sum_{\nu=1}^i G_{i\nu}^2 - \tilde{\lambda}^2|$ rigorously.

For off-diagonal elements this means

computing $G_{ij} := (\sum_{\nu=1}^j L_{i\nu}L_{j\nu} - \sum_{\nu=1}^{j-1} G_{i\nu}G_{j\nu})/G_{jj}$ approximatively and

estimating $|(LL^T - \tilde{\lambda}^2 I - GG^T)_{ij}| = |\sum_{\nu=1}^j L_{i\nu}L_{j\nu} - \sum_{\nu=1}^j G_{i\nu}G_{j\nu}|$ rigorously.

The computation and the estimation can essentially be done in one step. First the common part of both sums, resp. is evaluated *with* error estimation, then the midpoint is used for the floating-point component G_{ii}, G_{ij} of G , resp. and the interval part for the error estimation.

If only the four basic interval operations, that is IEEE 754 [19] arithmetic, is available that is the best we can do. If a precise scalar product [24], [25] is available then we can do better. For the diagonal elements we compute the exact value

$$dot := \sum_{\nu=1}^i L_{i\nu}^2 - \sum_{\nu=1}^{i-1} G_{i\nu}^2 - \tilde{\lambda}^2$$

and for S being the value of dot rounded to nearest we get $G_{ii} := fl(\sqrt{S})$, that is G_{ii} is the floating-point square root of S . Then we use the accumulating feature of the scalar product and compute the exact value of $dot - G_{ii}^2$. This value rounded to the smallest enclosing interval provides a very sharp bound for the error $(LL^T - \tilde{\lambda}^2 I - GG^T)_{ii}$. For the off-diagonal elements we proceed in a similar way. To avoid to formulate the algorithm twice we simply state in the diagonal case

Compute $S, \Delta S$ such that

$$\sum_{\nu=1}^i L_{i\nu}^2 - \sum_{\nu=1}^{i-1} G_{i\nu}^2 - \tilde{\lambda}^2 \in S \pm \Delta S.$$

For basic interval operations this means S being the midpoint, ΔS the radius of the left hand side computed in naive interval arithmetic. With the precise scalar product we proceed as described before. The off-diagonal elements are treated similarly.

Having an estimation on $E := LL^T - \tilde{\lambda}^2 I - GG^T$ and assuming the diagonal of G being nonnegative implies that $LL^T - \tilde{\lambda}^2 I - E$ is positive semidefinite. Hence perturbation theory tells us that the eigenvalues of $LL^T - \tilde{\lambda}^2 I$ are not smaller than $-\rho(E)$ (cf. [14], Corollary 8.1.3) and those of LL^T not smaller than $\tilde{\lambda}^2 - \rho(E)$. Now $\rho(E)$ can be estimated conveniently by $\|E\|_\infty$ which is done in the following algorithm. There the i th row sum is stored in e_i . When computing the ij -th component of G the error $(LL^T - \tilde{\lambda}^2 I - GG^T)_{ij}$ contributes to e_i and e_j due to symmetry. To obtain an upper bound on $\|E\|_\infty$ upward directed rounding is used in the computation of the e_i and e_{max} .

We give the algorithm for full matrix L . It can be altered for band matrices in a straightforward manner. Pivoting is omitted because $LL^T - \tilde{\lambda}^2 I$ is (hopefully) positive definite.

Given nonsingular lower triangular $L \in \mathbb{R}^{n \times n}$ and $\tilde{\lambda} \in \mathbb{R}$ do

$e_{\max} := 0$

for $i = 1 : n$ do $e_i := 0$;

for $i = 1 : n$ do

for $j = 1 : i - 1$ do

Compute $S, \Delta S$ such that

$$\sum_{\nu=1}^j L_{i\nu}L_{j\nu} - \sum_{\nu=1}^{j-1} G_{i\nu}G_{j\nu} \in S \pm \Delta S;$$

$$G_{ij} := fl(S/G_{jj});$$

Compute ΔT such that

$$|S - G_{ij}G_{jj}| \leq \Delta T;$$

$$d := \Delta S \triangle \Delta T; e_i := e_i \triangle d; e_j := e_j \triangle d;$$

Compute $S, \Delta S$ such that

$$\sum_{\nu=1}^i L_{i\nu}^2 - \sum_{\nu=1}^{i-1} G_{i\nu}^2 - \tilde{\lambda}^2 \in S \pm \Delta S;$$

$$G_{ii} := fl(\sqrt{S});$$

Compute ΔT such that

$$|S - G_{ii}^2| \leq \Delta T;$$

$$e_i := e_i \triangle \Delta S \triangle \Delta T;$$

$e_{\max} := \max_i e_i$;

Algorithm 2.1 Cholesky factorization of $LL^T - \tilde{\lambda}^2 I$ with lower bound for $\sigma_n(L)$

In precise computation, ΔS and ΔT as well as T would be zero according to (2.1). The main effort in the algorithm goes into the two inner products for computing S together with a validated bound. If L is a lower triangular band matrix of bandwidth p then the vector e needs only to be of length $p+1$ storing the values cyclically. Also, G needs only $(p+1)*(p+1)$ elements of storage.

It should be stressed that G is computed in floating-point arithmetic without presumptions on its accuracy. If the algorithm finishes successfully, i.e. the radicands are nonnegative, then the G_{ii} are nonnegative and therefore GG^T is positive semidefinite with

$$\|(LL^T - \tilde{\lambda}^2 I) - GG^T\|_{\infty} \leq e_{\max}. \quad (15)$$

The eigenvalues of LL^T are the squared singular values of L and are bounded from below by $\tilde{\lambda}^2 - e_{\max}$. This establishes the following theorem.

Theorem 2.1. If algorithm 2.1 finishes successfully (all square roots real) then $LL^T - (\tilde{\lambda}^2 - e_{\max})I$ is positive semidefinite. If $\tilde{\lambda}^2 \geq e_{\max}$ then

$$\sigma_n(L) \geq (\tilde{\lambda}^2 - e_{\max})^{1/2}.$$

The computing time for L with lower bandwidth p ($L_{ij} = 0$ for $i > j + p$) is less than $n \cdot p^2 + O(np)$ multiplications and additions plus $n(p + 1)$ divisions and n square roots.

Proof. The first part has been proved above, the computing time is a straightforward operation count. ■

In our applications we are particularly interested in sparse linear systems. This fact should be taken into account when implementing algorithm 2.1. For example, in case of a band matrix L the scalar products become very short compared to n .

Theorem 2.1 can be applied as follows. Consider some decomposition of A , for example $\tilde{L}\tilde{U} \approx A$ with $\tilde{A} := \tilde{L}\tilde{U}$. Then traditional norm estimates can be used to compute validated bounds for the solution together with theorem 2.1.

Theorem 2.2. Let $A \in \mathbb{R}^{n \times n}$, $b \in \mathbb{R}^n$ be given as well as nonsingular $\tilde{A} \in \mathbb{R}^{n \times n}$ and $\tilde{x} \in \mathbb{R}^n$. Define $\Delta A := \tilde{A} - A$ and suppose $\sigma_n(\tilde{A}) > n^{1/2} \cdot \|\Delta A\|_\infty$.

Then A is not singular and for $\hat{x} := A^{-1}b$ holds

$$\|\hat{x} - \tilde{x}\|_\infty \leq \frac{n^{1/2} \cdot \|b - A\tilde{x}\|_\infty}{\sigma_n(\tilde{A}) - n^{1/2} \cdot \|\Delta A\|_\infty}. \quad (16)$$

Proof. Since $\|\tilde{A}^{-1} \cdot \Delta A\|_2 \leq \sigma_n(\tilde{A})^{-1} \cdot \|\Delta A\|_2 \leq n^{1/2} \cdot \sigma_n(\tilde{A})^{-1} \cdot \|\Delta A\|_\infty < 1$ the matrix $I - \tilde{A}^{-1} \cdot \Delta A = \tilde{A}^{-1} \cdot A$ and hence A is invertible. Now

$$(I - \tilde{A}^{-1} \cdot \Delta A)(\hat{x} - \tilde{x}) = \tilde{A}^{-1} \cdot A \cdot (\hat{x} - \tilde{x}) = \tilde{A}^{-1} \cdot (b - A\tilde{x}).$$

Using $\|(I - F)^{-1}\| \leq (1 - \|F\|)^{-1}$ for convergent $F \in \mathbb{R}^{n \times n}$ this implies

$$\|\hat{x} - \tilde{x}\|_\infty \leq \frac{\|\tilde{A}^{-1} \cdot (b - A\tilde{x})\|_\infty}{1 - \|\tilde{A}^{-1} \cdot \Delta A\|_\infty} \quad (17)$$

and with $\|B\|_\infty \leq n^{1/2} \cdot \|B\|_2$ for $B \in \mathbb{R}^{n \times n}$

$$\|\hat{x} - \tilde{x}\|_\infty \leq \frac{n^{1/2} \cdot \sigma_n(\tilde{A})^{-1} \cdot \|b - A\tilde{x}\|_\infty}{1 - n^{1/2} \cdot \sigma_n(\tilde{A})^{-1} \cdot \|\Delta A\|_\infty}$$

proving the theorem. ■

In a practical application \tilde{A} is some floating-point decomposition of A , for example $\tilde{A} = \tilde{L}\tilde{U}$. Then the application of theorem 2.2 runs as follows. The nonsingularity of \tilde{A} is obvious. Compute an approximate solution \tilde{x} of $\tilde{A}x = b$ and a lower bound for $\sigma_n(\tilde{A})$ by

Hence $\|A\|_2 \leq (\|\Delta A\|_1 \cdot \|\Delta A\|_\infty)^{1/2}$ will in general be smaller than $n^{1/2} \cdot \|\Delta A\|_\infty$. Applying this to (2.4) we obtain the following result.

Theorem 2.3. Let $A \in \mathbb{R}^{n \times n}$, $b \in \mathbb{R}^n$ be given as well as nonsingular $\tilde{A} \in \mathbb{R}^{n \times n}$ and $\tilde{x} \in \mathbb{R}^n$. Define $\Delta A := \tilde{A} - A$ and suppose $\sigma_n(\tilde{A}) > (\|\Delta A\|_1 \cdot \|\Delta A\|_\infty)^{1/2}$.

Then A is not singular and for $\hat{x} := A^{-1}b$ holds

$$\|\hat{x} - \tilde{x}\|_\infty \leq \|\hat{x} - \tilde{x}\|_2 \leq \frac{\|b - A\tilde{x}\|_2}{\sigma_n(\tilde{A}) - (\|\Delta A\|_1 \cdot \|\Delta A\|_\infty)^{1/2}}. \quad (19)$$

Theorem 2.3 also follows by a fixed point argument. Using $\tilde{A} = LU$ and a disk of radius r instead of $[x]$ in (2.7) gives according to theorem 1.1

$$\begin{aligned} \tilde{A}^{-1} \cdot (b - A\tilde{x} + (\tilde{A} - A) \cdot U_r) \subseteq \text{int}(U_r) \Rightarrow \\ A \text{ is not singular and } A^{-1}b \in \tilde{x} + U_r. \end{aligned} \quad (20)$$

The inclusion in (2.7) is satisfied if

$$\sigma_n(\tilde{A})^{-1} \cdot (\|b - A\tilde{x}\|_2 + r \cdot \|\Delta A\|_2) < r.$$

This yields a bound on r and with a continuity argument (2.6).

The heuristic is that the elements of ΔA are roughly of the same size, namely $\epsilon \|A\|$. In the application of (2.6) we have to check $\sigma_n(\tilde{A}) > (\|\Delta A\|_1 \cdot \|\Delta A\|_\infty)^{1/2}$ to verify $\rho(\tilde{A}^{-1} \cdot \Delta A) < 1$ which is according to (2.5) more likely to happen than $\sigma_n(\tilde{A}) > n^{1/2} \cdot \|\Delta A\|_\infty$. Moreover, computing \tilde{x} by Gaussian elimination we know that the residual $\|b - A\tilde{x}\|$ will be of the order $\epsilon \cdot \|A\| \cdot \|\tilde{x}\|$ (cf. [14]).

In the following we add some computational hints for specific cases being relevant in practice.

I) A is M -matrix. Apply [32].

If A is symmetric positive definite we can use algorithm 2.1 to calculate a lower bound for $\sigma_n(A)$ directly. when replacing

$$\begin{aligned} \sum_{\nu=1}^j L_{i\nu} L_{j\nu} \quad \text{by } A_{ij} \quad \text{in row 7} \quad \text{and} \\ \sum_{\nu=1}^i L_{i\nu}^2 \quad \text{by } A_{ii} \quad \text{in row 13} \quad , \end{aligned}$$

then obviously

$$\sigma_n(A) \geq (\tilde{\lambda}^2 - e_{max})^{1/2}.$$

Replacing \tilde{A} by A in theorem 2.3 then yields

$$\|\hat{x} - \tilde{x}\|_\infty \leq \|\hat{x} - \tilde{x}\|_2 \leq (\tilde{\lambda}^2 - e_{max})^{-1/2} \cdot \|b - A\tilde{x}\|_2. \quad (21)$$

- II) *A is symmetric positive definite.* Compute a floating-point Cholesky decomposition $A \approx \tilde{G}\tilde{G}^T$ and an approximation $\tilde{\sigma}$ of the smallest singular value of A . Apply algorithm 2.1 altered as described above with $\tilde{\lambda} = 0.9 \cdot \tilde{\sigma}$ to compute a lower bound on $\sigma_n(A)$ and apply (2.8).

In case A is not symmetric positive definite one may use the following method. Having some approximate decomposition $A \approx \tilde{F} \cdot \tilde{G}$ compute an approximation $\tilde{\lambda}$ to the smallest singular value of A by inverse power method applied to $\tilde{F}\tilde{G} \cdot (\tilde{F}\tilde{G})^T$. If \tilde{F}, \tilde{G} are triangular this is inexpensive. Then apply theorem 2.1 with some obvious modifications to $A^T A - \tilde{\lambda}^2 I$ to bound $\sigma_n(A^T A)$.

This approach is working only for moderate condition numbers because the condition number of $A^T A$ is that of A squared. For working precision ϵ this limits the scope of application to $\text{cond}(A) < \epsilon^{-1/2}$ rather than $\text{cond}(A) < \epsilon^{-1}$.

In contrast we estimate the smallest singular value of the factors of the decomposition separately. We have to take provision that the condition numbers of the factors are of the same order, namely $\text{cond}(A)^{1/2}$. In this case the square of the condition number of the factors is still of the order of $\text{cond}(A)$ and no additional restrictions are imposed on A .

In the case A is symmetric we can do a little bit better than using LDL^T . Instead, let D_1, D_2 be diagonal such that $D_1 D_2 = D, |D_1| = |D_2| = |D|^{1/2}$. Then $LDL^T = (LD_1) \cdot (LD_2)^T$ and the usual LDL^T decomposition can be modified in an obvious way to compute $L_1 := L \cdot D_1$ and $L_2 := L \cdot D_2$ directly instead of L and D . Furthermore $D_1 = Q \cdot D_2$ with Q being a diagonal matrix with $+1$ or -1 in the diagonal thus being orthogonal. Therefore LD_1 and LD_2 have the same singular values and lower bound for $\sigma_n(LD_1)$ suffices for our purposes.

Despite saving computing time the heuristic is that $\sigma_n(LD_1)^2$ provides a better lower estimate for $\sigma_n(LDL^T)$ than $\sigma_n(L)^2 \cdot \sigma_n(D)$. Practical examples support this heuristic to a certain point. The same heuristic applies to general nonsymmetric matrices.

- III) *A is symmetric.* Compute an approximate $\tilde{L}_1 \cdot \tilde{L}_2^T$ decomposition as described above and an approximation σ of the smallest singular value of \tilde{L}_1 . Apply algorithm 2.1 with $\tilde{\lambda} = 0.9 \cdot \tilde{\sigma}$ to compute a lower bound σ on $\sigma_n(\tilde{L}_1)$ and apply theorem 2.2 or 2.3 with $\tilde{A} := \tilde{L}_1 \cdot \tilde{L}_2^T$ and $\sigma_n(\tilde{A}) \geq \sigma^2$.

In the general case we may apply an LU -decomposition. However, L tends to be fairly well-conditioned whereas the condition of A moves into U . Thus we may run into difficulties

trying to estimate $\sigma_n(U^T U)$. On the other hand the LDM^T -decomposition can be altered in an obvious way to distribute $D = D_1 \cdot D_2$, $|D_1| = |D_2| = |D|^{1/2}$ both in L and M as we did in the $L_1 L_2^T$ -decomposition in the symmetric case. This yields an LM -decomposition, L and M no longer being unit lower triangular. The heuristic is that then L and M are more or less equally conditioned, the condition number not being much bigger than the square root of the condition number of A .

- IV) *A is general nonsymmetric.* Compute an approximate $\tilde{L} \cdot \tilde{M}^T$ -decomposition of A and approximations $\tilde{\sigma}_1, \tilde{\sigma}_2$ for the smallest singular value of \tilde{L}, \tilde{M} , respectively. Apply algorithm 2.1 with $\tilde{\lambda}_1 = 0.9 \cdot \tilde{\sigma}_1, \tilde{\lambda}_2 = 0.9 \cdot \tilde{\sigma}_2$ to compute a lower bound σ_1, σ_2 on $\sigma_n(\tilde{L}), \sigma_n(\tilde{M})$ and apply theorem 2.2 or 2.3 with $\tilde{A} = \tilde{L} \cdot \tilde{M}^T$ and $\sigma_n(\tilde{A}) \geq \sigma_1 \cdot \sigma_2$.

It should be pointed out that the heuristic for cases III) and IV) works for many examples but also has its drawbacks. In the moment we do not know a general strategy for choosing a decomposition $A \approx \tilde{F}\tilde{G}$ which maximizes $\sigma_n(\tilde{F}) \cdot \sigma_n(\tilde{G})$. In case of symmetric positive definite A the method of choice is of course the Cholesky decomposition $A = GG^T$ with $\sigma_n(A) = \sigma_n(G)^2$.

Let $L \in \mathbb{R}^{n \times n}$ be of lower triangular of bandwidth p . Then approximations of the smallest singular value of L are either computed by

- inverse power iteration for LL^T at the cost of $2np$ ops per iteration or
- using some condition estimator at the cost of $c \cdot np$ ops, c small.

As has been pointed out before this is small against np^2 . Thus the total computing time for either of the algorithms for a linear system $Ax = b$ with A of lower, upper bandwidth p, q , respectively, $p \ll n, q \ll n$ is

I) <i>A is an M-matrix:</i>	$n \cdot pq$	ops
II) <i>A is symmetric positive definite:</i>	$n \cdot p^2$	ops
III) <i>A is symmetric indefinite:</i>	$\frac{3}{2}n \cdot p^2$	ops
IV) <i>A is general matrix:</i>	$n \cdot (pq + p^2 + q^2)$	ops.

Finally we want to mention how to use our methods in an interval setting, that is to solve $[A]x = [b]$, $[A] \in \mathbb{IIR}^{n \times n}$, $[b] \in \mathbb{IIR}^n$. Theorem 2.3 extends as follows.

Theorem 2.4. Let $[A] \in \mathbb{IIR}^{n \times n}$, $[b] \in \mathbb{IIR}^n$ be given as well as nonsingular $\tilde{A} \in \mathbb{R}^{n \times n}$ and $\tilde{x} \in \mathbb{R}^n$. Define $\Delta A := |[A] - \tilde{A}|$ and suppose $\sigma_n(\tilde{A}) > (\|\Delta A\|_1 \cdot \|\Delta A\|_\infty)^{1/2}$.

Then every $A \in [A]$ is nonsingular and for every $\hat{x} := A^{-1}b, A \in [A], b \in [b]$ holds

$$\|\hat{x} - \tilde{x}\|_\infty \leq \|\hat{x} - \tilde{x}\|_2 \leq \frac{\|[b] - [A] \cdot \tilde{x}\|_2}{\sigma_n(\tilde{A}) - (\|\Delta A\|_1 \cdot \|\Delta A\|_\infty)^{1/2}}. \quad (22)$$

The **proof** follows by applying theorem 2.3 to each $A \in [A], b \in [b]$. ■

We shortly describe an algorithm for solving a general interval linear system. We use the property $A \in [A] \Rightarrow \|A\|_p \leq \|\text{mid}([A])\|_p + \|\text{rad}([A])\|_p$ for $p \in \{1, \infty\}$.

Let $[A] \in \mathbb{IIR}^{n \times n}, [b] \in \mathbb{IIR}^n$ be given.

- 1) For $mA := \text{mid}([A])$ compute an approximate decomposition $\tilde{L} \cdot \tilde{M}^T := \tilde{A} \approx mA$ (see IV)) in floating-point arithmetic together with estimates ζ_1, ζ_∞ on $\|\tilde{A} - mA\|_1, \|\tilde{A} - mA\|_\infty$, resp.
- 2) Solve $\tilde{L} \cdot \tilde{M}^T \cdot \tilde{x} = mb, mb := \text{mid}([b])$ by floating-point backward and forward substitution to obtain \tilde{x}
- 3) Compute approximations for the smallest singular value $\tilde{\sigma}_1, \tilde{\sigma}_2$ of \tilde{L}, \tilde{M} by floating-point inverse power method applied to $\tilde{L}\tilde{L}^T, \tilde{M}\tilde{M}^T$, resp. or by some condition estimator
- 4) Apply algorithm 2.1 to compute lower bounds σ_1, σ_2 on $\sigma_n(\tilde{L}), \sigma_n(\tilde{M})$ using $\lambda := 0.9 \cdot \tilde{\lambda}_i$. If algorithm 2.1 does not finish successfully try smaller values for λ_i .
- 5) Calculate $z = \sup(\|[b] - [A] \cdot \tilde{x}\|)$ and upper bounds $\eta_1 \geq \|\text{rad}([A])\|_1, \eta_\infty \geq \|\text{rad}([A])\|_\infty$ using interval arithmetic (for η_1, η_∞ upwardly directed roundig suffices).
- 6) If $\mu := \sigma_1 \cdot \sigma_2 - ((\zeta_1 + \eta_1)(\zeta_\infty + \eta_\infty))^{1/2} > 0$ then every $A \in [A]$ is nonsingular and
$$\|\hat{x} - \tilde{x}\|_\infty \leq \|\hat{x} - \tilde{x}\|_2 \leq \mu^{-1} \cdot \|z\|_2$$
for every $\hat{x} = A^{-1}b$ with $A \in [A], b \in [b]$.

Algorithm 2.2 Inclusion of the solution of a general interval linear system

If very high accuracy of the inclusion is desired \tilde{x} may be stored in \tilde{x}_1 and \tilde{x}_2 with $\tilde{x} = \tilde{x}_1 + \tilde{x}_2$ (staggered correction, see [29], [34]). In this case $b - A\tilde{x}_1 - A\tilde{x}_2$ should be calculated in double the working precision. Using this method frequently very high or least significant bit accuracy is achieved. A simpler way is to perform a residual iteration

$$x^{k+1} := x^k + \tilde{M}^{-T} \tilde{L}^{-1} (b - Ax^k) \tag{23}$$

as usual. Only in the final step the addition is not executed but $\tilde{x}_1 := x^k$ and $\tilde{x}_2 := \tilde{M}^{-T} \tilde{L}^{-1} (b - Ax^k)$ are stored in separate vectors. This saves computing time and produces similar results to storing \tilde{x} in two parts \tilde{x}_1, \tilde{x}_2 from the beginning.

3 Computational results

In the following we give numerical results for three different types of our algorithm:

- (1) The *symmetric positive definite* case using a Cholesky-decomposition and proceeding as described in (II).
- (2) The *symmetric case* using a modified LDL^T -decomposition *without* pivoting as described in (III).
- (3) The *general case* using an LU -decomposition *with* pivoting from LAPACK.

In the following tables we display

n	dimension of the matrix
$\text{cond}(A)$	approximation of the $\ \cdot\ _\infty$ -condition number of A
iter	number of inverse power iterations to obtain an approximation for $\sigma_n(A)$
$\sigma_{\min}(A)$	lower bound for the smallest singular value of A
$\ \hat{x} - \tilde{x}\ _\infty / \ \tilde{x}\ _\infty$	upper bound for the relative error of the approximate solution \tilde{x} .

The condition number is estimated using the vector obtained by the inverse power iteration. Working accuracy is IEEE 754 double precision (approximately 17 decimals). As described in (2.10) we split \tilde{x} into \tilde{x}_1, \tilde{x}_2 and compute $b - A\tilde{x}_1 - A\tilde{x}_2$ in quadruple precision.

In all of the following examples the

- right hand side b is computed such that the solution \hat{x} of $Ax = b$ is $\hat{x}_i := (-1)^{i+1} \cdot 1/i$.

This introduces different magnitudes in the solution together with some roughness.

The first example, only displayed for reference purposes, is a discretisation of a Poisson equation

$$n := \begin{pmatrix} 4 & -1 & & \\ -1 & 4 & \ddots & \\ & & \ddots & \ddots \\ & & & \ddots & \ddots \end{pmatrix}; \quad A := \begin{pmatrix} M & -I & & \\ -I & M & \ddots & \\ & & \ddots & \ddots \\ & & & \ddots & \ddots \end{pmatrix} \quad (24)$$

with I being the identity matrix. We used three different bandwidths p .

Another example with high condition numbers are Hilbert matrices, $A_{ij} := 1/(i + j - 1)$. The following table shows the results.

n	cond	iter	$\sigma_{min}(A)$	$\ \hat{x} - \tilde{x}\ _\infty / \ \tilde{x}\ _\infty$
5	6.94E+05	3	1.81E-03	1.10E-22
6	2.26E+07	3	3.29E-04	4.44E-21
7	7.42E+08	3	5.91E-05	1.76E-19
8	2.45E+10	3	1.05E-05	1.88E-14
9	8.08E+11	3	1.87E-06	72.45E-16
10	2.68E+13	3	3.31E-07	1.86E-11
11	8.84E+14	3	5.83E-08	8.41E-10
12	2.60E+16	3	1.03E-08	2.38E-11
13	2.72E+17	2	1.21E-09	failed

Table 3.3. Hilbert matrices

Using Neumaier's example (1.13) we can show the behaviour for larger dimensions. We used $A = 10^{-1} \cdot LL^T$ producing a matrix of bandwidth 2. The factor 10^{-1} is introduced to make the factors of A not exactly representable. Otherwise a decomposition algorithm would rapidly produce the *exact* Cholesky factors. Here we observe decreasing precision of $\|\hat{x} - \tilde{x}\|_\infty / \|\tilde{x}\|_\infty$ with increasing condition number.

n	cond	iter	$\sigma_{min}(A)$	$\ \hat{x} - \tilde{x}\ _\infty / \ \tilde{x}\ _\infty$
100	1.26E+04	3	2.68E-02	3.49E-21
200	4.95E+04	3	1.35E-02	2.71E-20
500	3.06E+05	3	5.43E-03	8.50E-20
1000	1.22E+06	3	2.72E-03	3.40E-19
2000	4.87E+06	3	1.36E-03	1.36E-18
5000	3.04E+07	3	5.44E-04	8.47E-18
10000	1.22E+08	3	2.72E-04	3.39E-17
20000	4.87E+08	3	1.36E-04	1.35E-16
50000	3.04E+09	3	5.44E-05	8.47E-16
100000	1.22E+10	3	2.72E-05	3.39E-15
500000	3.04E+11	3	5.44E-06	8.47E-14
1000000	1.22E+12	3	2.72E-06	3.39E-13

Table 3.4. Neumaier's example with $A = 10^{-1}LL^T$, L from (1.13)

Next we go to the symmetric indefinite case. The first example is taken from [16], (4.20) with $a = 1$, bandwidth 2.

$$A := \begin{pmatrix} -1 & 2 & 1 & & & & \\ & 2 & 0 & 2 & 1 & & \\ & 1 & 2 & 0 & 2 & 1 & \\ & \dots & \dots & \dots & \dots & \dots & \\ & & & 1 & 2 & 0 & 2 \\ & & & & 1 & 2 & -1 \end{pmatrix}$$

Example (4.20) from [16]

The eigenvalues are $\lambda_k = \left(1 - 2 \cos \frac{k\pi}{n+1}\right)^2 - 3$, $1 \leq k \leq n$. We also display the computed upper bound on $\|A - \tilde{L}_1 \cdot \tilde{L}_2^T\|_2$. It is $\sigma_n(\tilde{L}_1) = \sigma_n(\tilde{L}_2)$.

n	cond	$\ A - \tilde{L}_1 * \tilde{L}_2^T\ _2$	iter	$\sigma_{min}(\tilde{L}_1)$	$\ \hat{x} - \tilde{x}\ _\infty / \ \tilde{x}\ _\infty$
100	5.37E+01	8.82E-14	3	1.01E-02	8.65E-21
200	9.02E+01	8.82E-14	3	7.66E-03	1.80E-20
500	3.29E+02	8.82E-14	3	2.62E-03	1.28E-19
1000	6.18E+02	2.59E-13	3	1.28E-03	1.02E-18
2000	1.29E+03	6.54E-12	3	3.27E-04	3.43E-18
5000	3.33E+03	6.54E-12	3	1.43E-04	1.82E-17
10000	5.52E+03	6.54E-12	3	9.57E-05	7.57E-17
20000	1.13E+04	6.54E-12	3	5.68E-05	2.62E-16
50000	3.24E+04	7.75E-12	3	2.04E-05	4.13E-15
100000	6.26E+04	7.75E-12	3	1.04E-05	7.62E-14

Table 3.5. Example (4.20) from [16], $a = 1$

The results show that, as before, few inverse power iterations are necessary to obtain an approximation for the smallest singular value of \tilde{L}_1 . The iteration is stopped when two successive iterates differ relatively less than 10^{-3} . Note that $\sigma_n(\tilde{L}_1)$ is fairly small in magnitude. This is due to the fact that the decomposition is performed without pivoting. Nevertheless sharp inclusions of the solution are achieved.

The next two tables show the behaviour for larger bandwidths. We use the abbreviation $M(a, b, c \dots)$ denoting a symmetric matrix with value a in the diagonal, b in the first sub-diagonal, c in the second and so forth.

n	cond	$\ A - \tilde{L}_1 * \tilde{L}_2^T\ _2$	iter	$\sigma_{min}(\tilde{L}_1)$	$\ \hat{x} - \tilde{x}\ _\infty / \ \tilde{x}\ _\infty$
100	2.07E+01	6.61E-14	6	2.43E-02	5.06E-20
200	2.82E+01	2.64E-13	6	1.55E-02	1.24E-19
500	1.16E+02	2.64E-13	4	6.78E-03	6.44E-19
1000	1.96E+02	2.64E-13	5	3.66E-03	2.18E-18
2000	4.18E+02	4.08E-12	5	1.12E-03	2.48E-17
5000	1.47E+03	4.46E-12	3	4.49E-04	1.42E-16
10000	1.98E+03	4.46E-12	6	2.93E-04	3.30E-16
20000	4.24E+03	2.15E-11	5	9.86E-05	3.65E-15
50000	1.14E+04	2.15E-11	5	5.06E-05	1.14E-14

Table 3.6. $M(1, -2, 3, 4, -5)$, bandwidth 4

n	cond	$\ A - \tilde{L}_1 * \tilde{L}_2^T\ _2$	iter	$\sigma_{min}(\tilde{L}_1)$	$\ \hat{x} - \tilde{x}\ _\infty / \ \tilde{x}\ _\infty$
100	1.08E+01	1.09E-12	4	2.05E-02	1.36E-19
200	4.38E+01	1.83E-12	3	7.91E-03	8.74E-19
500	7.46E+01	1.83E-12	4	5.26E-03	2.06E-18
1000	2.23E+02	6.48E-12	5	2.11E-03	1.24E-17
2000	2.54E+02	6.48E-12	6	1.89E-03	1.56E-17
5000	9.65E+02	6.48E-12	6	6.94E-04	1.28E-16
10000	2.23E+03	1.13E-11	5	2.41E-04	1.01E-15
20000	3.12E+03	1.13E-11	6	1.93E-04	1.54E-15

Table 3.7. $M(1, -2, 3, 4, -5, 5, 4, 3, 2, 1)$, bandwidth 9

Again the comparatively small values of $\sigma_n(\tilde{L}_1)$ are due to the lack of pivoting.

Finally we give some examples for the general case. First we show random matrices with upper and lower bandwidth 8 and uniformly distributed entries in the interval $[-1,1]$.

n	cond	$\ A - \tilde{L}\tilde{U}\ _2$	$\sigma_{min}(\tilde{L})$	$\sigma_{min}(\tilde{U})$	$\ \hat{x} - \tilde{x}\ _\infty / \ \tilde{x}\ _\infty$
100	1.2E+03	1.2E-16	7.1E-02	2.5E-02	5.5E-26
200	3.5E+03	1.4E-16	1.0E-01	5.6E-03	2.6E-25
500	9.3E+04	1.4E-16	2.6E-02	1.9E-03	4.1E-23
1000	6.3E+04	2.1E-16	1.3E-02	4.0E-04	2.1E-21
2000	4.4E+05	3.5E-16	1.4E-02	1.0E-04	7.7E-24
5000	4.9E+05	3.2E-16	7.7E-03	1.7E-04	1.3E-23
10000	8.7E+05	2.9E-16	1.8E-02	5.9E-05	6.0E-23
20000	2.2E+05	3.4E-16	2.4E-02	1.0E-04	5.5E-23
50000	1.6E+06	2.8E-16	1.0E-02	2.3E-05	7.6E-22

Table 3.8. Random matrices, upper and lower bandwidth 8

The LU -decomposition is performed using routines DGBTRF and DGBTRS from LAPACK with pivoting. The smallest singular values of \tilde{L} and \tilde{U} are not too near. This improves when distributing the diagonal of \tilde{U} among \tilde{L} and \tilde{U} . Finally we show an example being unsymmetric in upper and lower bandwidth.

n	cond	$\ A - \tilde{L}\tilde{U}\ _2$	$\sigma_{min}(\tilde{L})$	$\sigma_{min}(\tilde{U})$	$\ \hat{x} - \tilde{x}\ _\infty / \ \tilde{x}\ _\infty$
100	5.1E+02	1.1E-16	9.1E-02	2.0E-02	1.3E-26
200	1.4E+03	1.8E-16	4.5E-02	1.2E-03	8.1E-24
500	7.8E+06	2.3E-16	4.3E-02	2.7E-06	4.6E-23
1000	3.6E+07	2.2E-16	3.0E-02	3.8E-07	3.8E-22
2000	2.0E+05	5.8E-16	4.0E-02	5.9E-05	3.4E-24
5000	2.3E+06	4.3E-16	1.0E-02	7.7E-06	2.5E-22

Table 3.9. Random matrices, upper/lower bandwidth 8/6

Random matrices with symmetric upper and lower bandwidth are fairly well-conditioned. This changes when the bandwidth becomes unsymmetric. Then for moderate dimension we run into fairly ill-conditioned matrices. Again the numbers become better when distributing of \tilde{U} among \tilde{L} and \tilde{U} .

4 Conclusion

The presented algorithm in its different versions for symmetric positive definite, symmetric indefinite and general matrices works for high dimensions. Possible improvements and open questions are the following.

Using a condition estimator instead of inverse power iteration would eventually be cheaper but has not been tested yet. For symmetric indefinite and for general matrices it is not clear how to choose a proper decomposition $A \approx \tilde{F} \cdot \tilde{G}$ in order to minimize $\sigma_n(\tilde{F}) \cdot \sigma_n(\tilde{G})$. Using an LDL^T or LDM^T decomposition with D equally distributed among the other factors works fine in many cases but also has its drawbacks. The estimations given by the algorithm are ∞ -norm estimates on the relative error of an approximate solution \tilde{x} . Componentwise error estimates are not yet available. However, the results obtained up to now look promising.

References

- [1] Alefeld, G.; Herzberger, J.: Introduction to Interval Computations, Academic Press (1983)
- [2] Anderson, E.: Robust Triangular Solves for Use in Condition Estimation, Cray Research (1991)
- [3] High-Accuracy Arithmetic Subroutine Library, Program Description and User's Guide, Release 3, IBM Publications, Document Number SC 33-6164-3 (1986).
- [4] Arioli, M.; Demmel, J.W.; Duff, I.S.: Solving Sparse Linear Systems with Backward Error, SIAM J. Matrix Anal. Appl. 10, No. 2, 165–190 (1989)
- [5] Bauch, H.; Jahn, K.-U.; Oelschlägel, D.; Süsse, H.; Wiebigke, V.: Intervallmathematik, Theorie und Anwendungen; Mathematisch-naturwissenschaftliche Bibliothek, Bd. 72, B.G. Teubner, Leipzig (1987)
- [6] Bischof, Ch. H.; Tang, P.T.P.: Robust Incremental Condition Estimators, Argonne National Lab. (1992)
- [7] Böhm, H.: Berechnung von Polynomnullstellen und Auswertung arithmetischer Ausdrücke mit garantierter, maximaler Genauigkeit, Ph.D. dissertation, University of Karlsruhe (1983)
- [8] Böhm, H.; Rump, S.M.: Least Significant bit Evaluation for Arithmetic Expressions, Computing 30, 189–199 (1983)
- [9] Cline, A.K.; Moler, G.B.; Stewart, G.W.; Wilkinson, J.H.: An Estimate for the Condition Number of a Matrix, SIAM J. Num. Anal. 16, 368–375 (1979)
- [10] Cline, A.K.; Conn, A.R.; Van Loan, C.: Generalizing the LINPACK Condition Estimator, in: Numerical Analysis, ed., J. P. Hennart, Lecture Notes in Mathematics, No. 909, Springer-Verlag, New York (1982)
- [11] Cordes, D.; Kaucher, E.: Self-Validating Computation for Sparse Matrix Problems, in: E. Kaucher, U. Kulisch, Ch. Ullrich (Eds.): Computerarithmetic: Scientific Computation and Programming Languages. B.G. Teubner Verlag, Stuttgart (1987)
- [12] Cordes, D.: Spärlich besetzte Matrizen, in: U. Kulisch (Ed.): Wissenschaftliches Rechnen mit Ergebnisverifikation — Eine Einführung, ausgearbeitet von S. Geörg, R. Hammer und D. Ratz. Akademie Verlag, Berlin, und Vieweg Verlagsgesellschaft, Wiesbaden (1989)

- [13] Duff, I.S.; Erisman, A.M.; Reid, J.K.: *Direct Methods for Sparse Matrices*, Clarendon Press, Oxford (1986).
- [14] Golub, G. and v. Loan, C.: *Matrix Computations*, John Hopkins University Press, second edition (1989)
- [15] Grimes, R.G.; Lewis, J.G.: *Condition Number Estimation for Sparse Matrices*, *SIAM J. Sci. and Stat. Comp.* 2, 384–388 (1991)
- [16] Gregory, R.T.; Karney, D.L.: *A Collection of Matrices for Testing Computational Algorithms*, John Wiley (1969).
- [17] Hager, W.: *Condition Estimates*, *SIAM J. Sci. and Stat. Comp.* 5, 311–316 (1984)
- [18] Higham, N.J.: *Fortran Codes for Estimating the One-norm of a Real or Complex Matrix, with Applications to Condition Estimation*, *ACM Trans. Math. Soft.* 14, 381–396 (1987)
- [19] *IEEE Standard for Binary Floating-Point Arithmetic*, ANSI/IEEE Standard 754 (1985).
- [20] Jansson, Ch.: private communication.
- [21] Krawczyk, R.: *Newton-Algorithmen zur Bestimmung von Nullstellen mit Fehlerschranken*, *Computing* 4, 187–201 (1969)
- [22] Krawczyk, R.: *Fehlerabschätzung bei linearer Optimierung*, in "Interval Mathematics", edited by K. Nickel, *Lecture Notes in Computer Science* 29, Springer (1975).
- [23] Krämer, W.: *Verified Solution of Eigenvalue Problems with Sparse Matrices*. *Proceedings of 13th World Congress on Computation and Applied Mathematics*, IMACS '91, Dublin, 32–33 (1991)
- [24] Kulisch, U.: *Grundlagen des numerischen Rechnens*, Reihe Informatik 19, BI-Verlag (1976).
- [25] Kulisch, U.; Miranker, W.L.: *Computer Arithmetic in Theory and Practice*, Academic Press (1981)
- [26] Moore, R.E.: *A Test for Existence of Solutions for Non-Linear Systems*, *SIAM J. Numer. Anal.* 4 (1977)
- [27] Moore, R.E.: *Methods and Applications of Interval Analysis*, SIAM, Philadelphia (1979).

- [28] Neumaier, A.: Interval Methods for Systems of Equations, Cambridge University Press (1990)
- [29] Rump, S.M.: Kleine Fehlerschranken bei Matrixproblemen, Dissertation, Universität Karlsruhe (1980)
- [30] Rump, S.M.: Solving Algebraic Problems with High Accuracy, Habilitationsschrift, in: A New Approach to Scientific Computation, Hrsg. U.W. Kulisch und W.L. Miranker, Academic Press, 51–120 (1983)
- [31] Rump, S.M.: On the Solution of Interval Linear Systems, COMPUTING 47, p. 337-353 (1992).
- [32] Rump, S.M.: Inclusion of the Solution for Large Linear Systems with M-Matrix, Report of the Forschungsschwerpunkt Informations- und Kommunikationstechnik 91.3, Technical University Hamburg-Harburg (1991)
- [33] Schwandt, H.: An Interval Arithmetic Approach for the Construction of an almost Globally Convergent Method for the Solution of the Nonlinear Poisson Equation on the Unit Square, SIAM J. Sci. Stat. Comp., 5, No. 2 (1984)
- [34] Stetter, H.J.: Sequential Defect Correction in High-Accuracy Floating-Point Arithmetics, in: Numerical Analysis (Proceedings, Dundee 1983), Lecture Notes in Math. 1066, 186–202 (1984)