

ACCURATE SOLUTION OF DENSE LINEAR SYSTEMS PART II: ALGORITHMS USING DIRECTED ROUNDING

SIEGFRIED M. RUMP *

Abstract. In Part I and this Part II of our paper we investigate how extra-precise evaluation of dot products can be used to solve ill-conditioned linear systems rigorously and accurately. In Part I only rounding to nearest is used. In this Part II we improve the results significantly by permitting directed rounding. Linear systems with tolerances in the data are treated, and a comfortable way is described to compute error bounds for extremely ill-conditioned linear systems with condition number up to about \mathbf{u}^{-2}/n , where \mathbf{u} denotes the relative rounding error unit in a given working precision. We improve a method by Hansen/Bliek/Rohn/Ning/Kearfott/Neumaier. Of the known methods by Krawczyk, Rump, Hansen et al., Ogita and Nguyen we show that our presented Algorithm `LsErrBnd` seems the best compromise between accuracy and speed. Moreover, for input data with tolerances, a new method to compute componentwise inner bounds is presented. For not too wide input data they demonstrate that the computed inclusions are often almost optimal. All algorithms are given in executable Matlab code and are available from my homepage.

Key words. Linear systems, Matlab, rounding to nearest, data with tolerances, inner bounds, preconditioning, (extremely) ill-conditioned, extra-precise accumulation of dot products, Gaussian elimination, BLAS, LAPACK, error-free transformation, error analysis, rigorous error bounds.

AMS subject classifications. 65F05, 65G20

1. Introduction and notation. The paper divides into two parts: In Part I all algorithms use only the four basic floating-point operations in rounding to nearest, in the present Part II we use directed rounding and methods different from Part I to obtain superior results. All algorithms in both parts are presented in executable Matlab code.

The methods in Part I are based on norm estimates, verifying convergence of some residual matrix by approximating its Perron vector. In this Part II we verify the H -property of some matrix and demonstrate how this can be used to effectively compute verified error bounds of the solution of a linear system. Moreover, we show an efficient method to compute so-called “inner” inclusions of a linear system the data of which is afflicted with tolerances.

Dividing the paper into two parts serves also didactical purposes. In Part I we demonstrate that using only rounding to nearest allows to give simple algorithms to produce rigorous results. The algorithms presented in Part II can still be formulated in rounding to nearest, however, at the cost of easy readability.

All algorithms are given in executable Matlab code for which we reserve the “verbatim”-font. For instance, $\mathbf{C}=\mathbf{A}*\mathbf{B}$ means that \mathbf{C} is the result of the floating-point multiplication $\mathbf{A}*\mathbf{B}$, where \mathbf{A} and \mathbf{B} are compatible quantities (scalar, vector, matrix). For analyzing the error we use ordinary mathematical notation, for example in $P = \mathbf{A} \cdot \mathbf{B}$ the verbatim-font is used for floating-point quantities so that P is the exact (real) product of \mathbf{A} and \mathbf{B} . For $\mathbf{A}, \mathbf{B} \in \mathbb{F}^{n \times n}$ this implies $|P - \mathbf{C}| \sim \mathbf{u}|\mathbf{A}| \cdot |\mathbf{B}|$.

Comparison between vectors and matrices is always to be understood entrywise, for example $x \leq y$ for $x, y \in \mathbb{R}^n$ means $x_i \leq y_i$ for $1 \leq i \leq n$. For $A \in \mathbb{R}^{n \times n}$, the Ostrowski’s comparison matrix $\langle A \rangle \in \mathbb{R}^{n \times n}$ is

*Institute for Reliable Computing, Hamburg University of Technology, Schwarzenbergstraße 95, Hamburg 21071, Germany, and Visiting Professor at Waseda University, Faculty of Science and Engineering, 3-4-1 Okubo, Shinjuku-ku, Tokyo 169-8555, Japan (rump@tu-harburg.de).

defined by

$$(1.1) \quad \langle A \rangle_{ii} := |A_{ii}| \quad \text{and} \quad \langle A \rangle_{ij} := -|A_{ij}| \quad \text{for } i \neq j .$$

If $\langle A \rangle$ is an M -matrix, then A is called H -matrix. In that case A and $\langle A \rangle$ are non-singular, and $|A^{-1}| \leq \langle A \rangle^{-1}$. Finally, $\varrho(A)$ denotes the spectral radius of A .

This part II of the paper is organized as follows. In the next section we discuss several methods to obtain rigorous error bounds for linear systems based on H -matrices. In particular an efficient and best way is shown how to verify the H -property. In Section 3 we discuss how to implement these methods for rigorous error bounds using directed rounding.

Up to this point, standard Matlab suffices. From Section 4 on the algorithms become too involved and we use INTLAB [36], the Matlab toolbox for reliable computing. The toolbox INTLAB is entirely written in Matlab and thus portable on many environments. The only feature of INTLAB we need are the basic interval operations, so other libraries such as Intlib [9], Profil/Bias [10, 11] or b4m [39] may be used as well.

In Section 4 linear systems the data of which are afflicted with tolerances are treated. Using interval operations the code becomes easier to read without sacrificing performance and/or accuracy. Now the floor is prepared to discuss alternative approaches to compute rigorous error bounds in Section 5. In particular a method originated by Hansen is discussed and improved. In Section 6 we show how to obtain inclusions for extremely ill-conditioned matrices, i.e. with condition number up to \mathbf{u}^{-2} , and finally we show a new method to calculate inner inclusions, even if only one entry of the matrix and/or the right hand side is afflicted with a tolerance. This allows to judge the quality of outer inclusions. Detailed computational results and a conclusion finish the paper.

2. Rigorous error bounds for linear systems. Let a linear system $Ax = b$ with $A \in \mathbb{R}^{n \times n}$ and $b \in \mathbb{R}^n$ be given. As in Part I let $R \in \mathbb{R}^{n \times n}$ be an approximate inverse of A , for example computed by the Matlab command `inv`. Note that there are no a priori assumptions on A and R , in particular no accuracy requirement on R . For the following explanation assume the matrices A and R to be non-singular; in the following theorems this will be verified a posteriori by the methods.

Define $C := RA$. For an approximation \tilde{x} to the solution $A^{-1}b$ we show several ways how to estimate

$$(2.1) \quad \Delta := \tilde{x} - A^{-1}b = C^{-1}R(A\tilde{x} - b) = C^{-1}c \quad \text{with} \quad c := R(b - A\tilde{x}) .$$

In Part I of this paper we used normwise error estimates. Define $T := \text{diag}(t)$ for a positive vector $t \in \mathbb{R}^n$. Defining $F := I - C$ and exploring

$$(2.2) \quad C^{-1}c = c + C^{-1}Fc = c + T(I - T^{-1}FT)^{-1}T^{-1}Fc$$

and using $|Ex| \leq \|x\|_\infty \cdot |E|e$ for $E \in \mathbb{R}^{n \times n}$, $x \in \mathbb{R}^n$ and $e := (1, \dots, 1) \in \mathbb{R}^n$ we obtain

$$(2.3) \quad \begin{aligned} |\tilde{x} - A^{-1}b| &\leq |c| + \|(I - T^{-1}FT)^{-1}T^{-1}Fc\|_\infty \cdot Te \\ &\leq |c| + \frac{\|T^{-1}Fc\|_\infty}{1 - \|T^{-1}FT\|_\infty} \cdot t , \end{aligned}$$

provided $\|T^{-1}FT\|_\infty < 1$ (see Theorem 4.14 in Part I). Note that this true for any $0 < t \in \mathbb{R}^n$. Since $\|T^{-1}FT\|_\infty = \|T^{-1}|F|T\|_\infty = \|T^{-1}|F|Te\|_\infty$, the obvious choice for Te is the Perron vector t of $|F|$. Then $|F|t = \varrho t$ implies $T^{-1}|F|Te = \varrho e$, minimizing $\|T^{-1}FT\|_\infty$. For very ill-conditioned matrices, the choice $t = e$ may fail due to $\|F\|_\infty \geq 1$, see Figure 4.1 in Part I; otherwise, however, due to rounding errors in finite precision, sometimes the choice $t = e$ is superior to the Perron vector.

The previous result is based on the Neumann expansion $C^{-1} = (I - F)^{-1} = I + C^{-1}F$. If A is not too ill-conditioned, then $C = RA$ is not too far from the identity matrix and likely to be an H -matrix, which

means that $\langle C \rangle$ is an M -matrix. The matrix C is an H -matrix, also called generalized diagonally dominant, if and only if there exists some positive $v \in \mathbb{R}^n$ such that $u := \langle C \rangle v > 0$. For the moment assume such a vector v to be given.

Then $C = RA$ is an H -matrix, so that A and R are non-singular. Denote by

$$(2.4) \quad \langle C \rangle := D - E \quad \text{with } D > 0, E \geq 0$$

the splitting of $\langle C \rangle$ into diagonal and off-diagonal part. Note that $u = \langle C \rangle v > 0$ implies $D > 0$. Set

$$(2.5) \quad G := I - \langle C \rangle D^{-1} = ED^{-1} \geq 0 .$$

Following [19] define

$$(2.6) \quad w_k := \max_i \frac{G_{ik}}{u_i} \quad \text{for } 1 \leq k \leq n ,$$

so that $0 \leq w \in \mathbb{R}^n$ and $G \leq uw^T$. Multiplying $I - \langle C \rangle D^{-1} \leq uw^T$ from the left by $\langle C \rangle^{-1} \geq 0$ and using $|\langle C^{-1} \rangle| \leq \langle C \rangle^{-1}$, which is true [6] for any matrix C , implies

$$(2.7) \quad |\langle C^{-1} \rangle| \leq \langle C \rangle^{-1} \leq D^{-1} + vw^T .$$

Direct application of (2.7) to (2.1) gives

$$(2.8) \quad |\tilde{x} - A^{-1}b| \leq (D^{-1} + vw^T)|c| .$$

In contrast to the previous approach we use an implicit preconditioning of the matrix $\langle C \rangle$ to ensure $[D^{-1}\langle C \rangle]_{ii} \equiv 1$. For the splitting $C := \tilde{D} - \tilde{E}$ into diagonal and off-diagonal part we have $|\tilde{D}| = D$ and $|\tilde{E}| = E$, so that the identity $C^{-1} = \tilde{D}^{-1}(I + \tilde{E}C^{-1})$, $|\langle C^{-1} \rangle| \leq \langle C \rangle^{-1}$ and $\tilde{x} - A^{-1}b = C^{-1}c$ yield

$$(2.9) \quad |\tilde{x} - A^{-1}b| \leq \epsilon \quad \Rightarrow \quad |\tilde{x} - A^{-1}b| \leq D^{-1}(|c| + E\epsilon) .$$

As in (2.2) we may apply $C^{-1} = (I - F)^{-1} = I + C^{-1}F$ to $\Delta = C^{-1}c$ to obtain

$$(2.10) \quad |\tilde{x} - c - A^{-1}b| \leq (D^{-1} + vw^T)|Fc| .$$

Finally expanding the Neumann-series one more term gives $C^{-1} = I + F + C^{-1}F^2$, and therefore

$$(2.11) \quad |\tilde{x} - c - Fc - A^{-1}b| \leq (D^{-1} + vw^T)|F^2c| .$$

All the established bounds have to be estimated covering rounding errors which will be discussed in the next section. When relaxing the bounds (2.10) and (2.11) using $|Fc| \leq |F||c|$ etc., the estimates become particularly simple because the products are just computed in rounding to upwards. Extensive tests suggest that this weakens the computed bounds marginally.

We compared the normwise bounds (2.3) choosing $e = (1, \dots, 1)^T$ and the Perron vector of $|F|$ for t with (2.8), (2.9) applied to (2.8), (2.10) and (2.11) for various types of matrices with various right hand sides. We used in particular the matrices in Table 5.2 in Part I of this paper, and also ill-conditioned matrices up to dimension 1000.

In this test, (2.11) was almost always the best bound and (2.10) was a little weaker. The bound (2.8) was usually worse than (2.11) by a factor between 1.5 and 2, in few cases up to a factor 15. Hence the results are as expected.

The normwise bound (2.3) often failed for ill-conditioned matrices and $t = e$ because of $\|F\|_\infty \geq 1$. When choosing the Perron vector for t , the bound (2.3) was usually worse than (2.11) by a factor around 2,

sometimes up to a factor 20. The mentioned factors occur only for very ill-conditioned matrices, not too far from \mathbf{u}^{-1} . For matrices with moderate condition number all bounds are almost identical.

Besides (2.9) we may use $C^{-1} = I + FC^{-1}$ to deduce that $|\Delta| = |C^{-1}c| \leq \epsilon$ implies $|\Delta| \leq |c| + |F|\epsilon$, which means $|\Delta| \leq \min(\epsilon, |c| + |F|\epsilon)$. Applying this a few times to the discussed bounds, all final bounds were not too far apart. Since this improvement costs only some $\mathcal{O}(n^2)$ operations, a good choice seems to use (2.8), which computes in $\mathcal{O}(n)$ operations, and to improve it as mentioned.

THEOREM 2.1. *Let $A, R \in \mathbb{R}^{n \times n}$ and $b, \tilde{x} \in \mathbb{R}^n$ be given. Let $0 < v \in \mathbb{R}^n$ be such that $u := \langle RA \rangle v > 0$. Denote by $\langle RA \rangle := D - E$ the splitting of $\langle RA \rangle$ into diagonal and off-diagonal part, and define $w \in \mathbb{R}^n$ by (2.6). Then A and R are non-singular, and*

$$(2.12) \quad |A^{-1}b - \tilde{x}| \leq (D^{-1} + vw^T)|c| \quad \text{for } c := R(b - A\tilde{x}) .$$

Moreover, $|A^{-1}b - \tilde{x}| \leq \epsilon$ implies

$$(2.13) \quad |A^{-1}b - \tilde{x}| \leq D^{-1}(|c| + E\epsilon) \quad \text{and} \quad |A^{-1}b - \tilde{x}| \leq |c| + |I - RA|\epsilon .$$

2.1. Verification of H -property. It remains to find a positive vector v with $u := \langle C \rangle v > 0$, which means to prove that $C = RA$ is an H -matrix. There is a vast amount of literature on this problem, see, for example, [14] and the references therein. The optimal choice of v can be seen as follows.

Since RA is not too far from the identity matrix, the vector $e := (1, \dots, 1)^T \in \mathbb{R}^n$ seems a proper choice for v and is often used, for example in [19]. However, this means that $\langle C \rangle = D - E$ is strictly diagonally dominant, which need not be the case due to corruption of R by rounding errors. A better choice is $D^{-1}e$, but even then for ill-conditioned matrices RA it may happen that RA is an H -matrix, whereas this fact cannot be verified by both choices of the vector v . This is not uncommon, see Figure 3.1 in Subsection 3.2. After spending considerable effort to compute RA , it is particularly annoying to fail to prove RA to be an H -matrix only because of some standard choice of v .

Therefore it is worth to spend a few vector iterations to adapt the vector v to the particular circumstances. The matrix $\langle C \rangle = D - E$ is an M -matrix if and only if $r := \rho(D^{-1}E) < 1$. But $D^{-1}E$ is a nonnegative matrix, thus its spectral radius corresponds to an eigenvalue, the Perron root. Usually $D^{-1}E$ is at least nonnegative irreducible if not positive, so that the corresponding eigenvector is positive, and a power iteration

$$(2.14) \quad v^{(k+1)} := D^{-1}E v^{(k)} \quad \text{for } v^{(0)} := D^{-1}e$$

implies that $\max_i v_i^{(k+1)} / v_i^{(k)}$ decreases monotonically to the Perron root [38]. Assume $D^{-1}E v^{(k)} \approx r v^{(k)}$ for $r < 1$, then

$$(2.15) \quad \langle C \rangle v^{(k)} = D(I - D^{-1}E)v^{(k)} \approx (1 - r)Dv^{(k)} > 0 .$$

Therefore, starting with $v := D^{-1}e$, we perform a few power iterations to find a positive vector v satisfying $\langle C \rangle v > 0$. The computational cost for one iteration is one matrix-vector multiplication plus $\mathcal{O}(n)$ operations. The resulting vector v is required to be positive, which is not true if a row of E is identically zero. This is taken care of by adding some small positive constant ϵ to the iterates. The Matlab code needs directed rounding, therefore it is deferred to Section 3, see Algorithm 3.2 (`MVector`).

The vector u determines the size of w in (2.6), and hence the quality of the upper bounds of $\langle C \rangle^{-1}$ used in (2.11) and also in (2.3). Therefore we add in practice another one or two iterations if the vector u is too small in magnitude. It adds few operations, but may improve the bounds significantly. To the optimality of the choice of v consider the following lemma.

LEMMA 2.2. *Let a matrix $C \in \mathbb{C}^{n \times n}$ be given, and denote by $N \in \mathbb{R}^{n \times n}$ an arbitrary irreducible non-negative matrix. If $C_{\nu\nu} = 0$ for some $1 \leq \nu \leq n$, then C is not an H -matrix. Suppose $C_{\nu\nu} \neq 0$ for all $1 \leq \nu \leq n$. Then for small enough positive ε the following is true.*

For the splitting $D - E := \langle C \rangle - \varepsilon N$ into diagonal and off-diagonal part, the matrix $M := D^{-1}E$ has a (unique) positive eigenvector v to its spectral radius, and C is an H -matrix if and only if $\langle C \rangle v$ is positive. If the off-diagonal part of $\langle C \rangle$ is irreducible, then the above is true for $\varepsilon = 0$. In that case, for any positive vector w , $\langle C \rangle w > 0$ implies

$$(2.16) \quad \min_{1 \leq i \leq n} \frac{[\langle C \rangle w]_i}{w_i} < \min_{1 \leq i \leq n} \frac{[\langle C \rangle v]_i}{v_i} .$$

PROOF. If C is an H -matrix, then $\langle C \rangle w > 0$ for some positive w , and therefore $C_{\nu\nu} \neq 0$ for all $1 \leq \nu \leq n$. For small enough positive ε , the matrix C is an H -matrix if and only $\langle C \rangle - \varepsilon N$ is an M -matrix. Assume such an ε be given and $C_{\nu\nu} \neq 0$ for all $1 \leq \nu \leq n$, so that D is invertible.

The off-diagonal part of $\langle C \rangle$ is non-positive by definition, so that E is non-negative irreducible, and so is $M = D^{-1}E$. By Perron-Frobenius theory [38] there is a unique positive eigenvector v of M corresponding to its spectral radius. If $\langle C \rangle v > 0$, then C is an H -matrix. If conversely C is an H -matrix, then $D - E = \langle C \rangle - \varepsilon N$ is an M -matrix and $\lambda := \rho(D^{-1}E) < 1$, so that

$$(2.17) \quad \langle C \rangle v = D(I - D^{-1}E)v + \varepsilon Nv = D(1 - \lambda)v + \varepsilon Nv > 0 .$$

If E is irreducible, then the assertions are valid for $\varepsilon = 0$. In that case Perron-Frobenius Theory implies

$$\min_{1 \leq i \leq n} \frac{[D^{-1}Ew]_i}{w_i} < \lambda < \max_{1 \leq i \leq n} \frac{[D^{-1}Ew]_i}{w_i}$$

for any positive vector $w \neq v$. Hence there an index k with $[D^{-1}Ew]_k > \lambda w_k$. Using $\langle C \rangle = D(I - D^{-1}E)$ and (2.17) with $\varepsilon = 0$ we conclude

$$\frac{[\langle C \rangle w]_k}{w_k} < \frac{[D(1 - \lambda)w]_k}{w_k} = (1 - \lambda)D_{kk} = \frac{[\langle C \rangle v]_k}{v_k} = \frac{[\langle C \rangle v]_i}{v_i}$$

with the last equality being valid for all $i \in \{1, \dots, n\}$. The lemma is proved. \square

In practice, a choice for N is the matrix of all 1's. The ε -perturbation may be avoided by using the usual splitting $D - E := \langle C \rangle$ and applying the theorem to the irreducible normal form of $D^{-1}E$. For our purposes, however, Algorithm 3.2 (MVector) is sufficient.

3. Rigorous error bounds using directed rounding. In Part I of this paper we discussed the computation of rigorous bounds using solely floating-point arithmetic in rounding to nearest. Special care was necessary to cover all rounding errors, in particular for results in the underflow range.

This is simpler when using directed rounding. Denote by \mathbb{F} a set of p -bit binary floating-point numbers including infinity and NaN, i.e.

$$(3.1) \quad \mathbb{F} = \{ M \cdot 2^{e-p+1} \mid M, e \in \mathbb{Z}, |M| \leq 2^p - 1, e_{min} \leq e \leq e_{max} \} \cup \{-\infty, +\infty, \text{NaN}\} .$$

For details see [15] and Part I of this paper. We consider three different ways to round real numbers into floating-point numbers: Rounding to nearest fl_{\square} , rounding to downwards fl_{∇} and rounding to upwards fl_{Δ} . For $x \in \mathbb{R}$ they are defined by

$$(3.2) \quad \begin{aligned} |\text{fl}_{\square}(x) - x| &= \min\{|f - x| : f \in \mathbb{F}\} , \\ \text{fl}_{\nabla}(x) &:= \max\{f \leq x : f \in \mathbb{F}\} , \\ \text{fl}_{\Delta}(x) &:= \min\{f \geq x : f \in \mathbb{F}\} , \end{aligned}$$

where the ambiguity in the definition of fl_\square is resolved by rounding ties to even. For $a, b \in \mathbb{F}$ and $\circ \in \{+, -, \cdot, /\}$ these rounding functions define floating-point operations $\text{fl}_\square(a \circ b)$, $\text{fl}_\nabla(a \circ b)$ and $\text{fl}_\Delta(a \circ b)$, respectively, mapping $\mathbb{F} \times \mathbb{F}$ into \mathbb{F} . All those operations are mandatory in the IEEE 754 floating-point standard [7, 8], and most computers comply with this standard. For $a, b \in \mathbb{F}$ it follows

$$(3.3) \quad \text{fl}_\nabla(a \circ b) \leq a \circ b \leq \text{fl}_\Delta(a \circ b) \quad \text{and} \quad \text{fl}_\nabla(a \circ b) = \text{fl}_\Delta(a \circ b) \Leftrightarrow a \circ b \in \mathbb{F}.$$

The default in Matlab is rounding to nearest. Fortunately, directed rounding is accessible in Matlab through INTLAB [36], the Matlab toolbox for reliable computing. As has been mentioned, the toolbox INTLAB is entirely written in Matlab and thus portable on many environments. The command `setround(i)` changes the rounding mode according to the following table:

	Matlab command	effect
(3.4)	<code>setround(0)</code>	rounding to nearest
	<code>setround(-1)</code>	rounding to downwards
	<code>setround(1)</code>	rounding to upwards

The `setround` command changes the control word of the processor so that henceforth *all* floating-point operations are executed in the specified rounding mode until the next `setround` command. This implies that for given floating-point numbers $\mathbf{a}, \mathbf{b}, \mathbf{c}$ the code

$$(3.5) \quad \begin{array}{l} \text{setround}(-1) \\ \text{pinf} = \mathbf{a} \cdot \mathbf{b} - \mathbf{c}; \\ \text{setround}(1) \\ \text{psup} = \mathbf{a} \cdot \mathbf{b} - \mathbf{c}; \end{array}$$

computes floating-point numbers `pinf`, `psup` satisfying

$$(3.6) \quad \text{pinf} \leq \mathbf{a} \cdot \mathbf{b} - \mathbf{c} \leq \text{psup},$$

where $\mathbf{a} \cdot \mathbf{b} - \mathbf{c} \in \mathbb{R}$ denotes the true real result. Note that (3.5) always implies (3.6), also in the presence of overflow or underflow (if overflow occurs the bounds are, however, infinite or NaN). With a little thinking it becomes clear that the same code (3.5) is also working for vectors and matrices of compatible size, in which case `pinf`, `psup` are quantities of the corresponding dimension satisfying (3.6). Note that for $\mathbf{c} - \mathbf{a} \cdot \mathbf{b}$ this approach would not work correctly.

Directed rounding is accessible in standard Matlab on many architectures using the undocumented command `feature('setround',rnd)` with values 0.5, $-\infty$ and ∞ of `rnd` for switching the rounding to nearest, downwards and upwards, respectively.

3.1. Executable code implementing Theorem 2.1 using directed rounding. The following Algorithm 3.1 computes an approximate solution `xs` together with an error bound `err` such that $|\mathbf{A}^{-1}\mathbf{b} - \mathbf{xs}| \leq \mathbf{err}$, and proves non-singularity of the input matrix \mathbf{A} . It is based on Theorem 2.1. The exposition would become easier and more readable when using the interval operations as provided by INTLAB [36]. We avoid this dependency for the moment and present the algorithm using only directed rounding, i.e. pure Matlab code. From the next Section 4 on we use interval operations because the algorithms become more involved and would become difficult to read when using only directed rounding.

ALGORITHM 3.1. *Approximate solution `xs` of $Ax = b$ with error bound `err` such that $|\mathbf{A}^{-1}\mathbf{b} - \mathbf{xs}| \leq \mathbf{err}$ is rigorously satisfied (together with the proof of non-singularity of A).*

```

1 function [xs,err] = LssErrBndDirRdg(A,b)
2   setround(0)                                % rounding to nearest
```

```

3  err = repmat(NaN,size(b)); n = size(A,2); % initialization
4  R = inv(A); % approximate inverse of A
5  xs = ResidIter(A,b,R); % improved approximate solution
6  [rinf,rsup] = Dot2DirRdg([A b],[xs;-1],1); % inclusion of A*xs-b
7  setround(+1) % rounding to upwards
8  rmid = rinf + 0.5*(rsup-rinf); % rmid >= (rinf+rsup)/2
9  delta = abs(R)*(rmid-rinf); % delta >= |R|r with r:=(rsup-rinf)/2
10 csup = R*rmid + delta; % upper bound of R(A*xs-b)
11 Csup = R*A; % upper bound of RA
12 setround(-1) % rounding to downwards
13 cinf = R*rmid - delta; % [cinf,csup] inclusion of R(A*xs-b)
14 Cinf = R*A; % [Cinf,Csup] inclusion of RA
15 c = max(abs(cinf),abs(csup)); % upper bound of |R(A*xs-b)|
16 d = diag(Cinf); % lower bound of diag(RA)
17 if any(d<=0), return, end % C not H-matrix, no inclusion
18 E = max(abs(Cinf),abs(Csup)); E(1:n+1:n^2) = 0; % <C>=D-E with D,E>=0, D:=diag(d)
19 [u,v] = MVector(E,d); % vector iteration to verify M-property
20 if any(u<=0) || any(v<=0), return, end % verification of M-property failed
21 setround(1) % rounding to upwards
22 w = max(E./(u*d')); % upper bound u*w >= E*D^-1 (w row)
23 dinv = 1./d; % upper bound of diag(D^-1)
24 err = dinv.*c + v*(w*c); % error bound for xs
25 N = inf; % initialize constant
26 for iter=1:15 % at most 15 improvements of error bound
27   err = min(err, dinv.*(c + E*err) ); % improved error bound
28   Nold = N; N = max(err); % norm of correction
29   if N >= 0.99*Nold, break, end % stop iteration if no improvement
30 end
31 setround(0) % rounding to nearest

```

3.2. Proof of correctness of Algorithm 3.1 (LssErrBndDirRdg). Let an $n \times n$ matrix A and right hand side $\mathbf{b} \in \mathbb{F}^n$ with floating-point entries be given. Algorithm 3.1 is executable Matlab/INTLAB code, so all operations are floating-point operations with floating-point result. In contrast, *all of the operations in the following analysis are the exact real operations*. For example, using R as computed in line 4, $C := R \cdot A \in \mathbb{R}^{n \times n}$ denotes the true matrix product of the (floating-point) matrices R and A . In general, $C \notin \mathbb{F}^{n \times n}$.

After making sure that the rounding is set to nearest in line 2, an approximate inverse R of A and an approximate solution \mathbf{xs} is computed in lines 4 and 5, the latter improved by a residual iteration using Algorithm 4.13 in Part I of this paper. Nothing is required on the quality of R and \mathbf{xs} for the following theoretical analysis. Assume for the moment that the quantities $\mathbf{rinf}, \mathbf{rsup} \in \mathbb{F}^n$ computed in line 6 satisfy

$$(3.7) \quad \mathbf{rinf} \leq A \cdot \mathbf{xs} - \mathbf{b} \leq \mathbf{rsup} .$$

We come to that at the end of this section. We split the computation of \mathbf{delta} in line 9 into

$$\mathbf{rrad} = \mathbf{rmid} - \mathbf{rinf}; \quad \mathbf{delta} = \mathbf{abs}(R) * \mathbf{rrad};$$

both computed in rounding to upwards as specified in line 7. Then using line 8

$$(\mathbf{rinf} + \mathbf{rsup})/2 = \mathbf{rinf} + (\mathbf{rsup} - \mathbf{rinf})/2 \leq \mathbf{rmid} \quad \text{and} \quad \mathbf{rmid} - \mathbf{rinf} \leq \mathbf{rrad} ,$$

so that

$$(3.8) \quad \mathbf{rmid} - \mathbf{rrad} \leq \mathbf{rinf} \quad \text{and} \quad \mathbf{rsup} \leq 2 \cdot \mathbf{rmid} - \mathbf{rinf} \leq \mathbf{rmid} + \mathbf{rrad} .$$

Define $\mu := \mathbf{rmid}$ and $\rho := \mathbf{rrad}$, then (3.7) implies $\mathbf{A} \cdot \mathbf{xs} - \mathbf{b} = \mu + r$ for $|r| \leq \rho$. We use the idea in (3.5) and (3.6) and rounding to upwards to conclude

$$(3.9) \quad \begin{aligned} \mathbf{R} \cdot r &\leq |\mathbf{R} \cdot r| \leq |\mathbf{R}| \cdot \mathbf{rrad} \leq \mathbf{delta} \\ \mathbf{R} \cdot (\mathbf{A} \cdot \mathbf{xs} - \mathbf{b}) &= \mathbf{R} \cdot (\mu + r) \leq \mathbf{R} \cdot \mathbf{rmid} + \mathbf{delta} \leq \mathbf{csup} \\ \mathbf{R} \cdot \mathbf{A} &\leq \mathbf{Csup} . \end{aligned}$$

After switching the rounding to downwards in line 12 we obtain from lines 13 and 14

$$(3.10) \quad \begin{aligned} \mathbf{cinf} &\leq \mathbf{R} \cdot \mathbf{rmid} - \mathbf{delta} \leq \mathbf{R} \cdot (\mu + r) = \mathbf{R} \cdot (\mathbf{A} \cdot \mathbf{xs} - \mathbf{b}) \\ \mathbf{Cinf} &\leq \mathbf{R} \cdot \mathbf{A} . \end{aligned}$$

This nice implicit transformation of infimum-supremum into midpoint-radius is due to Oishi [27]. The operations in lines 15 and 16 cause no rounding error, so that after line 17

$$(3.11) \quad |\mathbf{R} \cdot (\mathbf{A} \cdot \mathbf{xs} - \mathbf{b})| \leq \mathbf{c} \quad \text{and} \quad \text{diag}(\mathbf{R} \cdot \mathbf{A}) \geq \mathbf{d} > 0$$

is satisfied. After extracting the diagonal of \mathbf{E} in line 18 we have

$$(3.12) \quad \langle \mathbf{R} \cdot \mathbf{A} \rangle \geq \text{diag}(\mathbf{d}) - \mathbf{E} \quad \text{with } \mathbf{d} > 0 \text{ and } \mathbf{E} \geq 0 .$$

For the true splitting $\langle \mathbf{R} \cdot \mathbf{A} \rangle = D - E$ this implies

$$(3.13) \quad D \geq \text{diag}(\mathbf{d}) \quad \text{and} \quad E \leq \mathbf{E} ,$$

so that the error bounds in Theorem 2.1 remain valid when replacing D and E by the computed quantities $\text{diag}(\mathbf{d})$ and \mathbf{E} , respectively.

The computation of the vectors \mathbf{u} and \mathbf{v} in line 19 is based on the described Perron iteration for $D^{-1}E$ using the computed quantities \mathbf{d} and \mathbf{E} . The corresponding Algorithm 3.2 (`MVector`) is given after the proof. It computes vectors \mathbf{u} and \mathbf{v} satisfying $(\text{diag}(\mathbf{d}) - \mathbf{E}) \cdot \mathbf{v} \geq \mathbf{u}$, so that $\langle \mathbf{R} \cdot \mathbf{A} \rangle = D - E$ implies

$$(3.14) \quad \langle \mathbf{R} \cdot \mathbf{A} \rangle \cdot \mathbf{v} \geq (\text{diag}(\mathbf{d}) - \mathbf{E}) \cdot \mathbf{v} \geq \mathbf{u} .$$

Thus arriving in line 21 means $\mathbf{u} > 0$ and $\mathbf{v} > 0$, so that $\mathbf{R} \cdot \mathbf{A}$ is an H -matrix, and \mathbf{A} (and \mathbf{R}) are proved to be non-singular. Using (3.13) we see that the (row) vector \mathbf{w} computed in rounding to upwards in line 22 satisfies

$$(3.15) \quad \mathbf{w}_k \geq \max_i \frac{\mathbf{E}_{ik}}{\mathbf{u}_i \mathbf{d}_k} \geq \max_i \frac{(\mathbf{E}D^{-1})_{ik}}{\mathbf{u}_i} \quad \text{for all } 1 \leq k \leq n ,$$

so that

$$(3.16) \quad \mathbf{u}_i \mathbf{w}_k \geq (\mathbf{E}D^{-1})_{ik} \quad \text{for all } 1 \leq i, k \leq n .$$

It follows that \mathbf{err} computed in line 24 in rounding to upwards is an upper bound of the right hand side of (2.12), so that indeed $|\mathbf{A}^{-1}\mathbf{b} - \mathbf{xs}| \leq \mathbf{err}$. The correctness of the final bound follows by (2.13) and the rounding to upwards.

The Perron iteration based on Subsection 2.1 is performed by the following algorithm `MVector`.

ALGORITHM 3.2. *Perron vector iteration to find the optimal vector $\mathbf{v} > 0$ proving $\langle C \rangle = D - E$ to be an M -matrix by $\langle C \rangle \mathbf{v} > 0$. To save memory, D is represented by its diagonal.*

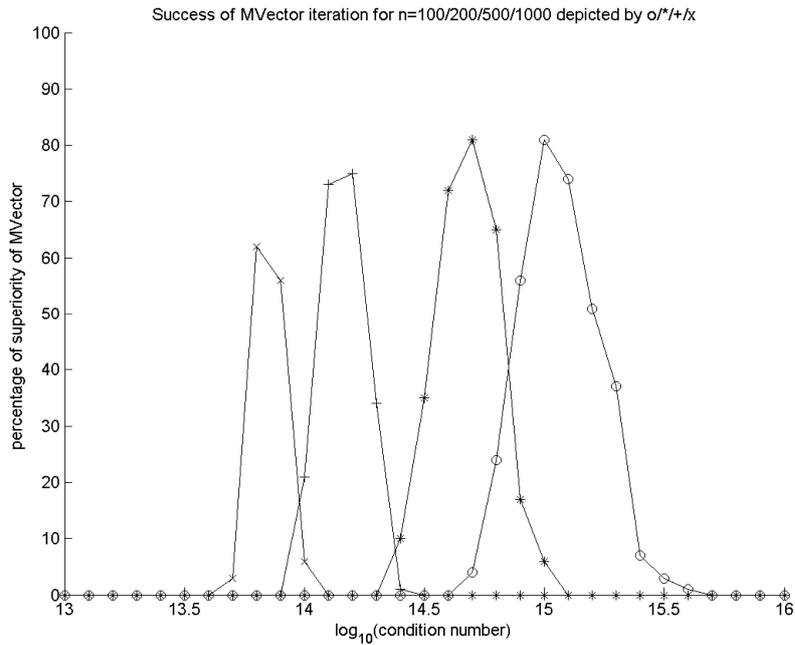


FIG. 3.1. Percentage that $\langle R \cdot A \rangle x > 0$ is satisfied for choosing the Perron vector computed by Algorithm 3.2 (Mvector) for x , but not satisfied for $x := (1, \dots, 1)^T$. Graphs for random matrices of dimensions $n = 100$ (o), $n = 200$ (*), $n = 500$ (+) and $n = 1000$ (x) and condition numbers ranging from 10^{13} to 10^{16} are shown.

```
function [u,v] = MVector(E,d)
% Vector iteration on M=D-E, D,E>=0, D:=diag(d), to verify M-property
setround(-1) % Vector iteration in rounding downwards
vnew = 1./d; % first guess
minu = -inf; % initialization
for iter=1:15 % at most 15 iterations
    v = vnew; minuold = minu; % update
    w = E*(-v);
    u = d.*v + w; % vector iteration, u <= M*v
    minu = min(u); % new minimum of M*v
    vnew = -w./d + eps; % update of guess of H-vector
    r = vnew./v; % ratio new and old Perron vector
    if ( max(r)<1.001*min(r) ) && ( ( minu>0 ) || ( minu<minuold ) )
        return % sufficiently accurate
    end
end
end
```

If successful, the output are two positive vectors u and v with $(D-E) \cdot v \geq u > 0$, so that $D-E$ is an H -matrix provided $E \geq 0$. This is true because the entire iteration is performed in rounding to downwards, so that always $w \leq E \cdot (-v) = -E \cdot v$ and

$$u \leq D \cdot v + w \leq (D - E) \cdot v \leq \langle R \cdot A \rangle v .$$

As has been mentioned, it may happen for ill-conditioned matrices that $\langle R \cdot A \rangle x > 0$ is satisfied for choosing an approximation of the Perron vector of $D^{-1}E$ for x , but is not satisfied for the standard choice $x := (1, \dots, 1)^T$. The percentage of such cases for random matrices of different dimensions and condition numbers is displayed in Figure 3.1.

It remains to specify Algorithm `Dot2DirRdg` and to prove the correctness of (3.7). Using directed rounding it is easy to rewrite Algorithm 3.4 (`Dot2Near`) from Part I:

ALGORITHM 3.3. *Approximation of the matrix product $A \cdot B$ for $A \in \mathbb{F}^{m \times k}$, $B \in \mathbb{F}^{k \times n}$ “as if” accumulated in twice the working precision and rounded into working precision with rigorous error term.*

```

1  function [res,ressup] = Dot2DirRdg(A,B,Incl)
2      setround(0) % rounding to nearest
3      [p,e] = TwoProduct(A(:,1),B(1,:)); % error-free transformation of first product
4      if Incl, einf = e; esup = e; end % (underflow is taken care of later)
5      k = size(A,2); % inner dimension
6      for i=2:k % matrix product by rank-1 updates
7          [h,r] = TwoProduct(A(:,i),B(i,:)); % error-free transformation of i-th product
8          [p,q] = TwoSum(p,h); % error-free transformation of accumulated sum
9          if Incl % rigorous error bounds to be computed
10             setround(-1); einf = einf+(q+r); % lower bound of sum of errors
11             setround(+1); esup = esup+(q+r); % upper bound of sum of errors
12             setround(0) % rounding to nearest
13         else % approximation to be computed
14             e = e + (q + r); % approximate accumulation of errors
15         end
16     end
17     if Incl % output inclusion [res,ressup]
18         setround(+1) % rounding to upwards
19         delta = max(1.5*k*eps,1)*realmin; % cover possible underflow
20         ressup = p + (esup + delta); % upper bound of A*B
21         setround(-1); % rounding to downwards
22         res = p + (einf - delta); % lower bound of A*B
23         setround(0) % rounding to nearest
24     else % output approximation res
25         res = p + e; % extra-precise approximation
26     end

```

We claim that the call `res = Dot2DirRdg(A,B,0)` computes an approximation `res` of the product $A \cdot B$ with extra-precise accumulation of the dot products, and after the call `[res,err] = Dot2DirRdg(A,B,1)` the computed quantities `res` and `ressup` satisfy

$$(3.17) \quad \text{res} \leq A \cdot B \leq \text{ressup} .$$

For `Incl = 0`, the approximation `res` is computed as in Algorithm 3.4 (`Dot2Near`). For the computation of the rigorous error bounds, assume for the moment that no underflow occurs. Then both `TwoProduct` and `TwoSum` in lines 3, 7 and 8 are error-free transformations, see (3.1) in Part I of this paper. If in line 14 there would be no rounding error, then after line ending the loop in line 16 the true sum $p + e$ would be equal to the true product $A \cdot B$. But in lines 10, 11, 20 and 22 the summations are performed with directed rounding, and the result (3.17) follows. Possible underflow is covered in line 19 according to (3.1) in Part I of this paper, so that (3.17) is always true. Note that setting the rounding to nearest in line 12 is necessary to ensure that `TwoProduct` and `TwoSum` are indeed error-free transformations.

We proved the following result.

THEOREM 3.4. *Let `xs` and `err` be the results of Algorithm 3.1 (`LssErrBndDirRdg`) applied to $A \in \mathbb{F}^{n \times n}$ and $\mathbf{b} \in \mathbb{F}^n$. If Algorithm 3.1 ends successfully, i.e. `err` is not a vector of NaN’s, then the input matrix A is*

non-singular and

$$(3.18) \quad |\mathbf{A}^{-1}\mathbf{b} - \mathbf{xs}| \leq \mathbf{err} .$$

3.3. Improvement of Algorithm 3.1. In this subsection we discuss an algorithmic improvement of Algorithm 3.1 which - for simplicity - was not mentioned before. Theorem 2.1, which is implemented by Algorithm 3.1, is based on the splitting $\langle RA \rangle = D - E$ so that error bounds for the matrix product RA are necessary. The standard error estimate [5] for a dot product $x^T y$ for $x, y \in \mathbb{R}^k$ evaluated in floating-point is

$$(3.19) \quad |\text{fl}(x^T y) - x^T y| \leq \gamma_k |x^T y| + k\eta/2$$

using $\gamma_k := k\mathbf{u}/(1 - k\mathbf{u})$, where the extra term $k\eta/2$ covers underflow; in IEEE 754 double precision (binary64) we have $\mathbf{u} = 2^{-53}$ and $\eta = 2^{-1074}$. For $R, A \in \mathbb{F}^{n \times n}$ it follows immediately

$$(3.20) \quad |\text{fl}(RA) - RA| \leq \gamma_n |R||A| + n\eta/2 \cdot ee^T .$$

However, two matrix multiplications are necessary to bound RA , the same as for calculating RA in rounding downwards and upwards. But the application of Theorem 2.1 requires only a lower bound of the diagonal of $|RA|$ and an upper bound of the off-diagonal elements of $|RA|$. More precisely, the off-diagonal elements in E are only needed for the Perron iteration in Algorithm 3.2 (`MVector`) and the rigorous computation of w satisfying (2.6). Both can be covered by the following code requiring only one matrix multiplication:

```

1  setround(-1);           % rounding to downwards
2  Cinf = R*A;            % lower bound Cinf <= RA
3  d = diag(Cinf);        % lower bound of diag(RA)
4  if all(d>0)            % C may be H-matrix
5      setround(1);       % rounding to upwards
6      e = ones(n,1);     % vector (1,...,1)^T
7      g = n*eps/(-(n*eps-1)); % upper bound for gamma_2n
8      Ee = ( abs(Cinf)*e-d ) + ( g*(abs(R)*(abs(A)*e)-d) + max(n^2*eps,1)*realmin*e );

```

The rounding mode in line 1 implies $\text{Cinf} \leq RA$, so that $0 < \mathbf{d} \leq \text{diag}(RA)$ in line 5. Rounding to downwards defines floating-point operations with relative rounding error unit $2\mathbf{u}$ and maximum error η in the underflow range rather than $\eta/2$. Therefore (3.20) and $\langle RA \rangle = D - E$ imply

$$(3.21) \quad E = |E| \leq |\text{Cinf}| - D + \frac{2n\mathbf{u}}{1 - 2n\mathbf{u}} (|R||A| - D) + n\eta \cdot ee^T .$$

Thus observing the definition $\mathbf{eps} := 2^{-52} = 2\mathbf{u}$ in Matlab, $\mathbf{realmin} = \frac{1}{2}\mathbf{u}^{-1}\eta = \mathbf{eps}^{-1}\eta$, the rounding mode and a little thinking proves

$$(3.22) \quad 0 < \mathbf{d} \leq D = \text{diag}(RA) \quad \text{and} \quad Ee \leq \mathbf{Ee} .$$

Up to now only one matrix multiplication $\text{Cinf} = R * A$ was necessary by avoiding to calculate $|R||A|$ explicitly. This is also possible in the further calculations. The only occurrence of E in Algorithm 3.2 (`MVector`) is $\mathbf{w} = \mathbf{E}*(-\mathbf{v})$ in line 8. Replacing this line by $\mathbf{w} = (-\max(\mathbf{v})) * \mathbf{Ee}$, setting $\mu := \max_i(v_i)$, observing $E, v \geq 0$ and the rounding to downwards ensures $\mathbf{w} \leq -\mu \cdot \mathbf{Ee} \leq -E \cdot (\mu e) \leq -E\mathbf{v}$ as desired. Furthermore E is only used in lines 22 and 27 of Algorithm 3.1 (`LssErrBndDirRdg`). Replace line 22 by

```

uinv = 1./u;                % upper bound of u_i^-1
w = max(uinv(:,ones(1,n)).*abs(Cinf)).*dinv' - uinv' ;
w = w + (g*((uinv'*abs(R))*abs(A)).*dinv') + max(n*eps*max(uinv)*max(dinv),1)*realmin);

```

Here $\mathbf{uinv}(:, \mathbf{ones}(1, n))$ is the $n \times n$ -matrix with columns \mathbf{uinv} , so that $\mathbf{uinv}(:, \mathbf{ones}(1, n)) \cdot \mathbf{*abs}(\mathbf{Cinf})$ is an upper bound of $\tilde{D}|\mathbf{Cinf}|$ for \tilde{D} denoting the diagonal matrix with entries \mathbf{u}_ν^{-1} . The maximum in the first computation of \mathbf{w} is the (row) vector of the maxima of columns, so that by $\mathbf{u}_\nu^{-1} \leq \mathbf{uinv}_\nu$ and $D_{\nu\nu}^{-1} \leq \mathbf{dinv}_\nu$ the first computed \mathbf{w} satisfies

$$\mathbf{w}_k \geq \max_i \frac{|\mathbf{Cinf}|_{ik}}{\mathbf{u}_i D_{kk}} - \mathbf{u}_k^{-1} = \max_i \frac{(|\mathbf{Cinf}| - D)_{ik}}{\mathbf{u}_i D_{kk}} = \max_i \frac{((|\mathbf{Cinf}| - D)D^{-1})_{ik}}{\mathbf{u}_i}$$

for $1 \leq k \leq n$. For the error term we have

$$\gamma_{2n} \cdot \max_i \frac{(|R||A| - D)_{ik}}{\mathbf{u}_i D_{kk}} + \max_i \frac{\mathit{neta}}{\mathbf{u}_i D_{kk}} \leq [\mathbf{uinv}^T |R||A|D^{-1}]_k + n(\max_\nu \mathbf{uinv}_\nu)(\max_\nu \mathbf{dinv}_\nu) \cdot \mathit{eta} .$$

Hence, putting things together, (3.21) implies

$$\mathbf{w}_k \geq \max_i \frac{(ED^{-1})_{ik}}{\mathbf{u}_i} \quad \text{for } 1 \leq k \leq n ,$$

so that $(ED^{-1})_{ik} \leq \mathbf{u}_i \mathbf{w}_k$ is satisfied as required in the analysis following (2.6). Finally, replacing $\mathbf{E*err}$ in line 27 by $\mathbf{Ee*max}(\mathbf{err})$, an upper bound of the estimate in (2.13) is computed. This shows that with only one matrix multiplication $\mathbf{Cinf}=\mathbf{R*A}$ verified bounds can be computed.

The quality of the bounds is almost identical to those by the original Algorithm 3.1 (`LssErrBndDirRdg`) provided that Algorithm 3.2 (`MVector`) succeeds to compute vectors $\mathbf{u}, \mathbf{v} \in \mathbb{F}^n$ with $\langle \mathbf{R}\mathbf{A} \rangle \mathbf{v} \geq \mathbf{u} > 0$. This may fail due to the relaxed estimate of the upper bound of E in (3.21). In that case $\mathbf{Csup}=\mathbf{R*A}$ is computed as in Algorithm 3.1, a new Perron vector is computed by `MVector` based on the improved E , and we continue as before.

The advantage is that for not too ill-conditioned problems one matrix multiplication can be saved. In any case, whether the approach succeeds or not, only few additional $\mathcal{O}(n^2)$ operations are necessary. Asymptotically, if successful, the computing time decreases by about one third. In practice, as has been mentioned, the performance of matrix multiplications is better than that of matrix inversion, and also some interpretation overhead is added. For timing results see Figure 8.2.

4. Linear systems with data afflicted with tolerances. Denote by \mathbb{IR} be the set of non-empty closed and bounded intervals

$$(4.1) \quad \mathbf{x} = [x_1, x_2] = \{x \in \mathbb{R} : x_1 \leq x \leq x_2\} .$$

Operations $\circ : \mathbb{IR} \times \mathbb{IR} \rightarrow \mathbb{IR}$ for $\circ \in \{+, -, \cdot, /\}$ are defined by

$$(4.2) \quad \mathbf{x}, \mathbf{y} \in \mathbb{IR} : \quad \mathbf{x} \circ \mathbf{y} := \bigcap \{ \mathbf{z} \in \mathbb{IR} : Z \subseteq \mathbf{z} \} \quad \text{for } Z := \{x \circ y : x \in \mathbf{x}, y \in \mathbf{y}\}$$

provided $0 \notin \mathbf{y}$ in case of division. For these scalar operations always $\mathbf{x} \circ \mathbf{y} = Z$, i.e. there is no difference to the power set operation, and obviously the inclusion property

$$(4.3) \quad x \in \mathbf{x}, y \in \mathbf{y} \quad \Rightarrow \quad x \circ y \in \mathbf{x} \circ \mathbf{y}$$

holds true. Intervals with equal endpoints define the natural embedding of \mathbb{R} into \mathbb{IR} . Using the partial ordering we define the sets \mathbb{IR}^n and $\mathbb{IR}^{m \times n}$ of non-empty and compact interval vectors and matrices similar to (4.1), respectively, with the natural embedding as before. The same definition (4.2) is used for operations between compatible vector and matrix quantities. The sets $\mathbb{IF}, \mathbb{IF}^n$ and $\mathbb{IF}^{m \times n}$ are defined in the same way except that the bounds are floating-point quantities, for example

$$(4.4) \quad \mathbf{A} = [A_1, A_2] = \{A \in \mathbb{R}^{m \times n} : A_1 \leq A \leq A_2\} \in \mathbb{IF}^{m \times n} \quad \text{for } A_1, A_2 \in \mathbb{F}^{m \times n} .$$

Note that although the bounds are floating-point quantities, the intervals comprise of all real quantities within the bounds. Obviously, $\mathbb{IF}^{m \times n} \subseteq \mathbb{IR}^{m \times n}$. Again, operations on \mathbb{IF} , \mathbb{IF}^n and $\mathbb{IF}^{m \times n}$ can be defined as in (4.2), and again the important inclusion property (4.3) is true. For more details see [18].

In view of practical applications it not necessary and usually not wise to insist on narrowest inclusions as in (4.2). In the vast majority of applications (as in this paper) only the inclusion property (4.3) is used. This allows the design of fast interval operations based on BLAS libraries (cf. [35]). Based on that, all interval operations are available in INTLAB [36]. As an example, for matrices $A, B \in \mathbb{F}^{n \times n}$, an inclusion $[C_{\text{inf}}, C_{\text{sup}}]$ of $A \cdot B$ is computed by

```
setround(-1); Cinf = A*B;
setround(+1); Csup = A*B;
```

This inclusion is not best possible, but satisfies the inclusion property (4.3). This approach was used in lines 14 and 11 of Algorithm 3.1 (`LssErrBndDirRdg`). Real quantities may be replaced by an including interval with floating-point endpoints, so that interval operations allow rigorous error bounds for operations with real numbers.

INTLAB makes use of the operator concept in Matlab by defining a data type `intval`. For a floating-point (scalar, vector or matrix) quantity x , the type cast `intval(x)` is a variable of type `intval` with left and right bound x . Moreover, `infsup(xinf, xsup)` checks $x_{\text{inf}} \leq x_{\text{sup}}$ and defines an interval with the left and right bounds x_{inf} and x_{sup} , respectively, whereas `x=midrad(xmid, xrad)` checks $x_{\text{rad}} \geq 0$ and defines an interval x containing $[x_{\text{mid}} - x_{\text{rad}}, x_{\text{mid}} + x_{\text{rad}}]$. In this case x may be a superset of $[x_{\text{mid}} - x_{\text{rad}}, x_{\text{mid}} + x_{\text{rad}}]$ because real intervals are stored in INTLAB by their left and right bound. The type cast and initializations of interval quantities may be applied to vector and matrix quantities x_{inf} , x_{sup} , x_{mid} and x_{rad} as well.

If at least one operand in a Matlab/INTLAB operation is of type `intval`, then the corresponding interval operation is executed, possibly using the natural embedding for the other operand. For example, after line 5 in Algorithm 3.1 (`LssErrBndDirRdg`) the statement `r=R*(A*intval(xs)-b)` would calculate a rigorous inclusion $r \in \mathbb{IF}^n$ of $R \cdot (A \cdot xs - b)$. This is not necessarily true for `rs=intval(R)*(A*xs-b)` because the residual `res=A*xs-b` is calculated in floating-point with $res \in \mathbb{F}^n$. For all INTLAB operators the rounding mode before and after execution is the same.

As a first example we rewrite Algorithm 3.3 (`Dot2DirRdg`) in INTLAB using interval operations:

ALGORITHM 4.1. *Matrix product $A \cdot B$ for $A \in \mathbb{F}^{m \times k}$, $B \in \mathbb{F}^{k \times n}$ "as if" accumulated in twice the working precision with approximate result $C \in \mathbb{F}^{m \times n}$ for $\text{Incl} = 0$, and rigorous inclusion $C \in \mathbb{IF}^{m \times n}$ for $\text{Incl} = 1$.*

```
function C = Dot2(A,B,Incl)
    if nargin==2, Incl=0; end           % no error term if "Incl" not specified
    [p,e] = TwoProduct(A(:,1),B(1,:)); % error-free transformation of first product
    k = size(A,2);                     % inner dimension
    for i=2:k                           % matrix product by rank-1 updates
        [h,r] = TwoProduct(A(:,i),B(i,:)); % error-free transformation of i-th product
        [p,q] = TwoSum(p,h);           % error-free transformation of accumulated sum
        if Incl                         % rigorous error bounds to be computed
            t = q + intval(r);         % inclusion of sum of errors
        else                             % approximation to be computed
            t = q + r;                 % approximation of sum of errors
        end
        e = e + t;                     % approximate accumulation of errors
    end
    if Incl, e=e+midrad(0,max(1.5*k*eps,1)*realmin); end % cover possible underflow
```

```
C = p + e; % extra-precise result
```

In view of `Dot2DirRdg` the algorithm seems self-explaining, only shorter and simpler. Note that only \mathbf{t} , \mathbf{e} and \mathbf{C} are interval quantities for `Incl = 1`; for `Incl = 0` only floating-point quantities occur in the algorithm. Also note that the multiplications $1.5 \cdot k \cdot \text{eps} \cdot \text{realmin}$ do not cause a rounding error and may safely be executed in rounding to nearest.

Next we rewrite Algorithm 3.1 (`LssErrBndDirRdg`) for linear systems the data of which may be afflicted with tolerances. In this case the set of all $A^{-1}b$ for A and b within the tolerances is included together with the proof of non-singularity of all such matrices A .

ALGORITHM 4.2. *Rigorous inclusion $\mathbf{x} \in \mathbb{IF}^n$ of the solution of the linear systems $A\mathbf{x} = \mathbf{b}$. If the matrix and/or right hand side are afflicted with tolerances, i.e. $A \in \mathbb{IF}^{n \times n}$ and/or $\mathbf{b} \in \mathbb{IF}^n$, then the inclusion covers all linear systems within the given data (together with the proof of non-singularity of all matrices in A).*

```
1 function x = LssErrBnd(A,b,Illco)
2   setround(0) % rounding to nearest
3   x = intval(NaN(size(b))); n = size(A,2); % initialization
4   R = inv(mid(A)); % approximate inverse of mid(A)
5   xs = ResidIter(mid(A),mid(b),R); % approximate solution of midpoint system
6   tol = ( any(rad(A(:))>0) ) || ( any(rad(b)>0) ); % check tolerances
7   if tol % input data with tolerances
8     c = mag(R*(A*intval(xs)-b)); % upper bound of |R(A*xs-b)|
9   else % extra-precise accumulation
10    c = mag(R*Dot2([A b],[xs;-1],1)); % upper bound of |R(A*xs-b)|
11  end
12  C = R*intval(A); % inclusion of RA
13  d = diag(C.inf); % lower bound of diag(RA)
14  if all(d>0) % C may be H-matrix
15    E = mag(C); E(1:n+1:n^2) = 0; % <C>=D-E with D,E>=0, D:=diag(d)
16    [u,v] = MVector(E,d); % vector iteration to verify M-property
17    if all(u>0) && all(v>0) % C proved to be H-matrix
18      setround(1) % rounding to upwards
19      w = max(E./(u*d')); % upper bound (u*w)_ik >= (E*D^-1)_ik
20      dinv = 1./d; % upper bound of diag(D^-1)
21      err = dinv.*c + v*(w*c); % error bound for xs
22      N = inf; % initialize constant
23      for iter=1:15 % at most 15 improvements of error bound
24        err = min(err, dinv.*(c + E*err) ); % improved error bound
25        Nold = N; N = max(err); % norm of correction
26        if N >= 0.99*Nold, break, end % stop iteration if no improvement
27        if max(abs(err)./abs(xs))<=eps, break, end % already maximally accurate
28      end
29      x = xs + midrad(0,err); % final inclusion
30      setround(0) % rounding to nearest
31      return
32    end
33  end
34  if ( nargin>2 ) && Illco && (~tol) % solution by phase II
35    x = LssErrBnd(Dot2(R,A,1),Dot2(R,b,1),0);
```

```

36     end
37     setround(0)                                % rounding to nearest

```

In the remainder of this section Algorithm `LssErrBnd` is analyzed and proved to be correct for the call

$$(4.5) \quad \mathbf{x} = \text{LssErrBnd}(\mathbf{A}, \mathbf{b}, \text{Illco}) \quad \text{with Illco} = 0 ,$$

in particular for input data with tolerances. Despite the type cast `intval` the functions `mid`, `rad` and `mag` are used. For a floating-point number $x \in \mathbb{F}$ or interval $\mathbf{x} \in \mathbb{IF}$ they are defined as in Table 4.1.

TABLE 4.1
Definition of mid, rad and mag for non-interval and interval input x.

	$x \in \mathbb{F}$	$\mathbf{x} \in \mathbb{IF}$	
<code>mid(x)</code>	x	$\mathbf{mx} := \text{mid}(\mathbf{x}) \in \mathbb{F}$	} such that $\mathbf{x} \subseteq [\mathbf{mx} - \mathbf{rx}, \mathbf{mx} + \mathbf{rx}]$
<code>rad(x)</code>	0	$\mathbf{rx} := \text{rad}(\mathbf{x}) \in \mathbb{F}$	
<code>mag(x)</code>	$ x $	$\max\{ x : x \in \mathbf{x}\} \in \mathbb{F}$	

Note that the quantities $\mathbf{mx} := \text{mid}(\mathbf{x})$ and $\mathbf{rx} := \text{rad}(\mathbf{x})$ are computed such that $\mathbf{x} \subseteq [\mathbf{mx} - \mathbf{rx}, \mathbf{mx} + \mathbf{rx}]$; they are not necessarily equal to the true midpoint and radius, respectively, because those do not need to be in \mathbb{F} . All definitions apply to vector and matrix quantities entrywise.

First assume $\mathbf{A} \in \mathbb{F}^{n \times n}$ and $\mathbf{b} \in \mathbb{F}^n$ is given. Then mathematically there is no difference to Algorithm 3.1 (`LssErrBndDirRdg`): By definition $\text{mid}(\mathbf{A}) = \mathbf{A}$, $\text{mid}(\mathbf{b}) = \mathbf{b}$, and $\text{tol} = 0$, so that \mathbf{R} and \mathbf{xs} are identical, and \mathbf{c} as computed in line 10 as well. No matter how the inclusion \mathbf{C} of $\mathbf{R} \cdot \mathbf{A}$ in line 12 is computed, all assertions remain valid as long as $\mathbf{R} \cdot \mathbf{A} \subseteq \mathbf{C}$ (in fact, `INTLAB` computes \mathbf{C} in line 10 exactly as in Algorithm `LssErrBndDirRdg`). The rest of Algorithm `LssErrBnd` up to line 33 computes the same quantities as Algorithm `LssErrBndDirRdg`, only the final result is stored in an interval vector $\mathbf{x} \in \mathbb{F}^n$ rather than in an approximation with separate error term. The simple check in line 27 for already achieved maximum accuracy is fast and proves to be useful. Since `Illco` = 0, line 35 is not executed.

Next assume $\mathbf{A} \in \mathbb{IF}^{n \times n}$ or $\mathbf{b} \in \mathbb{IF}^n$, or both. In lines 4 and 5 an approximate inverse of $\text{mid}(\mathbf{A})$ and an approximate solution \mathbf{xs} of the midpoint linear system is computed. Mathematically the concrete values are irrelevant, all assertions are valid for any $\mathbf{R} \in \mathbb{F}^{n \times n}$ and $\mathbf{xs} \in \mathbb{F}^n$. The extra-precise evaluation of the correction $R(\mathbf{A}\tilde{\mathbf{x}} - \mathbf{b})$ does not gain accuracy because interval quantities are involved in \mathbf{A} or \mathbf{b} . This would change if only few entries in \mathbf{A} or \mathbf{b} have nonzero diameter, however, this seems untypical. Therefore an upper bound $\mathbf{c} \in \mathbb{F}^n$ of $|\mathbf{R} \cdot (\mathbf{A} \cdot \mathbf{xs} - \mathbf{b})|$ is computed by standard interval arithmetic in line 8 using `mag` as in Table 4.1.

Now let fixed but arbitrary $A \in \mathbb{R}^{n \times n}$ and $b \in \mathbb{R}^n$ with $A \in \mathbf{A}$ and $b \in \mathbf{b}$ be given (in case \mathbf{A} or \mathbf{b} is a non-interval quantity, we use the natural embedding). Note that A and b are arbitrary real quantities within the interval bounds. Then the inclusion property (4.3) implies $C := \mathbf{R} \cdot A \in \mathbf{C}$ for \mathbf{C} as computed in line 12.

Let $C = D - E$ the splitting of C into diagonal and off-diagonal part. The left, i.e. lower bound of the diagonal of \mathbf{C} is stored in \mathbf{d} , and $\mathbf{d} > 0$ implies $D \geq \text{diag}(\mathbf{d}) > 0$. Moreover, $C \in \mathbf{C}$ implies $|C| \leq \text{mag}(\mathbf{C})$, so that $|E| \leq E$ after execution of line 15. When arriving in line 18, vectors $\mathbf{u}, \mathbf{v} > 0$ have been computed by `MVector(E, d)` in line 16 satisfying $(\text{diag}(\mathbf{d}) - E) \cdot \mathbf{v} \geq \mathbf{u} > 0$. Using $D \geq 0$ this implies

$$\langle C \rangle \cdot \mathbf{v} = (D - |E|) \cdot \mathbf{v} \geq (\text{diag}(\mathbf{d}) - E) \cdot \mathbf{v} > 0 ,$$

so that $\langle C \rangle$ is an M -matrix. Next an upper bound of $\langle \text{diag}(\mathbf{d}) - E \rangle^{-1}$ is computed as in Algorithm 3.1 (`LssErrBndDirRdg`). Hence a standard result from the theory of M -matrices (Theorem 2.5.4 in [6]) implies

$$C^{-1} \leq \langle C \rangle^{-1} \leq \langle \text{diag}(\mathbf{d}) - E \rangle^{-1} ,$$

so that $|\mathbf{R} \cdot (\mathbf{A} \cdot \mathbf{x}\mathbf{s} - \mathbf{b})| \leq \mathbf{c}$ and the rounding to upwards as set in line 18 prove that A is non-singular and that \mathbf{x} is an inclusion of $A^{-1}b$. Since A and b were chosen arbitrarily within \mathbf{A} and \mathbf{b} , this proves the following result.

THEOREM 4.3. *Let a matrix $\mathbf{A} \in \mathbb{F}^{n \times n}$ and a right hand side $\mathbf{b} \in \mathbb{F}^n$ be given, and let $\mathbf{x} \in \mathbb{F}^n$ be the result of Algorithm 4.2 (`LssErrBnd`) computed by the call $\mathbf{x} = \text{LssErrBnd}(\mathbf{A}, \mathbf{b}, 0)$. Then*

$$(4.6) \quad A \text{ is non-singular and } A^{-1}b \in \mathbf{x}$$

for $A = \mathbf{A}$ and $b = \mathbf{b}$. If $\mathbf{A} \in \mathbb{IF}^{n \times n}$ and/or $\mathbf{b} \in \mathbb{IF}^n$, then (4.6) is true for all $A \in \mathbf{A}$ and $b \in \mathbf{b}$. The computational effort for $\mathbf{A} \in \mathbb{F}^{n \times n}$ is $3n^3 + \mathcal{O}(n^2)$ operations, and for $\mathbf{A} \in \mathbb{IF}^{n \times n}$ it is $4n^3 + \mathcal{O}(n^2)$ operations.

5. Alternative approaches to calculate rigorous error bounds for linear systems. Given $A \in \mathbb{R}^{n \times n}$ and $b \in \mathbb{R}^n$, we discuss in this section alternative approaches to prove A to be non-singular and to calculate rigorous error bounds for $A^{-1}b$.

5.1. The Krawczyk-operator. Let $A \in \mathbb{R}^{n \times n}$ and $b \in \mathbb{R}^n$ be given together with some approximate inverse R of A , an approximate solution \tilde{x} , and some interval vector $\mathbf{x} \in \mathbb{IR}^n$. Suppose

$$(5.1) \quad \tilde{x} + R(b - A\tilde{x}) + (I - RA)(\mathbf{x} - \tilde{x}) \subseteq \mathbf{x} ,$$

where for simplicity of exposition all operations are power set operations. Then Krawczyk [12] showed that, if $\|I - RA\| < 1$ for some norm, A is non-singular and $A^{-1}b \in \mathbf{x}$. The left hand side of (5.1) is called the Krawczyk-operator. This has been simplified and improved by Rump [32] as follows. Suppose

$$(5.2) \quad R(b - A\tilde{x}) + (I - RA)\mathbf{x} \subseteq \text{int}(\mathbf{x}) ,$$

where $\text{int}(\cdot)$ denotes the topological interior. Then A is non-singular and $A^{-1}b \in \tilde{x} + \mathbf{x}$. No assumption on a norm of $I - RA$ is necessary. It is often superior to compute an inclusion of the error with respect to the approximation \tilde{x} .

Moreover, Krawczyk did not provide a method how to find a suitable \mathbf{x} except the obvious choice $\tilde{x} \pm \varepsilon$. In [32] we introduced and analyzed the “so-called” epsilon-inflation to construct a suitable \mathbf{x} . The method is implemented as algorithm `verifylss(A, b)` in INTLAB. Although the method is more than 30 years old, it seems to be the state of the art for computing rigorous inclusions for general dense matrices [4].

A successful application of (5.2) requires the necessary condition $\varrho(|I - RA|) < 1$ to hold true. Under reasonable assumptions and using the epsilon-inflation, this condition is also sufficient [32], i.e. the interval-iteration ends successfully. For small enough positive δ and $E = ee^T$ this condition implies $r := \varrho(\delta E + |I - RA|) < 1$, and by Perron-Frobenius Theory the existence of some $0 < u \in \mathbb{R}^n$ with

$$(5.3) \quad (\delta E + |I - RA|)u = ru \quad \text{with } r < 1 .$$

Hence, adapting the proof of Proposition 3.7.2 in [18],

$$(5.4) \quad ru > |I - RA|u \geq (\langle I \rangle - \langle RA \rangle)u = u - \langle RA \rangle u ,$$

so that $\langle RA \rangle u > (1 - r)u > 0$ proves RA to be an H -matrix. Conversely, such as for $RA = -I$, RA may be an H -matrix without $\varrho(|I - RA|) < 1$ being true.

As a consequence, Theorem 2.1 is always applicable whenever the Krawczyk-operator or (5.2) is. If R is scaled so that all diagonal elements of RA are equal to 1, then $D = I$ and $|I - RA| = E$, so that $\varrho(|I - RA|) < 1$ is equivalent to RA being an H -matrix, and (5.2) and Algorithm 4.2 (`LssErrBnd`) have the same scope of applicability.

The mentioned interval-iteration has the same effect as a Perron iteration. Therefore we did not find a case where Theorem 2.1 succeeds to compute an inclusion but (5.2) fails, see the computational results presented in Section 8.

5.2. The theorem by Hansen/Blied/Rohn/Ning/Kearfott/Neumaier. The theoretical background for linear systems with tolerances, which we used implicitly in the previous section, is as follows. The mignitude and magnitude of a (scalar) interval $\mathbf{x} \in \mathbb{IR}$ is defined by

$$(5.5) \quad \text{mig}(\mathbf{x}) := \min\{|x| : x \in \mathbf{x}\} \in \mathbb{R} \quad \text{and} \quad \text{mag}(\mathbf{x}) := \max\{|x| : x \in \mathbf{x}\} \in \mathbb{R} .$$

The definition, like for midpoint and radius, extends to interval vectors and matrices entrywise. For the magnitude we use also the intuitive notation $|\mathbf{x}|$. For an interval matrix $\mathbf{A} \in \mathbb{IR}^{n \times n}$, Ostrowski's comparison matrix $\langle \mathbf{A} \rangle \in \mathbb{R}^{n \times n}$ is, similar to (1.1), defined by [18]

$$(5.6) \quad \langle \mathbf{A} \rangle_{ii} := \text{mig}(\mathbf{A}_{ii}) \quad \text{and} \quad \langle \mathbf{A} \rangle_{ij} := -\text{mag}(\mathbf{A}_{ij}) \text{ for } i \neq j ,$$

If $C := \langle \mathbf{A} \rangle \in \mathbb{R}^{n \times n}$ is an M -matrix, then \mathbf{A} is called H -matrix. In that case every $A \in \mathbf{A}$ is non-singular, and $|A^{-1}| \leq C^{-1}$ (cf. [18]). The hull of a set $X \subseteq \mathbb{R}^n$ is defined to be the smallest interval vector containing X , i.e.

$$(5.7) \quad \text{hull}(X) := \bigcap \{\mathbf{x} \in \mathbb{IR}^n : X \subseteq \mathbf{x}\} .$$

We use interval operations as defined in (4.2). For given $\mathbf{A} \in \mathbb{IR}^{n \times n}$, $\mathbf{b} \in \mathbb{IR}^n$ define

$$(5.8) \quad \Sigma(\mathbf{A}, \mathbf{b}) := \{x \in \mathbb{R}^n : \exists A \in \mathbf{A}, \exists b \in \mathbf{b}, Ax = b\} .$$

Theorem 4.3 shows that Algorithm 4.2 (`LssErrBnd`) computes some \mathbf{x} with $\text{hull}(\Sigma(\mathbf{A}, \mathbf{b})) \subseteq \mathbf{x}$. This inclusion, however, may be an overestimation due to data dependencies. In Section 7 we will discuss how to measure this. Much work has been invested to reduce this overestimation [18]. It was hoped that a clue would be the theorem to be discussed in this subsection: For certain matrices the *exact hull* of the solution set $\Sigma(\mathbf{A}, \mathbf{b})$ can be computed efficiently. In general, the problem is NP-hard.

For given $\mathbf{A} \in \mathbb{IR}^{n \times n}$ one may use, as in Algorithm 4.2 (`LssErrBnd`), an approximate inverse R of the midpoint matrix as a preconditioner. Define $\mathbf{C} := R \cdot \mathbf{A} \in \mathbb{IR}^{n \times n}$ and $C := \langle \mathbf{C} \rangle \in \mathbb{R}^{n \times n}$. If the tolerances in \mathbf{A} are not too wide and the matrices in \mathbf{A} not too ill-conditioned, it is not unlikely that \mathbf{C} is an H -matrix. In that case the following theorem is applicable. As the number of authors suggests, the theorem has a history [3, 1, 30, 24, 19, 20]. We state the version by Neumaier [19, 20], who presented a short and nice proof. Over there also historical remarks appear.

THEOREM 5.1. *Let $\mathbf{C} \in \mathbb{IR}^{n \times n}$ be an H -matrix, $\mathbf{c} \in \mathbb{IR}^n$ a right hand side,*

$$(5.9) \quad r = \langle \mathbf{C} \rangle^{-1} |\mathbf{c}|, \quad \delta_i = (\langle \mathbf{C} \rangle^{-1})_{ii}, \quad \alpha_i = \langle \mathbf{C} \rangle_{ii} - 1/\delta_i, \quad \beta_i = r_i/\delta_i - |\mathbf{c}_i| .$$

Then $\text{hull}(\Sigma(\mathbf{C}, \mathbf{c}))$ is contained in the vector $\mathbf{x} \in \mathbb{IR}^n$ with the components

$$(5.10) \quad \mathbf{x}_i = \frac{\mathbf{c}_i + [-\beta_i, \beta_i]}{\mathbf{C}_{ii} + [-\alpha_i, \alpha_i]} .$$

Moreover, if the midpoint of \mathbf{C} is diagonal, then $\text{hull}(\Sigma(\mathbf{C}, \mathbf{c})) = \mathbf{x}$.

The astonishing part of this theorem is the last sentence: If all off-diagonal entries of \mathbf{C} are centered around zero, then \mathbf{x} is equal to the true hull of the solution set $\text{hull}(\Sigma(\mathbf{A}, \mathbf{b}))$.

For the practical application of Theorem 5.1, the quantities in (5.9) have to be estimated. Following we discuss and improve Neumaier's approach [19]. The critical part is an upper bound for the inverse of $\langle \mathbf{C} \rangle$, where for the diagonal δ a lower bound is required as well.

Since $\langle \mathbf{C} \rangle \in \mathbb{R}^{n \times n}$, we split $\langle \mathbf{C} \rangle = D - E$ into diagonal and off-diagonal part as in (2.4). Then the estimate (2.7) of $\langle \mathbf{C} \rangle^{-1}$ requires a positive vector $v \in \mathbb{R}^n$ with $u = \langle \mathbf{C} \rangle v > 0$. Neumaier [19] proposes to compute an approximate inverse B of the midpoint of $\langle \mathbf{C} \rangle$ and uses the natural choice $v := Be$ with $e := (1, \dots, 1)^T$, for which one can expect that $\langle \mathbf{C} \rangle v$ is not too far from e . Then he computes a nonnegative vector w similar to in (2.6) satisfying $I - \langle \mathbf{C} \rangle B \leq ww^T$. By $uw^T \geq 0$ also $\langle \mathbf{C} \rangle B - I \leq uw^T$, and multiplying from the left by $\langle \mathbf{C} \rangle^{-1} \geq 0$ and using $u = \langle \mathbf{C} \rangle v$ he concludes

$$(5.11) \quad B - vw^T \leq \langle \mathbf{C} \rangle^{-1} \leq B + vw^T .$$

This is used by Neumaier [19, 20] to estimate α and β in (5.10).

Neumaier [21] was not interested in giving a most efficient algorithm but in a nice and simple proof of Theorem 5.1. Choosing B to be an approximate inverse of (the midpoint of) $\langle \mathbf{C} \rangle$ can be interpreted as an extra preconditioning¹ which requires extra matrix multiplications, so that the total computational effort to compute bounds for $A^{-1}b$ is more than doubled.

Another common choice for v is $v = e$ resulting in $B = I$. However this may result in failure to verify $\langle \mathbf{C} \rangle$ to be an M -matrix as has been discussed earlier.

But the choices of B and v are independent: $B \in \mathbb{R}^{n \times n}$ can be chosen arbitrarily, only $u = \langle \mathbf{C} \rangle v > 0$ must be satisfied for some positive v . As has been discussed in Subsection 2.1, the Perron vector computed by Algorithm 3.2 (`MVector`) is the best choice for v , and a simple preconditioning matrix B for $\langle \mathbf{C} \rangle$ suffices.

THEOREM 5.2. *Let $\mathbf{C} \in \mathbb{IR}^{n \times n}$ and $\mathbf{c} \in \mathbb{IR}^n$ together with a positive vector v with $u := \langle \mathbf{C} \rangle v > 0$ be given. Let $D - E = \langle \mathbf{C} \rangle \in \mathbb{R}^{n \times n}$ be the splitting of $\langle \mathbf{C} \rangle$ into diagonal and off-diagonal part, and let $w \in \mathbb{R}^n$ be such that $I - \langle \mathbf{C} \rangle B \leq ww^T$. Then \mathbf{C} is an H -matrix and D is invertible. Abbreviating $d_i := D_{ii} = \langle \mathbf{C} \rangle_{ii}$, the quantities $\alpha, \beta \in \mathbb{R}^n$ defined in Theorem 5.1 satisfy*

$$(5.12) \quad \alpha_i \leq \frac{d_i}{1 + 1/((Ev)w^T)_{ii}} \quad \text{and} \quad \beta_i \leq w^T |\mathbf{c}| \cdot d_i v_i ,$$

and $\text{hull}(\Sigma(\mathbf{C}, \mathbf{c}))$ is contained in the vector $\mathbf{x} \in \mathbb{IR}^n$ defined by (5.10). Moreover,

$$(5.13) \quad \beta_i \leq \frac{[E(D^{-1} + vw^T)|\mathbf{c}]_i - N_i |\mathbf{c}|_i}{1 + N_i} ,$$

where $N_i := [(ED^{-1})^2]_{ii}$.

REMARK 1. Note that only the diagonal of $(ED^{-1})^2$ is needed for N_i , so the computational effort is $\mathcal{O}(n^2)$.

REMARK 2. If $\mathbf{C} := R \cdot \mathbf{A}$ for some given interval matrix $\mathbf{A} \in \mathbb{IR}^{n \times n}$ and an approximate midpoint inverse R as preconditioner, then the off-diagonal elements E in $\langle \mathbf{C} \rangle$ can be expected to be small.

PROOF. First we proceed as in (2.5) and (2.6) to see

$$(5.14) \quad D^{-1} \leq \langle \mathbf{C} \rangle^{-1} \leq D^{-1} + vw^T ,$$

where the lower bound follows by $(D - E)^{-1} = \langle \mathbf{C} \rangle^{-1} = D^{-1}(I - ED^{-1})^{-1}$, the Neumann expansion and $D, E \geq 0$. As has been mentioned in Section 2, the expansion $D^{-1}(I - ED^{-1})^{-1} = D^{-1}(I + ED^{-1}(I - ED^{-1})^{-1})$ yields

$$(5.15) \quad \langle \mathbf{C} \rangle^{-1} \leq D^{-1}(I + E(D^{-1} + vw^T)) .$$

Expanding the Neumann series one more term and observing that the diagonal of E is zero, we obtain

$$(5.16) \quad \text{diag}(D^{-1}(I + (ED^{-1})^2)) \leq \text{diag}(\langle \mathbf{C} \rangle^{-1}) \leq \text{diag}(D^{-1}(I + (Ev)w^T)) ,$$

¹Note that a right inverse is necessary or `B=inv(C)'` in Matlab notation.

and using (5.16) and $\delta_i = (\langle \mathbf{C} \rangle^{-1})_{ii}$ yields

$$(5.17) \quad \alpha_i = d_i - 1/\delta_i \leq d_i \left(1 - \frac{1}{1 + ((Ev)w^T)_{ii}} \right) = \frac{d_i}{1 + 1/((Ev)w^T)_{ii}}.$$

The upper bound for β uses $\delta_i \geq d_i^{-1}$ as by (5.14) to conclude

$$(5.18) \quad \beta_i = (\langle \mathbf{C} \rangle^{-1} |\mathbf{c}|)_i / \delta_i - |\mathbf{c}|_i \leq \frac{d_i^{-1} |\mathbf{c}|_i + (vw^T |\mathbf{c}|)_i}{d_i^{-1}} - |\mathbf{c}|_i,$$

so that (5.12) follows. A little computation using (5.16) shows (5.13) and completes the proof. \square

Unfortunately, according to our numerical evidence, also these improved estimates on the quantities used in Theorem 5.1 do not lead to better bounds than our methods presented in the previous section, see the computational results in Section 8. This seems surprising because the bound given in Theorem 5.1 is provably optimal under reasonable assumptions.

However, the optimality is only true if the entries in \mathbf{C} are independent, which is not true for $\mathbf{C} = R \cdot \mathbf{A}$. There is hardly an efficient remedy for this because the computation of the hull of the solution set $\Sigma(\mathbf{A}, \mathbf{b})$, even of an approximation to a certain degree, is an NP-hard problem (see [28, 13, 29, 2]).

5.3. Ogita's method. Let $A \in \mathbb{R}^{n \times n}$ and $b \in \mathbb{R}^n$ together with an approximate inverse R of A be given. In [26] Ogita et al. assumes $\|I - RA\| < 1$ and uses the standard estimate

$$\|A^{-1}\| = \|(I - (I - RA))^{-1}R\| \leq \frac{\|R\|}{1 - \|I - RA\|}$$

to conclude²

$$(5.19) \quad \|\tilde{x} - A^{-1}b\| = \|\tilde{z} - A^{-1}(A\tilde{z} - r)\| \leq \|\tilde{z}\| + \frac{\|R(A\tilde{z} - r)\|}{1 - \|I - RA\|},$$

where $r = A\tilde{x} - b$ is the residual of some approximate solution \tilde{x} of the original system $Ax = b$, and \tilde{z} is an approximate solution of the residual equation $Az = r$. The authors use $\|\cdot\|_\infty$. Since $\|I - RA\|_\infty = \| |I - RA| e \|$, the scope of applicability could be increased to that of `LssErrBnd` by choosing the Perron vector v instead of e to verify $\|D_v^{-1} |I - RA| D_v\| < 1$.

The residual is computed using some extra-precise dot product accumulation, so there is some similarity to (4.1) in Part I. The idea is that the residual r is small, so that the residual $A\tilde{z} - r$ of the residual equation $Az = r$ is extra small.

The verification in [26] is carried out in rounding to nearest and estimation of rounding errors. In the Matlab code we received by Ogita he uses interval operations and INTLAB, so this method may safely be compared with our algorithms. Computational results are presented in Section 8.

Moreover, a method for inverse monotone matrices, i.e. A is non-singular and $A^{-1} \geq 0$, is presented in [25]. Given a monotone matrix C and an approximate solution \tilde{y} of $Cy = e$ with $e = (1, \dots, 1)^T$, they use

$$\|A^{-1}\|_\infty = \|A^{-1}e\|_\infty = \|A^{-1}(e - Ay) + y\|_\infty \leq \|A^{-1}\|_\infty \|Ay - e\|_\infty + \|y\|_\infty$$

to conclude

$$(5.20) \quad \|A^{-1}\|_\infty \leq \frac{\|y\|_\infty}{1 - \|Ay - e\|_\infty}.$$

If inverse monotonicity is known a priori from some given application, then (5.20) gives a fast estimate of $\|A^{-1}\|_\infty$. This method is also applicable to the matrix $C = RA$ using $C^{-1} \leq \langle C \rangle^{-1}$ after verifying inverse

²In [26] the numerator in the last term in (5.19) was given as $\|R\| \cdot \|A\tilde{z} - r\|$, which can obviously be improved into $\|R(A\tilde{z} - r)\|$; in the numerical tests, (5.19) was used.

monotonicity of $\langle C \rangle$ by $\langle C \rangle v > 0$ for positive v . Note, however, that (2.7) yields a componentwise estimate of C^{-1} in contrast to the normwise estimate (5.20).

Numerical tests suggest that usually (2.7) is better than (5.20), but sometimes worse. However, we are interested in estimating $|C^{-1}c|$ for c denoting the correction $R(A\tilde{x} - b)$. In (2.8) we could estimate the absolute value of this product, whereas using (5.20) only $\|C^{-1}\|_\infty \|c\|_\infty$ is estimated. The latter is always worse than (2.8).

5.4. Nguyen's method. Let $A \in \mathbb{R}^{n \times n}$ together with an approximate inverse R of A be given. Then Nguyen [23, 22] discusses another method based on verifying that RA is an H -matrix. Rather than the obvious choice $v = e := (1, \dots, 1)^T$ for checking $\langle RA \rangle v > 0$ he also uses an iteration on v . His iteration is based on some Jacobi-iteration. In our notation $D - E = C := \langle RA \rangle$ he uses

ALGORITHM 5.3. *Nguyen's vector iteration to find v with $\langle RA \rangle v > 0$.*

```

v = e
for i = 1 : i_max
    v = D-1(e + Ev)
    u = Cv
    if u > 0, break, end if
end for

```

This iteration is almost as good as the Perron iteration presented in Subsection 2.1; however, it happens for ill-conditioned matrices that it fails to verify RA to be an H -matrix, whereas Algorithm 3.2 (**MVector**) based on the Perron iteration succeeds, see Section 8. This confirms the analysis in Subsection 2.1.

Suppose $v > 0$ is given with $u := \langle RA \rangle v > 0$. Then for $b, \tilde{x} \in \mathbb{R}^n$, $r := A\tilde{x} - b$, $c := R(A\tilde{x} - b)$ and $D_u := \text{diag}(u)$, Proposition 4.1.9 in [18] gives

$$(5.21) \quad |\tilde{x} - A^{-1}b| \leq \|D_u^{-1}c\|_\infty v.$$

This is true by observing $C^{-1} = \langle RA \rangle^{-1} \geq 0$ and $|\tilde{x} - A^{-1}b| = |C^{-1}c| \leq C^{-1} \cdot \|D_u^{-1}c\|_\infty u = \|D_u^{-1}c\|_\infty v$. Nguyen uses this to obtain an error bound similar to (2.8). But rather than improving the initial value $\tilde{x} := Rb$ by an a priori residual iteration, he combines this with an inclusion by (5.21) as follows:

```

Compute an inclusion r of  $A\tilde{x} - b$  using extra-precise dot product accumulation
Compute an inclusion of z of  $R\mathbf{r}$ 
Initial error bound  $e := \|D_u^{-1}\mathbf{r}\|_\infty v$ 
While (not converged)
    Apply five Gauss-Seidel iterations
    If accurate enough, break
    Update  $\tilde{x}$  by  $\tilde{x} + e$ 
    Update  $e$  by  $e - \text{mid}(e)$ 
    Recompute r and z
End while

```

Nguyen discusses three basic ideas to speed up the inclusion method. First, the iteration can be stopped if the already achieved inclusion is sufficiently accurate. To decide this he uses two stopping criteria, namely

$$-\log_2 \left(\max_i \frac{\text{rad}(e)}{|\tilde{x}_i|} \right) \geq 52 \quad \text{or} \quad \max_i \frac{|\text{rad}(e'_i) - \text{rad}(e_i)|}{|\tilde{x}_i|} < \mathbf{u},$$

where e' denotes the error bound in the previous iteration. This trick may avoid unnecessary iterations by utilizing the knowledge of an inclusion. Since Gauss-Seidel iterations suffer significantly from interpretation overhead, he uses a Jacobi-iteration instead.

Second, rather than improving an approximate solution \tilde{x} a priori by a (floating-point) residual iteration and then to start the verification process, Nguyen uses an initial approximation and combines the verification process with the interval residual iteration process. If the computed inclusion is not accurate enough, he uses the midpoint of the interval correction \mathbf{z} to improve the approximate solution \tilde{x} . Then he proceeds to compute the next inclusion.

This trick is working fine if the matrix A is not too ill-conditioned; otherwise, however, the inclusion \mathbf{r} of the residual is multiplied by R . Although \mathbf{r} is very narrow due to the extra-precise evaluation, the (small) radii are amplified by the ill-conditioned matrix R resulting in relatively wide intervals of the product \mathbf{z} . Hence the midpoint of \mathbf{z} , which is an inclusion of the correction $R(A\tilde{x} - b)$, is often not as good a correction for \tilde{x} as $R(A\tilde{x} - b)$ computed in floating-point. As a consequence, more residual and/or Gauss/Seidel or Jacobi iterations are necessary, see the computational results in Section 8.

Third, he uses that $\langle RA \rangle$ is close to the identity matrix and proposes a relaxed version, where certain interval operations are replaced by floating-point operations in rounding to upwards. We received the Matlab code from Nguyen using INTLAB for his method, so this method may also safely be compared with our algorithms. Computational results are presented in Section 8.

5.5. Sparse matrices. We briefly mention yet another approach [34] for symmetric positive definite matrices. Let two $n \times n$ matrices $A^T = A$ and G be given and assume $\|A - \alpha I - G^T G\|_2 \leq \tau$ for some positive α and τ . Then for all i ,

$$(5.22) \quad \lambda_i(A) = \alpha + \lambda_i(A - \alpha I) \geq \alpha + \lambda_i(G^T G) - \tau \geq \alpha - \tau ,$$

so that $\alpha > \tau$ implies that A is symmetric positive definite and $\sigma_{\min}(A) \geq \alpha - \tau$. Note that A is not assumed a priori to be positive definite. An error bound for an approximate solution \tilde{x} follows by

$$(5.23) \quad \|\tilde{x} - A^{-1}b\|_2 \leq \|A^{-1}\|_2 \cdot \|A\tilde{x} - b\|_2 \leq (\alpha - \tau)^{-1} \|A\tilde{x} - b\|_2 .$$

Apparently this still the only efficient method to compute rigorous error bounds for sparse matrices; an efficient method for general or even symmetric indefinite matrices is still to be found. We mention this method for completeness; it is designed for sparse matrices and not competitive for full matrices to the methods presented in this paper. On the other hand, the latter use an approximate inverse and are thus not suited for sparse matrices.

6. Inclusion of extremely ill-conditioned linear systems. In Part I of this paper we presented a method to compute verified error bounds for the solution of a linear system $Ax = b$ with extremely ill-conditioned matrix, i.e. $\text{cond}(A) \gtrsim \mathbf{u}^{-1}$. We gave Algorithm 4.20 (`LssIllcoErrBndNear`) in Part I in executable code using only rounding to nearest.

If directed rounding is available and the possibility to compute an inclusion of the solution set $\Sigma(\mathbf{A}, \mathbf{b})$ of a linear system the data of which is afflicted with tolerances, then we may proceed in a much simpler way.

Let $A \in \mathbb{F}^{n \times n}$ and $b \in \mathbb{F}^n$ be given. For $R \in \mathbb{F}^{n \times n}$ let $\mathbf{C} \in \mathbb{IF}^{n \times n}$ and $\mathbf{c} \in \mathbb{IF}^n$ be such that

$$(6.1) \quad RA \in \mathbf{C} \quad \text{and} \quad Rb \in \mathbf{c} .$$

If one of our algorithms applied to \mathbf{C} and \mathbf{c} succeeds to compute an inclusion \mathbf{x} of $\Sigma(\mathbf{C}, \mathbf{c})$, then we know that 1) all matrices $C \in \mathbf{C}$ are non-singular, and 2) that $C^{-1}c \in \mathbf{x}$ is true for all $C \in \mathbf{C}$ and $c \in \mathbf{c}$. This is in particular true for $C := RA \in \mathbf{C}$ and $c := Rb \in \mathbf{c}$, which implies that A (and R) is non-singular, and that

$$(6.2) \quad (RA)^{-1}Rb = A^{-1}b \in \mathbf{x} .$$

Thus we may apply Algorithm 4.2 (`LssErrBnd`) to an interval matrix and vector including RA and Rb , respectively. This is done in line 35 of Algorithm `LssErrBnd` by the call `LssErrBnd(A, b, 1)`.

This approach is not applicable if the original matrix A is already afflicted with tolerances and the midpoint matrix, say, is extremely ill-conditioned. The weighted 2-norm distance to the nearest singular matrix is $\text{cond}(A)$, so that for $\text{cond}(A) \gtrsim \mathbf{u}^{-1}$ it is likely that small tolerances on A produce a singular matrix.

It is important to calculate RA and Rb using extra-precise dot products. Suppose an inclusion \mathbf{C} of RA is computed in working precision. Then standard rounding error analysis tells that the radius of \mathbf{C} is of the order $\mathbf{u}|R||A|$. Even for a perfect approximate inverse R this is of the order $\mathbf{u} \cdot \text{cond}(A)$ or larger, so that for extremely ill-conditioned matrices \mathbf{C} contains singular matrices. We obtain the following result.

THEOREM 6.1. *Let a matrix $\mathbf{A} \in \mathbb{F}^{n \times n}$ and a right hand side $\mathbf{b} \in \mathbb{F}^n$ be given, and let $\mathbf{x} \in \mathbb{F}^n$ be the result of Algorithm 4.2 (`LssErrBnd`) computed by the call $\mathbf{x} = \text{LssErrBnd}(\mathbf{A}, \mathbf{b}, 1)$. Then*

$$(6.3) \quad \mathbf{A} \text{ is non-singular and } \mathbf{A}^{-1}\mathbf{b} \in \mathbf{x} .$$

The computational effort is $7n^3 + \mathcal{O}(n^2)$ operations plus $4n^2$ extra-precise dot product evaluations. Numerical evidence suggests that an inclusion is obtained up to $\text{cond}(\mathbf{A}) \lesssim \mathbf{u}^{-2}$.

7. Inner inclusions for linear systems with tolerances. Consider the interval matrix \mathbf{A} and interval vector \mathbf{b}

$$(7.1) \quad \mathbf{A} = \begin{pmatrix} 2 & [0.25, 1] \\ [1, 2] & [-3, -2.5] \end{pmatrix} \quad \text{and} \quad \mathbf{b} = \begin{pmatrix} [-1, 0.5] \\ [0, 0.5] \end{pmatrix} .$$

Application of Algorithm 4.2 (`LssErrBnd`) computes the inclusion (rounded to two decimals)

$$(7.2) \quad \mathbf{x} = \begin{pmatrix} [-0.56, 0.40] \\ [-0.57, 0.30] \end{pmatrix}$$

together with the proof that each matrix $A \in \mathbf{A}$ is non-singular. The question arises how conservative these error bounds are. An answer to that can be provided by a so-called inner inclusion. For given $\mathbf{A} \in \mathbb{IR}^{n \times n}$ and $\mathbf{b} \in \mathbb{IR}^n$, Algorithm `LssErrBnd` computes $\mathbf{x} \in \mathbb{IF}^n$ with $\Sigma(\mathbf{A}, \mathbf{b}) \subseteq \mathbf{x}$. An interval vector $\mathbf{y} \in \mathbb{R}^n$ is called “inner inclusion” if $\mathbf{y} \subseteq \text{hull}(\Sigma(\mathbf{A}, \mathbf{b}))$. It follows that

$$(7.3) \quad \forall i \in \{1, \dots, n\} \exists A_1, A_2 \in \mathbf{A} \exists b_1, b_2 \in \mathbf{b} : (A_1^{-1}b_1)_i \leq (\inf(\mathbf{y}))_i \quad \text{and} \quad (\sup(\mathbf{y}))_i \leq (A_2^{-1}b_2)_i .$$

This does not imply that for $y \in \mathbf{y}$ there exist $A \in \mathbf{A}$ and $b \in \mathbf{b}$ with $Ay = b$ (see Figure 7.1). Nevertheless the quality of an (outer) inclusion can be judged to be close to optimal if there is not too much difference between this outer and an inner inclusion.

Inner inclusions were introduced by Neumaier [16, 17]. In [33] I showed how to compute inner inclusions with practically no additional computational effort. There is the drawback that the inner inclusion was often empty if not all entries of the matrix and right hand side were intervals of nonzero diameter. However, even if only one entry of the right hand side is an interval of nonzero diameter, the solution set $\Sigma(\mathbf{A}, \mathbf{b})$ has nonzero diameter in all entries except when A^{-1} has certain zero entries. If entries of \mathbf{A} are intervals with nonzero diameter, then this is true except when all matrices A^{-1} for $A \in \mathbf{A}$ have certain zero entries. Both is highly unlikely. An example with only one entry of the matrix afflicted with a tolerance is given at the end of this section, results for only \mathbf{b} being afflicted with tolerances are given in Section 8.

To show how to compute inner inclusions we use two representations of intervals simultaneously. For an interval $\mathbf{x} \in \mathbb{IR}$ we use both the infimum-supremum and midpoint-radius representation

$$(7.4) \quad \mathbf{x} = [x_{inf}, x_{sup}] = \{x \in \mathbb{R} : x_{inf} \leq x \leq x_{sup}\} = \langle mx, rx \rangle = \{x \in \mathbb{R} : |x - mx| \leq rx\} .$$

The same notation applies to interval vectors and matrices using entrywise comparison and absolute value. Note that the notation implies $x_{inf} \leq x_{sup}$ in case of $[x_{inf}, x_{sup}]$, and $rx \geq 0$ in case of $\langle mx, rx \rangle$. Define

$\mathcal{S}_n := \{S \in \mathbb{R}^{n \times n} : |S| \leq I\}$ and $\mathcal{E}_{mn} := \{E \in \mathbb{R}^{m \times n} : |E_{ij}| \leq 1 \text{ for all } i, j\}$. We omit the indices when the dimension of \mathcal{S} and \mathcal{E} is clear from the context. Note that the signature matrices $\text{diag}(\pm 1, \dots, \pm 1)$ are the vertices of \mathcal{S} . Then

$$(7.5) \quad \begin{aligned} \mathbf{x} = \langle mx, rx \rangle \in \mathbb{R}^n &\Leftrightarrow \mathbf{x} = \{mx + S \cdot rx : S \in \mathcal{S}\} = \{mx + E \circ rx : E \in \mathcal{E}_{n1}\} \\ \mathbf{A} = \langle mA, rA \rangle \in \mathbb{R}^{m \times n} &\Leftrightarrow \mathbf{A} = \{mA + E \circ rA : E \in \mathcal{E}_{mn}\}, \end{aligned}$$

where \circ denotes the Kronecker product (entrywise multiplication).

LEMMA 7.1. *Let $\mathbf{y} = \langle my, ry \rangle \in \mathbb{R}^{m \times n}$ be given. Suppose $\mu, \rho \in \mathbb{R}^{m \times n}$ satisfy*

$$(7.6) \quad my - ry \leq \mu - \rho \quad \text{and} \quad \mu + \rho \leq my + ry,$$

where entries ρ_{ij} may be negative. Define the splitting $\rho^+ := (|\rho| + \rho)/2$ and $\rho^- := (|\rho| - \rho)/2$, so that $\rho = \rho^+ - \rho^-$ with $0 \leq \rho^-, \rho^+ \in \mathbb{R}^{m \times n}$ and $\rho^+ \circ \rho^- = 0$. Then there exists $\sigma \in \mathbb{R}^{m \times n}$ with $|\sigma| \leq \rho^-$ and

$$(7.7) \quad \mu - E \circ \rho^+ - \sigma \in \langle my, ry \rangle \quad \text{for all } E \in \mathcal{E}.$$

PROOF. Define $\sigma \in \mathbb{R}^{m \times n}$ entrywise by

$$(7.8) \quad \sigma_{ij} := \begin{cases} \mu_{ij} - (my + ry)_{ij} & \text{if } \rho_{ij}^- > 0 \text{ and } \mu_{ij} > (my + ry)_{ij}, \\ 0 & \text{if } \rho_{ij}^- = 0 \text{ or } (my - ry)_{ij} \leq \mu_{ij} \leq (my + ry)_{ij}, \\ \mu_{ij} - (my - ry)_{ij} & \text{if } \rho_{ij}^- > 0 \text{ and } \mu_{ij} < (my - ry)_{ij}. \end{cases}$$

Let $i \in \{1, \dots, m\}$, $j \in \{1, \dots, n\}$ and $s \in \mathbb{R}$ with $|s| \leq 1$ be fixed but arbitrary. Then (7.6) and $\rho_{ij}^+ \rho_{ij}^- = 0 = \rho_{ij}^+ \sigma_{ij}$ imply

$$(7.9) \quad \mu_{ij} + s \cdot \rho_{ij}^+ - \sigma_{ij} := \begin{cases} s \cdot \rho_{ij}^+ + (my + ry)_{ij} = (my + ry)_{ij} \\ \mu_{ij} + s \cdot \rho_{ij}^+ = \mu_{ij} + s \cdot \rho_{ij} \\ s \cdot \rho_{ij}^+ + (my - ry)_{ij} = (my - ry)_{ij} \end{cases}$$

for the three cases as in (7.8). With (7.6) always $\mu_{ij} + s \cdot \rho_{ij}^+ - \sigma_{ij} \in \langle my_{ij}, ry_{ij} \rangle$, and therefore (7.7). Moreover (7.6) yields

$$|\sigma_{ij}| = \begin{cases} \mu_{ij} - (my + ry)_{ij} \leq -\rho_{ij} = \rho_{ij}^- & \text{if } \rho_{ij}^- > 0 \text{ and } \mu_{ij} > (my + ry)_{ij}, \\ 0 & \text{if } \rho_{ij}^- = 0 \text{ or } (my - ry)_{ij} \leq \mu_{ij} \leq (my + ry)_{ij}, \\ (my - ry)_{ij} - \mu_{ij} \leq -\rho_{ij} = \rho_{ij}^- & \text{if } \rho_{ij}^- > 0 \text{ and } \mu_{ij} < (my - ry)_{ij}. \end{cases}$$

This completes the proof. \square

Let $\mathbf{A} \in \mathbb{R}^{n \times n}$ and $\mathbf{b} \in \mathbb{R}^n$ be given with

$$(7.10) \quad \mathbf{A} = [A_{inf}, A_{sup}] = \langle mA, rA \rangle \quad \text{and} \quad \mathbf{b} = [b_{inf}, b_{sup}] = \langle mb, rb \rangle$$

for suitable $A_{inf}, A_{sup}, mA, rA \in \mathbb{R}^{n \times n}$ and $b_{inf}, b_{sup}, mb, rb \in \mathbb{R}^n$, where $rA, rb \geq 0$. Given $\tilde{x} \in \mathbb{R}^n$ define

$$(7.11) \quad \mathbf{y} = \langle my, ry \rangle := [b_{inf} - mA \cdot \tilde{x} - rA \cdot |\tilde{x}|, b_{sup} - mA \cdot \tilde{x} + rA \cdot |\tilde{x}|] \in \mathbb{R}^n.$$

Let fixed but arbitrary $y \in \mathbf{y}$ be given, and define a signature matrix T such that $T\tilde{x} = |\tilde{x}|$. Using (7.5) implies $y = my + S \cdot ry$ for some $S \in \mathcal{S}$, so that

$$y = my + S \cdot ry = mb - mA \cdot \tilde{x} + S \cdot (rb + rA \cdot |\tilde{x}|) = mb + S \cdot rb - (mA - S \cdot rA \cdot T)\tilde{x}.$$

Hence $mb + S \cdot rb \in \mathbf{b}$ and $mA - S \cdot rA \cdot T \in \mathbf{A}$ yield

$$(7.12) \quad my - ry \leq y \leq my + ry \quad \Leftrightarrow \quad \exists A \in \mathbf{A} \exists b \in \mathbf{b} : y = b - A\tilde{x},$$

where the reverse implication follows by the definition (7.11). This is the basis to compute an inner inclusion by the following code.

ALGORITHM 7.2. *Rigorous inner bounds $\mathbf{yinf}, \mathbf{ysup} \in \mathbb{F}^n$ of the solution of the linear systems $Ax = b$ the data of which are afflicted with tolerances.*

```

1  if tol && ( ( nargin==2 ) || ( ~Illco ) )      % compute inner inclusion
2    mA = -(inf_(A)-sup(A))/2; rA = -(inf_(A)-mA); % inner inclusion of A
3    resinf = inf_(b) + mA*(-xs) + rA*(-abs(xs)); % lower inner bound
4    setround(-1)                                % rounding to downwards
5    ressup = sup(b) + mA*(-xs) + rA*abs(xs);    % upper inner bound
6    mu = resinf + 0.5*(ressup-resinf);          % inner midpoint
7    rho = mu - resinf;                          % inner radius, maybe negative
8    csup = R*mu + abs(R)*rho;                   % upper inner bound correction
9    setround(1)                                  % rounding to upwards
10   cinf = R*mu + abs(R)*(-rho);                % lower inner bound correction
11   e = mag(speye(n)-C)*err;                    % inner correction
12   yinf = xs + cinf + e;                       % inner lower bound
13   ysup = -(e - csup - xs);                   % inner upper bound
14   end

```

THEOREM 7.3. *Let $A \in \mathbb{IF}^{n \times n}$ and $\mathbf{b} \in \mathbb{IF}^n$ be given. Suppose Algorithm 4.2 (LssErrBnd) arrives at line 30 after the call $\mathbf{x} = \text{LssErrBnd}(A, \mathbf{b}, 0)$. Then, as by Theorem 4.3, $\mathbf{x} \in \mathbb{IF}^n$ computed in line 29 is an inclusion of the solution set $\Sigma(A, \mathbf{b})$, and all $A \in \mathbf{A}$ are non-singular. Suppose the code in Algorithm 7.2 is inserted and executed between lines 29 and 30 in Algorithm 4.2 (LssErrBnd). Then for all $i \in \{1, \dots, n\}$ there exist $A_1, A_2 \in \mathbf{A}$ and $b_1, b_2 \in \mathbf{b}$ with*

$$(7.13) \quad (A_1^{-1}b_1)_i \leq \mathbf{yinf}_i \quad \text{and} \quad \mathbf{ysup}_i \leq (A_2^{-1}b_2)_i .$$

The additional computational effort to compute inner inclusions is $17n^2 + \mathcal{O}(n)$ floating-point operations.

REMARK 1. Note that the code in Algorithm 7.2 is only executed if the matrix and/or the right hand side contain data with tolerances, and not the option `Illco=1` for extremely ill-conditioned matrices was chosen. This is because otherwise the solution set has empty interior, or the sensitivity does not permit the computation of an inner inclusion, respectively.

REMARK 2. Also note that (7.13) is valid even if $\mathbf{yinf}_i > \mathbf{ysup}_i$. In that case the inner inclusion is empty for the i -th component of $\Sigma(A, \mathbf{b})$. However, at least some information is provided on the lower and upper bound of any \hat{x}_i with $\hat{x} \in \Sigma(A, \mathbf{b})$. An example is given in (7.30).

PROOF of Theorem 7.3. We refer to the line numbering in Algorithm 7.2, and we use the midpoint-radius representation $\mathbf{A} =: \mathbf{A} = \langle mA, rA \rangle \in \mathbb{IR}^{n \times n}$ and $\mathbf{b} =: \mathbf{b} = \langle mb, rb \rangle \in \mathbb{IR}^n$ as in (7.10). Note that $\mathbf{inf}_-(\mathbf{x})$ ³ and $\mathbf{sup}(\mathbf{x})$ for a (real) interval quantity \mathbf{x} in INTLAB is only accessing the lower and upper bound, respectively, so that $\mathbf{inf}(\mathbf{A}) = A_{inf}$ and so forth without rounding error. We split the computations in line 2 into

$$\begin{aligned} mA1 &= (-\mathbf{inf}_-(\mathbf{A})-\mathbf{sup}(\mathbf{A}))/2; \quad mA = -mA1; \\ rA1 &= \mathbf{inf}_-(\mathbf{A})-mA; \quad rA = -rA1; \end{aligned}$$

The rounding mode in line 2 is still to upwards as set in line 18 in Algorithm 4.2 (LssErrBnd), so that

$$mA1 \geq (-A_{inf} - A_{sup})/2 \quad \text{and} \quad rA1 \geq A_{inf} - mA .$$

³For syntactical reasons, $\mathbf{inf}(\mathbf{x})$ cannot be used in Matlab/INTLAB because `inf` is reserved for the constant ∞ .

Negation is error-free, so that

$$(7.14) \quad \mathbf{mA} \leq (A_{inf} + A_{sup})/2 = mA \quad \text{and} \quad \mathbf{rA} \leq \mathbf{mA} - A_{inf} \leq mA - A_{inf} = rA .$$

Hence

$$(7.15) \quad mA - rA = A_{inf} \leq \mathbf{mA} - \mathbf{rA} \quad \text{and} \quad \mathbf{mA} + \mathbf{rA} \leq mA + rA .$$

Therefore (7.6) is satisfied for $\mathbf{y} := \langle mA, rA \rangle$, $\mu := \mathbf{mA}$ and $\rho := \mathbf{rA}$. Splitting $\mathbf{rA} = \mathbf{rA}^+ - \mathbf{rA}^-$ and using Lemma 7.1 there exists $\Sigma \in \mathbb{R}^{n \times n}$ with $|\Sigma| \leq \mathbf{rA}^-$ and

$$(7.16) \quad \mathbf{mA} - \mathbf{rA}^+ \cdot S - \Sigma \in \langle mA, rA \rangle \quad \text{for all } S \in \mathcal{S} .$$

Setting $\tilde{x} := \mathbf{xs}$ and using $S \in \mathcal{S}$ such that $S\tilde{x} = -|\tilde{x}|$, rounding to upwards and line 3 in Algorithm 7.2 imply

$$(7.17) \quad \mathbf{resinf} \geq b_{inf} + \mathbf{mA} \cdot (-\tilde{x}) + \mathbf{rA} \cdot (-|\tilde{x}|) = b_{inf} - \mathbf{mA} \cdot \tilde{x} + \mathbf{rA}^+ S\tilde{x} + \mathbf{rA}^- |\tilde{x}| .$$

By (7.16) and $\Sigma \leq |\Sigma| \leq \mathbf{rA}^-$ there exist $A_1 \in \mathbf{A}$ and $b_1 \in \mathbf{b}$ with

$$(7.18) \quad \mathbf{resinf} \geq b_{inf} - (\mathbf{mA} - \mathbf{rA}^+ S - \Sigma) \cdot \tilde{x} = b_1 - A_1 \tilde{x} .$$

Similarly, using rounding to downwards, we conclude for $S \in \mathcal{S}$ with $S\tilde{x} = |\tilde{x}|$ and line 5 in Algorithm 7.2 that

$$(7.19) \quad \mathbf{ressup} \leq b_{sup} + \mathbf{mA} \cdot (-\tilde{x}) + \mathbf{rA} \cdot |\tilde{x}| = b_{sup} - \mathbf{mA} \cdot \tilde{x} + \mathbf{rA}^+ S\tilde{x} - \mathbf{rA}^- |\tilde{x}| ,$$

so that using $-\mathbf{rA}^- \leq -|\Sigma| \leq \Sigma$ there exist $A_2 \in \mathbf{A}$ and $b_2 \in \mathbf{b}$ with

$$(7.20) \quad \mathbf{ressup} \leq b_{sup} - (\mathbf{mA} - \mathbf{rA}^+ S - \Sigma) \cdot \tilde{x} = b_2 - A_2 \tilde{x} .$$

Lines 6 and 7 and rounding to downwards imply

$$\begin{aligned} \mathbf{mu} &\leq \mathbf{resinf} + (\mathbf{ressup} - \mathbf{resinf})/2 = (\mathbf{resinf} + \mathbf{ressup})/2 \\ \mathbf{rho} &\leq \mathbf{mu} - \mathbf{resinf} , \end{aligned}$$

so that (7.18) and (7.20) yield

$$(7.21) \quad b_1 - A_1 \tilde{x} \leq \mathbf{resinf} \leq \mathbf{mu} - \mathbf{rho} \quad \text{and} \quad \mathbf{mu} + \mathbf{rho} \leq 2\mathbf{mu} - \mathbf{resinf} \leq \mathbf{ressup} \leq b_2 - A_2 \tilde{x} .$$

Hence for $\mathbf{y} = \langle my, ry \rangle$ defined in (7.11) we use (7.12) to conclude

$$(7.22) \quad my - ry \leq b_1 - A_1 \tilde{x} \leq \mathbf{mu} - \mathbf{rho} \quad \text{and} \quad \mathbf{mu} + \mathbf{rho} \leq b_2 - A_2 \tilde{x} \leq my + ry .$$

Abbreviating $\mu := \mathbf{mu}$ and $\rho := \mathbf{rho}$, Lemma 7.1 implies for the splitting $\rho = \rho^+ - \rho^-$ that there exist some $\sigma \in \mathbb{R}^n$ with $|\sigma| \leq \rho^-$ and $\mu - S \cdot \rho^+ - \sigma \in \langle my, ry \rangle$ for all $S \in \mathcal{S}$, and (7.12) gives

$$(7.23) \quad \forall S \in \mathcal{S} \exists A \in \mathbf{A} \exists b \in \mathbf{b} : b - A\tilde{x} = \mu - S \cdot \rho^+ - \sigma .$$

For the remainder of the proof let $i \in \{1, \dots, n\}$ be fixed but arbitrary. For $R \in \mathbb{R}^{n \times n}$ define signature matrices $|S_1| = |S_2| = I$ such that $(RS_1)_{ik} \geq 0$ and $(RS_2)_{ik} \leq 0$ for all $k \in \{1, \dots, n\}$. By (7.23) there exist $A_1, A_2 \in \mathbf{A}$ and $b_1, b_2 \in \mathbf{b}$ with $b_\nu - A_\nu \tilde{x} = \mu - S_\nu \rho^+ - \sigma$ for $\nu \in \{1, 2\}$. Then

$$(7.24) \quad \begin{aligned} [R(b_1 - A_1 \tilde{x})]_i &= [R\mu - |R|\rho^+ - R\sigma]_i \leq [R\mu - |R|\rho^+ + |R|\rho^-]_i = [R\mu - |R|\rho]_i \quad \text{and} \\ [R(b_2 - A_2 \tilde{x})]_i &= [R\mu + |R|\rho^+ - R\sigma]_i \geq [R\mu + |R|\rho^+ - |R|\rho^-]_i = [R\mu + |R|\rho]_i . \end{aligned}$$

For $R := \mathbf{R}$ the computation of `csup` and `cinf` in lines 8 and 10 and the rounding modes in use imply

$$\mathbf{R} \cdot \mu - |\mathbf{R}| \cdot \rho \leq \text{cinf} \quad \text{and} \quad \text{csup} \leq \mathbf{R} \cdot \mu + |\mathbf{R}| \cdot \rho ,$$

so that by (7.24) and the definition of μ and ρ

$$(7.25) \quad [R(b_1 - A_1 \tilde{x})]_i \leq \text{cinf}_i \quad \text{and} \quad \text{csup}_i \leq [R(b_2 - A_2 \tilde{x})]_i .$$

For fixed but arbitrary $A \in \mathbf{A}$, $b \in \mathbf{b}$ define $C := RA$ and $c := R(A\tilde{x} - b)$, so that, as in (2.1) and (2.2),

$$(7.26) \quad A^{-1}b = \tilde{x} - C^{-1}c = \tilde{x} - c - (I - C)C^{-1}c .$$

For the quantities `xs` and `err` computed before line 29 in Algorithm 4.2 (`LssErrBnd`) we know by Theorem 4.3 that $|A^{-1}b - \mathbf{xs}| \leq \mathbf{err}$ for all $A \in \mathbf{A}$ and $b \in \mathbf{b}$. This is in particular true for $A_1, A_2 \in \mathbf{A}$ and $b_1, b_2 \in \mathbf{b}$ as in (7.25). Moreover, $RA \in \mathbf{C}$ for all $A \in \mathbf{A}$ for \mathbf{C} as computed in line 12 in Algorithm 4.2, so that the computation of `e` in line 11 of Algorithm 7.2 in rounding to upwards and `mag` as defined in Table 4.1 imply $|I - RA||C^{-1}c| \leq \mathbf{e}$. By (7.26) and $\mathbf{xs} = \tilde{x}$ it follows that `yinf` computed in line 12 of Algorithm 7.2 satisfies

$$(7.27) \quad [A_1^{-1}b_1]_i = \tilde{x}_i + [R(b_1 - A_1 \tilde{x})]_i - [(I - RA_1)(\tilde{x} - A_1^{-1}b_1)]_i \leq \mathbf{xs}_i + \text{cinf}_i + \mathbf{e} \leq \text{yinf}_i .$$

Now `ysup` is computed by `z=e-csup-xs` and `ysup=-z` in rounding to upwards. Therefore $\mathbf{z} \geq \mathbf{e} - \text{csup} - \mathbf{xs}$ and $\text{ysup} \leq \mathbf{xs} + \text{csup} - \mathbf{e}$, so that similar to the previous conclusion we obtain

$$(7.28) \quad [A_2^{-1}b_2]_i = \tilde{x}_i + [R(b_2 - A_2 \tilde{x})]_i - [(I - RA_2)(\tilde{x} - A_2^{-1}b_2)]_i \geq \mathbf{xs}_i + \text{csup}_i - \mathbf{e} \geq \text{ysup}_i .$$

This completes the proof because i was chosen fixed but arbitrary in $\{1, \dots, n\}$. □

Coming back to our toy example (7.1), Algorithm 4.2 (`LssErrBnd`) together with Algorithm 7.2 compute outer and inner bounds as follows, rounded to 2 decimals:

$$(7.29) \quad \begin{pmatrix} [-0.35, 0.18] \\ [-0.27, 0.00] \end{pmatrix} \subseteq \text{hull}(\Sigma(\mathbf{A}, \mathbf{b})) \subseteq \begin{pmatrix} [-0.56, 0.40] \\ [-0.57, 0.30] \end{pmatrix}$$

The situation is depicted in Figure 7.1. As can be seen there is some distance between the inner and outer inclusion. In this case the latter is in fact not far from the optimal inclusion:

$$\text{hull}(\Sigma(\mathbf{A}, \mathbf{b})) = \begin{pmatrix} [-0.48, 0.30] \\ [-0.55, 0.19] \end{pmatrix} \quad \text{to 2 decimal places.}$$

In this toy example it is not difficult to compute the (optimal) interval hull of the solution set. Note that in general this is an NP-hard problem, and it is NP-complete to decide that an interval matrix contains a singular matrix [28].

For wider input data, some or all of the entries of the inner inclusion may become empty. If, for instance, \mathbf{A}_{11} in (7.1) is changed from $[2, 2]$ to $[1, 1.75]$, then the second entry of the inner inclusion becomes empty (all data rounded to 2 decimal places):

$$(7.30) \quad \begin{pmatrix} [-0.15, -0.07] \\ \text{---} \end{pmatrix} \subseteq \text{hull}(\Sigma(\mathbf{A}, \mathbf{b})) = \begin{pmatrix} [-0.93, 0.50] \\ [-0.84, 0.34] \end{pmatrix} \subseteq \begin{pmatrix} [-1.2, 0.93] \\ [-0.96, 0.65] \end{pmatrix}$$

Although the inner inclusion for the second component of $\Sigma(\mathbf{A}, \mathbf{b})$ is empty because $\text{yinf}_2 = -0.09$ and $\text{ysup}_2 = -0.21$, at least the information

$$(7.31) \quad \inf(\text{hull}(\Sigma(\mathbf{A}, \mathbf{b}))_2) \leq -0.09 \quad \text{and} \quad -0.21 \leq \sup(\text{hull}(\Sigma(\mathbf{A}, \mathbf{b}))_2)$$

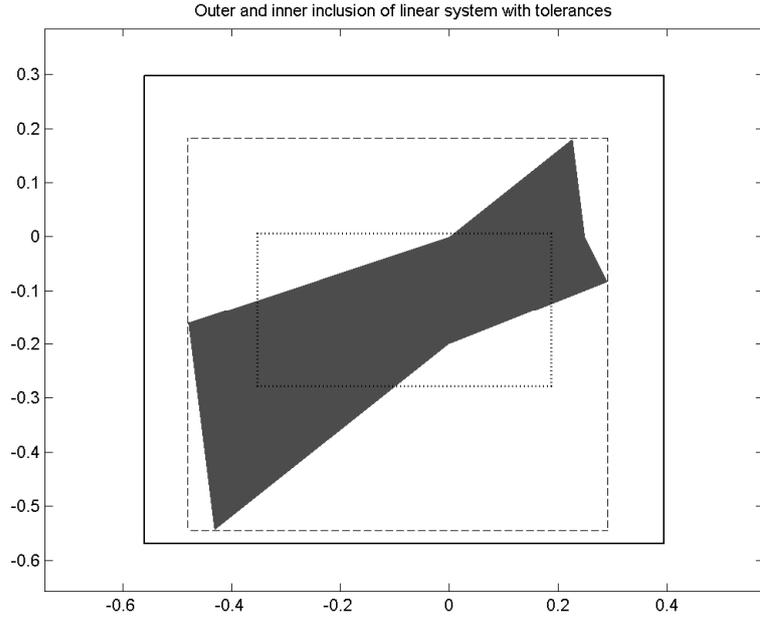


FIG. 7.1. Outer and inner inclusions for \mathbf{A} and \mathbf{b} as in (7.1).

is obtained. Note that nevertheless there is not too much difference between the optimal inclusion and the bounds computed by Algorithm 4.2 (`LssErrBnd`). For narrow input intervals this is typical and can be explained as follows. Up to rounding errors, $[\mathbf{resinf}, \mathbf{ressup}]$ is equal to \mathbf{y} as defined in (7.11), and also up to rounding errors $[\mathbf{cinf}, \mathbf{csup}]$ is equal to the hull of $\{\mathbf{R} \cdot r : r \in [\mathbf{resinf}, \mathbf{ressup}]\}$. Hence the difference of the inner and outer inclusion is determined by the size of \mathbf{e} , which estimates $|(I - C)C^{-1}\mathbf{c}|$ as by (7.26). This in turn is the product of $I - R \cdot \mathbf{A}$ and the maximum difference of $A^{-1}\mathbf{b}$ and $\tilde{\mathbf{x}}$ for $A \in \mathbf{A}$ and $\mathbf{b} \in \mathbf{b}$. In particular it almost vanishes if the input matrix \mathbf{A} has no tolerances at all.

This implies that the inner and outer inclusion are almost identical if there are only tolerances in the right hand side. As an example, we generate random matrices of different dimensions and condition number, and some random right hand side by \mathbf{b} with entries from a normal distribution with mean zero and standard deviation one. Then each entry b_i is afflicted with a tolerance $b_i \pm r_i$, where r_i is taken from a uniform distribution in $[0, 1]$. Note that this implies that the radius of the solution set is roughly of the order of the condition number. The ratio of the median of the radii of the inner and outer inclusion for different dimensions and condition numbers is displayed in Figure 7.2. A ratio close to 1 means that the inner and outer inclusions are almost identical. As can be seen the ratio becomes less than one only for large condition numbers, reflecting the increasing influence of accumulated rounding errors.

As has been mentioned, the previous approaches in [16, 17, 33] basically compute inner inclusions only if all entries in the matrix \mathbf{A} and the right hand side \mathbf{b} are intervals with non-empty interior. Theorem 7.3 is applicable if only one entry in the matrix or the right hand side is a proper interval. A typical example is as follows. We generate a random matrix $A \in \mathbb{R}^{1000 \times 1000}$ with random right hand side and replace some random entry like $A_{529,386}$ by $[-|A_{529,386}|, |A_{529,386}|]$. Then a typical outer and inner inclusion of any solution component may be as follows:

$$[0.5158, 0.5564] \subseteq \text{hull}(\Sigma(\mathbf{A}, \mathbf{b}))_{817} \subseteq [0.5143, 0.5579].$$

In many test examples the result were similar or better, i.e. less difference between the inner and outer inclusion. For smaller uncertainties like replacing $A_{913,209}$ by $A_{913,209} \cdot (1 \pm 0.1)$ instead, there is almost no

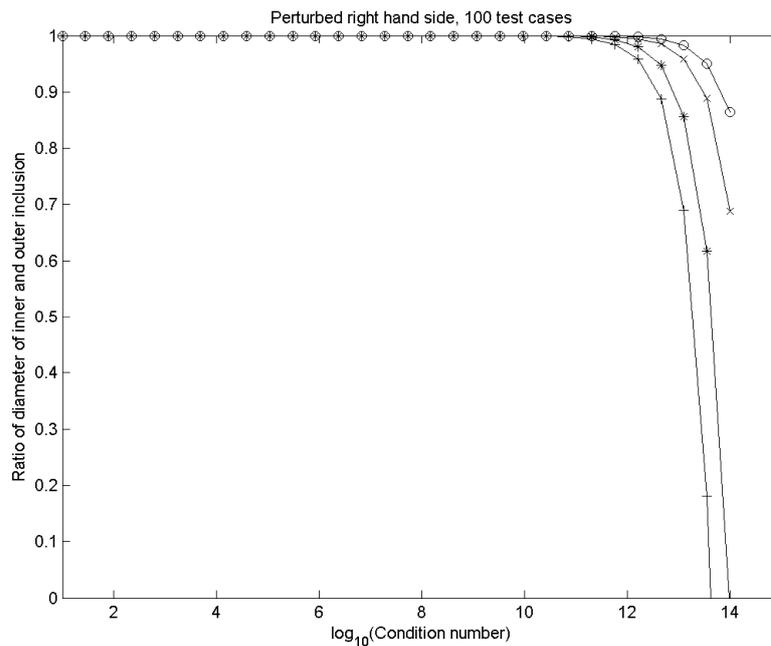


FIG. 7.2. Median ratio of radii of inner and outer inclusion for perturbations in the right hand side, dimensions $n=100$ (\circ), $n=200$ (\times), $n=500$ ($*$) and $n=1000$ ($+$).

difference between the outer and inner inclusion. This allows, for example, to prove

$$\text{hull}(\Sigma(\mathbf{A}, \mathbf{b}))_{308} = [-0.9731, -0.9687]$$

up to four decimal digits. The results for the other components are completely similar.

8. Computational results. Following we report computational results. All algorithms are tested in Matlab version 7.11.0.584 (R2010b) on an Intel Core i7 CPU M640 with 2.8 GHz, INTLAB version 6 and Windows 7 operating system. We discuss the numerical behavior for the following algorithms:

	method
I	Algorithm 4.2 (<code>LssErrBnd</code>) with the improvement described in Subsection 3.3
II	Algorithm 4.16 (<code>LssErrBndNear</code>) from Part I of this paper
(8.1) III	<code>verifylss</code> in INTLAB (see Subsection 5.1)
IV	Nguyen's method (see Subsection 5.4)
V	Ogita's method (see Subsection 5.3)
VI	Theorem 5.2

To allow a fair comparison, all algorithms are written in Matlab/INTLAB, and all algorithm use the same extra-precise dot product accumulation. For the vector residuals all algorithms use `Dot_` as provided in INTLAB rather than `Dot2`. Mathematically `Dot_` is identical to Algorithm 4.1 (`Dot2`) but suffers less from interpretation overhead due to a tricky vectorization of operations. `Dot_` is only suited for vector residuals; for the matrix residual in `LssErrBnd` for extremely ill-conditioned problems, `Dot2` is used.

We use the source code of `LssErrBnd` and `LssErrBndNear` as presented in Part II and Part I of this paper, respectively; only `LssErrBnd` is improved as described in Subsection 3.3. This gains some computing time for not too ill-conditioned problems. Nguyen and Ogita kindly provided the Matlab/INTLAB source code of their routines `certifylss_relaxed` [22] and `vlssr` [26], respectively. The improved version of Neumaier's algorithm based on Hansen's method in Theorem 5.2 using (5.12) and (5.13) we implemented ourselves. Algorithm 3.1 (`LssErrBndDirRdg`) was presented to display an algorithm using only rounding to nearest

TABLE 8.1

Median relative error of the bounds computed by Algorithm 4.2 (**LssErrBnd**) [I], Algorithm 4.16 (**LssErrBndNear**) from Part I of this paper [II], **verifylss** in INTLAB [III], Nguyen’s method [IV], Ogita’s method [V], and Theorem 5.2 [VI] for the matrices from Table 5.1 in Part I and right hand side $\mathbf{b}=\mathbf{A}\cdot\mathbf{randn}(n,1)$.

matrix	n	$\text{cond}(A)$	I	II	III	IV	V	VI
Pascal	17	$2.2 \cdot 10^{16}$	$1.5 \cdot 10^{-16}$	$2.8 \cdot 10^{-16}$	$1.5 \cdot 10^{-16}$	$1.5 \cdot 10^{-16}$	$7.0 \cdot 10^{-16}$	$1.5 \cdot 10^{-16}$
	18	$2.5 \cdot 10^{17}$	$2.0 \cdot 10^{-16}$	-	$2.0 \cdot 10^{-16}$	$2.0 \cdot 10^{-16}$	-	$2.0 \cdot 10^{-16}$
Hilbert	11	$2.5 \cdot 10^{14}$	$1.8 \cdot 10^{-16}$	$1.8 \cdot 10^{-16}$	$1.6 \cdot 10^{-16}$	$1.6 \cdot 10^{-16}$	$4.0 \cdot 10^{-16}$	$1.8 \cdot 10^{-16}$
	12	$8.8 \cdot 10^{15}$	$1.0 \cdot 10^{-14}$	-	$1.8 \cdot 10^{-12}$	$4.5 \cdot 10^{-14}$	-	$4.9 \cdot 10^{-15}$
ScHilbert	11	$3.0 \cdot 10^{14}$	$1.8 \cdot 10^{-16}$	$1.8 \cdot 10^{-16}$	$1.8 \cdot 10^{-16}$	$1.8 \cdot 10^{-16}$	$1.1 \cdot 10^{-15}$	$1.8 \cdot 10^{-16}$
	12	$7.8 \cdot 10^{15}$	$3.2 \cdot 10^{-16}$	-	$2.0 \cdot 10^{-16}$	$3.6 \cdot 10^{-14}$	-	$2.6 \cdot 10^{-16}$
InvHilbert	11	$2.9 \cdot 10^{14}$	$1.4 \cdot 10^{-16}$	$1.4 \cdot 10^{-16}$	$1.4 \cdot 10^{-16}$	$1.4 \cdot 10^{-16}$	$6.3 \cdot 10^{-16}$	$1.4 \cdot 10^{-16}$
	12	$1.0 \cdot 10^{16}$	$2.0 \cdot 10^{-16}$	-	$2.0 \cdot 10^{-16}$	$3.5 \cdot 10^{-13}$	$1.7 \cdot 10^{-13}$	$2.0 \cdot 10^{-16}$
Boothroyd	11	$1.9 \cdot 10^{15}$	$1.5 \cdot 10^{-16}$	$1.9 \cdot 10^{-16}$	$1.5 \cdot 10^{-16}$	$1.5 \cdot 10^{-16}$	$3.5 \cdot 10^{-15}$	$1.5 \cdot 10^{-16}$
	12	$8.5 \cdot 10^{16}$	$2.7 \cdot 10^{-15}$	-	$1.7 \cdot 10^{-16}$	$8.9 \cdot 10^{-14}$	-	$2.2 \cdot 10^{-15}$
Vander	12	$3.0 \cdot 10^{14}$	$1.5 \cdot 10^{-16}$	$1.5 \cdot 10^{-16}$	$1.5 \cdot 10^{-16}$	$1.5 \cdot 10^{-16}$	$2.0 \cdot 10^{-16}$	$1.5 \cdot 10^{-16}$
	13	$1.3 \cdot 10^{16}$	$1.9 \cdot 10^{-16}$	$3.1 \cdot 10^{-16}$	$1.9 \cdot 10^{-16}$	$1.6 \cdot 10^{-16}$	$2.1 \cdot 10^{-16}$	$1.9 \cdot 10^{-16}$

and avoiding interval operations. Since the results are identical to those of **LssErrBnd**, they need not to be discussed separately.

To our knowledge the only other Matlab codes for computing rigorous bounds for the solution of linear systems are by Hargreaves [4] and Rohn [31]. The first author admits that **verifylss** is superior to his routine; the second author provides a library for computing rigorous error bounds for a number of basic problems in linear algebra. However, for square linear systems he uses **verifylss**.

To begin with, we display in Table 8.1 the results for the “usual suspects” of ill-conditioned matrices as listed in Table 5.1 in Part I⁴. Almost always in the following the results for right hand side $\mathbf{b}=\mathbf{randn}(n,1)$ and $\mathbf{b}=\mathbf{A}\cdot\mathbf{randn}(n,1)$ are practically identical, thus we mostly display only one of them. Recall that the Matlab function **rand** generates pseudo-random values drawn from a uniform distribution on the unit interval, whereas **randn** produces pseudo-random values drawn from a normal distribution with mean zero and standard deviation one. In Table 8.1 we use $\mathbf{b}=\mathbf{A}\cdot\mathbf{randn}(n,1)$. The largest dimensions are displayed for which at least one of the Algorithms I-VI succeeds. Here and in the following always the median of the relative error of the inclusions is displayed.

As can be seen, **LssErrBnd** [I] and Theorem 5.2 [VI] are almost always maximally accurate, **verifylss** [III] and Nguyen’s method [IV] are sometimes slightly less accurate. Those four methods have all the same scope of applicability (which is clear for methods [I] and [VI] because they use the same Perron vector) of about $\text{cond}(A) \lesssim \mathbf{u}^{-1}$. **LssErrBndNear** [II] and Ogita’s method [V] fail for some ill-conditioned test matrices. For **LssErrBndNear** this is because all error estimations are performed in rounding to nearest, whereas Ogita’s method could be improved by using a Perron vector to scale $I - RA$. We did not do this but used the original algorithm provided by the author.

For condition numbers up to \mathbf{u}^{-1} , matrices are generated by **randsvd(n,cnd)** from the Matlab matrix gallery. Here some random orthogonal transformation is applied from the left and right to a diagonal matrix with specified singular values. In Figure 8.1 the median of the relative errors of the inclusion computed

⁴For a discussion of those “well-known suspects” of ill-conditioned matrices A and alternative ways to solve $Ax = b$ see the beginning of Section 5 in Part I.

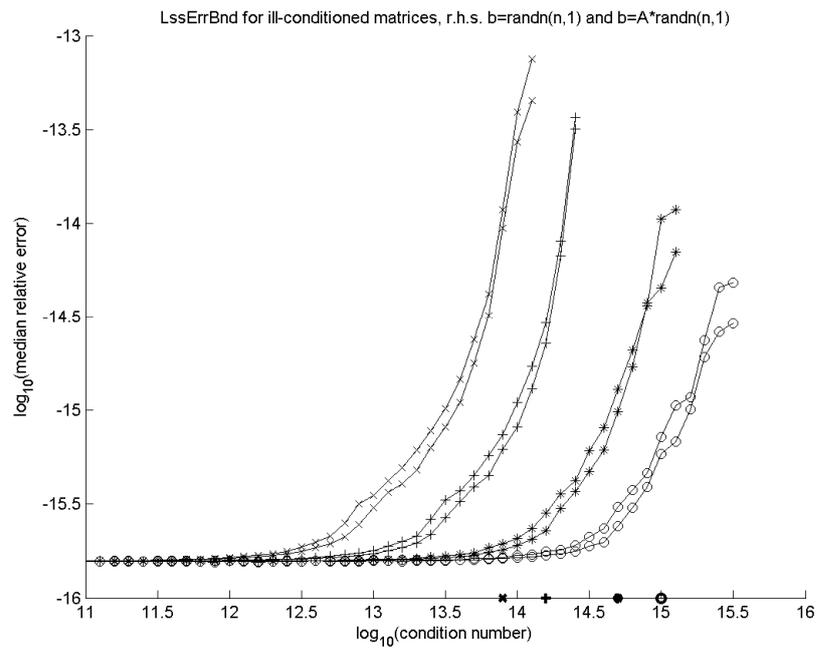


FIG. 8.1. Results of Algorithm 4.2 (**LssErrBnd**) for ill-conditioned random matrices of dimension $n=100$ (\circ), $n=200$ ($*$), $n=500$ ($+$) and $n=1000$ (\times).

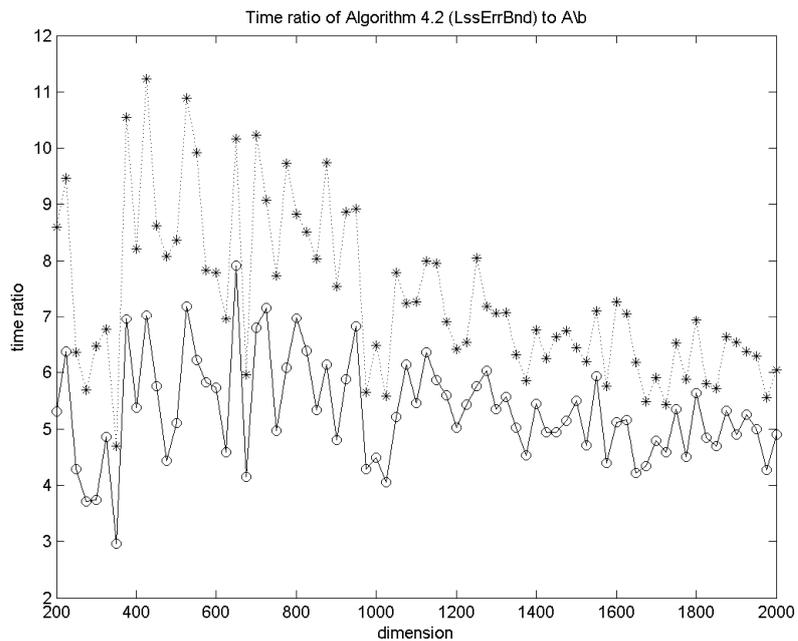


FIG. 8.2. Time ratio between Algorithm 4.2 (**LssErrBnd**) and the standard Matlab call $A \setminus b$, the former with extra-precise evaluation of dot products (dotted line, $*$) and with dot products in working precision (solid line, \circ).

by Algorithm 4.2 (**LssErrBnd**) [I] is displayed for 100 test cases for each dimension, for different condition numbers, and for right hand sides $b=\text{randn}(n,1)$ and $b=A*\text{randn}(n,1)$. Only results for condition number larger than 10^{11} are displayed; for smaller condition numbers the results are always of maximum accuracy. Compared to Figure 5.4 in Part I for **LssErrBndNear** the curves are shifted to the right and show the advantage of using directed rounding.

For $n = 100$, $n = 200$, $n = 500$ and $n = 1000$, **LssErrBnd** did not fail in all test cases for $\text{cond}(A) \leq 1.0 \cdot 10^{15}$,

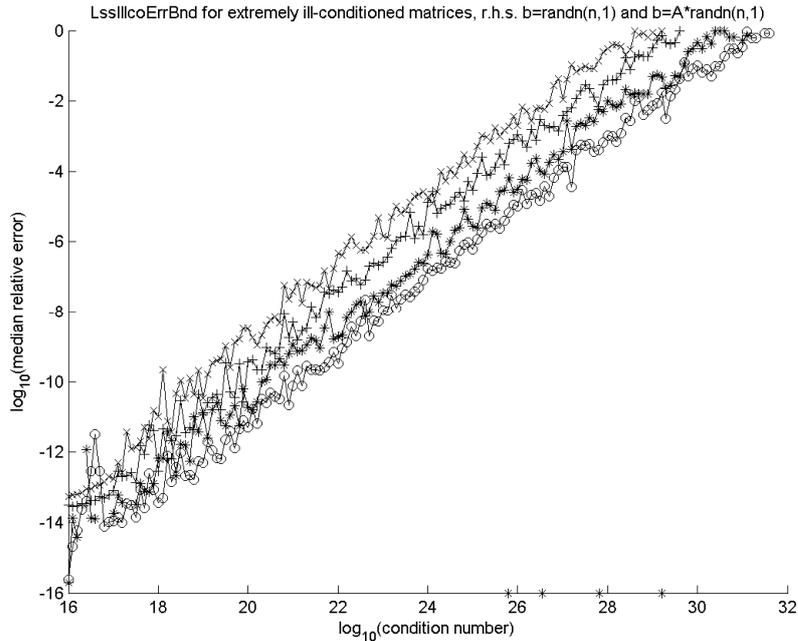


FIG. 8.3. Results of Algorithm 4.2 (**LssErrBnd**) for extremely ill-conditioned random matrices of dimension $n=100$ (o), $n=200$ (*), $n=500$ (+) and $n=1000$ (x).

$1.0 \cdot 10^{14}$, $1.6 \cdot 10^{14}$ and $7.9 \cdot 10^{13}$, which is better than \mathbf{u}^{-1}/n . Recall that the applicability of Algorithm 4.16 (**LssErrBndNear**) in Part I was roughly limited to condition numbers up to \mathbf{u}^{-1}/n^2 .

The computing time of **LssErrBnd** [I] is similar to **LssErrBndNear** [II] (see subsequent tables). In Figure 8.2 the time ratio between Algorithm 4.2 (**LssErrBnd**) and the Matlab command $A \setminus b$ is displayed for 100 random matrices of dimensions from 100 to 2000. Here apples and oranges are compared for two reasons: First, $A \setminus b$ computes an approximation compared to rigorous error bounds by **LssErrBnd**, and second $A \setminus b$ uses plain Gaussian elimination without residual iteration in contrast to **LssErrBnd**. To make it a little more fair, Figure 8.2 shows the ratio for **LssErrBnd** as specified in Algorithm 4.2 with extra-precise residual correction in the upper curve (*), and with only one residual correction in working precision (ensuring backward stability [37]) in the lower curve (o). The latter inclusions are less accurate, as is $A \setminus b$.

With extra-precise residual corrections, the factor in computing time between **LssErrBnd** and $A \setminus b$ is about 6 to 7, whereas without the factor improves to about 5. Note that by Theorem 4.3 the theoretical ratio is 9 for **LssErrBnd** as stated, and 6 with the improvement discussed in Subsection 3.3.

In Figure 8.3 results of Algorithm 4.2 (**LssErrBnd**) for extremely ill-conditioned matrices are shown. We use the methods as described in Part I to generate extremely ill-conditioned matrices. Comparing Figure 5.6 in Part I and Figure 8.3 in this part, an improvement of applicability of **LssErrBnd** over Algorithm 4.16 (**LssI11coErrBndNear**) in Part I, similar to Table 8.1 for small dimension, is observed. The latter succeeds to compute rigorous inclusions for condition numbers up to about \mathbf{u}^{-2}/n^2 , whereas **LssErrBnd** succeeds up to about \mathbf{u}^{-2}/n . Otherwise the graphs in Figure 8.3 are similar to those in Figure 5.6 in Part I but shifted to the right. Again it shows the improvement by using directed rounding.

In our examples, **LssErrBnd** succeeded for dimension $n = 100$, $n = 200$, $n = 500$ and $n = 1000$ in all test cases for $\text{cond}(A) \leq 1.6 \cdot 10^{29}$, $6.6 \cdot 10^{27}$, $6.1 \cdot 10^{25}$ and $3.5 \cdot 10^{26}$, respectively (the value for $n = 500$ is smaller than for $n = 1000$ by the randomness of the examples).

Next we give detailed comparisons between the six methods listed in (8.1). First, for dimension $n = 100$ matrices with different condition numbers are generated 100 each. Then the methods [I] to [VI] are applied

TABLE 8.2

Results of Algorithm 4.2 (`LssErrBnd`) [I], Algorithm 4.16 (`LssErrBndNear`) from Part I of this paper [II], `verifylss` in INTLAB [III], Nguyen's method [IV], Ogita's method [V], and Theorem 5.2 [VI], dimension $n = 100$, r.h.s. `randn(n,1)`, $K = 100$ test cases.

cond(A)	median(median rel. error)						percentage of success					
	I	II/I	III/I	IV/I	V/I	VI/I	I	II	III	IV	V	VI
10^1	$1.6 \cdot 10^{-16}$	1.00	1.00	1.00	12.82	1.00	100	100	100	100	100	100
10^4	$1.6 \cdot 10^{-16}$	1.00	1.00	1.00	1.00	1.00	100	100	100	100	100	100
10^7	$1.6 \cdot 10^{-16}$	1.00	1.00	1.00	1.00	1.00	100	100	100	100	100	100
10^{10}	$1.6 \cdot 10^{-16}$	1.00	1.00	1.00	18.06	1.00	100	100	100	100	100	100
10^{12}	$1.6 \cdot 10^{-16}$	1.01	1.00	1.00	1.64	1.00	100	100	100	100	100	100
10^{13}	$1.6 \cdot 10^{-16}$	1.09	1.00	1.00	15.26	1.00	100	100	100	100	100	100
10^{14}	$1.6 \cdot 10^{-16}$	92.38	1.00	0.96	20.24	1.00	100	33	100	100	100	100
$5 \cdot 10^{14}$	$2.4 \cdot 10^{-16}$	-	0.99	0.69	13.24	1.05	100	0	100	100	92	100
$6 \cdot 10^{14}$	$3.0 \cdot 10^{-16}$	-	1.00	0.57	12.75	1.09	100	0	100	100	84	100
$7 \cdot 10^{14}$	$3.6 \cdot 10^{-16}$	-	1.00	1.00	18.08	1.15	100	0	100	100	69	100
$8 \cdot 10^{14}$	$4.9 \cdot 10^{-16}$	-	1.00	$9 \cdot 10^{16}$	34.33	1.23	100	0	100	100	34	100
$9 \cdot 10^{14}$	$7.0 \cdot 10^{-16}$	-	0.99	$9 \cdot 10^{16}$	19.71	1.38	98	0	98	98	26	98
10^{15}	$1.1 \cdot 10^{-15}$	-	0.98	$9 \cdot 10^{16}$	45.47	1.42	96	0	96	94	21	96
$2 \cdot 10^{15}$	$1.1 \cdot 10^{-14}$	-	0.93	$9 \cdot 10^{16}$	-	1.76	30	0	30	23	0	30
$3 \cdot 10^{15}$	$1.4 \cdot 10^{-14}$	-	0.83	$9 \cdot 10^{16}$	-	1.98	4	0	4	2	0	4
$4 \cdot 10^{15}$	$2.1 \cdot 10^{-15}$	-	0.93	$9 \cdot 10^{16}$	-	1.98	2	0	2	2	0	2

to a linear system with right hand side $\mathbf{b}=\mathbf{randn}(n,1)$. Denote by μ_k the median over the 100 test cases of the median relative errors of the result of algorithm k . Then in Table 8.2 we display μ_1 for `LssErrBnd` [I] in column 2, and in columns 3 to 7 the ratio μ_k/μ_1 for $2 \leq k \leq 6$. Furthermore, in columns 8 to 13 the percentage of success out of the 100 test cases is displayed.

`LssErrBnd` [I] shows a slight decrease of accuracy with increasing condition number, but results guarantee at least 14 correct digits. For larger condition numbers, the inclusions of `verifylss` [III] are a little better (up to 17%), those of Theorem 5.2 [VI] a little weaker, and those by Ogita's method [V] provide about 1 decimal digit less accuracy. Nguyen's method [IV] is better for condition numbers around 10^{14} and significantly weaker for condition numbers up to \mathbf{u}^{-1} . Here the ratio $9 \cdot 10^{16}$ indicates that at least half of the inclusion components are intervals containing 0.

As expected, `LssErrBndNear` fails for condition numbers approaching \mathbf{u}^{-1}/n . The scope of applicability of methods [I] (and of course of [VI]) and [III] is the same, that of Nguyen's and Ogita's a little less. Although in favor of `LssErrBnd`, we do not display computing times for small dimensions because the measurement error is too high.

For $n = 100$ we also display the results for right hand side $\mathbf{A} \cdot \mathbf{randn}(n,1)$ in Table 8.3 because the behavior of `verifylss` is a little different. Now for condition numbers close to \mathbf{u}^{-1} the result is less accurate than `LssErrBnd`, the other results are similar. However, the difference is on a high level: both routines verify some 14 correct figures of the solution. For all accuracy ratios μ_k/μ_1 between 0.5 and 2 the results can be considered to be of similar quality.

Next we come to dimension $n = 1000$ in Table 8.4. Here the results for right hand side `randn(n,1)` are displayed; those for $\mathbf{A} \cdot \mathbf{randn}(n,1)$ are very similar. Now we display the computing time of Algorithms [II]-[VI] relative to that of Algorithm [I] in columns 8 to 12. As expected the maximally treatable condition number decreases compared to $n = 100$. Algorithm `LssErrBnd` verifies at least 13 correct digits of the

TABLE 8.3

Results of Algorithm 4.2 (LssErrBnd) [I], Algorithm 4.16 (LssErrBndNear) from Part I of this paper [II], `verifylss` in INTLAB [III], Nguyen’s method [IV], Ogita’s method [V], and Theorem 5.2 [VI], dimension $n = 100$, r.h.s. `A*randn(n,1)`, $K = 100$ test cases.

cond(A)	median(median rel. error)						percentage of success					
	I	II/I	III/I	IV/I	V/I	VI/I	I	II	III	IV	V	VI
10^1	$1.6 \cdot 10^{-16}$	1.00	1.00	1.00	16.43	1.00	100	100	100	100	100	100
10^4	$1.6 \cdot 10^{-16}$	1.00	1.00	1.00	1.00	1.00	100	100	100	100	100	100
10^7	$1.6 \cdot 10^{-16}$	1.00	1.00	1.00	1.00	1.00	100	100	100	100	100	100
10^{10}	$1.6 \cdot 10^{-16}$	1.00	1.00	1.00	26.34	1.00	100	100	100	100	100	100
10^{12}	$1.6 \cdot 10^{-16}$	1.01	1.00	1.00	5.11	1.00	100	100	100	100	100	100
10^{13}	$1.6 \cdot 10^{-16}$	1.19	1.00	0.99	2.87	1.00	100	100	100	100	100	100
10^{14}	$1.7 \cdot 10^{-16}$	110	1.03	0.94	4.13	1.00	100	30	100	100	100	100
$5 \cdot 10^{14}$	$2.9 \cdot 10^{-16}$	-	1.25	0.60	9.91	1.06	100	0	100	100	96	100
$6 \cdot 10^{14}$	$3.9 \cdot 10^{-16}$	-	1.24	3.04	14.85	1.11	100	0	100	100	81	100
$7 \cdot 10^{14}$	$4.2 \cdot 10^{-16}$	-	1.30	0.42	12.34	1.14	100	0	100	100	65	100
$8 \cdot 10^{14}$	$5.6 \cdot 10^{-16}$	-	1.34	$1.8 \cdot 10^6$	20.01	1.23	100	0	100	100	37	100
$9 \cdot 10^{14}$	$8.2 \cdot 10^{-16}$	-	1.34	$2.3 \cdot 10^{13}$	16.13	1.31	100	0	100	99	31	100
10^{15}	$1.1 \cdot 10^{-15}$	-	1.34	$2.5 \cdot 10^{14}$	22.09	1.44	96	0	96	96	17	96
$2 \cdot 10^{15}$	$6.9 \cdot 10^{-15}$	-	1.33	$1.1 \cdot 10^{15}$	-	1.85	37	0	37	28	0	37
$3 \cdot 10^{15}$	$3.4 \cdot 10^{-15}$	-	1.42	$2.4 \cdot 10^{14}$	-	2.01	6	0	6	5	0	6
$4 \cdot 10^{15}$	$1.1 \cdot 10^{-14}$	-	1.56	-	-	2.08	1	0	1	0	0	1

TABLE 8.4

Results of Algorithm 4.2 (LssErrBnd) [I], Algorithm 4.16 (LssErrBndNear) from Part I of this paper [II], `verifylss` in INTLAB [III], Nguyen’s method [IV], Ogita’s method [V], and Theorem 5.2 [VI], dimension $n = 1000$, r.h.s. `randn(n,1)`, $K = 100$ test cases. An asterisk at the time ratio indicates that not for all 100 test cases inclusions were computed successfully.

cond(A)	median(median rel. error)						relative comp. time				
	I	II/I	III/I	IV/I	V/I	VI/I	II/I	III/I	IV/I	V/I	VI/I
10^1	$1.6 \cdot 10^{-16}$	1.00	1.00	1.00	1.00	1.00	1.21	1.17	1.75	1.57	1.57
10^4	$1.6 \cdot 10^{-16}$	1.00	1.00	1.00	1.00	1.00	1.21	1.17	1.79	1.59	1.58
10^7	$1.6 \cdot 10^{-16}$	1.00	1.00	1.00	1.00	1.00	1.19	1.16	2.25	1.48	1.57
10^{10}	$1.6 \cdot 10^{-16}$	1.01	1.02	0.98	0.98	0.98	1.17	1.05	2.66	1.58	1.42
10^{12}	$1.6 \cdot 10^{-16}$	160	10.5	0.96	1.01	1.00	0.90*	0.88	2.75	1.55	1.22
$4 \cdot 10^{13}$	$1.4 \cdot 10^{-15}$	-	75.1	0.11	1.65	1.29	-	1.00	4.26	2.26	1.21
$5 \cdot 10^{13}$	$2.1 \cdot 10^{-15}$	-	77.4	0.08	2.31	1.50	-	1.04	4.51	2.44*	1.20
$6 \cdot 10^{13}$	$3.1 \cdot 10^{-15}$	-	76.6	0.05	5.17	1.80	-	1.18	4.92	2.50*	1.20
$7 \cdot 10^{13}$	$5.4 \cdot 10^{-15}$	-	73.8	0.05	7.73	2.21	-	1.51	4.96*	2.59*	1.19
$8 \cdot 10^{13}$	$1.0 \cdot 10^{-14}$	-	67.5	$3 \cdot 10^8$	34.5	2.53	-	1.54*	4.13*	2.48*	1.17*
$9 \cdot 10^{13}$	$2.5 \cdot 10^{-14}$	-	61.1	$2 \cdot 10^{14}$	-	2.59	-	1.62*	3.09*	-	1.17*
10^{14}	$5.1 \cdot 10^{-14}$	-	50.4	$3 \cdot 10^{14}$	-	2.45	-	1.58*	2.76*	-	1.17*
$2 \cdot 10^{14}$	-	-	-	-	-	-	-	-	-	-	-

solution, `verifylss` is weaker. Again Theorem 5.2 [VI] is a little weaker, Ogita’s method [V] about one decimal digit less accurate, whereas, similar to dimension $n = 100$, Nguyen’s method [IV] is superior for condition numbers around $4 \cdot 10^{13}$, but computes a number of zero intervals for larger condition numbers.

The performance of the algorithms depends on the computing time spend to improve the accuracy of the final result. It is not surprising that `LssErrBndNear` is sometimes faster than `LssErrBnd` because all computations

TABLE 8.5

Results of Algorithm 4.2 (`LssErrBnd`) [I], Algorithm 4.16 (`LssErrBndNear`) from Part I of this paper [II], `verifylss` in INTLAB [III], Nguyen's method [IV], Ogita's method [V], and Theorem 5.2 [VI], dimension $n = 1000$, r.h.s. `randn(n,1)`, $K = 100$ test cases.

cond(A)	median(median rel. error)						percentage of success					
	I	II/I	III/I	IV/I	V/I	VI/I	I	II	III	IV	V	VI
10^1	$1.6 \cdot 10^{-16}$	1.00	1.00	1.00	1.00	1.00	100	100	100	100	100	100
10^4	$1.6 \cdot 10^{-16}$	1.00	1.00	1.00	1.00	1.00	100	100	100	100	100	100
10^7	$1.6 \cdot 10^{-16}$	1.00	1.00	1.00	1.00	1.00	100	100	100	100	100	100
10^{10}	$1.6 \cdot 10^{-16}$	1.01	1.02	0.98	0.98	0.98	100	100	100	100	100	100
10^{12}	$1.6 \cdot 10^{-16}$	160	10.50	0.96	1.01	1.00	100	91	100	100	100	100
$4 \cdot 10^{13}$	$1.4 \cdot 10^{-15}$	-	75.1	0.11	1.65	1.29	100	0	100	100	100	100
$5 \cdot 10^{13}$	$2.1 \cdot 10^{-15}$	-	77.4	0.08	2.31	1.50	100	0	100	100	97	100
$6 \cdot 10^{13}$	$3.1 \cdot 10^{-15}$	-	76.6	0.05	5.17	1.80	100	0	100	100	55	100
$7 \cdot 10^{13}$	$5.4 \cdot 10^{-15}$	-	73.8	0.05	7.73	2.21	100	0	100	99	9	100
$8 \cdot 10^{13}$	$1.0 \cdot 10^{-14}$	-	67.5	$3 \cdot 10^8$	34.5	2.53	87	0	87	80	5	87
$9 \cdot 10^{13}$	$2.5 \cdot 10^{-14}$	-	61.1	$2 \cdot 10^{14}$	-	2.59	54	0	54	43	0	54
10^{14}	$5.1 \cdot 10^{-14}$	-	50.4	$2 \cdot 10^{14}$	-	2.45	26	0	26	15	0	26
$2 \cdot 10^{14}$	-	-	-	-	-	-	0	0	0	0	0	0

are performed in rounding to nearest and no interpretation overhead for operations with interval data types occur. But this happens only for moderately ill-conditioned examples: For well-conditioned matrices the improvements discussed in Subsection 3.3 accelerate `LssErrBnd` to be faster than `LssErrBndNear`, and for more ill-conditioned examples `LssErrBndNear` fails.

In the design of `LssErrBnd` we compromised between speed and accuracy, where we did not iterate further if an accuracy of about 13 correct decimal digits is achieved. In particular we tried to maintain this compromise over a large range of dimensions and condition numbers. Other methods like Nguyen's [IV] compute better inclusions for certain condition numbers at the price of increasing computing time; for large condition numbers the inclusions are much weaker.

We did not compromise on the applicability of `LssErrBnd` but tried with high priority to compute an inclusion whenever possible. In our examples, the range of applicability of `LssErrBnd` [I], `verifylss` [III] and Theorem 5.2 [VI] is the same, whereas the other methods failed in cases where the afore mentioned succeeded. This is indicated in Table 8.4 with an asterisk at the relative computing times in columns 8 to 12, i.e. the corresponding method was successful in a smaller number of cases than methods [I], [III] and [VI]. The exact percentage of success for the same set of examples as in Table 8.4 is displayed in Table 8.5.

Finally we consider data with tolerances. Only `LssErrBnd` [I], `verifylss` [III] and our implementation of Theorem 5.2 [VI] accept input data with tolerances. Again there is not much difference between a right hand side `randn(n,1)` and `A*randn(n)`, so one of them is chosen in the following. Results for linear systems with only the right hand side afflicted with tolerances were already shown in Section 7, in particular Figure 7.2. Matrices with condition number beyond \mathbf{u}^{-1} the data of which is afflicted with tolerances likely contain a singular matrix. To save space we consider only condition numbers 10^{13} and 10^3 .

We first generate a random matrix A of dimension $n = 1000$ with specified condition number and a right hand side $b = \text{randn}(n, 1)$. Then, using midpoint-radius notation, algorithms [I], [III] and [VI] are applied to the matrix \mathbf{A} and right hand side \mathbf{b} with $\mathbf{A}_{ij} := \langle A_{ij}, r|A_{ij}| \rangle$ and $\mathbf{b}_j := \langle b_j, r|b_j| \rangle$ for $1 \leq i, j \leq n$. It means that all entries in \mathbf{A} and \mathbf{b} are afflicted with a relative tolerance r . It can be expected that \mathbf{A} contains singular matrices for $r \gtrsim \text{cond}(A)^{-1}$.

TABLE 8.6

Results of Algorithm 4.2 (**LssErrBnd**) [I], **verifylss** in INTLAB [III], and Theorem 5.2 [VI], matrix and right hand side with relative radius $r/\text{cond}(A)$, dimension $n = 1000$, r.h.s. **A*randn(n,1)**, condition number 10^{13} , $K = 100$ test cases.

r	median(median radius)		relative comp. time		percentage of success		
	III/I	VI/I	III/I	VI/I	I	III	VI
$5 \cdot 10^{-16}$	1.01	1.03	0.96	1.05	100	100	100
$6 \cdot 10^{-16}$	1.01	1.03	0.96	1.05	100	100	100
$7 \cdot 10^{-16}$	1.01	1.03	0.96	1.04	100	100	100
$8 \cdot 10^{-16}$	1.01	1.04	0.96	1.05	100	100	100
$9 \cdot 10^{-16}$	1.01	1.04	0.96	1.04	100	100	100
10^{-15}	1.01	1.04	0.96	1.05	100	100	100
$2 \cdot 10^{-15}$	1.01	1.08	0.96	1.04	100	100	100
$3 \cdot 10^{-15}$	1.01	1.16	0.99	1.06	100	100	100
$4 \cdot 10^{-15}$	1.00	1.27	1.01	1.04	100	100	100
$5 \cdot 10^{-15}$	0.99	1.41	1.13	1.05	100	100	100
$6 \cdot 10^{-15}$	0.96	1.55	1.24	1.05	100	100	100
$7 \cdot 10^{-15}$	0.92	1.72	1.30	1.06	85	85	85
$8 \cdot 10^{-15}$	0.91	1.80	1.33	1.07	53	53	53
$9 \cdot 10^{-15}$	0.91	1.82	1.32	1.06	16	16	16

TABLE 8.7

Results of Algorithm 4.2 (**LssErrBnd**) [I], **verifylss** in INTLAB [III], and Theorem 5.2 [VI], matrix and right hand side with relative radius $r/\text{cond}(A)$, dimension $n = 1000$, r.h.s. **randn(n,1)**, condition number 10^3 , $K = 100$ test cases.

r	median(median radius)		relative comp. time		percentage of success		
	III/I	VI/I	III/I	VI/I	I	III	VI
$5 \cdot 10^{-6}$	1.01	1.06	1.15	1.08	100	100	100
$6 \cdot 10^{-6}$	1.00	1.08	1.15	1.07	100	100	100
$7 \cdot 10^{-6}$	1.02	1.11	1.19	1.10	100	100	100
$8 \cdot 10^{-6}$	1.01	1.14	1.18	1.07	100	100	100
$9 \cdot 10^{-6}$	1.01	1.18	1.18	1.08	100	100	100
10^{-5}	1.01	1.22	1.21	1.10	100	100	100
$2 \cdot 10^{-5}$	0.96	1.84	1.71	1.08	95	95	95
$3 \cdot 10^{-5}$	-	-	-	-	0	0	0

In Table 8.6 results are shown for 100 test cases each with $\text{cond}(A) = 10^{13}$ and different values of r . For μ_k denoting the median over the 100 test cases of the median relative errors of the result of algorithm k , the ratio μ_k/μ_1 for $k \in \{3, 6\}$ is displayed in columns 2 and 3. Next the ratio of computing times between Algorithm k and **LssErrBnd** is shown, and in the last three columns the percentage of success.

As can be seen the range of applicability is the same for all examples. The quality of **verifylss** [III] is better for larger radii requiring more computing time; for small radius better inclusions are achieved in less time. The inclusions by Theorem 5.2 is always weaker than **LssErrBnd** requiring more computing time. This is in particular interesting because of the optimality of the inclusion if there would be no dependencies in RA as by Theorem 5.1.

Finally we show In Table 8.7 similar results for condition number 10^3 allowing for larger radii. Again the scope of applicability of all three methods Algorithm 4.2 (**LssErrBnd**) [I], **verifylss** in INTLAB [III], and Theorem 5.2 [VI] is the same in the examples. Now **LssErrBnd** delivers almost always the best bounds requiring less computing time than the other methods. The results for other dimensions are similar and

therefore omitted.

Note that both for condition number 10^{13} and 10^3 the maximal relative radius for which inclusions are achieved is not too far from $\text{cond}(A)^{-1}$. Also note that success of an algorithm proves all matrices within \mathbf{A} to be non-singular, but failure does not mean that there exists a singular matrix within the tolerances; it means that the problem is too difficult to solve by the given means. It is possible to prove that an interval matrix contains a singular matrix by using some heuristic to “move” into the right direction within \mathbf{A} to find matrices $A_1, A_2 \in \mathbf{A}$ with $\det(A_1)\det(A_2) \leq 0$. A gap remains where nothing can be said by an efficient algorithm because of the NP-hardness of the problem.

9. Conclusion. In this Part II of the paper directed rounding is used to obtain algorithms with a wider range of applicability. Interval notation increases readability substantially without sacrificing accuracy and/or speed. In particular extremely ill-conditioned linear systems with condition number up to \mathbf{u}^{-2} are easily treated by solving the corresponding preconditioned interval linear system. The only additional tool other than directed rounding is the extra-precise accumulation of dot products, which in turn is implemented using solely floating-point operations in working precision.

A new method to compute so-called “inner inclusion” verifies the quality of the computed bounds in case of input data afflicted with tolerances. It is applicable even if only a single entry in the matrix and/or right hand side is afflicted with a tolerance. If tolerances are not too wide, there is not too much difference between the outer and inner inclusion, indicating that the inclusions are almost optimal.

Acknowledgement. My thanks to Florian Bünger and an anonymous referee for helpful comments.

REFERENCES

- [1] C. Blied. *Computer methods for design automation*. Ph.D. dissertation, Massachusetts Institute of Technology MIT, 1992.
- [2] G.E. Coxson. Computing Exact Bounds on Elements of an Inverse Interval Matrix is NP-Hard. *Reliable Computing*, 5(2):137–142, 1999.
- [3] E.R. Hansen. Bounding the Solution Set of Interval Linear Systems. *SIAM Journal on Numerical Analysis (SINUM)*, 29:1493–1503, 1992.
- [4] G. Hargreaves. Interval Analysis in MATLAB. Master’s thesis, University of Manchester, 2002. <http://www.manchester.ac.uk/mims/eprints>.
- [5] N. J. Higham. *Accuracy and stability of numerical algorithms*. SIAM Publications, Philadelphia, 2nd edition, 2002.
- [6] R.A. Horn and Ch. Johnson. *Topics in Matrix Analysis*. Cambridge University Press, 1991.
- [7] *ANSI/IEEE 754-1985: IEEE Standard for Binary Floating-Point Arithmetic*. New York, 1985.
- [8] *ANSI/IEEE 754-2008: IEEE Standard for Floating-Point Arithmetic*. New York, 2008.
- [9] R.B. Kearfott, M. Dawande, and C. Hu. Intlib: A portable Fortran-77 interval standard function library. *ACM Trans. Math. Software*, 20:447–459, 1994.
- [10] O. Knüppel. PROFIL / BIAS — A Fast Interval Library. *Computing*, 53:277–287, 1994.
- [11] O. Knüppel and T. Simenec. PROFIL/BIAS extensions. Technical Report 93.5, Forschungsschwerpunkt Informations- und Kommunikationstechnik, Technische Universität Hamburg-Harburg, 1993.
- [12] R. Krawczyk. Newton-Algorithmen zur Bestimmung von Nullstellen mit Fehlerschranken. *Computing*, 4:187–201, 1969.
- [13] V. Kreinovich, A.V. Lakeyev, and S.I. Noskov. Optimal Solution of Interval Linear Systems Is Intractable (NP-Hard). *Interval Computations*, 1:6–14, 1993.
- [14] Lei Li. On the iterative criterion for generalized diagonally dominant matrices. *SIAM Journal Matrix Anal. Appl. (SIMAX)*, 24(1):17–24, 2002.
- [15] J.M. Muller, N. Brisebarre, F. de Dinechin, C.P. Jeannerod, V. Lefèvre, G. Melquiond, N. Revol, D. Stehlé, and S. Torres. *Handbook of Floating-Point Arithmetic*. Birkhäuser Boston, 2010.
- [16] A. Neumaier. Overestimation in Linear Interval Equations. *SIAM J. Numer. Anal. (SINUM)*, 24(1):207–214, 1987.
- [17] A. Neumaier. *Rigorous Sensitivity Analysis for Parameter-Dependent Systems of Equations*. J. Math. Anal. Appl. 144, 1989.
- [18] A. Neumaier. *Interval Methods for Systems of Equations*. Encyclopedia of Mathematics and its Applications. Cambridge University Press, 1990.
- [19] A. Neumaier. A simple derivation of the Hansen-Blied-Rohn-Ning-Kearfott enclosure for linear interval equations. *Reliable Computing*, 5:131–136, 1999.

- [20] A. Neumaier. Erratum to: A simple derivation of the Hansen-Bliek-Rohn-Ning-Kearfott enclosure for linear interval equations. *Reliable Computing*, 6:227, 1999.
- [21] A. Neumaier. private communication, 2010.
- [22] H.D. Nguyen. *Efficient algorithms for verified scientific computing: numerical linear algebra using interval arithmetic*. Dissertation, Laboratoire LIP, INRIA ENS Lyon, 2011.
- [23] H.D. Nguyen and N. Revol. Accuracy issues in linear algebra using interval arithmetic. *SCAN conference Lyon*, 2010.
- [24] S. Ning and R.B. Kearfott. A Comparison of some Methods for Solving Linear Interval Equations. *SIAM J. Numer. Anal. (SINUM)*, 34(4):1289–1305, 1997.
- [25] T. Ogita, S. Oishi, and Y. Ushiro. Fast verification of solutions for sparse monotone matrix equations. volume 15 of *Comput. Suppl.*, pages 175–187. Springer, Wien, 2001.
- [26] T. Ogita, S. Oishi, and Y. Ushiro. Fast inclusion and residual iteration for solutions of matrix equations. 16:171–184, 2002.
- [27] S. Oishi and S.M. Rump. Fast verification of solutions of matrix equations. *Numer. Math.*, 90(4):755–773, 2002.
- [28] S. Poljak and J. Rohn. Checking Robust Nonsingularity Is NP-Hard. *Math. of Control, Signals, and Systems* 6, pages 1–9, 1993.
- [29] J. Rohn. Linear Interval Equations: Computing Enclosures with Bounded Relative or Absolute Overestimation is NP-hard. In R.B. Kearfott and V. Kreinovich, editors, *Applications of Interval Computations*, pages 81–89. Kluwer Academic Publisher, 1991.
- [30] J. Rohn. Stability of the Optimal Basis of a Linear Program under Uncertainty. *Operations Research Letters* 13, pages 9–12, 1993.
- [31] J. Rohn. VERSOFT: Verification software in MATLAB/INTLAB. 2009.
- [32] S.M. Rump. *Kleine Fehlerschranken bei Matrixproblemen*. PhD thesis, Universität Karlsruhe, 1980.
- [33] S.M. Rump. Rigorous Sensitivity Analysis for Systems of Linear and Nonlinear Equations. *Math. Comput.*, 54(10):721–736, 1990.
- [34] S.M. Rump. Verification Methods for Dense and Sparse Systems of Equations. In J. Herzberger, editor, *Topics in Validated Computations — Studies in Computational Mathematics*, pages 63–136, Elsevier, Amsterdam, 1994.
- [35] S.M. Rump. Fast and parallel interval arithmetic. *BIT Numerical Mathematics*, 39(3):539–560, 1999.
- [36] S.M. Rump. INTLAB - INTerval LABoratory. In Tibor Csendes, editor, *Developments in Reliable Computing*, pages 77–104. Kluwer Academic Publishers, Dordrecht, 1999.
- [37] R. Skeel. Iterative Refinement Implies Numerical Stability for Gaussian Elimination. *Math. Comp.*, 35(151):817–832, 1980.
- [38] R.S. Varga. *Matrix Iterative Analysis*. Prentice-Hall, Englewood Cliffs, N.J., 1962.
- [39] J. Zemke. b4m - BIAS for Matlab. Technical report, Inst. f. Informatik III, Technische Universität Hamburg-Harburg, 1998.