

# ACCURATE SOLUTION OF DENSE LINEAR SYSTEMS PART I: ALGORITHMS IN ROUNDING TO NEAREST

SIEGFRIED M. RUMP \*

**Abstract.** We investigate how extra-precise accumulation of dot products can be used to solve ill-conditioned linear systems accurately. For a given  $p$ -bit working precision, extra-precise evaluation of a dot product means that the products and summation are executed in  $2p$ -bit precision, and that the final result is rounded into the  $p$ -bit working precision. Denote by  $\mathbf{u} = 2^{-p}$  the relative rounding error unit in a given working precision. We treat two types of matrices: First up to condition number  $\mathbf{u}^{-1}$ , and second up to condition number  $\mathbf{u}^{-2}$ . For both types of matrices we present two types of methods: First for calculating an approximate solution, and second for calculating rigorous error bounds for the solution together with the proof of non-singularity of the matrix of the linear system. In the first part of this paper we present algorithms using only rounding to nearest, in Part II we use directed rounding to obtain better results. All algorithms are given in executable Matlab code and are available from my homepage.

**Key words.** Linear systems, Matlab, rounding to nearest, 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 solution of a linear system  $Ax = b$  is a ubiquitous task in numerical computations. In Part I and II of this paper we present different methods to compute guaranteed error bounds for the solution of a linear system, i.e. with a *certified* accuracy. The methods in this Part I are based on norm estimates, in particular verifying convergence of some residual matrix by approximating its Perron vector, whereas the methods in Part II are based on the verification of the  $H$ -property of some matrix. Moreover, in the present Part I of the paper we 1) present a method to compute an approximation for *extremely* ill-conditioned linear systems which is *likely* to be accurate.

In the present Part I all algorithms use only the four basic floating-point operations in rounding to nearest, in Part II directed rounding is used as well. The challenge for the first part is to use only standard Matlab code in rounding to nearest without additional mex-files and to derive simple and fast algorithms. All algorithms in both parts are presented in executable Matlab-code.

Let a floating-point format with relative rounding error unit  $\mathbf{u}$  be given. The forward error of an approximation  $\tilde{x}$  computed by a standard algorithm like Gaussian elimination is of the order  $\mathbf{u} \cdot \text{cond}(A)$  [20]. This naturally bounds the applicability to matrices  $A$  with  $\text{cond}(A) \lesssim \mathbf{u}^{-1}$ , which means  $\text{cond}(A) \lesssim 10^{16}$  in IEEE 754 double precision (binary64). For larger condition numbers,  $\tilde{x}$  is expected to have no correct digit.

Skeel [48] showed that one step of the classical residual iteration with the residual computed in *working precision* produces a backward stable result. It is also known that, for condition numbers up to  $\text{cond}(A) \lesssim \mathbf{u}^{-1}$ , an almost maximally accurate result (of order  $\mathbf{u}$ ) is achieved if the dot products in the residual iteration are accumulated in twice the working precision (see also [12]). But for condition numbers of the order  $\mathbf{u}^{-1}$  and larger, again no correct digit can be expected. So basically we face the dichotomy of either high accuracy or no accuracy at all.

---

\*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](mailto:rump@tu-harburg.de)).

TABLE 1.1

Constants for single (binary32) and double (binary64) precision in the IEEE 754 floating-point standard .

	$p$	$e_{min}$	$e_{max}$
single precision	24	-126	127
double precision	53	-1022	1023

We will show that the accumulation of dot products in twice the working precision (with result rounded into working precision) suffices to compute an accurate approximation of the solution of a linear system for condition numbers up to  $\text{cond}(A) \lesssim \mathbf{u}^{-2}$ . Moreover, we show how rigorous error bounds can be computed including the proof on non-singularity of the input matrix  $A$ . Practical experience suggests that this approach works successfully up to condition numbers  $c \cdot \text{cond}(A) \lesssim \mathbf{u}^{-2}$  with  $c \approx n^2$  in IEEE 754 binary64 (double precision). This factor will be improved to about  $c \approx n$  in Part II of this paper.

We want to stress that our bounds are mathematically completely rigorous including all possible sources of errors (provided the compiler and operating system work to their specification). Although it is in principle known how to compute rigorous error bounds in rounding to nearest, the corresponding algorithms are involved, and taking care of underflow they become unwieldy. One reason to divide the paper in two parts is to clearly distinguish between algorithms using solely rounding to nearest (Part I), and those using directed rounding (Part II).

There are other, very good but not completely rigorous approaches. For example, an upper bound of the condition number  $\text{cond}(A)$  implies an error bound of an approximate solution of  $Ax = b$ . There are many  $\mathcal{O}(n^2)$  condition number estimators (cf. [18, 20]), usually providing good approximations. By the principle of the methods these are *lower bounds* for the condition number, and for every estimator counterexamples are known where the condition number is grossly underestimated.

In another approach [26, 51] the authors use a statistic way for estimating rounding errors. Using a so-called stochastic arithmetic they propose a method to determine the number of significant digits of a computed result. Also those results are correct with a high degree of certainty, but not with complete rigor.

Yet another approach [12] uses the vast experience in solving linear systems very thoughtfully to produce approximations with “likely correct error terms” [12]. It seems that no counterexample is known where the claimed accuracy is not valid, but it is not *proved* to be correct.

To repeat it, beyond accurate approximations for very ill-conditioned linear systems, we are also interested in mathematically rigorous error bounds. Such rigorous bounds are, for example, mandatory in so-called “computer-assisted proofs” [16], which recently gain interest. For example, Tucker [50] received the 2004 EMS prize awarded by the European Mathematical Society for “giving a rigorous proof that the Lorenz attractor exists for the parameter values provided by Lorenz. This was a long standing challenge to the dynamical system community, and was included by Smale in his list of problems for the new millennium. The proof uses computer estimates with rigorous bounds based on higher dimensional interval arithmetics.”

Concerning notation denote by  $\mathbb{F}$  a set of  $p$ -bit binary floating-point numbers including  $\infty$  and NaN, i.e. [30]

$$(1.1) \quad \mathbb{F} = \{ M \cdot 2^{e-p+1} \mid M, e \in \mathbb{Z}, |M| \leq 2^p - 1, e_{min} \leq e \leq e_{max} \} \cup \{-\infty, +\infty, \text{NaN}\} .$$

For single (binary32) and double (binary64) precision in the IEEE 754 floating-point standard [22, 23] the constants are as in Table 1.1. We assume floating-point operations in rounding to nearest, tie to even, as in the IEEE 754 standard. That means there is a mapping  $\text{fl} : \mathbb{R} \rightarrow \mathbb{F}$  such that  $|\text{fl}(x) - x| = \min_{f \in \mathbb{F}} |f - x|$  for all  $x \in \mathbb{R}$ , and for  $a, b \in \mathbb{F}$  floating-point operations  $\circ_{\text{fl}} : \mathbb{F} \times \mathbb{F} \rightarrow \mathbb{F}$  with  $\circ \in \{+, -, \cdot, /\}$  are defined by

$$(1.2) \quad a \circ_{\text{fl}} b := \text{fl}(a \circ b) .$$

Therefore the result  $a \circ_{\text{fl}} b \in \mathbb{F}$  is a best approximation of  $a \circ b \in \mathbb{R}$ . The relative rounding error unit is defined by  $\mathbf{u} = 2^{-p}$ . A floating-point number  $M \cdot 2^{e-p+1}$  is normalized if  $|M| \geq 2^{p-1}$ , the smallest normalized positive floating-point number is  $\mathbf{realmin} = 2^{e_{\min}}$ , and the smallest unnormalized positive floating-point number is  $\mathbf{eta} = 2^{e_{\min}-p+1}$ . All algorithms are given in executable Matlab code using IEEE 754 double precision, but the results are valid in any floating-point arithmetic complying with the IEEE 754 standard.

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  $1 \leq i \leq n$ . Executable Matlab-code is written using 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  $|P - \mathbf{C}| \sim \mathbf{u}|\mathbf{A}| \cdot |\mathbf{B}|$ .

In principle, it is not difficult to transform an algorithm using directed rounding into an algorithm using only rounding to nearest [36, 38]. However, often the code becomes unwieldy. In this paper we use new estimates on floating-point summation and dot products [45] to derive particularly simple algorithms to compute verified error bounds using only rounding to nearest.

The paper is organized as follows. In the next section we describe an algorithm for calculating an approximation of the solution of extremely ill-conditioned linear systems, i.e. up to condition number  $\mathbf{u}^{-2} \approx 10^{32}$  in IEEE 754 double precision. An algorithm for the necessary extra-precise evaluation of dot products is presented in Section 3. It uses only basic floating-point operations in rounding to nearest. The following Section 4 splits in five subsections including the computation of rigorous error bounds in rounding to nearest for ill-conditioned linear systems, for extremely ill-conditioned linear systems as well as for extra-precise dot product accumulation. This suggests a hybrid algorithm addressing both ill-conditioned and extremely ill-conditioned systems which will be discussed in Part II of this paper. Computational results and a conclusion finish the paper.

**2. Accurate approximations for very ill-conditioned linear systems.** The aim of this section is an algorithm to compute an accurate approximation of a linear system  $Ax = b$  with  $\mathbf{u}^{-1} \leq \text{cond}(A) \lesssim \mathbf{u}^{-2}$ . Besides the basic floating-point operations  $\{+, -, \cdot, /\}$  the algorithm requires only the accumulation of dot products in twice the working precision. For  $x, y \in \mathbb{F}^n$  this is denoted by

$$(2.1) \quad \mathbf{res} = \text{Dot2Near}(x', y); \quad \% \text{ accumulation of dot product } x^T y \text{ in twice the working precision .}$$

For given p-precision  $x, y \in \mathbb{F}^n$  this means that  $\mathbf{res} \in \mathbb{F}$  is the result obtained when calculating the products  $x_i y_i$  in 2p-precision, accumulating the sum  $\sum x_i y_i$  in 2p-precision and rounding the result of the sum into working precision (p-precision). Therefore the error of  $\mathbf{res}$  is bounded by  $\mathbf{u}|\mathbf{res}| + c\mathbf{u}^2(|x|^T|y|)$  for a small constant  $c$ , the second term addressing the accumulation error and the first term the rounding into working precision. In [12] this is called “extra-precise” accumulation of dot products. We use similarly

$$(2.2) \quad \begin{aligned} \mathbf{res} &= \text{Dot2Near}(\mathbf{R}, \mathbf{b}); & \% \text{ accumulation of dot products in } \mathbf{R} \cdot \mathbf{b} \text{ in twice the working precision} \\ \mathbf{res} &= \text{Dot2Near}(\mathbf{R}, \mathbf{A}); & \% \text{ accumulation of dot products in } \mathbf{R} \cdot \mathbf{A} \text{ in twice the working precision.} \end{aligned}$$

There is a vast amount of literature devoted to the accurate computation of sums and dot products, among them [10, 19, 29, 31, 40, 39, 46, 44, 47, 55, 56, 57]. One way to implement `Dot2Near` using only basic floating-point operations in working precision is XBLAS [28, 53]. Here two IEEE 754 double precision (binary64) floating-point numbers are used to represent a quadruple precision number, and accurate arithmetical operations on pairs are defined. This does more than necessary for our purposes because the result is a pair, but for `Dot2Near` we need only the first part. Another way is to use some multiple-precision package like [3, 15].

Since the accurate computation of residuals, possibly with rigorous error bounds, is of central importance

TABLE 2.1  
*Results for n=10; A=hilb(10) and x=A\b and y=inv(A)\*b.*

	$\ A\tilde{x} - b\ _2$	$\ A\tilde{y} - b\ _2$	$\ A^{-1}b - \tilde{x}\ _2/\ A^{-1}b\ _2$	$\ A^{-1}b - \tilde{y}\ _2/\ A^{-1}b\ _2$
$b = \text{randn}(n,1);$	$2.1 \cdot 10^{-4}$	$1.7 \cdot 10^{-3}$	$1.2 \cdot 10^{-4}$	$1.5 \cdot 10^{-5}$
$b = A*\text{randn}(n,1);$	$1.2 \cdot 10^{-16}$	$3.9 \cdot 10^{-5}$	$5.3 \cdot 10^{-5}$	$1.4 \cdot 10^{-4}$

for this paper, we discuss several methods and give executable code in Section 3. For the moment assume an algorithm `Dot2Near` as described to be given.

Let  $A \in \mathbb{F}^{n \times n}$  with  $\text{cond}(A) \geq \mathbf{u}^{-1}$  be given. Then an approximation  $\tilde{x}$  of the solution of a linear system with matrix  $A$  computed by Gaussian elimination in double precision is expected to be heavily corrupted by rounding errors, and no digit of  $\tilde{x}$  can be expected to be correct.

It is well-known [14, 21] that “in the vast majority of practical computational problems, it is unnecessary and inadvisable to actually compute  $A^{-1}$ .” In particular an approximation  $\tilde{y} := Rb$  using some approximate inverse  $R$  of  $A$  requires not only three times as much operations than Gaussian elimination, it is also, in general, less accurate and less stable [21]. Nevertheless an approximate inverse is our key to solve such extremely ill-conditioned linear systems.

More precisely, numerical evidence suggests that on the one hand, depending on the right hand side  $b$ , the residual  $\|A\tilde{x} - b\|_2$  for  $\tilde{x} = A \setminus b$  is sometimes much smaller than  $\|A\tilde{y} - b\|_2$  with  $\tilde{y} = Rb$ . However, there seems to be on the other hand, in general, not much difference between  $\|A^{-1}b - \tilde{x}\|_2/\|A^{-1}b\|_2$  and  $\|A^{-1}b - \tilde{y}\|_2/\|A^{-1}b\|_2$ . Results for a not untypical example are given in Table 2.1.

Let  $R$  be an approximate inverse computed in working precision (e.g. by the Matlab command `R=inv(A)`), and assume  $\text{cond}(A) \geq \mathbf{u}^{-1}$ . Then  $R$  is also expected to be entirely corrupted by rounding errors with no correct digit. Nevertheless the approximate inverse  $R$  contains useful information. This corresponds to the fact that the rounding errors in Gaussian elimination are by no means random, see [49, 52].

In about 1984 I derived an algorithm squeezing out this information. Because of lack of analysis, I did not publish it. In [37] Oishi et al. analyzed a modification of this algorithm, and in [43] I analyzed the original algorithm. It requires the accumulation of dot products in some  $K$ -fold precision and storing the result in an unevaluated sum of  $L$  floating-point numbers; otherwise floating-point operations in working precision are used.

For this algorithm it is important to store an approximate inverse in an unevaluated sum of matrices  $R_1 + \dots + R_k$ , where the summands are computed recursively starting with the first approximate inverse  $R$ . Under reasonable assumptions it is shown in [43] that with  $k$  summands matrices  $A$  up to condition number  $\mathbf{u}^{-k}$  can be treated, i.e. the spectral radius of  $I - (\sum_{\nu=1}^k R_\nu)A$  becomes less than 1.

The first step of this algorithm is given by the following executable Matlab code. We abbreviate “accumulation of dot products in twice the working precision and rounding into working precision” by “extra-precise accumulation” in the comments.

```

1  R = inv(A);                               % approximate inverse
2  while any(isinf(R(:))) || any(isnan(R(:)))
3      R = inv(A.*(1+randn(n)*eps));         % inversion of perturbed matrix
4  end
5  C = Dot2Near(R,A);                         % extra-precise accumulation
6  Cinv = inv(C);                             % multiplicative correction for R

```

For  $A \in \mathbb{F}^{n \times n}$  with  $\text{cond}(A) \geq \mathbf{u}^{-1}$  it may happen that the “approximate inverse”  $R$  computed in the first line contains infinity- and NaN-components. In that case the matrix  $A$  is slightly perturbed and inverted

again. Note that mathematically, due to the large condition number of  $\mathbf{A}$ , this may change the entries of  $\mathbf{R}$  in the first digit. In any case,  $\mathbf{R}$  is completely corrupted by rounding errors.

Nevertheless it can be observed that

$$(2.3) \quad \text{cond}(\mathbf{C}) \sim \mathbf{u} \cdot \text{cond}(\mathbf{A}), \quad \text{even for } \text{cond}(\mathbf{A}) \gg \mathbf{u}^{-2}.$$

We cannot expect a mathematically rigorous analysis, but in [43] arguments are given for that. Examples with condition number up to  $10^{300}$ , just before overflow, confirm this observation.

For  $\text{cond}(\mathbf{A}) \sim \beta \mathbf{u}^{-1}$  this means  $\text{cond}(\mathbf{C}) \sim \beta$ , so that for  $\beta \lesssim \mathbf{u}^{-1}$  we can expect some accuracy in  $\mathbf{C}\text{inv}$ . The next step in [43] is to compute  $\mathbf{C}\text{inv} \cdot \mathbf{R}$  in twice the working precision but to store the result in an unevaluated sum  $\mathbf{R}_1 + \mathbf{R}_2$ . Then it is shown that  $I - (\mathbf{R}_1 + \mathbf{R}_2)\mathbf{A}$  is convergent for  $\text{cond}(\mathbf{A}) \lesssim \mathbf{u}^{-2}$ .

One might use the single matrix  $\text{Dot2Near}(\mathbf{C}\text{inv}, \mathbf{R}) \in \mathbb{F}^{n \times n}$ , certainly a good if not best approximation to  $\mathbf{C}\text{inv} \cdot \mathbf{R}$ , directly as a preconditioner for  $\mathbf{A}$ . However, for  $\text{cond}(\mathbf{A}) > \mathbf{u}^{-1}$  a single preconditioning matrix  $\mathbf{B} \in \mathbb{F}^{n \times n}$  cannot, in general, force  $I - \mathbf{B}\mathbf{A}$  to be convergent. Even for  $\mathbf{B}$  being the nearest floating-point matrix to  $\mathbf{A}^{-1}$ , that is  $\mathbf{B} = \mathbf{A}^{-1} + \Delta$  with  $\|\Delta\| \sim \mathbf{u}\|\mathbf{A}^{-1}\|$  for some norm, we have  $\|I - \mathbf{B}\mathbf{A}\| = \|\Delta \cdot \mathbf{A}\| \sim \mathbf{u} \cdot \text{cond}(\mathbf{A}) > 1$ . Thus the unevaluated sum  $\mathbf{R}_1 + \mathbf{R}_2$  rather than  $\mathbf{C}\text{inv} \cdot \mathbf{R}$  is necessary in [43] to serve as a preconditioning matrix.

For computing an accurate approximation of the solution of a linear system, even for  $\text{cond}(\mathbf{A}) > \mathbf{u}^{-1}$ , no unevaluated sum as a preconditioning matrix is necessary. In the following algorithm we use only  $\text{Dot2Near}$  as specified in (2.2), i.e. accumulation of dot products in twice the working precision with the result stored in working precision, and only floating-point operations in rounding to nearest.

ALGORITHM 2.1. *Accurate approximate solution  $\mathbf{x}_s$  of  $A\mathbf{x} = \mathbf{b}$  for extremely ill-conditioned  $A$ .*

```

1 function  $\mathbf{x}_s = \text{LssIllcoApprox}(\mathbf{A}, \mathbf{b})$ 
2    $n = \text{size}(\mathbf{A}, 1);$  % dimension of linear system
3    $\mathbf{R} = \text{inv}(\mathbf{A});$  % approximate inverse
4   while any(isinf(R(:))) || any(isnan(R(:)))
5      $\mathbf{R} = \text{inv}(\mathbf{A} \cdot (1 + \text{randn}(n) \cdot \text{eps}));$  % inversion of perturbed matrix
6   end
7    $\mathbf{C} = \text{Dot2Near}(\mathbf{R}, \mathbf{A});$  % extra-precise accumulation
8    $\mathbf{C}\text{inv} = \text{inv}(\mathbf{C});$  % multiplicative correction for  $\mathbf{R}$ 
9    $\mathbf{x}_s = \mathbf{C}\text{inv} \cdot \text{Dot2Near}(\mathbf{R}, \mathbf{b});$  % first approximate solution
10   $\mathbf{N} = \text{inf}; \text{iter} = 0;$  % initialization of constants
11  while iter < 5 % at most 5 residual iterations
12    iter = iter + 1;  $\mathbf{N}\text{old} = \mathbf{N};$  % update constants
13     $\text{res} = \text{Dot2Near}([\mathbf{A} \ \mathbf{b}], [\mathbf{x}_s; -1]);$  % residual  $\mathbf{A} \cdot \mathbf{x}_s - \mathbf{b}$  (extra-precise acc.)
14     $\mathbf{d} = \mathbf{C}\text{inv} \cdot \text{Dot2Near}(\mathbf{R}, \text{res});$  % correction for  $\mathbf{x}_s$ 
15     $\mathbf{N} = \text{norm}(\mathbf{d}, 1);$  % norm of correction
16    if  $\mathbf{N} < \mathbf{N}\text{old}$ ,  $\mathbf{x}_s = \mathbf{x}_s - \mathbf{d};$  end % correction acceptable
17    if  $\mathbf{N} \geq 0.1 \cdot \mathbf{N}\text{old}$ , break, end % stop iteration if no improvement
18  end
    
```

OBSERVATION 2.2. *Let  $\mathbb{F}$  be a set of floating-point numbers with relative rounding error unit  $\mathbf{u}$  together with floating-point operations complying with the IEEE 754 arithmetic standard [22, 23]. Let  $\mathbf{x}_s$  be the result of Algorithm 2.1 ( $\text{LssIllcoApprox}$ ) applied to a matrix  $\mathbf{A} \in \mathbb{F}^{n \times n}$  with  $\text{cond}(\mathbf{A}) \lesssim \mathbf{u}^{-2}$  and to a right hand side  $\mathbf{b} \in \mathbb{F}^n$ . Then numerical evidence suggests that the relative error of  $\mathbf{x}_s$  to the exact solution  $\mathbf{A}^{-1}\mathbf{b}$  is of the order  $\mathbf{u} + \mathbf{u}^2 \text{cond}(\mathbf{A})$ .*

To start with, we do not fully understand why Algorithm 2.1 works that good for extremely ill-conditioned matrices. Obviously  $(\mathbf{R}\mathbf{A})^{-1}\mathbf{R} = \mathbf{A}^{-1}$ , so that without the presence of rounding errors  $\mathbf{C}\text{inv} \cdot \mathbf{R} \cdot \mathbf{b} = \mathbf{A}^{-1}\mathbf{b}$ .

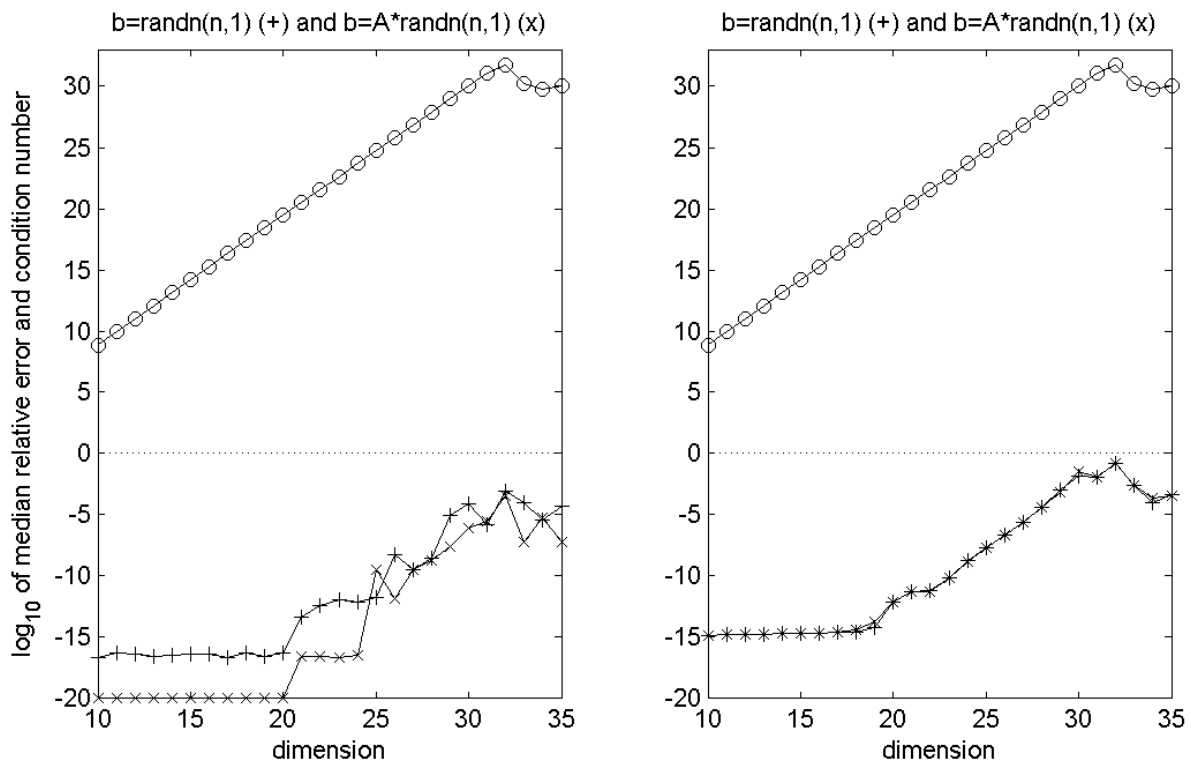


FIG. 2.1. Results of Algorithm 2.1 (`LssIllcoApprox`) (left) and of Algorithm 4.20 (`LssIllcoErrBndNear`) (right) for Pascal matrices with right hand sides  $\mathbf{b}=\text{randn}(n,1)$  and  $\mathbf{b}=\mathbf{A}*\text{randn}(n,1)$ . The upper parts show the condition number, the lower parts the relative error of the results.

This is the main idea explored in [43]. Over there, we identified  $\mathbf{C}\text{inv} \cdot \mathbf{R}$  as a suitable preconditioning matrix for  $\mathbf{A}$ ; here we are interested in solving the linear system and compute  $\mathbf{C}\text{inv} \cdot (\mathbf{R} \cdot \mathbf{b})$ .

An approximate solution  $\mathbf{x}_s$  is computed in line 9, and in lines 11 – 18 some residual iteration is applied to it. An approximate solution  $\mathbf{d}$  of the residual system is computed in the same way. Since  $\mathbf{d}$  is the correction to  $\mathbf{x}_s$ , the pair  $(\mathbf{x}_s, \mathbf{d})$  might be regarded as an unevaluated sum. However, the pair is not proceeded but added (in working precision) into the new  $\mathbf{x}_s$  in line 16.

The extra-precise accumulation of dot products is used for the preconditioning  $\mathbf{R} \cdot \mathbf{A}$  in line 7, for the residual  $\mathbf{A} \cdot \mathbf{x}_s - \mathbf{b}$  in line 13, and for the multiplication of a right hand side by  $\mathbf{R}$  in lines 9 and 14. If one of these dot product operations is replaced by the usual computation in working precision, the results deteriorate significantly or the algorithm fails completely. One may be inclined to use `Dot2Near` for the multiplication by  $\mathbf{C}\text{inv}$  in lines 9 and 14 as well; however, numerical evidence suggests that usually this does not improve the results (sometimes to the contrary): only in some extreme situations it is advantageous.

To this end we give one typical example for the performance of Algorithm 2.1 (`LssIllcoApprox`). It is not obvious how to construct an extremely ill-conditioned matrix with floating-point entries<sup>1</sup>. Ways to construct ill-conditioned floating-point matrices are introduced in [42, 33, 34]. Here we use Pascal matrices defined by  $A_{ij} := \binom{i+j-2}{j-1}$ . Up to dimension  $n = 31$  the entries are double precision floating-point numbers; for  $n > 31$  the entries are rounded to the nearest floating-point number. To increase numerical stability we use always equilibrated matrices as in [12], see Section 5.

As can be seen in Figure 2.1, the condition number (upper line, computed by the symbolic toolbox in Matlab)

<sup>1</sup>For a discussion of “well-known suspects” of ill-conditioned matrices  $\mathbf{A}$  and alternative ways to solve  $\mathbf{A}\mathbf{x} = \mathbf{b}$  see the beginning of Section 5

increases monotonically up to dimension  $n = 32$ , then the matrix entries become corrupted by the conversion into floating-point. Although the Pascal matrix is not exactly representable for  $n = 32$ , the condition number is accidentally larger than for  $n = 31$ .

The lower lines in the left graph show the median relative error of the approximation computed by Algorithm 2.1 against the exact solution computed by the Matlab symbolic toolbox: For given  $\mathbf{A} \in \mathbb{F}^{n \times n}$  and  $\mathbf{b} \in \mathbb{F}^n$ , the solution  $x \in \mathbb{Q}^n$  of the linear system as a vector of rational numbers is computed and compared to the computed approximation  $\mathbf{xs} \in \mathbb{F}^n$ . For two different right hand sides  $\mathbf{b} = \text{randn}(n,1)$  (+) and  $\mathbf{b} = \mathbf{A} * \text{randn}(n,1)$  (x) the median relative error of  $\mathbf{xs}$  to  $x$  is displayed. The right graph shows the quality of rigorous error bounds computed by Algorithm 4.20 (`LssI11coErrBndNear`) and will be discussed later.

It can be seen that up to condition number  $10^{16}$  the approximate solution is of full accuracy, whereas for condition number  $10^k$  with  $k > 16$  about  $32 - k$  digits are correct. This acknowledges Observation 2.2. Extensive numerical tests with various types of matrices and right hand sides confirm that up to condition numbers of about  $10^{30}$  the computed approximations can be expected to have some accuracy, see Section 5.

To this end, we can deduce the accuracy of the computed approximation by *knowing* the exact solution by some oracle or by computing it in rational arithmetic. In the next sections we show how mathematically rigorous error bounds can be *computed* in floating-point rounding to nearest.

**3. Accurate dot products in rounding to nearest.** Following we describe an algorithm realizing `Dot2Near` as used in Section 2. Although Algorithm 2.1 (`LssI11coApprox`) needs only an accurate *approximation*, the following algorithm also computes rigorous error bounds for a dot product. Therefore it is also suitable for our algorithms computing rigorous error bounds to be discussed in the next section.

The algorithm to be presented requires only the basic floating-point operations  $+$ ,  $-$ ,  $\cdot$ ,  $/$  in working precision and no additional features such as access to mantissa and/or exponent, assembly language routines or alike, and the algorithm is also free of branches. This improves, as analyzed by Langlois [27], the computational performance on today's architectures significantly.

In [31] Neumaier, as a young student, developed a fast algorithm with provably improved accuracy. It computes an approximation of a dot product with a quality "as if" computed in twice the working precision. His paper is written in German and did not receive wide attention. A modern formulation of this algorithm was presented in [35] and is based on error-free transformations.

We first need an error-free transformation to split a floating-point number into a high and low order part. This is done by Dekker's method as in Algorithm 3.1.

ALGORITHM 3.1. *Error-free splitting of a floating-point number  $a$  into two parts  $x, y$  such that  $a = x + y$ .*

```
function [x,y] = Split(a,s)
    c = (2^s+1)*a;           % splitting in (53-s)-bit and (s-1)-bit part
    x = c - (c-a);         % high order part
    y = a - x;             % low order part
```

As a result,  $a = x + y$  for all  $a \in \mathbb{F}$ , also in the presence of underflow. Moreover, for 53-bit precision input  $a$ , the summands  $x$  and  $y$  have at most  $53 - s$  and  $s - 1$  significant bits, respectively. So for  $s = 27$ , both summands have at most 26 significant bits adding to a 53-bit number, an apparent contradiction. This is possible because Dekker's ingenious and fast algorithm uses the sign bit as an extra bit of information. Note that Algorithm 3.1 works correctly for vector and matrix input as well, and possible overflow may be avoided by some scaling.

For completeness we repeat algorithms `TwoSum` and `TwoProduct`, error-free transformations of the sum and product of two floating-point numbers into the nearest floating-point approximation and the true error,

respectively.

ALGORITHM 3.2. *Error-free transformation of  $a + b$  into  $x + y$ .*

```
function [x,y] = TwoSum(a,b)
    x = a + b;                % floating-point approximation of a+b
    z = x - a;
    y = ( a - (x-z)) + (b-z); % exact error of x
```

ALGORITHM 3.3. *Error-free transformation of  $a \cdot b$  into  $x + y$ .*

```
function [x,y] = TwoProduct(a,b)
    x = a * b;                % floating-point approximation of a*b
    [a1,a2] = Split(a,27);   % error-free splitting a = a1+a2
    [b1,b2] = Split(b,27);   % error-free splitting b = b1+b2
    y = a2*b2 - (((x-a1*b1)-a2*b1)-a1*b2); % exact error of x (if no underflow)
```

Algorithm 3.2 (`TwoSum`) is due to Knuth [24], and Algorithm 3.3 (`TwoProduct`) is due to G.W. Veltkamp (see [8]). The two algorithms satisfy for all  $a, b \in \mathbb{F}$  the rigorous error estimates [30]

$$(3.1) \quad \begin{aligned} [x, y] = \text{TwoSum}(a, b) &\Rightarrow a + b = x + y \quad \text{and} \\ [x, y] = \text{TwoProduct}(a, b) &\Rightarrow a \cdot b = x + y + \eta \quad \text{with } |\eta| \leq 3\text{eta} . \end{aligned}$$

Recall that `eta` denotes the smallest positive (unnormalized) floating-point number. In IEEE 754 double precision  $\text{eta} = 2^{-1074}$ , so that the smallest positive normalized floating-point number is  $\text{realmin} = \frac{1}{2}\mathbf{u}^{-1}\text{eta}$ .

Note that the results of `TwoSum` always satisfy  $x + y = a + b$ , also in the presence of underflow, whereas the results of `TwoProduct` satisfy  $x + y = a \cdot b$  if no underflow occurs.

The routine `TwoSum` requires 17 floating-point operations. The new IEEE 754 floating-point standard [23] requires an FMA (Fused Multiply and Add) operation, which is already available on some processors. It computes  $a \cdot b + c$  with one final rounding to nearest. With FMA, `TwoProduct` can be replaced by

```
x = a*b;
y = FMA(a,b,-x);
```

thus requiring only two floating-point operations instead of 17. Based on those routines a summation algorithm was presented in [35], which is almost identical with Neumaier's [31]. We formulate it directly for matrix multiplication using rank-1 updates.

ALGORITHM 3.4. *Approximation of the matrix product  $A \cdot B$  "as if" accumulated in twice the working precision and rounded into working precision with rigorous error term.*

```
function [res,err] = Dot2Near(A,B)
    [p,e] = TwoProduct(A(:,1),B(1,:)); % error-free transformation of first product
*   E = abs(e); % for error term
    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
        t = q + r; % sum of errors
        e = e + t; % accumulation of errors
*   E = E + abs(t); % accumulation for error term
    end
    res = p + e; % extra-precise approximation
```



TABLE 3.1

Number of floating-point operations of `Dot2Near` without and with error bound, and without and with FMA.

	without FMA	with FMA
<code>Dot2Near</code> without error bound	$25n$	$10n$
<code>Dot2Near</code> with error bound	$27n$	$12n$
<code>DotXBLAS</code> without error bound	$37n$	$22n$

```
*      epss = 0.5*eps;                % relative rounding error unit 2^(-53)
*      err0 = max(6*k*epss,1)*realmin + (k+2)*epss*ufp(E) + epss*ufp(res);
*      err = err0 + 3*epss*ufp(err0);  % rigorous error bound for res
```

Apparently Neumaier did not know about the error-free transformations `TwoSum` and `TwoProduct`. He developed his algorithm as a sequence of floating-point operations with some reminiscence to the Kahan-Babuška algorithm [2].

Here we added a simplified computation of a rigorous error term which we need and explain later. It is based on the new analysis of floating-point summation in [45]. The error term is computed in the lines marked with an asterisk. If no error term is needed, all those lines can be omitted.

**THEOREM 3.5.** *Let  $A \in \mathbb{F}^{m \times k}$  and  $B \in \mathbb{F}^{k \times n}$  with  $(k + 2)\mathbf{u} \leq 1$  be given, and let `res` and `err` be the quantities computed by Algorithm 3.4 (`Dot2Near`). Then, also in the presence of underflow,*

$$(3.2) \quad |A \cdot B - \text{res}| \leq \mathbf{u}|A \cdot B| + \gamma_k^2 |A| |B| + 5k\mathbf{eta} ,$$

where `eta` denotes the smallest positive (unnormalized) floating-point number. Moreover,

$$(3.3) \quad |A \cdot B - \text{res}| \leq \text{err} .$$

The factor  $k + 2$  in the second last line of Algorithm 3.4 cannot be replaced by  $k + 1$ .

The first estimate (3.2) follows by the corresponding estimate in [35] for dot products, the correctness of the error bound (3.3) will be shown in Subsection 4.2. The first term in the right hand side of (3.2) reflects the unavoidable error of rounding the final sum `res=p+e` into working precision, the second term reflects the accuracy of the unevaluated sum `p + e` before this addition, and the third term covers possible underflow. This means that the quality is “as if” computed in twice the working precision and rounded into working precision as in (2.2).

The number of floating-point operations up to  $\mathcal{O}(1)$  for `Dot2Near` is as in Table 3.1. For comparison, the data for XBLAS [28] are displayed as well.

**4. Rigorous error bounds for linear systems including extremely ill-conditioned matrices with  $\text{cond}(A) \lesssim \mathbf{u}^{-2}$ .** For a matrix  $E \in \mathbb{R}^{n \times n}$  with  $\|E\| < 1$  in some matrix norm it is well-known [18, 32] that  $I \pm E$  is non-singular and  $\|(I \pm E)^{-1}\| \leq (1 - \|E\|)^{-1}$ . For any given  $R, A \in \mathbb{R}^{n \times n}$  and  $\tilde{x}, b \in \mathbb{R}^n$ ,  $\|I - RA\| < 1$  implies  $A$  (and  $R$ ) to be non-singular and

$$(4.1) \quad \|\tilde{x} - A^{-1}b\|_\infty = \|(I - (I - RA))^{-1}R(A\tilde{x} - b)\|_\infty \leq \frac{\|R(A\tilde{x} - b)\|_\infty}{1 - \|I - RA\|_\infty} .$$

Moreover, we may use  $(I - E)^{-1} = (I - E^2)^{-1}(I + E)$  for  $E := I - RA$  to deduce

$$(4.2) \quad \|\tilde{x} - A^{-1}b\|_\infty \leq \frac{\|(I + E)R(A\tilde{x} - b)\|_\infty}{1 - \|I - RA\|_\infty^2} .$$

If  $\|I - RA\|_\infty$  is not close to one, then the accuracy of the bounds is determined by the size of the residual  $\|A\tilde{x} - b\|_\infty$ . The computation of the residual is, of course, subject to heavy cancellation. However, when

computed with Algorithm 3.4 (`Dot2Near`) presented in the previous section we may expect accurate error bounds for the maximum error of  $\tilde{x}$  to the exact solution  $A^{-1}b$ .

Both (4.1) and (4.2) are uniform, normwise error bounds for all entries of the approximation  $\tilde{x}$ . If the entries of  $\tilde{x}$  differ largely in magnitude, it is superior to use the following entrywise error estimate by T. Yamamoto [54].

**THEOREM 4.1.** *Let  $A, R \in \mathbb{R}^{n \times n}$  and  $b, \tilde{x} \in \mathbb{R}^n$  be given. Define  $E := I - RA$  and  $\delta := R(A\tilde{x} - b)$ , and assume  $\|E\|_\infty < 1$ . Then  $A$  is non-singular and*

$$(4.3) \quad |\tilde{x} - A^{-1}b| \leq |\delta| + \frac{\|\delta\|_\infty}{1 - \|E\|_\infty} \cdot |E|e ,$$

where  $e := (1, \dots, 1)^T \in \mathbb{R}^n$ .

**PROOF.** For  $u, v, x \in \mathbb{R}^n$  it holds  $|u^T v| \leq \|u\|_1 \|v\|_\infty$ , and therefore  $|Ex| \leq \|x\|_\infty \cdot |E|e \in \mathbb{R}^n$ . Hence  $(I - E)^{-1} = I + E(I - E)^{-1}$  yields

$$\begin{aligned} |\tilde{x} - A^{-1}b| &= |(I - E)^{-1}R(A\tilde{x} - b)| = |(I + E(I - E)^{-1})\delta| \leq |\delta| + \|(I - E)^{-1}\delta\|_\infty \cdot |E|e \\ &\leq |\delta| + \frac{\|\delta\|_\infty}{1 - \|E\|_\infty} \cdot |E|e . \end{aligned}$$

□

Next we describe how to compute a rigorous upper bound of the right hand side in (4.3) in rounding to nearest. It is applicable up to condition numbers of about  $\mathbf{u}^{-1}$ , and it is valid under all circumstances, also in the presence of underflow.

A result in overflow is rounded in IEEE 754 to  $\pm\infty$ , which means for all of the following algorithms that the final result contains either  $\infty$  or NaN. Therefore we may safely assume that no intermediate overflow occurs because this is monitored in the final result.

**4.1. Bounds for summation and dot product in rounding to nearest (accumulation in working precision).** Algorithm 3.4 (`Dot2Near`) estimates the error of a dot product in rounding to nearest. However, often a less accurate dot product, only accumulated in working precision, is sufficient. Let  $x, y \in \mathbb{F}^n$  be given, and denote by  $\tilde{s}$  the result of the floating-point approximation of  $x^T y$  computed by a standard for-loop. Then, provided no underflow occurs, the classical Wilkinson estimate [20] is

$$(4.4) \quad |\tilde{s} - x^T y| \leq \gamma_n |x^T| |y| \quad \text{for } n\mathbf{u} < 1 ,$$

where  $\gamma_n := n\mathbf{u}/(1 - n\mathbf{u})$ . Since the computation of  $\gamma_n$  in floating-point causes again rounding errors to be controlled, the code for rigorous estimates is unwieldy, in particular if underflow is allowed. Moreover, the right hand side is not known.

Fortunately there is a simple way to avoid the nasty  $\gamma_n$  terms but to obtain nevertheless rigorous error estimates *including possible underflow*. For given  $A \in \mathbb{F}^{m \times k}$  and  $B \in \mathbb{F}^{k \times n}$  consider the following Algorithm 4.2, where `realmin` is the Matlab constant denoting the smallest positive normalized floating-point number.<sup>2</sup>

**ALGORITHM 4.2.** *Rigorous error bound  $|A \cdot B - \text{res}| \leq \text{err}$ , also in the presence of underflow.*

```
function [res,err] = DotErr(A,B)
    res = A*B;                % approximation of product
    D = abs(A)*abs(B);        % product of absolute values
    U = ufp(D);                % unit in the first place of D
    k = size(A,2);            % inner dimension
    err = (k+2)*(eps/2*U) + 1.5*realmin; % error bound for res
```

<sup>2</sup>The rounding error unit `eps` in Matlab is  $2\mathbf{u}$ .

We assume that the products  $\mathbf{A}*\mathbf{B}$  and  $\text{abs}(\mathbf{A})*\text{abs}(\mathbf{B})$  is executed in the same order of evaluation. The algorithm uses the “unit in the first place”  $\text{ufp}(r)$  of  $r \in \mathbb{R}$  defined by

$$(4.5) \quad \text{ufp}(r) := 2^{\lceil \log_2 |r| \rceil},$$

with  $\text{ufp}(0) := 0$ , which may be smaller than  $|r|$  by up to a factor of 2. Algorithm 3.5 in [44], which we repeat for convenience, computes the unit in the first place in floating-point rounding to nearest without branch.

ALGORITHM 4.3. *Unit in the first place of a floating-point number, vector or matrix.*

```
function res = ufp(x)
    q = (1/eps+1)*x;           % eps = 2^(-52) in double precision
    res = abs(q-(1-eps/2)*q); % unit in the first place of x
```

A possible overflow in the computation of  $\mathbf{q}$  may be avoided by some scaling. The algorithm works for vectors and matrices as well. Concerning Algorithm 4.2, we proved in [45] the following.

THEOREM 4.4. *Let  $\mathbf{A} \in \mathbb{F}^{m \times k}$  and  $\mathbf{B} \in \mathbb{F}^{k \times n}$  with  $(k+2)\mathbf{u} \leq 1$  be given, and let  $\mathbf{res}$  and  $\mathbf{err}$  be the quantities computed by Algorithm 4.2 (DotErr). Then, also in the presence of underflow,*

$$(4.6) \quad |\mathbf{A} \cdot \mathbf{B} - \mathbf{res}| \leq \mathbf{err}.$$

The factor  $k+2$  in Algorithm 4.2 cannot be replaced by  $k+1$ .

The  $\text{ufp}$ -concept proved to be very useful in verifying the validity of floating-point estimates, and to obtain sharp estimates. I introduced this concept in [46] to prove the delicate estimations in there. Among the many properties of the unit in the first place (cf. [46]) we only need the following. For  $a, b \in \mathbb{F}$  and  $\circ \in \{+, -, \cdot, /\}$ , as in (1.2),  $\text{fl}(a \circ b)$  is the result of the floating-point approximation of  $a \circ b$ . Then the standard estimate of the error of  $\text{fl}(a \circ b)$  is improved into

$$(4.7) \quad f = \text{fl}(a \circ b) \Rightarrow f = a \circ b + \delta \quad \text{with} \quad |\delta| \leq \mathbf{u} \cdot \text{ufp}(a \circ b) \leq \mathbf{u} \cdot \text{ufp}(f) \leq \mathbf{u}|f|$$

for  $\circ \in \{+, -\}$ . If  $\circ \in \{\cdot, /\}$  and  $a \circ b$  is not in the underflow range, then (4.7) is true as well. If  $\circ \in \{\cdot, /\}$  and  $a \circ b$  is in the underflow range, then  $|\delta| \leq \frac{1}{2}\mathbf{eta}$ , where  $\mathbf{eta}$  denotes the smallest positive unnormalized floating-point number. Using  $\text{ufp}(x) \leq |x|$  for  $x \in \mathbb{R}$  this implies, for example,

$$(4.8) \quad a \circ b = (1 + \varepsilon_1) \cdot \text{fl}(a \circ b) = \text{fl}(a \circ b)/(1 + \varepsilon_2) \quad \text{with} \quad |\varepsilon_\nu| \leq \mathbf{u}, \quad 1 \leq \nu \leq 2$$

if  $\circ \in \{+, -\}$ , or if  $\circ \in \{\cdot, /\}$  and  $a \circ b$  is not in the underflow range.

This means that the error bound (4.6) computed by Algorithm 4.2 (DotErr) is not only rigorous and simple and valid if underflow occurs, but using the unit in the first place ( $\text{ufp}$ ) it may also be sharper than the classical Wilkinson-estimate (4.4) by up to a factor 2. For numerical evidence we compared in [45] the error estimate by DotErr and (4.4) for a matrix product  $\mathbf{R}*\mathbf{A}$ , where  $\mathbf{A}$  is randomly generated with specified condition number and  $\mathbf{R}=\text{inv}(\mathbf{A})$ . For dimensions from 10 to 1000 and condition numbers from 1 to  $\mathbf{u}^{-1}$ , the value of the bound by DotErr is uniformly about 0.7 times the classical Wilkinson-estimate (4.4).

For the sum of floating-point numbers rigorous bounds are computed by the following Algorithm 4.5.

ALGORITHM 4.5. *Rigorous error bound  $|\sum_{i=1}^n \mathbf{p}(i) - \mathbf{res}| \leq \mathbf{err}$ , also in the presence of underflow.*

```
function [res,err] = SumErr(p)
    n = length(p);           % number of summands
    res = sum(p);           % approximation of sum
    D = sum(abs(p));        % sum of absolute values
    U = ufp(D);             % unit in the first place of D
    err = (n-1)*(eps/2*U); % error bound for d
```

**THEOREM 4.6.** [45] *Let  $p \in \mathbb{F}^n$  with  $n\mathbf{u} \leq 1$  be given, and let  $\mathbf{res}$  and  $\mathbf{err}$  be the quantities computed by Algorithm 4.5 (SumErr). Then, also in the presence of underflow,*

$$(4.9) \quad \left| \sum_{i=1}^n p(i) - \mathbf{res} \right| \leq \mathbf{err} .$$

*The estimation is sharp for all  $n \leq \mathbf{u}^{-1}$ .*

Underflow is no problem for summation since a floating-point sum with a result in the underflow range is exact. For later usage we note that for  $\mathbf{U}$  as computed in SumErr

$$(4.10) \quad \left| \sum_{i=1}^n p(i) - \mathbf{res} \right| \leq (n-1) \cdot \mathbf{u} \cdot \mathbf{U}$$

is always true, also for  $n > \mathbf{u}^{-1}$ . If  $\mathbf{U}$  is in the underflow range, the right hand side can be replaced by zero. Besides, we need an upper bound for the sum of nonnegative floating-point numbers as computed by the following algorithm.

**ALGORITHM 4.7.** *Rigorous error bound  $\sum_{i=1}^n p(i) \leq \mathbf{ubnd}$ , also in the presence of underflow, for nonnegative  $p \in \mathbb{F}^n$ .*

```
function ubnd = SumPosBnd(p)
    n = length(p);           % number of summands
    S = sum(p);              % approximation of sum
    ubnd = S + (n+1)*(0.5*eps*ufp(S)); % upper bound for sum
```

**THEOREM 4.8.** *Let nonnegative  $p \in \mathbb{F}^n$  with  $(n+1)\mathbf{u} \leq 1$  be given, and let  $\mathbf{ubnd}$  be the quantity computed by Algorithm 4.7 (SumPosBnd). Then, also in the presence of underflow,*

$$(4.11) \quad \sum_{i=1}^n p(i) \leq \mathbf{ubnd} .$$

**PROOF.** If  $\mathbf{U} := \mathbf{ufp}(S)$  is in the underflow range, then, because the summands are nonnegative, all intermediate sums are in the underflow range and no error occurs at all. If  $\mathbf{U}$  is not in the underflow range, then there is no error in the floating-point computation of  $\mathbf{delta} := (n+1) * (0.5 * \mathbf{eps} * \mathbf{ufp}(S))$  because  $\mathbf{U}$  is a power of 2,  $\mathbf{eps}/2 = \mathbf{u}$  and  $(n+1)\mathbf{u} \leq 1$ . Hence  $\mathbf{delta} = (n+1) \cdot \mathbf{u} \cdot \mathbf{U} \leq \mathbf{U} \leq S$  and  $\mathbf{ufp}(S + \mathbf{delta}) \leq \mathbf{ufp}(2S) = 2\mathbf{ufp}(S) = 2\mathbf{U}$ , so that (4.7) implies

$$\mathbf{ubnd} \geq S + \mathbf{delta} - \mathbf{u} \cdot \mathbf{ufp}(S + \mathbf{delta}) \geq S + (n-1) \cdot \mathbf{u} \cdot \mathbf{U} \geq \sum_{i=1}^n p(i)$$

by (4.10) and the nonnegativity of the summands. □

Before we come to the computation of a rigorous upper bound of the right hand side in (4.3) in rounding to nearest, we have to prove the correctness of the error bound of Algorithm 3.4 (Dot2Near). We need the following auxiliary lemma.

**LEMMA 4.9.** *Let  $a, b \in \mathbb{F}$  be given, and let  $\mathbf{s}$  and  $\mathbf{sabs}$  be computed by the following Matlab commands:*

```
s = a+b;           % floating-point approximation
signal = 1+eps;    % successor of 1
sabs = signal*abs(s); % upper bound for |a+b|
```

*Then, also in the presence of underflow,*

$$(4.12) \quad |a + b| \leq \mathbf{sabs} .$$

PROOF. If  $\mathbf{s}$  is in the underflow range, then  $\mathbf{s} = \mathbf{a} + \mathbf{b}$  and the result follows. Otherwise  $|\mathbf{a} + \mathbf{b}| \leq |\mathbf{s}| + \mathbf{u} \cdot \text{ufp}(\mathbf{s}) < |\mathbf{s}| + 2\mathbf{u} \cdot \text{ufp}(\mathbf{s})$  by (4.7). But  $|\mathbf{s}| + 2\mathbf{u} \cdot \text{ufp}(\mathbf{s})$  is the successor of  $\mathbf{s}$  and therefore a floating-point number, so that  $|\mathbf{s}| + 2\mathbf{u} \cdot \text{ufp}(\mathbf{s}) \leq (1 + 2\mathbf{u})|\mathbf{s}| = \mathbf{sigma1} \cdot |\mathbf{s}|$  and the monotonicity of the rounding proves the result.  $\square$

**4.2. Correctness of the error bound of Algorithm 3.4 (Dot2Near).** For the proof of correctness of the error term, we need the following error estimation of recursive floating-point summation. As in (1.2) denote by  $\text{fl} : \mathbb{R} \rightarrow \mathbb{F}$  the rounding such that  $\text{fl}(a \circ b) \in \mathbb{F}$  is nearest to  $a \circ b \in \mathbb{R}$  for  $a, b \in \mathbb{R}$  and  $\circ \in \{+, -, \cdot, / \}$ . This is the definition in the IEEE 754 floating-point standard. For a given vector  $p \in \mathbb{F}^n$  of floating-point numbers, let  $\tilde{s}_n$  and  $\tilde{S}_n$  be computed by the following algorithm:

ALGORITHM 4.10. *Recursive summation with error estimation.*

$$\begin{aligned} \tilde{s}_1 &= p_1; \quad \tilde{S}_1 = |p_1| \\ \text{for } k &= 2 : n \\ s_k &= \tilde{s}_{k-1} + p_k; \quad \tilde{s}_k = \text{fl}(s_k) \\ S_k &= \tilde{S}_{k-1} + |p_k|; \quad \tilde{S}_k = \text{fl}(S_k) \end{aligned}$$

Then (4.7) implies the following rigorous error estimate, also in the presence of underflow:

$$(4.13) \quad \left| \tilde{s}_n - \sum_{i=1}^n p_i \right| \leq (n-1)\mathbf{u} \cdot \text{ufp}(\tilde{S}_n) \quad \left[ \leq (n-1)\mathbf{u}\tilde{S}_n \right].$$

To analyze Algorithm 3.4 (Dot2Near), we first rephrase it to distinguish intermediate results in the loop and to identify the true and the rounded result of the intermediate operations. It suffices to analyze a single dot product  $x^T y$  for  $x, y \in \mathbb{F}^n$  with  $(n+2)\mathbf{u} \leq 1$ . Then, using (3.1), Algorithm 3.4 is equivalent to the following:

$$\begin{aligned} p_1 + \tilde{e}_1 &= x_1 \cdot y_1 + \eta_1; \quad E_1 = |\tilde{e}_1| \\ 2 \leq i \leq n & \left\{ \begin{array}{l} h_i + r_i = x_i \cdot y_i + \eta_i \\ p_i + q_i = p_{i-1} + h_i \\ t_i = q_i + r_i; \quad \tilde{t}_i = \text{fl}(t_i) \\ e_i = \tilde{e}_{i-1} + \tilde{t}_i; \quad \tilde{e}_i = \text{fl}(e_i) \\ E_i = E_{i-1} + |\tilde{t}_i|; \quad \tilde{E}_i = \text{fl}(E_i) \end{array} \right. \\ \rho &= p_n + \tilde{e}_n; \quad \mathbf{res} = \text{fl}(\rho) \end{aligned}$$

It follows

$$(4.14) \quad \begin{aligned} x^T y + \sum_{i=1}^n \eta_i &= p_1 + \tilde{e}_1 + \sum_{i=2}^n (h_i + r_i) = p_n + \tilde{e}_1 + \sum_{i=2}^n (q_i + r_i) \\ &= p_n + \tilde{e}_1 + \sum_{i=2}^n \tilde{t}_i + \sum_{i=2}^n (t_i - \tilde{t}_i). \end{aligned}$$

Now  $\tilde{e}_n$  is the floating-point sum of  $\tilde{e}_1 + \sum_{i=2}^n \tilde{t}_i$  and  $\tilde{E}_n$  is the floating-point sum of the absolute values, so (4.10) gives

$$(4.15) \quad \left| \tilde{e}_1 + \sum_{i=2}^n \tilde{t}_i - \tilde{e}_n \right| \leq (n-1)\mathbf{u} \cdot \text{ufp}(\tilde{E}_n).$$

Furthermore,  $|t_i - \tilde{t}_i| \leq \mathbf{u}|\tilde{t}_i|$ , so again (4.10) and  $(n+2)\mathbf{u} \leq 1$  imply

$$(4.16) \quad \sum_{i=2}^n |t_i - \tilde{t}_i| \leq \mathbf{u} \cdot \sum_{i=2}^n |\tilde{t}_i| \leq \mathbf{u} \cdot [\tilde{E}_n + (n-2)\mathbf{u} \cdot \text{ufp}(\tilde{E}_n)] \leq 3\mathbf{u} \cdot \text{ufp}(\tilde{E}_n).$$

Combining (4.14), (4.15) and (4.16) with (3.1) yields

$$(4.17) \quad |x^T y - (p_n + \tilde{e}_n)| \leq 3\mathbf{neta} + (n+2)\mathbf{u} \cdot \text{ufp}(\tilde{E}_n).$$

Now  $|\mathbf{res} - (p_n + \tilde{e}_n)| = |\text{fl}(p_n + \tilde{e}_n) - (p_n + \tilde{e}_n)| \leq \mathbf{u} \cdot \text{ufp}(\mathbf{res})$  and  $\mathbf{realmin} = \frac{1}{2}\mathbf{u}^{-1}\mathbf{eta}$  give

$$(4.18) \quad |x^T y - \mathbf{res}| \leq \mathbf{u} \cdot \text{ufp}(\mathbf{res}) + \max(6n\mathbf{u}, 1) \cdot \mathbf{realmin} + (n+2)\mathbf{u} \cdot \text{ufp}(\tilde{E}_n).$$

If  $\mathbf{res}$  is in the underflow range, then there is no rounding error in the addition  $p_n + \tilde{e}_n$ , so that  $\mathbf{res} = p_n + \tilde{e}_n$ , and  $\mathbf{u} \cdot \text{ufp}(\mathbf{res})$  in (4.18) can be omitted. If  $\tilde{E}_n$  as a sum of nonnegative numbers is in the underflow range, then all  $\tilde{t}_i$  are in the underflow range, and there is no error in the sums  $q_i + r_i$ . Hence  $t_i = \tilde{t}_i$  for  $2 \leq i \leq n$ , the right hand side of (4.16) can be replaced by zero so that  $(n+2)\mathbf{u} \cdot \text{ufp}(\tilde{E}_n)$  in (4.18) can be omitted.

If neither  $\mathbf{res}$  nor  $\tilde{E}_n$  is in the underflow range, then the floating-point computation of the three summands in (4.18) does not cause a rounding error, so that in any case only the rounding errors in the two additions in the right hand side of (4.18) have to be taken care of. Again applying (4.10) to this floating-point sum of three non-negative summands proves<sup>3</sup>

$$(4.19) \quad |x^T y - \mathbf{res}| \leq \mathbf{err0} + 2\mathbf{u} \cdot \text{ufp}(\mathbf{err0}) \leq \text{fl}(\mathbf{err0} + 3\mathbf{u} \cdot \text{ufp}(\mathbf{err0})) = \mathbf{err}.$$

This proves the error estimate (3.3) in Theorem 3.5.

The presented algorithms for computing error bounds are suitable for a compiled programming language such as C or Fortran or a Matlab mex-file; a pure Matlab implementation suffers significantly from interpretation overhead.

### 4.3. Rigorous error bounds for linear systems in rounding to nearest (up to $\text{cond}(A) \lesssim \mathbf{u}^{-1}$ ).

In the remaining of this subsection we will prove that the quantity  $\mathbf{err}$  computed by the following Algorithm 4.11 (`LssErrBndNear0`) is an upper bound of the right hand side in (4.3) in Theorem 4.1 and thus an error bound for the approximate solution  $\mathbf{xs}$ . For didactical reasons we first state this preliminary version of the final Algorithm 4.16 (`LssErrBndNear`).

ALGORITHM 4.11. *Rigorous error bound for the solution of a linear system  $Ax = b$ , preliminary version.*

```
function [xs,err] = LssErrBndNear0(A,b)
    err = NaN(size(b));           % initialize result
    R = inv(A);                   % approximate inverse of A
    xs = ResidIter(A,b,R);        % approximate solution
    n = size(A,2);                % dimension of the linear system
    [D,eD] = Dot2Near([A b],[xs;-1]); % error bound D+/-eD for residual
    [aRD,eRD] = DotErr(R,D);      % error bound aRD+/-eRD of R*D
    [aReD,eReD] = DotErr(abs(R),eD); % error bound aReD+/-eReD of |R|*eD
    dd = abs(aRD) + eRD + aReD + eReD; % not yet upper bound of |R*(A*xs-b)|
    delta = dd + 2.5*(eps*ufp(dd)); % upper bound of |R*(A*xs-b)|
    [aRA,eRA] = DotErr(R,A);      % bounds for R*A
    RA_I = (1+eps)*abs(aRA-eye(n)); % upper bound of |aRA-I|
    E = (1+eps)*(RA_I+eRA);        % upper bound of |R*A-I|
    aE1 = sum(E,2);                % approximation of |R*A-I|*ones(n,1)
    uE1 = aE1 + (n+1)*(0.5*eps*ufp(aE1)); % upper bound of |R*A-I|*ones(n,1)
    Den = (1-max(uE1)) - 1.5*eps;  % lower bound of 1-||E||_inf
    if Den>0                       % algorithm successful
        err0 = (max(delta)/Den)*uE1 + realmin; % almost final error bound
        err = (1+eps)*(delta+err0);         % final error bound
    end
```

<sup>3</sup>The Matlab constant `eps` is  $2\mathbf{u}$ , so `eps` =  $0.5 * \mathbf{u}$  is used in Algorithm 3.4.

The command `Dot2Near([A b],[xs;-1])` in line 5 can obviously be replaced by a specialized algorithm `Dot2Near` taking into account the special structure. For ease of exhibition we refrain from that in this article.

In all of the following analysis we use the “verbatim”-font for the computed quantities, and all operations in the following analysis are the *exact, real* operations. For example,  $\mathbf{R} \cdot (\mathbf{A} \cdot \mathbf{xs} - \mathbf{b}) \in \mathbb{R}^n$  is the *true* correction of  $\mathbf{xs} \in \mathbb{F}^n$ . Note that taking the maximum and the absolute value of a vector does not cause any rounding error, so we may use  $\max(\mathbf{x})$  or  $\text{abs}(\mathbf{x})$  without causing any ambiguity, and similarly for the absolute value.

For the moment assume that the command `xs = ResidIter(A,b,R)` in line 3 computes some vector  $\mathbf{xs} \in \mathbb{F}^n$  of floating-point numbers. Of course, a good approximation to  $A^{-1}b$  is preferable, but for the proof of correctness of the bound this is irrelevant. Also assume  $(n+2)\mathbf{u} \leq 1$ . Then (3.3) in Theorem 3.5 implies

$$(4.20) \quad |\mathbf{D} - (\mathbf{A} \cdot \mathbf{xs} - \mathbf{b})| \leq \mathbf{eD} ,$$

and (4.6) in Theorem 4.4 gives

$$(4.21) \quad |\mathbf{aRD} - \mathbf{R} \cdot \mathbf{D}| \leq \mathbf{eRD} \quad \text{and} \quad |\mathbf{aReD} - |\mathbf{R}| \cdot \mathbf{eD}| \leq \mathbf{eReD} .$$

Define  $\delta := \mathbf{R} \cdot (\mathbf{A} \cdot \mathbf{xs} - \mathbf{b})$ , observe that `dd` is the sum of four nonnegative summands and that the computation of `delta` implements Algorithm 4.7 (`SumPosBnd`) for those four summands. Hence Theorem 4.8 is applicable and implies

$$(4.22) \quad |\delta| = |\mathbf{R} \cdot (\mathbf{A} \cdot \mathbf{xs} - \mathbf{b})| \leq |\mathbf{R} \cdot \mathbf{D}| + |\mathbf{R}| \cdot \mathbf{eD} \leq |\mathbf{aRD}| + \mathbf{eRD} + \mathbf{aReD} + \mathbf{eReD} \leq \mathbf{delta} .$$

Next

$$(4.23) \quad |\mathbf{aRA} - \mathbf{R} \cdot \mathbf{A}| \leq \mathbf{eRA} ,$$

and Lemma 4.9 implies

$$|\mathbf{aRA} - \mathbf{I}| \leq \mathbf{RA\_I} \quad \text{and} \quad \mathbf{RA\_I} + \mathbf{eRA} \leq \mathbf{E} .$$

Putting things together yields

$$(4.24) \quad |\mathbf{I} - \mathbf{R} \cdot \mathbf{A}| \leq |\mathbf{I} - \mathbf{aRA}| + |\mathbf{aRA} - \mathbf{R} \cdot \mathbf{A}| \leq \mathbf{E} .$$

Again using the code in Algorithm 4.7 (`SumPosBnd`) and Theorem 4.8 gives

$$(4.25) \quad |\mathbf{I} - \mathbf{R} \cdot \mathbf{A}| \cdot \mathbf{e} \leq \mathbf{E} \cdot \mathbf{e} \leq \mathbf{uE1} \quad \text{and} \quad \|\mathbf{I} - \mathbf{R} \cdot \mathbf{A}\|_\infty \leq \max(\mathbf{uE1}) .$$

To continue we first split the computation of `Den` into two parts by the Matlab statements

$$(4.26) \quad \mathbf{Den0} = 1 - \max(\mathbf{uE1}); \quad \mathbf{Den} = \mathbf{Den0} - 1.5 * \mathbf{eps};$$

Note that  $\mathbf{eps} = 2\mathbf{u}$ . If Algorithm 4.11 (`LssErrBndNear0`) finishes successfully, then  $1 \geq \mathbf{Den0} \geq \mathbf{Den} > 0$ , so that

$$(4.27) \quad \|\mathbf{I} - \mathbf{R} \cdot \mathbf{A}\|_\infty \leq \max(\mathbf{uE1}) < 1$$

and Theorem 4.1 is applicable. Abbreviate  $\mathbf{NE} := \max(\mathbf{uE1}) \geq \|\mathbf{I} - \mathbf{R} \cdot \mathbf{A}\|_\infty$  and suppose for the moment  $\mathbf{NE} > 0$ . Then  $\text{ufp}(1 - \mathbf{NE}) \leq \frac{1}{2}$  and (4.7) give

$$(4.28) \quad \mathbf{Den0} \leq 1 - \mathbf{NE} + \mathbf{u} \cdot \text{ufp}(1 - \mathbf{NE}) \leq 1 - \mathbf{NE} + \frac{1}{2}\mathbf{u} ,$$

so that again by (4.7) and  $\mathbf{eps} = 2\mathbf{u}$

$$(4.29) \quad \mathbf{Den} \leq \mathbf{Den0} - 3\mathbf{u} + \mathbf{u} \cdot \text{ufp}(\mathbf{Den0} - 3\mathbf{u}) \leq 1 - \mathbf{NE} + \frac{1}{2}\mathbf{u} - 3\mathbf{u} + \frac{1}{2}\mathbf{u} = 1 - \mathbf{NE} - 2\mathbf{u} .$$

Consider the Matlab statements

$$(4.30) \quad \mathbf{M} = \max(\mathbf{delta}); \quad \mathbf{F1} = \mathbf{M}/\mathbf{Den}; \quad \mathbf{F2} = \mathbf{F1} \cdot \mathbf{uE1}; \quad \mathbf{err0} = \mathbf{F2} + \mathbf{realmin};$$

Note that (4.22) implies  $\|\mathbf{R} \cdot (\mathbf{A} \cdot \mathbf{xs} - \mathbf{b})\|_\infty \leq \mathbf{M}$  and that the computations of  $\mathbf{err0}$  in (4.30) and in Algorithm 4.11 (`LssErrBndNear0`) are identical. The computation of  $\mathbf{delta}$  and Algorithm 4.2 (`DotErr`) imply  $\mathbf{delta} \geq \mathbf{eRD} \geq \mathbf{realmin}$ , and therefore  $\mathbf{M}/\mathbf{Den} \geq \mathbf{realmin}$  because  $\mathbf{Den} < 1$ . Hence the quotient  $\mathbf{F1}$  is not in the underflow range, and (4.8) implies

$$(4.31) \quad \mathbf{F1} \geq (1 - \mathbf{u}) \frac{\mathbf{M}}{\mathbf{Den}} \geq \frac{(1 - \mathbf{u})\mathbf{M}}{1 - \mathbf{NE} - 2\mathbf{u}} \geq \frac{(1 + \mathbf{u})\mathbf{M}}{1 - \mathbf{NE}}$$

with a little computation for the last inequality. The product  $\mathbf{F2}$  may be in the underflow range. If so, then

$$\mathbf{err0} \geq \mathbf{realmin} \geq \mathbf{F1} \cdot \mathbf{uE1},$$

and otherwise by (4.8),

$$\mathbf{err0} \geq \mathbf{F2} \geq \frac{\mathbf{F1} \cdot \mathbf{uE1}}{1 + \mathbf{u}}.$$

In any case (4.31) and (4.25) imply

$$(4.32) \quad \mathbf{err0} \geq \frac{\mathbf{F1} \cdot \mathbf{uE1}}{1 + \mathbf{u}} \geq \frac{\mathbf{M} \cdot \mathbf{uE1}}{1 - \mathbf{NE}} \geq \frac{\|\delta\|_\infty}{1 - \|\mathbf{I} - \mathbf{R} \cdot \mathbf{A}\|_\infty} \cdot |\mathbf{I} - \mathbf{R} \cdot \mathbf{A}| \cdot e.$$

Finally Lemma 4.9, (4.22) and (4.32) yield

$$(4.33) \quad \mathbf{err} \geq \mathbf{delta} + \mathbf{err0} \geq |\delta| + \frac{\|\delta\|_\infty}{1 - \|\mathbf{I} - \mathbf{R} \cdot \mathbf{A}\|_\infty} \cdot |\mathbf{I} - \mathbf{R} \cdot \mathbf{A}| \cdot e.$$

This finishes the proof for the case  $\mathbf{NE} > 0$ . If  $\mathbf{NE} = 0$ , then  $\mathbf{NE} = \max(\mathbf{uE1}) \geq \|\mathbf{I} - \mathbf{R} \cdot \mathbf{A}\|_\infty$  implies that  $\mathbf{R} = \mathbf{A}^{-1}$ , so that the right hand side in (4.3) reduces to  $|\delta|$ . Observing  $\mathbf{err} \geq |\delta|$  proves the following theorem.

**THEOREM 4.12.** *Let a set  $\mathbb{F}$  of floating-point numbers with relative rounding error unit  $\mathbf{u}$  together with floating-point operations complying with the IEEE 754 arithmetic standard [22, 23] be given.*

*Let  $\mathbf{res}$  and  $\mathbf{err}$  be the results of Algorithm 4.11 (`LssErrBndNear0`) applied to a matrix  $\mathbf{A} \in \mathbb{F}^{n \times n}$  and a vector  $\mathbf{b} \in \mathbb{F}^n$  with  $(n + 2)\mathbf{u} \leq 1$ . If the algorithm ends successfully, then  $\mathbf{A}$  is non-singular and*

$$(4.34) \quad \|\mathbf{xs} - \mathbf{A}^{-1}\mathbf{b}\| \leq \mathbf{err}.$$

It remains to explain  $\mathbf{xs} = \mathbf{ResidIter}(\mathbf{A}, \mathbf{b}, \mathbf{R})$  in line 3 of Algorithm 4.11 (`LssErrBndNear0`). The simplest is to perform just one residual iteration in working precision to ensure backward stability [48] of the computed  $\mathbf{xs}$  by the Matlab statements

$$(4.35) \quad \mathbf{xs} = \mathbf{R} \cdot \mathbf{b}; \quad \mathbf{xs} = \mathbf{xs} - \mathbf{R} \cdot (\mathbf{A} \cdot \mathbf{xs} - \mathbf{b});$$

and also to use `DotErr([A b], [xs; -1])` rather than `Dot2Near([A b], [xs; -1])` two lines later. In that case no extra-precise dot products are used at all and the quality of the inclusion is of the order  $\mathbf{u} \cdot \text{cond}(\mathbf{A})$ . A much better result and often inclusions of almost maximum accuracy, also for ill-conditioned matrices, is obtained by a residual iteration using `Dot2Near`. The stopping criterion implements numerical experience heuristically.

**ALGORITHM 4.13.** *Improvement of  $\mathbf{xs}$  by extra-precise residual iteration.*

```
function xs = ResidIter(A,b,R)
    xs = R*b; % first approximate solution
    normxs = norm(xs,inf); N = inf; % initialization of constants
```



TABLE 4.1

Ratio of measured computing times between Algorithm 4.11 (`LssErrBndNear0`) and the Matlab command `A\b`. Rigorous inclusions by `LssErrBndNear0` are computed using the Matlab implementation of Algorithm 3.4 (`Dot2Near`) and a C-implementation using mex-files.

	$n = 100$	$n = 200$	$n = 500$	$n = 1000$
Matlab <code>Dot2Near</code>	35.8	19.7	11.8	9.2
C-program <code>Dot2Near</code>	3.5	4.3	7.5	7.3

```

for iter=1:10                                % at most 10 residual iterations
    Nold = N;                                  % update constants
    d = R*Dot2([A b],[xs;-1]);                 % correction for xs
    N = norm(d,inf);                           % norm of correction
    if N<Nold, xs = xs-d; end                  % correction acceptable
    if ( ( iter==1 ) && ( N<1e-9*normxs ) ) || ( N<eps*normxs ) || ( N>=0.3*Nold )
        break                                  % stop iteration if well-conditioned
    end                                        % or no improvement
end
end
    
```

As has been mentioned, the implementation of `Dot2Near` in Matlab suffers severely from interpretation overhead. Nevertheless the main computational effort in Algorithm 4.11 (`LssErrBndNear0`) is the computation of the approximate inverse  $R = \text{inv}(A)$  and the bounds for the residual  $I - R * A$ , the latter requiring two matrix multiplications  $R * A$  and  $\text{abs}(R) * \text{abs}(A)$  in Algorithm 4.2 (`DotErr`). Therefore the theoretical time ratio between `LssErrBndNear0` and the Gaussian elimination by the Matlab command `xs = A\b` is 9.

In practice the ratio is better because matrix multiplication performs usually better than Gaussian elimination. Table 4.1 shows the ratio between the computing times of Algorithm 4.11 (`LssErrBndNear0`) and the Matlab command `xs = A\b` for different dimensions, the former first with the Matlab implementation of Algorithm 3.4 (`Dot2Near`), and second using a C-program and mex-file for `Dot2Near`.<sup>4</sup> The matrix and right hand side are chosen randomly. For ill-conditioned matrices, where more residual iterations are performed, the ratio increases because Matlab uses no residual iteration, apparently even not a single one in working precision.

As can be seen, using the Matlab implementation the ratio decreases because of the decreasing interpretation overhead; using the C-program there is an increase, possibly due to a (relatively) better performance of the Matlab command `A\b` and because the C-program is written straightforwardly without blocking.

Needless to say that Table 4.1 compares apples with oranges because Algorithm 4.11 (`LssErrBndNear0`) computes rigorous error bounds which are almost accurate to working precision, whereas the Matlab command `A\b` delivers only an approximation supposedly of quality  $\mathbf{u} \cdot \text{cond}(A)$ .

Replacing the computation `[aRA,eRA]=DotErr(R,A)` of the bounds for  $R \cdot A$  by `[aRA,eRA]=Dot2Near(R,A)` extends the range of applicability of Algorithm 4.11 (`LssErrBndNear0`), i.e. error bounds are computed for more ill-conditioned matrices. Moreover, the quality of the error bound improves for  $\text{cond}(A) \lesssim \mathbf{u}^{-1}$  at the price of increasing computing time. This will be shown in the next subsection.

**4.4. Improved bounds for ill-conditioned linear systems.** The successful computation of an error bound by Algorithm 4.11 (`LssErrBndNear0`) depends on  $\|I - RA\|_\infty < 1$ . For ill-conditioned  $A$ , i.e.  $\text{cond}(A) \approx \mathbf{u}^{-1}$ , it is not uncommon that  $\|I - RA\|_\infty \geq 1$  for a computed approximate inverse  $R$ , but some diagonal scaling rescues the situation by  $\|D^{-1}(I - RA)D\|_\infty < 1$ . Since  $\|I - RA\|_\infty = \| |I - RA| \|_\infty$ ,

<sup>4</sup>Many thanks to Prof. Ogita from Tokyo Woman's Christian University for providing the C- and mex-programs

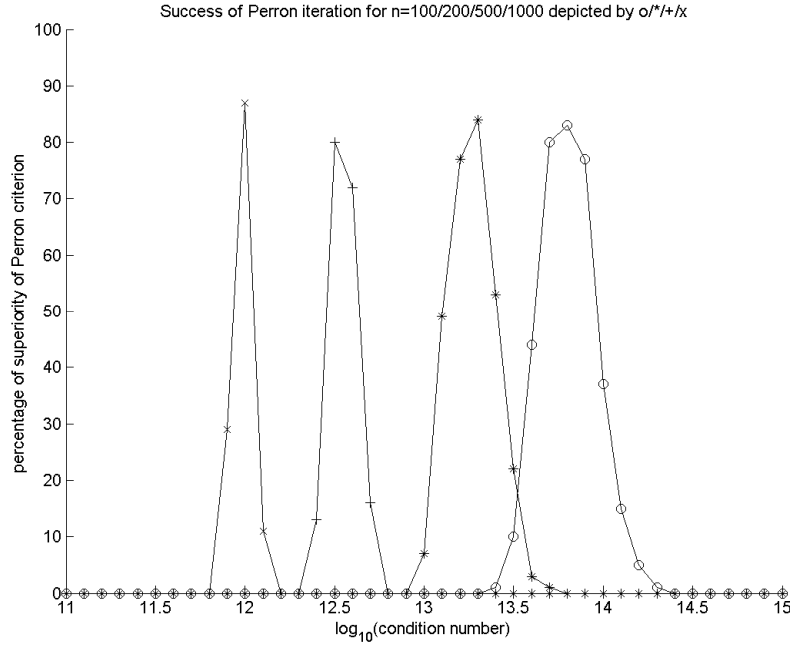


FIG. 4.1. Percentage that  $\|D^{-1}(I - RA)D\|_\infty < 1$  but  $\|I - RA\|_\infty \geq 1$  using  $D := \text{diag}(u)$  for  $u$  denoting the Perron vector of  $|I - RA|$  for dimensions  $n = 100$  (o),  $n = 200$  (\*),  $n = 500$  (+) and  $n = 1000$  (x) and condition numbers ranging from  $10^{13}$  to  $10^{15}$ .

the Perron vector of  $|I - RA|$  is the best choice.<sup>5</sup>

In Figure 4.1 the percentage of cases where  $\|D^{-1}(I - RA)D\|_\infty < 1$  but  $\|I - RA\|_\infty \geq 1$  using  $D := \text{diag}(u)$  for  $u$  denoting the Perron vector of  $|I - RA|$  for different dimensions and condition numbers is displayed. It shows that  $\|I - RA\|_\infty \geq 1$  begins to happen for condition numbers roughly starting with  $\mathbf{u}^{-1}/n$ , thus matching theory, whereas  $\|D^{-1}(I - RA)D\|_\infty < 1$  is still true for a little larger condition numbers.

For corresponding improved error bounds consider the following refinement of Theorem 4.1. It approaches the theoretically necessary condition that the spectral radius of  $|I - RA|$  is strictly less than one and shows the range of applicability of our approach.

**THEOREM 4.14.** *Let  $A, R \in \mathbb{R}^{n \times n}$  and  $b, \tilde{x}, u \in \mathbb{R}^n$  be given. Assume  $u > 0$ , and let  $D \in \mathbb{R}^{n \times n}$  be the diagonal matrix with  $u$  on the diagonal. Define  $E := I - RA$  and  $\delta := R(A\tilde{x} - b)$ , and assume  $\|D^{-1}|E|u\|_\infty < 1$ . Then  $A$  is non-singular and*

$$(4.36) \quad |\tilde{x} - A^{-1}b| \leq |\delta| + \frac{\|D^{-1}\delta\|_\infty}{1 - \|D^{-1}|E|u\|_\infty} \cdot |E|u.$$

**PROOF.** Define  $\hat{A} := AD$ ,  $\hat{R} := D^{-1}R$ , and  $\hat{x} := D^{-1}\tilde{x}$ . Then

$$\hat{\delta} := \hat{R}(\hat{A}\hat{x} - b) = D^{-1}\delta \quad \text{and} \quad \hat{E} := I - \hat{R}\hat{A} = D^{-1}ED.$$

Therefore (4.3) in Theorem 4.1 implies

$$|\tilde{x} - A^{-1}b| = D|\hat{x} - \hat{A}^{-1}b| \leq D|\hat{\delta}| + \frac{\|\hat{\delta}\|_\infty}{1 - \|\hat{E}\|_\infty} \cdot D|\hat{E}|e,$$

so that  $De = u$  and  $\|\hat{E}\|_\infty = \|\hat{E}|e\|_\infty$  finishes the proof.  $\square$

<sup>5</sup>In the rare case that  $|I - RA|$  is not irreducible, the Perron vector  $u$  of  $|I - RA| + \epsilon$  for small  $\epsilon > 0$  may be used to ensure  $u > 0$ , so that  $D = \text{diag}(u)$  is non-singular.

Next we show how this is used to compute rigorous bounds for linear systems in rounding to nearest. Suppose that some positive vector  $\mathbf{u} \in \mathbb{F}^n$  is given and that the quantities  $\delta$  and  $\mathbf{E}$  in Algorithm 4.11 (`LssErrBndNear0`) have been computed. Then consider the following Matlab commands.

```

Dd = delta./u; % approximation of D^-1*delta
uDd = max((1+eps)*Dd) + realmin; % upper bound of ||D^-1*delta||_inf
aEu = E*u; % approximation of |I-RA|u
e = (n+2)*(eps/2*ufp(aEu)) + 1.5*realmin; % bound |I-RA|u <= aEu+e
uEu = (1+eps)*(aEu+e); % upper bound of |I-RA|u
N0 = max(uEu./u); % approximation of ||D^-1|I-RA|u||_inf
N1 = (1+eps)*N0; % not yet bound of ||D^-1|I-RA|u||_inf
N2 = 1 - N1; % approximation of 1-||D^-1|I-RA|u||_inf
Den = N2 - 1.5*eps; % suitable denominator
if Den>0 % algorithm successful
    N3 = uDd./Den;
    err0 = N3*uEu + realmin; % upper bound of second term
    err = (1+eps)*(delta+err0); % final error term
end
    
```

To analyze the code we first need a version of Lemma 4.9 for multiplication and division.

LEMMA 4.15. *Let  $a, b \in \mathbb{F}$  be given, and let  $\mathbf{r}$  and  $\mathbf{rabs}$  be computed by the following Matlab commands:*

```

r = a*b; % floating-point approximation
rabs = (1+eps)*abs(r) + realmin; % upper bound for |a*b|
    
```

Then, also in the presence of underflow,

$$(4.37) \quad |\mathbf{a} \cdot \mathbf{b}| \leq \mathbf{rabs} .$$

The statement is also true when replacing multiplication by division in the first line of the Matlab code and in (4.37).

PROOF. In case of underflow,  $|\mathbf{a} \cdot \mathbf{b}|$  is bounded by  $\frac{1}{2}\mathbf{eta} < \mathbf{realmin}$ , and otherwise the result follows as in the proof of Lemma 4.9 because  $\mathbf{eps} = 2\mathbf{u}$  implies that  $1+\mathbf{eps}$  is the successor of 1.  $\square$

First note that (4.22) and Lemma 4.15 imply for positive  $\mathbf{u}$

$$(4.38) \quad \|D^{-1}\delta\|_{\infty} = \max_i \frac{\delta_i}{u_i} \leq \max_i \frac{\delta_i}{u_i} \leq \mathbf{uDd} .$$

Furthermore, (4.24), Theorem 4.4 and Lemma 4.9 imply

$$(4.39) \quad |I - \mathbf{R} \cdot \mathbf{A}| \cdot \mathbf{u} \leq \mathbf{E} \cdot \mathbf{u} \leq \mathbf{aEu} + \mathbf{e} \leq \mathbf{uEu} .$$

Assume for the moment that  $\max(\mathbf{N0}) \neq 0$ , and that in the Matlab statement  $\mathbf{N0} = \max(\mathbf{uEu}./\mathbf{u})$  no underflow occurs. Then also  $\mathbf{N1} \neq 0$ , and again using Lemma 4.15 yields

$$(4.40) \quad \mu := \|D^{-1}|I - \mathbf{R} \cdot \mathbf{A}| \cdot \mathbf{u}\|_{\infty} \leq \mathbf{N1} .$$

Assume an error bound is computed, i.e.  $\mathbf{Den} > 0$ . Then, as in the proof of Theorem 4.12,  $1 > \mathbf{N2} > \mathbf{Den} > 0$  so that  $\max\{\mathbf{ufp}(\mathbf{Den}), \mathbf{ufp}(1 - \mathbf{N1})\} \leq \frac{1}{2}$ ,  $\mathbf{eps} = 2\mathbf{u}$ , (4.7) and (4.40) show

$$(4.41) \quad \mathbf{Den} \leq \mathbf{N2} - 3\mathbf{u} + \mathbf{u} \cdot \mathbf{ufp}(\mathbf{Den}) \leq 1 - \mathbf{N1} + \mathbf{u} \cdot \mathbf{ufp}(1 - \mathbf{N1}) - \frac{5}{2}\mathbf{u} \leq 1 - \mu - 2\mathbf{u} ,$$

and  $\mathbf{uDd} \geq \mathbf{realmin}$  and (4.7) yield

$$(4.42) \quad \mathbf{N3} \geq (1 - \mathbf{u}) \frac{\mathbf{uDd}}{\mathbf{Den}} \geq (1 - \mathbf{u}) \frac{\mathbf{uDd}}{1 - \mu - 2\mathbf{u}} \geq (1 + \mathbf{u}) \cdot \frac{\mathbf{uDd}}{1 - \mu} .$$

If the multiplication  $N3 \cdot uEu$  causes underflow, then  $err0 \geq \text{realmin} \geq N3 \cdot uEu$ , and otherwise (4.38), (4.39) and (4.40) give

$$(4.43) \quad err0 \geq \frac{N3 \cdot uEu}{1 + \mathbf{u}} \geq \frac{uDd \cdot uEu}{1 - \mu} \geq \frac{\|D^{-1}\delta\|_\infty}{1 - \|D^{-1}|I - R \cdot A| \cdot \mathbf{u}\|_\infty} \cdot |(I - R \cdot A)| \cdot \mathbf{u}.$$

Finally (4.22) and Lemma 4.9 show that  $err$  is an upper bound of the right hand side in (4.36). It remains the case  $\max(N0) = 0$  or that in the Matlab statement  $N0 = \max(uEu./u)$  an underflow occurs. In that case (4.39) implies

$$(4.44) \quad \mu = \|D^{-1}|I - R \cdot A| \cdot \mathbf{u}\|_\infty = \max_i \frac{[|I - R \cdot A| \cdot \mathbf{u}]_i}{u_i} \leq \max_i \frac{uEu(i)}{u(i)} < \text{realmin},$$

so that  $N1 \leq \text{realmin}$ ,  $N2 = 1$  and  $Den = 1 - 3\mathbf{u}$ . As in (4.42) we conclude

$$(4.45) \quad N3 \geq (1 - \mathbf{u}) \frac{uDd}{Den} = \frac{1 - \mathbf{u}}{1 - 3\mathbf{u}} \cdot uDd \geq \frac{1 + \mathbf{u}}{1 - \text{realmin}} \cdot uDd > (1 + \mathbf{u}) \frac{uDd}{1 - \mu},$$

so that the lower bound of  $N3$  as in (4.42) is satisfied and we can continue as before. This proves that the computed vector  $err$  is indeed an upper bound for the right hand side in (4.36).

Both bounds (4.3) in Theorem 4.1 and (4.36) in Theorem 4.14 estimate the componentwise error of  $|\delta|$ , so the componentwise minimum of the bounds does as well. The additional effort to compute the bound in (4.36) is few  $\mathcal{O}(n^2)$  operations. Nevertheless it can be saved if the first bound in (4.3) is already accurate enough, i.e. if the relative error  $\max_i err_i/|xs|_i$  is small enough.

Putting things together, the following Algorithm 4.16 (`LssErrBndNear`) computes an approximation of the solution of a linear system together with a rigorous error bound. Note that we avoid to calculate the second bound based on Theorem 4.14 if  $\max_i err0_i/|xs|_i < 2\text{eps}$  because only  $err0$  can be improved and the offset  $\delta$  appears also in Theorem 4.1.

ALGORITHM 4.16. *Approximation  $xs$  and rigorous error bound  $err$  for the solution of a linear system  $Ax = b$ .*

```
function [xs,err] = LssErrBndNear(A,b)
... lines in LssErrBndNear0 before "if Den>0" ...
if Den>0
    % algorithm successful
    err0 = (max(delta)/Den)*uE1 + realmin; % almost final error bound
    err = (1+eps)*(delta+err0);          % final error bound
    if max(err0./abs(xs))<2*eps, return, end % sufficiently accurate bound
end
for i=1:2
    if i==1, u = PerronIter(E); end      % Perron vector for E
    if i==2, u = delta; end              % residual correction
    if u>0
        % vector u is suitable
        uDd = max((1+eps)*(delta./u)) + realmin; % upper bound of ||D^-1*delta||_inf
        aEu = E*u;                          % approximation of |I-RA|u
        e = (n+2)*(eps/2*ufp(aEu)) + 1.5*realmin; % bound |I-RA|u <= aEu+e
        uEu = (1+eps)*(aEu+e);              % upper bound of |I-RA|u
        Den = (1-(1+eps)*max(uEu./u))-1.5*eps; % suitable denominator
        if Den>0
            % choice of u successful
            err0 = (uDd./Den)*uEu + realmin; % upper bound of second term
            err = min(err,(1+eps)*(delta+err0)); % improved final error bound
        end
    end
end
```

```
end
end
```

The computed error bound is correct for any positive vector  $\mathbf{u}$ . To achieve accurate error bounds an approximation of the Perron vector of  $\mathbf{E}$ , the upper bound of  $|I - \mathbf{R} \cdot \mathbf{A}|$ , is preferable. To equilibrate the error terms, also  $\mathbf{u}=\mathbf{delta}$  proved sometimes to be a good choice. Since the additional effort is small, we use both. An approximation of the Perron vector is computed by the following algorithm.

ALGORITHM 4.17. *Approximation of the Perron vector of the nonnegative matrix  $E$ .*

```
function u = PerronIter(E)
    u = ones(size(E,2),1);           % first guess (1,...,1)^T
    for iter=1:10                   % at most 10 iterations
        v = E*u;                   % Perron iteration
        rho = v./u;                % ratio of iterates
        if min(rho)>=1, u=-1; return, end % matrix not convergent
        if max(rho)/min(rho)<1.05, break, end % sufficiently accurate
        u = v/norm(v,inf) + eps;    % update taking care of 0
    end
```

Note that for a nonnegative matrix  $E \in \mathbb{R}^{n \times n}$  and positive vector  $u \in \mathbb{R}^n$  Collatz [7] proved

$$(4.46) \quad \min_i \frac{(Eu)_i}{u_i} \leq \rho(E) \leq \max_i \frac{(Eu)_i}{u_i}$$

for  $\rho(\cdot)$  denoting the spectral radius, so that for  $\min(v./u) \geq 1$  the input matrix is not convergent and the iteration is stopped. Otherwise, adding  $\mathbf{eps}$  in the second last line ensures that the updated  $\mathbf{u}$  is positive.

**4.5. Rigorous error bounds for extremely ill-conditioned linear systems in rounding to nearest ( $\mathbf{up}$  to  $\text{cond}(A) \lesssim \mathbf{u}^{-2}$ ).** The obvious approach to calculate rigorous error bounds for extremely ill-conditioned linear systems is to combine Algorithm 2.1 (`LssIllcoApprox`) with an error bound by Theorems 4.1 or 4.14. However, this is not working.

Let a linear system  $Ax = b$  with  $\text{cond}(A) \gtrsim \mathbf{u}^{-1}$  be given. It seems reasonable to assume that the distance  $d := A^{-1}b - \mathbf{xs}$  of the exact solution  $A^{-1}b$  to its nearest floating-point vector  $\mathbf{xs}$  is of the order<sup>6</sup>  $\|d\| \sim \mathbf{u}\|A^{-1}b\|$  (indeed there are probabilistic arguments for that, see [43]). Moreover, for a matrix  $M \in \mathbb{R}^{n \times n}$  and a vector  $x \in \mathbb{R}^n$  which are not correlated we can expect  $\|Mx\|$  to be of the order  $n^{-1/2}\|M\|\|x\|$ , so that

$$(4.47) \quad \|A \cdot \mathbf{xs} - b\| = \|A \cdot d\| \sim \varphi' \cdot \mathbf{u}\|A\| \cdot \|A^{-1}b\| \sim \varphi \cdot \mathbf{u} \cdot \text{cond}(A) \cdot \|b\| \gtrsim \varphi \cdot \|b\|$$

for  $\varphi', \varphi$  not too far from 1. Note this is true without the presence of rounding errors in the computation of the residual. Let  $\mathbf{R}$  be an approximate inverse of  $A$ . Of course, mathematically  $\text{cond}(A) = \text{cond}(A^{-1})$ ; however,  $\mathbf{R}$  is a floating-point approximation, and rounding into floating-point has some smoothing effect [43], comparable to regularization, so that  $\text{cond}(\mathbf{R})$  can be expected to be not much larger than  $\mathbf{u}^{-1}$ . Hence

$$(4.48) \quad \mathbf{delta} := \|\mathbf{R} \cdot (A \cdot \mathbf{xs} - b)\| \gtrsim \varphi \cdot \mathbf{u}^{-1} \|b\| .$$

Note that this is true in general, no matter how accurate  $\mathbf{R}$  and  $\mathbf{xs}$  are. The reason is that both  $\mathbf{R}$  and  $\mathbf{xs}$  have floating-point entries and are thus of limited precision. Following Algorithm 2.1 (`LssIllcoApprox`),  $\mathbf{delta}$  is multiplied by  $\mathbf{Cinv}$ , the approximate inverse of  $\mathbf{R} \cdot \mathbf{A}$ . Under ideal circumstances,  $\text{cond}(\mathbf{Cinv}) = 1$ , but the approaches in Theorems 4.1 or 4.14 are expected to deliver useless error bounds of the size  $\mathbf{u}^{-1} \|b\|$ .

<sup>6</sup>Since matrix norms are equivalent, we do not to specify the norm for the following heuristic arguments.

To avoid this we construct an error bound depending directly on  $\mathbf{R} \cdot \mathbf{b}$ , without using an approximate solution  $\mathbf{xs}$ . Another way would be to store  $\mathbf{xs}$  in two parts as in [12]. However, we wanted to avoid that and to use strictly only (2.1) and (2.2) beyond ordinary floating-point arithmetic. Consider the following theorem.

**THEOREM 4.18.** *Let  $A, S \in \mathbb{R}^{n \times n}$  and  $b, u \in \mathbb{R}^n$  be given. Assume  $u > 0$  and let  $D \in \mathbb{R}^{n \times n}$  be the diagonal matrix with  $u$  on the diagonal. Define  $E := I - SA$  and assume  $\|D^{-1}|E|u\|_\infty < 1$ . Then  $A$  is non-singular and*

$$(4.49) \quad |A^{-1}b - S \cdot b| \leq \frac{\|D^{-1}S \cdot b\|_\infty}{1 - \|D^{-1}|E|u\|_\infty} \cdot |E|u .$$

For  $e := (1, \dots, 1)^T$  it follows in particular

$$(4.50) \quad |A^{-1}b - S \cdot b| \leq \frac{\|S \cdot b\|_\infty}{1 - \|E\|_\infty} \cdot |E|e .$$

**PROOF.** Using  $(I - E)^{-1} = D(I - D^{-1}ED)^{-1}D^{-1}$  and  $|Ex| \leq \|x\|_\infty \cdot |E|e \in \mathbb{R}^n$  for  $x \in \mathbb{R}^n$  as in the proof of Theorem 4.1, and  $De = u > 0$  and  $\|E\|_\infty = \||E|e\|_\infty$  it follows

$$\begin{aligned} |A^{-1}b - S \cdot b| &= |E(I - E)^{-1}Sb| = |ED \cdot (I - D^{-1}ED)^{-1}D^{-1}Sb| \\ &\leq \|(I - D^{-1}ED)^{-1}D^{-1}Sb\|_\infty \cdot |ED|e \\ &\leq \frac{\|D^{-1}S \cdot b\|_\infty}{1 - \|D^{-1}|E|u\|_\infty} \cdot |E|u . \end{aligned}$$

Setting  $u := e$  finishes the proof. □

Note that there is no restriction on  $S$ . Now the trick is, as explained following Observation 2.2, to define  $S := \mathbf{Cinv} \cdot \mathbf{R}$ , but rather than computing  $Sb$  to use  $\mathbf{Cinv} \cdot (\mathbf{R} \cdot \mathbf{b})$ . This has the appealing advantage to be faster and more accurate than  $(\mathbf{Cinv} \cdot \mathbf{R}) \cdot \mathbf{b}$ .

Corresponding error bounds in rounding to nearest based on Theorem 4.18 are not difficult to compute using the results of the previous subsections. First consider the following algorithm to compute error bounds for the product of two matrices where the second factor is afflicted with an error term.

**ALGORITHM 4.19.** *Rigorous bounds  $R \pm E$  of matrix products  $Q * \tilde{P}$  for  $P - eP \leq \tilde{P} \leq P + eP$ .*

```
function [R,E] = Prod2Bnd(Q,P,eP)
    [R,eR] = DotErr(Q,P);           % error bound R+/-eR of Q*P
    [aD,eD] = DotErr(abs(Q),eP);    % error bound aD+/-eD of |Q|*eP
    aE = eR + aD + eD;             % not yet upper bound of |R-Q*(P+/-eP)|
    E = aE + eps*ufp(aE);          % upper bound of |R-Q*(P+/-eP)|
```

Let  $Q \in \mathbb{F}^{m \times k}$  and  $P, eP \in \mathbb{F}^{k \times n}$  with  $eP \geq 0$  be given, and assume  $(n+2)\mathbf{u} \leq 1$ . Let  $\tilde{P} \in \mathbb{R}^{k \times n}$  with

$$P - eP \leq \tilde{P} \leq P + eP$$

be given. Then

$$|Q \cdot \tilde{P} - Q \cdot P| \leq |Q| \cdot eP .$$

The analysis of `DotErr` in Theorem 4.4 yields

$$|Q \cdot P - R| \leq eR \quad \text{and} \quad |Q| \cdot eP \leq aD + eD ,$$

hence

$$|Q \cdot \tilde{P} - R| \leq eR + aD + eD .$$

Therefore the analysis of `SumPosBnd` in Theorem 4.8 and  $\text{eps} = 2\mathbf{u}$  prove

$$(4.51) \quad |Q \cdot \tilde{P} - R| \leq E.$$

With these preliminaries we can state the algorithm to compute error bounds for extremely ill-conditioned linear systems in rounding to nearest.

ALGORITHM 4.20. *Approximation `xs` and rigorous error bound `err` for the solution of a linear system  $Ax = b$  for extremely ill-conditioned matrix  $A$ .*

```
function [xs,err] = LssIllcoErrBndNear(A,b)
    err = NaN(size(b)); % initialize result
    n = size(A,2); % dimension of the linear system
    R = inv(A); % approximate inverse
    while any(isinf(R(:))) || any(isnan(R(:)))
        R = inv(A.*(1+randn(n)*eps)); % inversion of perturbed matrix
    end
    [P,eP] = Dot2Near(R,A); % error bound P+/-eP for R*A
    Q = inv(P); % approximate inverse of P
    [aSA,eSA] = Prod2Bnd(Q,P,eP); % error bound aSA+/-eSA for Q*(R*A)
    [y,ey] = Dot2Near(R,b); % error bound y+/-ey for R*b
    [xs,eSb] = Prod2Bnd(Q,y,ey); % error bound xs+/-eSb for Q*(R*b)
    delta = (1+eps)*(abs(xs)+eSb); % upper bound of |Q*(R*b)|
    SA_I = (1+eps)*abs(aSA-eye(n)); % upper bound of |aSA-I|
    E = (1+eps)*(SA_I+eSA); % upper bound of |Q*(R*A)-I|
    aE1 = sum(E,2); % approximation of |Q*(R*A)-I|*ones(n,1)
    uE1 = aE1 + (n+1)*(0.5*eps*ufp(aE1)); % upper bound of |Q*(R*A)-I|*ones(n,1)
    Den = (1-max(uE1)) - 1.5*eps; % lower bound of 1-||E||_inf
    if Den>0 % algorithm successful
        err = (max(delta)/Den)*uE1 + realmin; % final error bound
        if max(err./abs(xs))<2*eps, return, end % bound sufficiently accurate
    end
    for i=1:2
        if i==1, u = PerronIter(E); end % Perron vector for E
        if i==2, u = delta; end % residual correction
        if u>0 % vector u is suitable
            uDd = max((1+eps)*(delta./u)) + realmin; % upper bound of ||D^-1*delta||_inf
            aEu = E*u; % approximation of |I-RA|u
            e = (n+2)*(eps/2*ufp(aEu)) + 1.5*realmin; % bound |I-RA|u <= aEu+e
            uEu = (1+eps)*(aEu+e); % upper bound of |I-RA|u
            Den = (1-(1+eps)*max(uEu./u))-1.5*eps; % suitable denominator
            if Den>0 % choice of u successful
                err = min(err, (uDd./Den)*uEu + realmin); % final error bound
            end
        end
    end
end
```

In the first lines the matrices  $R, P, eP, Q \in \mathbb{F}^{n \times n}$  are computed with

$$(4.52) \quad |P - R \cdot A| \leq eP.$$

Otherwise there are no assumptions on  $\mathbf{A}$ ,  $\mathbf{R}$  or  $\mathbf{Q}$ . Abbreviate  $\mathbf{S} := \mathbf{Q} \cdot \mathbf{R}$ , then (4.51) implies

$$(4.53) \quad |\mathbf{S} \cdot \mathbf{A} - \mathbf{aSA}| \leq \mathbf{eSA} .$$

Similarly,  $|\mathbf{y} - \mathbf{R} \cdot \mathbf{b}| \leq \mathbf{ey}$  and (4.51) imply

$$(4.54) \quad |\mathbf{S} \cdot \mathbf{b} - \mathbf{xs}| \leq \mathbf{eSb} ,$$

and Lemma 4.9 yields

$$(4.55) \quad |\mathbf{S} \cdot \mathbf{b}| \leq \mathbf{delta} \quad \text{and therefore} \quad \|\mathbf{S} \cdot \mathbf{b}\|_\infty \leq \max(\mathbf{delta}) .$$

As in (4.23)ff. in the analysis of Algorithm 4.11 (`LssErrBndNear0`) it follows

$$(4.56) \quad |I - \mathbf{S} \cdot \mathbf{A}| \leq \mathbf{E}, \quad \|I - \mathbf{S} \cdot \mathbf{A}\|_\infty \leq \max(\mathbf{uE1}) \quad \text{and finally} \quad \frac{\|\mathbf{S} \cdot \mathbf{b}\|_\infty}{1 - \|\mathbf{E}\|_\infty} \cdot |I - \mathbf{S} \cdot \mathbf{A}| \cdot \mathbf{e} \leq \mathbf{err}$$

provided  $\mathbf{Den} > 0$ , so that in this case  $\mathbf{err}$  is an upper bound of the right hand side in (4.50).<sup>7</sup> The remaining of the analysis is analogous to the proof of correctness of Algorithm 4.16 (`LssErrBndNear`) in (4.38)ff. with replacing  $\mathbf{R}$  by  $\mathbf{S}$ . Both bounds (4.49) and (4.50) are componentwise upper bounds for the error of  $\mathbf{S} \cdot \mathbf{b}$ , so the minimum of both is a valid bound as well. As before we use for  $\mathbf{u}$  both an approximation of the Perron vector of  $\mathbf{E}$  as well as  $\mathbf{delta}$ . If the computation of the first bound was not successful it is set to `NaN`, and observing that Matlab ignores `NaN`'s when computing a minimum, we proved the following theorem.

**THEOREM 4.21.** *Let a set  $\mathbb{F}$  of floating-point numbers with relative rounding error unit  $\mathbf{u}$  together with floating-point operations complying with the IEEE 754 arithmetic standard [22, 23] be given.*

*Let  $\mathbf{xs}$  and  $\mathbf{err}$  be the results of Algorithm 4.20 (`LssIllcoErrBndNear`) applied to a matrix  $\mathbf{A} \in \mathbb{F}^{n \times n}$  and a vector  $\mathbf{b} \in \mathbb{F}^n$  with  $(n+2)\mathbf{u} \leq 1$ . If the algorithm ends successfully, then  $\mathbf{A}$  is non-singular and*

$$(4.57) \quad |\mathbf{xs} - \mathbf{A}^{-1}\mathbf{b}| \leq \mathbf{err} .$$

*Computational evidence suggests that the algorithm ends successfully for condition numbers up to about  $\mathbf{u}^{-2}$ .*

As in Algorithm 2.1 (`LssIllcoApprox`) the precondition matrix  $\mathbf{R}$  is replaced by  $\mathbf{S} = \mathbf{Q} \cdot \mathbf{R}$ , and as before it is important to compute  $\mathbf{S} \cdot \mathbf{A}$  and  $\mathbf{S} \cdot \mathbf{b}$  in the order  $\mathbf{Q} \cdot (\mathbf{R} \cdot \mathbf{A})$  and  $\mathbf{Q} \cdot (\mathbf{R} \cdot \mathbf{b})$ , respectively. In the latter case this also reduces the computing time to  $\mathcal{O}(n^2)$  operations.

**5. 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. INTLAB [41] is the Matlab toolbox for reliable computing written by the author of this paper. To our knowledge there is no other publicly available Matlab code to compare with for rigorously solving linear systems using only rounding to nearest. If not stated otherwise, we always use the 2-norm condition number  $\|A^{-1}\|_2 \|A\|_2$ .

It is not obvious how to generate matrices with floating-point entries with condition number much larger than  $\mathbf{u}^{-1}$ . Of course, one may try the ‘‘usual suspects’’ like Hilbert matrices or the notoriously ill-conditioned Vandermonde matrices. However, as has been mentioned before, extremely ill-conditioned matrices hardly have a condition number much larger than  $\mathbf{u}^{-1}$  when rounded into floating-point. The situation is shown in Figure 5.1, where the condition number (computed by the symbolic toolbox in Matlab) of the (Matlab) floating-point approximation `hilb(n)` and the true Hilbert matrix is shown. This behavior is typical, with some exceptions mentioned below, that the condition number of non-singular floating-point matrices is roughly bounded by  $\mathbf{u}^{-1}$ .

<sup>7</sup>Note that unlike (4.3) there is no additive term in (4.50).



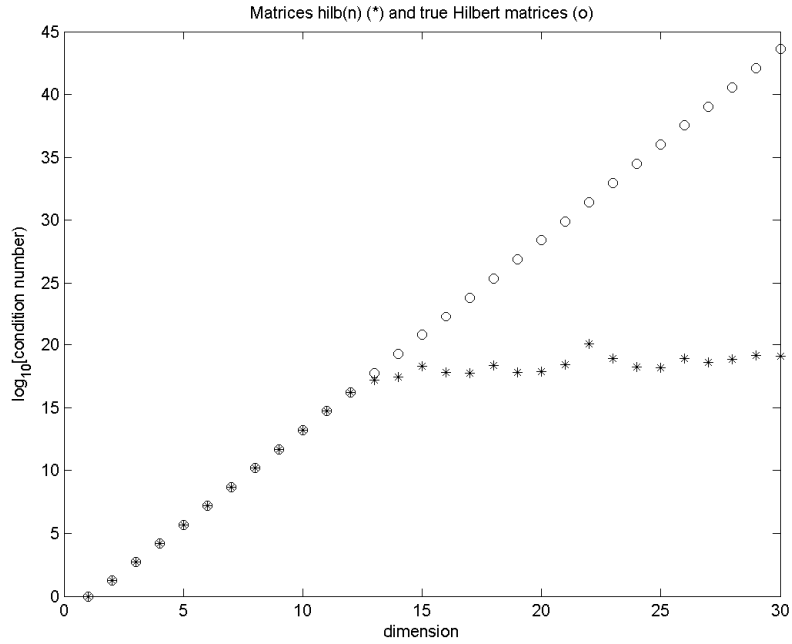


FIG. 5.1. *Logarithm of the condition number of  $\text{hilb}(n)$  (\*) and of the true Hilbert matrix  $H_{ij} = 1/(i + j - 1)$  (o).*

TABLE 5.1  
*Test matrices of small dimension (equilibrated).*

matrix	Matlab command	definition	$n_{max}$	$\text{cond}_{1 \leq n \leq 40}^{max}$
Pascal	<code>pascal(n)</code>	$\binom{i+j-1}{j-1}$	31	$6.0 \cdot 10^{31}$
Hilbert	<code>hilb(n)</code>	$1/(i+j-1)$	1	$6.3 \cdot 10^{19}$
scaled Hilbert		$\text{lcm}(1, \dots, 2n-1)/(i+j-1)$	21	$1.1 \cdot 10^{28}$
inverse Hilbert	<code>invhilb(n)</code>	defined by explicit formula	12	$4.8 \cdot 10^{37}$
Boothroyd		$\frac{\binom{n+i-1}{i-1} \cdot n \cdot \binom{n-1}{n-j}}{i+j-1}$	20	$1.1 \cdot 10^{30}$
Vandermonde	<code>vander(1:n)</code>	$i^{n-j}$	14	$4.7 \cdot 10^{61}$

Larger condition numbers are possible for some of the usual test matrices as long as they are exactly representable in floating-point. Well-known examples are listed in Table 5.1. Some of them are directly available in Matlab. The matrix entries like  $i^{n-j}$  for the Vandermonde matrix are computed in floating-point and coincide with the mathematical definition until dimension  $n \leq n_{max}$ . For larger dimensions the entries are corrupted by rounding errors and the matrix does not coincide with the mathematical definition. Usually this has a smoothing effect so that the condition number does not increase any more with increasing dimension. The maximum condition number for dimensions  $1 \leq n \leq 40$  is displayed in the last column of Table 5.1.

Usually it is preferable to equilibrate the input matrix to improve the condition number. This is also done in [12]. To avoid rounding errors we use the nearest power of 2 to equilibrate the input matrix. This works well except for Vandermonde matrices which show a very special and in some way strange behavior, see Section 5.4. Henceforth all matrices, also in Table 5.1, are equilibrated.

TABLE 5.2

Absolute and relative computing time in seconds for Algorithm 2.1 (`LssIllcoApprox`), the variable precision arithmetic package (`vpa`), Algorithm 4.20 (`LssIllcoErrBndNear`) and `sym(A,'f')\sym(b,'f')`. Based on an algorithm given in [42], matrices of condition number  $10^{25}$  are generated in INTLAB [41] by `randmat(n,1e25)` with random right hand side.

dimension	absolute computing time [sec]				relative computing time	
	Alg. 2.1	vpa	Alg. 4.20	sym	vpa/Approx	sym/ErrBnd
10	0.010	0.031	0.0040	0.039	3.2	9.7
20	0.019	0.089	0.0073	0.12	4.6	16.6
50	0.051	0.38	0.019	0.71	7.6	36.4
100	0.13	1.73	0.063	7.49	13.2	117.7
200	0.39	10.3	0.28	105.8	26.6	375.2
500	8.8	149.1	8.6	4140	16.9	479.1

Our algorithms are designed to solve general linear systems, not taking advantage of any structure of the matrix. Some of the test matrices mentioned so far do have a special structure, for example being totally nonnegative. Taking into account this property, many numerical problems can be solved with high relative accuracy of the result [25, 1, 13, 6]. A famous example is that the smallest singular value of the  $100 \times 100$  Hilbert matrix, which is of size  $10^{-151}$ , can be computed in IEEE 754 double precision to almost full accuracy. This was noted in [9], possibly the starting point for an extensive research on this topic.

A reason for this unexpected behavior, seemingly contradicting common perturbation analysis, is that it can be shown that some algorithms applied to such structured matrices perform only certain arithmetic operations, prohibiting in particular catastrophic cancellation [11].

Although this important research allows to solve a number of very ill-conditioned problems, certain structural properties of the matrix are mandatory. In particular a number of the “usual suspects” satisfy those properties. However, in contrast to our algorithms, those methods do not apply to general matrices.

**5.1. Results for `LssIllcoApprox`.** We start with some timing comparison. Note that all our algorithms are completely implemented in Matlab, and in particular the extra-precise accumulation of dot products suffers from interpretation overhead. Matlab offers in the symbolic toolbox some multi-precision arithmetic (`vpa`) for approximate calculations, and a rational arithmetic to compute the exact solution of linear systems. We compare the computing times for Algorithm 2.1 (`LssIllcoApprox`) (producing an approximation) with the variable precision arithmetic package (`vpa`), and Algorithm 4.20 (`LssIllcoErrBndNear`) (producing a rigorous error bound) with the precise result of the solution of the linear system  $Ax = b$  computed by `sym(A,'f')\sym(b,'f')`.<sup>8</sup>

As can be seen in Table 5.2 our algorithms are faster than the competitors from the symbolic toolbox. We note, however, that the comparison is not fair because `sym(A,'f')\sym(b,'f')` computes the exact, rational solution whenever the input matrix is non-singular, whereas `LssIllcoErrBnd` computes only error bounds and may fail. Apparently using the symbolic toolbox seems the only possibility in Matlab to compute rigorous results for extremely ill-conditioned matrices. Note that `LssIllcoApprox` improves the initial approximation by an extra-precise residual iteration whereas `LssIllcoErrBndNear` does not. This explains why for smaller dimension the former algorithm is slower than the latter.

Concerning the accuracy of the results, the right graph in Figure 2.1 shows already the performance of Algorithm 2.1 (`LssIllcoApprox`) for Pascal matrices for right hand sides  $b = \text{randn}(n, 1)$  and  $b = A * \text{randn}(n, 1)$ . 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

<sup>8</sup>The extra parameter ‘f’ ensures that the input  $A$  and  $b$  is converted into long format without error, respectively.

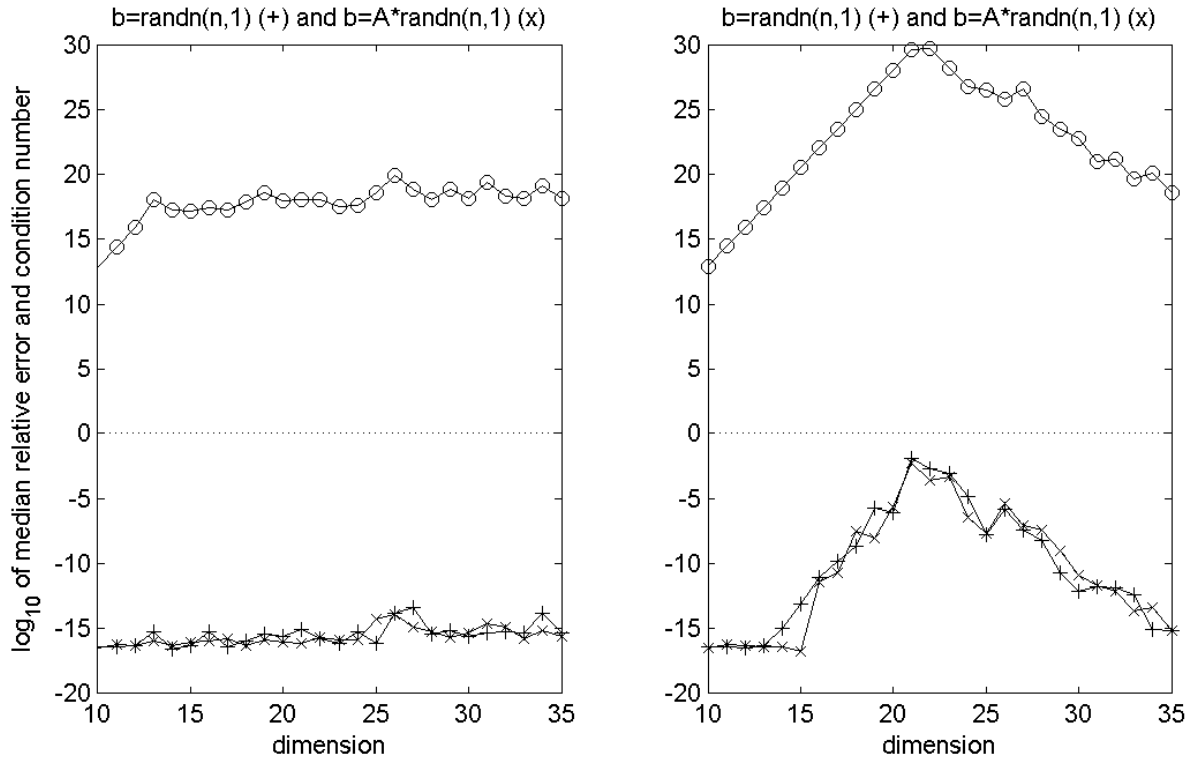


FIG. 5.2. Accuracy of Algorithm 2.1 (`LssI11coApprox`) for Hilbert and scaled Hilbert matrices as defined in Table 5.1. The upper parts show the condition number, the lower parts the relative error of the results.

mean zero and standard deviation one. As can be seen, the accuracy of the bounds decreases with increasing condition number. For all examples with maximum condition number up to  $6.0 \cdot 10^{31}$ , on the median at least 3 to 4 digits of the solution are correct.

For this and all of the following examples for Algorithm `LssI11coApprox` the number of residual iterations is mostly 1 or 2, in very few cases 3 iterations, but never more.

Next we use  $A = \text{hilb}(n)$  as defined in Matlab with entries approximating  $1/(i + j - 1)$ , and second the scaled Hilbert matrix with entries  $lcm(1, \dots, 2n - 1)/(i + j - 1)$ . Up to dimension  $n = 21$  the entries are an integer multiple of the original Hilbert matrix, for  $n > 21$  the entries are corrupted by rounding errors. Therefore, the scaled Hilbert matrices achieve much larger condition numbers, as the original Hilbert matrix, than  $\text{hilb}(n)$ , see Figure 5.1. In Figure 5.2 the median relative error of Algorithm 2.1 for the right hand sides  $b = \text{randn}(n, 1)$  and  $b = A * \text{randn}(n, 1)$  are printed in one graph each. Again the accuracy of the result is inverse proportional to the condition number.

In the left of Figure 5.3 we show the same results of Algorithm 2.1 (`LssI11coApprox`) for  $A = \text{invhilb}(n)$ , the approximation of the inverse Hilbert matrix for right hand sides  $b = \text{randn}(n, 1)$  and  $b = A * \text{randn}(n, 1)$ . Now the relative error is constantly less than  $10^{-14}$  although the condition number rises to almost  $10^{40}$ . A similar behavior is observed for Vandermonde matrices and will be discussed in Section 5.4.

As has been mentioned, in lines 9 and 14 in Algorithm 2.1 one might use extra-precise dot products by `Dot2Near` in the multiplication of the residual by `Cinv`. The results for matrices proposed by Boothroyd [5] with integer entries, where a checkerboard sign distribution produces its inverse, are displayed in the right of Figure 5.3. The results for right hand side  $b = \text{randn}(n, 1)$  are displayed; the results for  $b = A * \text{randn}(n, 1)$  are completely similar.

As can be seen there is not too much difference whether `Dot2Near` is used in lines 9 and 14 or not, and for

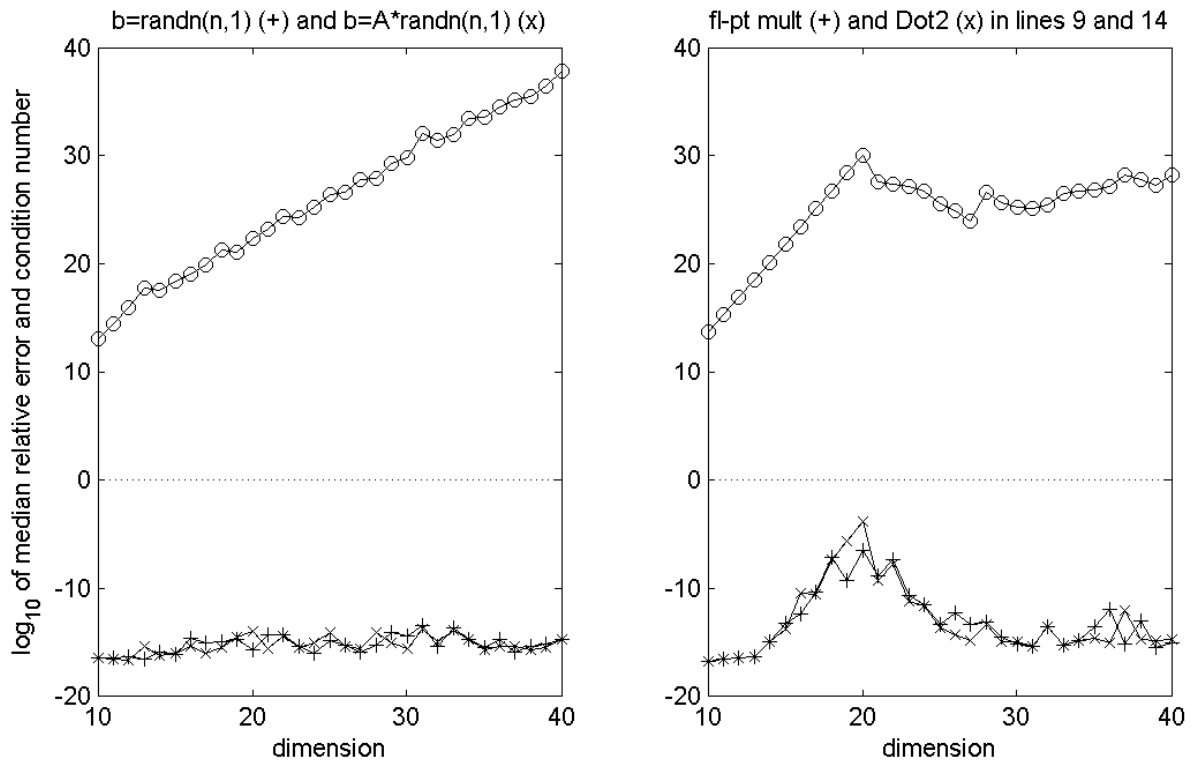


FIG. 5.3. Accuracy of Algorithm 2.1 (`LssIllcoApprox`) for inverse Hilbert and Boothroyd matrices as defined in Table 5.1. The upper parts show the condition number, the lower parts the relative error of the results.

other matrices and other right hand sides a similar behavior is observed. Therefore, this extra computing time in Algorithm 2.1 (`LssIllcoApprox`) is saved. Note again how the accuracy of the result corresponds to the condition number.

For larger dimensions it is very time consuming to compute an accurate solution to test against. But the results of Algorithm 4.20 (`LssIllcoErrBndNear`) are based on an approximate solution computed along the lines of Algorithm 2.1 (`LssIllcoApprox`); so from the following results for `LssIllcoErrBndNear` we can deduce that the approximations by Algorithm 2.1 (`LssIllcoApprox`) are at least as good as the computed bounds. Computational results for the latter are given in Section 5.3.

**5.2. Results for `LssErrBndNear`.** Next we discuss the results of Algorithm 4.16 (`LssErrBndNear`). Computing times are already given in Table 4.1. Those are for Algorithm 4.11 (`LssErrBndNear0`), but the difference to Algorithm 4.16 (`LssErrBndNear`) is only few  $\mathcal{O}(n^2)$  operations and negligible.

Concerning accuracy, we first consider the matrices in Table 5.1. Again we tested linear systems with right hand sides  $\mathbf{b}=\text{randn}(n,1)$  and  $\mathbf{b}=\mathbf{A}*\text{randn}(n,1)$ . As a typical example we display in Table 5.3 the results for Pascal matrices for dimensions 14 to 18. To save space we display for the other matrices from Table 5.1 only the results of Algorithm `LssErrBndNear` for the highest dimension it succeeded.

As can be seen the relative error is of the order  $\mathbf{u}$  and better. So the extra-precise residual iteration performs as expected, not only for approximations but also for the rigorous error bound. Also the maximum condition number for which `LssErrBndNear` successfully computes an error bound is not too far from  $\mathbf{u}^{-1}$ , which is due to the Perron iteration (see also Figure 4.1). However, this is partly due to the fact that the dimensions are small. Better results in this respect (using directed rounding) are obtained by Algorithm 4.2 (`LssErrBnd`) to be described in Section II of this paper. For all matrices except Vandermonde, this algorithm can handle one dimension larger.

TABLE 5.3

Median relative error of the bounds computed by Algorithm 4.16 (*LssErrBndNear*) for the matrices in Table 5.1.

matrix	$n$	$\text{cond}(A)$	$\mathbf{b}=\text{randn}(n,1)$	$\mathbf{b}=\mathbf{A}*\text{randn}(n,1)$
Pascal	14	$1.4 \cdot 10^{13}$	$5.1 \cdot 10^{-17}$	$1.2 \cdot 10^{-19}$
	15	$1.6 \cdot 10^{14}$	$3.3 \cdot 10^{-17}$	$1.2 \cdot 10^{-17}$
	16	$1.8 \cdot 10^{15}$	$4.8 \cdot 10^{-17}$	$1.0 \cdot 10^{-17}$
	17	$2.2 \cdot 10^{16}$	$2.0 \cdot 10^{-16}$	$6.3 \cdot 10^{-17}$
	18	$2.5 \cdot 10^{17}$	failed	failed
Hilbert	11	$7.4 \cdot 10^{12}$	$4.9 \cdot 10^{-17}$	$4.5 \cdot 10^{-17}$
inverse Hilbert	11	$1.1 \cdot 10^{13}$	$4.3 \cdot 10^{-17}$	$5.1 \cdot 10^{-17}$
scaled Hilbert	11	$8.7 \cdot 10^{12}$	$4.3 \cdot 10^{-17}$	$5.2 \cdot 10^{-17}$
Boothroyd	11	$5.0 \cdot 10^{13}$	$6.1 \cdot 10^{-17}$	$1.1 \cdot 10^{-19}$
Vandermonde	13	$3.0 \cdot 10^{14}$	$4.4 \cdot 10^{-17}$	$7.8 \cdot 10^{-17}$

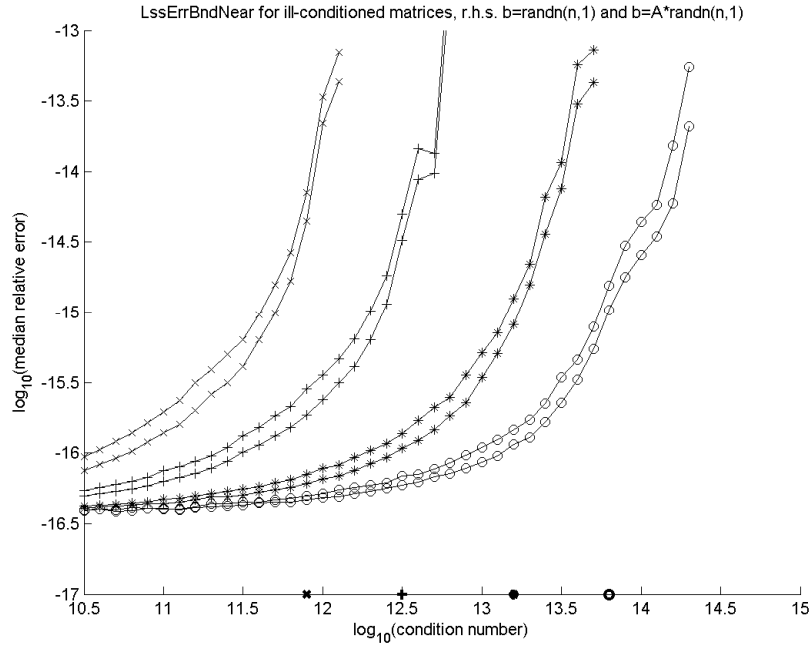


FIG. 5.4. Results of Algorithm 4.16 (*LssErrBndNear*) for ill-conditioned random matrices of dimension  $n=100$  ( $\circ$ ),  $n=200$  ( $*$ ),  $n=500$  ( $+$ ) and  $n=1000$  ( $\times$ ).

Next we treat ill-conditioned matrices of larger dimension. For dimensions up to  $\mathbf{u}^{-1}$  we may safely use `randsvd(n,cnd)` from the Matlab matrix gallery applying some random orthogonal transformation from the left and right to a diagonal matrix with specified singular values. As mentioned before, this approach is not applicable for condition numbers beyond  $\mathbf{u}^{-1}$  because of the inevitable presence of rounding errors.

In Figure 5.4 the median relative errors of all solution components of the result of Algorithm *LssErrBndNear* for right hand side  $\mathbf{b}=\text{randn}(n,1)$  and  $\mathbf{b}=\mathbf{A}*\text{randn}(n,1)$  over 100 samples is displayed for dimensions  $n = 100$  ( $\circ$ ),  $n = 200$  ( $*$ ),  $n = 500$  ( $+$ ) and  $n = 1000$  ( $\times$ ). Note that the maximum relative error over all samples is below  $10^{-13}$ , so that until failure for a condition number of roughly  $\mathbf{u}^{-1}/n$  the accuracy of the bounds is not too far from  $\mathbf{u}$ . For dimensions  $n = 100$ ,  $n = 200$ ,  $n = 500$  and  $n = 1000$  there is no failure in the 100 samples until condition numbers  $7.9 \cdot 10^{13}$ ,  $2.5 \cdot 10^{13}$ ,  $4.0 \cdot 10^{12}$  and  $1.6 \cdot 10^{12}$ , respectively, as depicted on the

x-axis in Figure 5.4. Again more ill-conditioned matrices can be treated in the same computing time using directed rounding as by Algorithm 4.2 (`LssErrBnd`) presented in Part II of this paper.

Finally we mention that one may compute the error bound of  $R \cdot A$  by `[aRA,eRA] = Dot2Near(R,A)` rather than `[aRA,eRA] = DotErr(R,A)`. Due to the interpretation overhead this would be costly; however, in all test examples we did not see a significant difference in accuracy.

Next we discuss the quality of rigorous error bounds computed by Algorithm 4.20 (`LssIllcoErrBndNear`) for extremely ill-conditioned matrices.

**5.3. Results for `LssIllcoErrBndNear`.** A timing comparison has already been shown in Table 5.2. Concerning accuracy we first treat the matrices in Table 5.1. The results for Pascal matrices are already shown in the right graph of Figure 2.1. As expected the rigorous error bound is weaker than the approximate solution shown in the left graph. In both cases the quality is nicely inverse proportional to the condition number.

The results for Hilbert and scaled Hilbert matrices do not reveal new information, so the graphs are omitted to save space. Similar to Figure 5.2 the rigorous error bound is weaker, but the behavior is similar. The same holds true for inverse Hilbert and Boothroyd matrices as shown in Figure 5.5. The results are similar to the approximate solutions computed by Algorithm `LssIllcoApprox` shown in Figure 5.3; the results for right hand sides `randn(n,1)` and `A*randn(n,1)` are practically identical.

Again the results are surprisingly good for inverse Hilbert matrices: For example, for dimension  $n = 40$  the matrix has a condition number larger than  $10^{37}$ , but nevertheless the error bound guarantees more than 10 correct digits of the solution for both right hand sides `b=randn(n,1)` and `b=A*randn(n,1)`. We have no conclusive explanation for that; a similar behavior is observed for Vandermonde matrices and will be discussed in Section 5.4.

Finally we consider extremely ill-conditioned matrices of dimensions up to  $n = 1000$ . A method how to construct extremely ill-conditioned matrices being exactly representable in floating-point of higher dimension is described in [42]. Moreover, interesting methods can be found in [33, 34]. Yet another way, which is used in INTLAB [41] for condition numbers up to about  $10^{100}$ , is to multiply a couple of sparse unit lower triangular matrices with small integer entries and to form  $A^T A$  until the desired condition number is achieved. For the following data we tried all methods with similar results.

In Figure 5.6 the results of Algorithm 4.20 (`LssIllcoErrBndNear`) for dimensions  $n \in \{100, 200, 500, 1000\}$  and for right hand sides `b=randn(n,1)` and `b=A*randn(n,1)` are shown in one graph. With increasing condition number the quality of the error bounds decreases. The maximum treatable condition number is of the order  $u^{-2}/n^2$ . In contrast, the results of Algorithm 4.16 (`LssErrBndNear`) as shown in Figure 5.4 are always of high accuracy - until the algorithm fails. This is due to the extra-precise residual iteration in `LssErrBndNear`, which could not be used in `LssIllcoErrBndNear`. For dimensions  $n = 100$ ,  $n = 200$ ,  $n = 500$  and  $n = 1000$  there is no failure of `LssIllcoErrBndNear` in the 100 samples until condition numbers  $6.2 \cdot 10^{25}$ ,  $1.5 \cdot 10^{26}$ ,  $4.7 \cdot 10^{25}$  and  $3.3 \cdot 10^{24}$ , respectively, as depicted on the x-axis in Figure 5.6. Accidentally the number for dimension  $n = 200$  is larger than for  $n = 100$  due to the randomness of the examples.

**5.4. Vandermonde matrices.** We observed in Figure 5.5 that the quality of the error bounds for the inverse Hilbert matrices was considerably better than one would expect from the condition number. For example, the condition number for  $n = 40$  is  $2.3 \cdot 10^{37}$ , but nevertheless the rigorous bounds by Algorithm 4.20 (`LssIllcoErrBndNear`) are accurate to 10 decimal digits.

A similar behavior can be observed for the notoriously ill-conditioned Vandermonde matrices [4, 17]. In Figure 5.7 we display the results of `LssIllcoErrBndNear` for the original Vandermonde matrix as defined

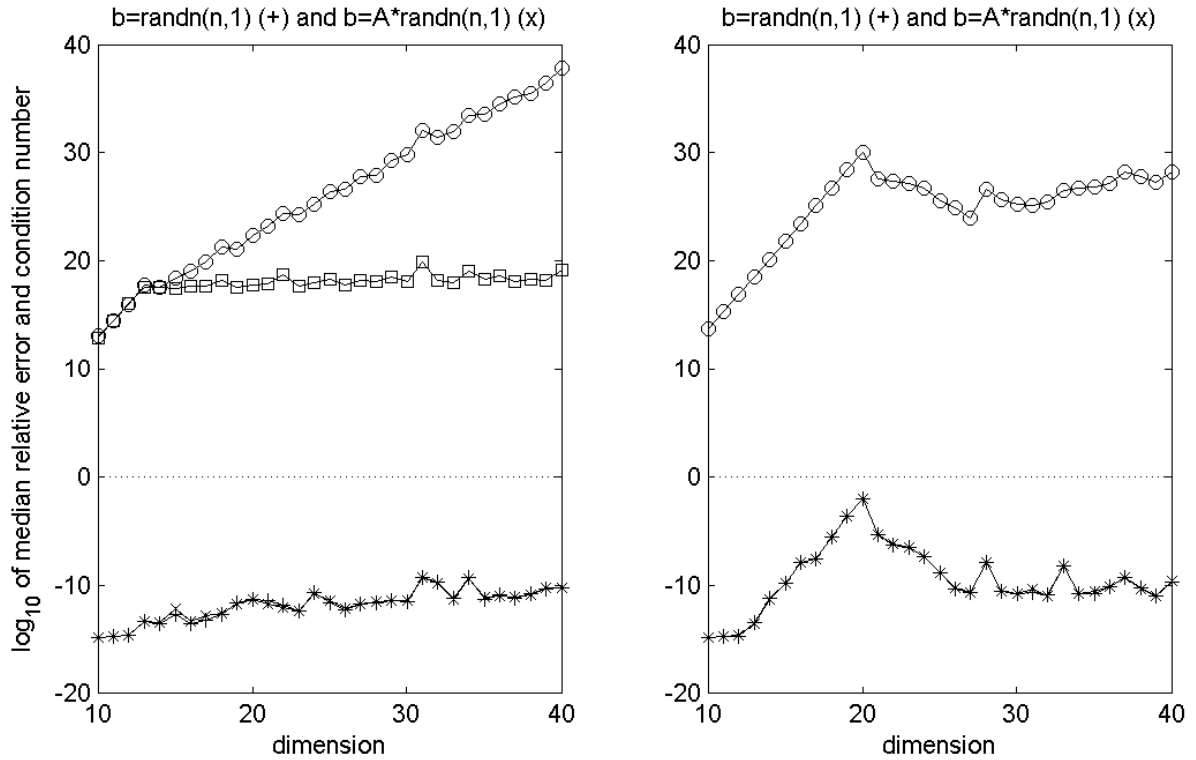


FIG. 5.5. Accuracy of Algorithm 4.20 (`LssIllcoErrBndNear`) for inverse Hilbert and Boothroyd matrices as defined in Table 5.1. The upper parts show the condition number, the lower parts the relative error of the results. In the left graph, “o” corresponds to the traditional condition number  $\|A^{-1}\|_2\|A\|_2$  and “□” to the Bauer-Skeel condition number  $\| |A^{-1}| \cdot |A| \|$ .

in Table 5.1 and of its equilibrated version. In both cases the right hand side is  $\mathbf{b}=\mathbf{randn}(n,1)$ ; the results for  $\mathbf{b}=\mathbf{A}*\mathbf{randn}(n,1)$  are completely similar. The circles in the upper halves denote the traditional condition number  $\|A^{-1}\|_2\|A\|_2$ , increasing rapidly beyond  $10^{40}$ . In the lower halves the median relative error over all solution components for the built-in Matlab call  $A \setminus b$  is given by “+”. For comparison also the median relative error of `LssIllcoErrBndNear` is depicted by “\*”. For small dimension sometimes the solution is exactly representable; in that case an error zero is replaced by  $10^{-20}$ .

In the left graph of Figure 5.7 we see that for condition numbers way beyond  $10^{20}$  the results of the built-in Matlab approximation  $A \setminus b$  still maintain some accuracy, apparently contradicting the well-accepted rule of thumb that in IEEE 754 double precision for condition number  $10^k$  about  $16 - k$  correct digits can be expected. In contrast, the linear system with the equilibrated matrix in the right graph shows the expected behavior: The relative error surpasses 1 at the condition number  $10^{16}$  (the dotted line in the upper half). Here is yet another example that equilibration need not to improve the quality of the result.

We do not have a convincing explanation for that. Similar contradictions to the mentioned rule of thumb have been observed in [12]. A reason might be the following. The Bauer-Skeel condition number  $\kappa := \| |A^{-1}| \cdot |A| \|_\infty$  is displayed by “□” in the left graph in Figures 5.5 and 5.7. This is the optimal normwise condition number achievable by left diagonal scaling, and it is also equal to the componentwise condition number with respect to relative perturbations of the matrix entries.

The accuracy of the Matlab-approximation  $A \setminus b$  satisfies the well-known rule of thumb with respect to the Bauer-Skeel condition number for the original Vandermonde matrix, see Figure 5.7; the same observation applies to inverse Hilbert matrices as in Figure 5.5. In any case the computational results for inverse Hilbert and Vandermonde matrices are nice but seem to be a kind of artefact.

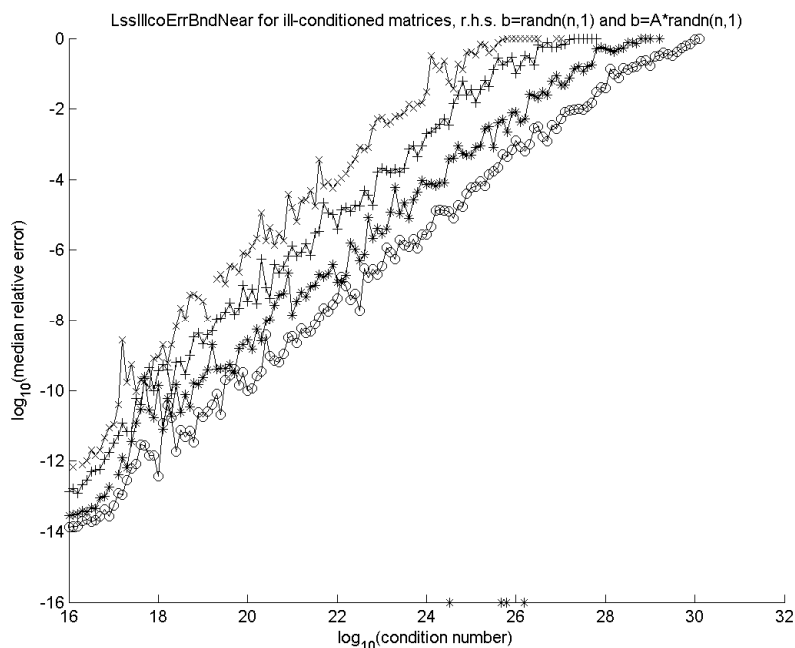


FIG. 5.6. Results of Algorithm 4.16 (`LssErrBndNear`) for ill-conditioned random matrices of dimension  $n=100$  ( $\circ$ ),  $n=200$  ( $*$ ),  $n=500$  ( $+$ ) and  $n=1000$  ( $\times$ ).

**6. Conclusion.** In this Part I of the article all algorithms use solely ordinary floating-point arithmetic in rounding to nearest. One purpose of the paper is to show that it is fairly simple to obtain rigorous error bounds subject to that constraint.

An algorithm for computing approximations of reasonable quality for extremely ill-conditioned matrices with condition number  $\gg \mathbf{u}^{-1}$  was given. Moreover, algorithms computing rigorous error bounds (including possible underflow), also for extremely ill-conditioned matrices, have been presented. All algorithms are given in executable Matlab-code. All algorithms use only the basic floating-point operations in rounding to nearest. For the extra-precise dot product, i.e. products accumulated in double the working precision with result rounded into working precision, an algorithm using only the basic floating-point operations in rounding to nearest was given as well.

The error bounds are of high quality, however, certain estimates in rounding to nearest seem improvable. Considerably sharper error bounds for ill-conditioned matrices and also for extremely ill-conditioned matrices are presented in Part II of this paper. They are a little bit more involved; of course, they can be computed in rounding to nearest, but are better and easier to discuss using directed rounding.

**Acknowledgement.** The author wishes to thank an anonymous referee for helpful comments, in particular for pointing to accurate algorithms to solve extremely ill-conditioned linear systems with special matrices.

#### REFERENCES

- [1] P. Alonso, J. Delgado, R. Gallego, and J.M. Peña. Growth factors of pivoting strategies associated to Neville Elimination. *J. Comp. Appl. Math.*, 235(7):1775–1762, 2011.
- [2] I. Babuška. Numerical Stability in Mathematical Analysis. *Information Processing*, 68:11–23, 1969.
- [3] D.H. Bailey, H. Yozo, X.S. Li, and B. Thompson. ARPREC: An Arbitrary Precision Computation Package. Technical Report LBNL-53651, Lawrence Berkeley National Laboratory, Berkeley, CA, 2002.
- [4] B. Beckermann. The condition number of real Vandermode, Krylov and positive definite Hankel matrices. *Numer. Math.*, 85(4):553–577, 2000.



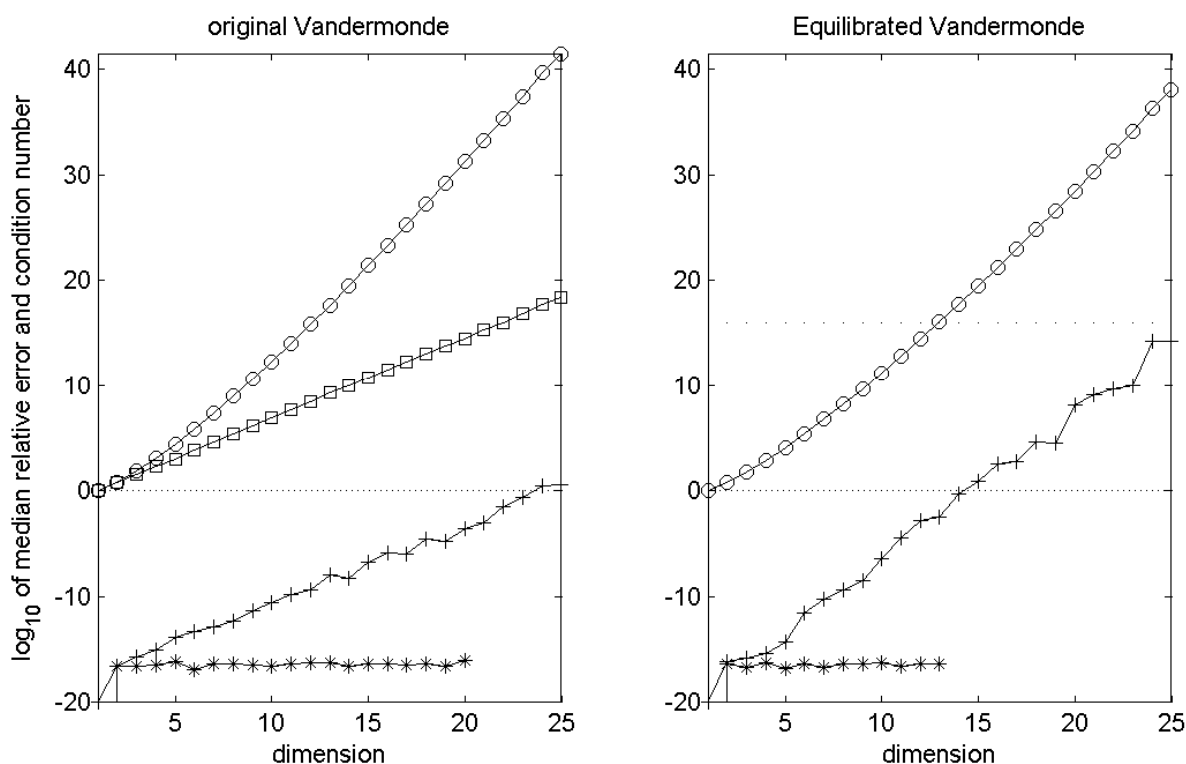


FIG. 5.7. Accuracy of Algorithm 4.20 (`LssIllcoErrBndNear`) (\*) and Matlab built-in `A\b` (+) for original (left graph) and equilibrated Vandermonde matrices (right graph) as defined in Table 5.1. The upper parts show the condition number, the lower parts the relative error of the results. In both cases the right hand side is `b=randn(n,1)`. Condition numbers of original and equilibrated Vandermonde matrices (o) and of optimally scaled matrices ( $\square$ ) are shown.

- [5] J. Boothroyd. Algorithm 274: Generation of Hilbert Derived Test Matrix. *Communications of the ACM*, 9(1):11, 1966.
- [6] N. Castro-González, J. Ceballos, F. Dopico, and J. Molera. Multiplicative Perturbation Theory and Accurate Solution of Least Squares Problems. submitted for publication, 2012.
- [7] L. Collatz. Einschließungssatz für die charakteristischen Zahlen von Matrizen. *Math. Z.*, 48:221–226, 1942.
- [8] T.J. Dekker. A floating-point technique for extending the available precision. *Numerische Mathematik*, 18:224–242, 1971.
- [9] J. Demmel. Accurate SVDs of structured matrices. *SIAM J. Math. Anal. (SIMA)*, 21(2):562–580, 1999.
- [10] J. Demmel and Y. Hida. Accurate and efficient floating point summation. *SIAM J. Sci. Comput. (SISC)*, 25:1214–1248, 2003.
- [11] J. B. Demmel, I. Dumitriu, O. Holtz, and P. Koev. Accurate and efficient expression evaluation and linear algebra. *Acta Numerica*, 2008:87–145, 2008.
- [12] J.B. Demmel, Y. Hida, W. Kahan, X.S. Li, S. Mukherjee, and E.J. Riedy. Error Bounds from Extra Precise Iterative Refinement. *ACM Transactions on Mathematical Software (TOMS)*, 32(2):325–351, 2006.
- [13] F.M. Dopico and J.M. Molera. Accurate solution of structured linear systems via rank-revealing decompositions. *IMA J. Numer. Anal.*, 32:1096–1116, 2012.
- [14] G.E. Forsythe, M.A. Malcolm, and C.B. Moler. *Computer Methods for Mathematical Computations*. Prentice Hall, Englewood Cliffs, NJ, 1977.
- [15] L. Fousse, G. Hanrot, V. Lefèvre, P. Pélicier, and P. Zimmermann. MPFR: A multiple-precision binary floating-point library with correct rounding. *ACM Trans. Math. Software*, 33(2), 2007. article 13.
- [16] A. Frommer. Proving Conjectures by Use of Interval Arithmetic. In U. Kulisch et al., editor, *Perspectives on enclosure methods. SCAN 2000, GAMM-IMACS international symposium on scientific computing, computer arithmetic and validated numerics, Univ. Karlsruhe, Germany, September 19-22, 2000*, Wien, 2001. Springer.
- [17] W. Gautschi. Optimally scaled and optimally conditioned Vandermonde and Vandermonde-like matrices. *BIT*, 51(1):103–125, 2011.
- [18] G.H. Golub and Ch. Van Loan. *Matrix Computations*. Johns Hopkins University Press, Baltimore, third edition, 1996.
- [19] N. J. Higham. The accuracy of floating point summation. *SIAM J. Sci. Comput.*, 14:783–799, 1993.
- [20] N. J. Higham. *Accuracy and stability of numerical algorithms*. SIAM Publications, Philadelphia, 2nd edition, 2002.
- [21] N.J. Higham. Experience with a matrix norm estimator. *SIAM J. Sci. Statist. Comput. (SISC)*, 11:804–809, 1990.
- [22] *ANSI/IEEE 754-1985: IEEE Standard for Binary Floating-Point Arithmetic*. New York, 1985.

- [23] *ANSI/IEEE 754-2008: IEEE Standard for Floating-Point Arithmetic*. New York, 2008.
- [24] D.E. Knuth. *The Art of Computer Programming: Seminumerical Algorithms*, volume 2. Addison Wesley, Reading, Massachusetts, 1969.
- [25] P. Koev. Accurate computations with totally nonnegative matrices. *SIAM J. Matrix Anal. Appl.*, 29:731–751, 2007.
- [26] M. La Porte and J. Vignes. Étude statistique des erreurs dans l'arithmétique des ordinateurs; application au contrôle des résultats d'algorithmes numériques. *Numer. Math.*, 23:63–72, 1974.
- [27] P. Langlois. Accurate algorithms in floating-point arithmetic. In *Lecture at the 12th GAMM-IMACS International Symposium on Scientific Computing (SCAN), Computer Arithmetic and Validated Numerics*, Duisburg, 2006. IEEE.
- [28] X. Li, J. Demmel, D. Bailey, G. Henry, Y. Hida, J. Iskandar, W. Kahan, S. Kang, A. Kapur, M. Martin, B. Thompson, T. Tung, and D. Yoo. Design, implementation and testing of extended and mixed precision BLAS. *ACM Trans. Math. Software*, 28(2):152–205, 2002.
- [29] M. Malcolm. On accurate floating-point summation. *Comm. ACM*, 14(11):731–736, 1971.
- [30] 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.
- [31] A. Neumaier. Rundungsfehleranalyse einiger Verfahren zur Summation endlicher Summen. *Zeitschrift für Angew. Math. Mech. (ZAMM)*, 54:39–51, 1974.
- [32] A. Neumaier. *Introduction to Numerical Analysis*. Cambridge University Press, 2001.
- [33] T. Nishi, T. Ogita, S. Oishi, and S. M. Rump. A Method for the Generation of a Class of Ill-conditioned Matrices. In *2008 International Symposium on Nonlinear Theory and its Applications, NOLTA'08, Budapest, Hungary, September 7-10*, pages 53–56, 2008.
- [34] T. Nishi, S.M. Rump, and S. Oishi. On the generation of very ill-conditioned integer matrices. *Nonlinear Theory and Its Applications (NOLTA), IEICE*, 2(2):226–245, 2011.
- [35] T. Ogita, S.M. Rump, and S. Oishi. Accurate sum and dot product. *SIAM Journal on Scientific Computing (SISC)*, 26(6):1955–1988, 2005.
- [36] T. Ogita, S.M. Rump, and S. Oishi. Verified solution of linear systems without directed rounding. Technical Report 2005-04, Advanced Research Institute for Science and Engineering, Waseda University, Tokyo, Japan, 2005.
- [37] S. Oishi, K. Tanabe, T. Ogita, and S.M. Rump. Convergence of Rump's method for inverting arbitrarily ill-conditioned matrices. *J. Comput. Appl. Math.*, 205(1):533–544, 2007.
- [38] K. Ozaki, T. Ogita, S. Miyajima, S. Oishi, and S.M. Rump. A method of obtaining verified solutions for linear systems suited for Java. *Journal of Computational and Applied Mathematics (JCAM)*, 199(2):337–344, 2006. Special issue on Scientific Computing, Computer Arithmetic, and Validated Numerics (SCAN 2004).
- [39] K. Ozaki, T. Ogita, S. M. Rump, and S. Oishi. Accurate matrix multiplication by using level 3 BLAS operation. In *Proceedings of the 2008 International Symposium on Nonlinear Theory and its Applications, NOLTA'08, Budapest, Hungary*, pages 508–511. IEICE, 2008.
- [40] D.M. Priest. *On Properties of Floating-Point Arithmetics: Numerical Stability and the Cost of Accurate Computations*. PhD thesis, Mathematics Department, University of California at Berkeley, CA, 1992. <ftp://ftp.icsi.berkeley.edu/pub/theory/priest-thesis.ps.Z>.
- [41] S.M. Rump. INTLAB - Interval Laboratory, Version 6. <http://www.ti3.tu-harburg.de/rump>, 1998–2011.
- [42] S.M. Rump. A Class of Arbitrarily Ill-conditioned Floating-Point Matrices. *SIAM J. Matrix Anal. Appl. (SIMAX)*, 12(4):645–653, 1991.
- [43] S.M. Rump. Inversion of extremely ill-conditioned matrices in floating-point. *Japan J. Indust. Appl. Math. (JJIAM)*, 26:1–29, 2009.
- [44] S.M. Rump. Ultimately Fast Accurate Summation. *SIAM Journal on Scientific Computing (SISC)*, 31(5):3466–3502, 2009.
- [45] S.M. Rump. Error estimation of floating-point summation and dot product. *BIT*, 51(1):201–220, 2012.
- [46] S.M. Rump, T. Ogita, and S. Oishi. Accurate floating-point summation part I: Faithful rounding. *SIAM J. Sci. Comput.*, 31(1):189–224, 2008.
- [47] J.R. Shewchuk. Adaptive precision floating-point arithmetic and fast robust geometric predicates. *Discrete Comput. Geom.*, 18(3):305–363, 1997.
- [48] R. Skeel. Iterative Refinement Implies Numerical Stability for Gaussian Elimination. *Math. Comp.*, 35(151):817–832, 1980.
- [49] L.N. Trefethen and R. Schreiber. Average-case stability of gaussian elimination. *SIAM J. Matrix Anal. Appl.*, 11(3):335–360, 1990.
- [50] W. Tucker. The Lorenz attractor exists. *C. R. Acad. Sci., Paris, Sér. I, Math.*, 328(12):1197–1202, 1999.
- [51] Jean Vignes. *Algorithmes numériques, analyse et mise en œuvre. 2*. Éditions Technip, Paris, 1980. Équations et systèmes non linéaires. [Nonlinear equations and systems], With the collaboration of René Alt and Michèle Pichat, Collection Langages et Algorithmes de l'Informatique.
- [52] D. Viswanath and L.N. Trefethen. Condition numbers of random triangular matrices. *SIAM J. Matrix Anal. Appl.*, 19(2):564–581, 1998.
- [53] XBLAS: A Reference Implementation for Extended and Mixed Precision BLAS. <http://crd.lbl.gov/~xiaoye/XBLAS/>.
- [54] T. Yamamoto. Error Bounds for Approximate Solutions of Systems of Equations. *Japan J. Appl. Math.*, 1:157–171, 1984.

- [55] Y.-K. Zhu and W. Hayes. Fast, guaranteed-accurate sums of many floating-point numbers. In G. Hanrot and P. Zimmermann, editors, *Proceedings of the RNC7 Conference on Real Numbers and Computers, Nancy, France*, pages 11–22. Loria, 2006.
- [56] Y.-K. Zhu, J.-H. Yong, and G.-Q. Zheng. A new distillation algorithm for floating-point summation. *SIAM J. Sci. Comput.*, 26(6):2066–2078, 2005.
- [57] G. Zielke and V. Drygalla. Genaue Lösung linearer Gleichungssysteme. *GAMM Mitt. Ges. Angew. Math. Mech.*, 26:7–108, 2003.