

1 VERIFIED ERROR BOUNDS FOR MATRIX DECOMPOSITIONS*

2 SIEGFRIED M. RUMP[†] AND TAKESHI OGITA[‡]

3 **Abstract.** In this note we consider common matrix factorizations such as LU decomposition of
4 a square and rectangular matrix, Cholesky and QR decomposition, singular value decomposition for
5 square and rectangular matrices, eigen-, Schur and Takagi decomposition. We first note that well-
6 conditioned factors tend to be sensitive to perturbations of the input matrix, while ill-conditioned
7 factors tend to be insensitive. It seems that this behaviour has not been recognized in numerical
8 analysis. We develop a formula for the relation between condition number of the factor and its
9 sensitivity with respect to input perturbations, and give reasons for that.

10 Our main focus is to describe verification methods for the factors of the mentioned decompo-
11 sitions. That means to prove existence of the factorization together with rigorous entrywise error
12 bounds for the factors. Our goal is to develop algorithms requiring $\mathcal{O}(Pp^2)$ operations for an $m \times n$
13 matrix with $P := \max(m, n)$ and $p := \min(m, n)$. Moreover, bounds of high quality are aimed for,
14 often not far from maximal accuracy. A main tool to achieve that is accurate dot products based on
15 error-free transformations. Since preconditioning based on approximate inverses is used, our methods
16 are restricted to full matrices.

17 **Key words.** Sensitivity of matrix factors, verified inclusions, error-free transformations, LU
18 decomposition, Cholesky decomposition, QR decomposition, singular value decomposition, eigende-
19 composition, Schur decomposition, polar decomposition, Takagi decomposition, INTLAB

20 **MSC codes.** 65G20, 65F99

21 **1. Introduction.** Verification methods are mathematical theorems the assump-
22 tions of which can be verified on a digital computer. The assumptions are verified
23 with mathematical rigor including all procedural, rounding and other sources of error,
24 thus the assertions are true with mathematical rigor. The error bounds are computed
25 together with the proof of existence and often uniqueness of the solution. Problems
26 cover systems of linear and nonlinear equations, eigenproblems or ordinary and partial
27 differential equations. For the theoretical foundation and algorithms see [26, 32, 28].

28 Verification methods aim to formulate the assumptions in such a way that they
29 can be rigorously verified on a computer, and that it is likely that they are satisfied
30 for not too ill-conditioned problems. The computing time should be of the same order
31 as that of a standard numerical algorithm. The bounds should be narrow.

32 There is a general limit to verification methods, namely, they are not applicable to
33 ill-posed problems. That is the price we have to pay by using floating-point operations
34 combined with error estimates rather than computing exactly like in computer algebra.
35 For example, it is possible to verify that a matrix is nonsingular, even for very large
36 condition numbers. However, it is not possible to verify that a matrix is singular
37 because that problem is ill-posed in the sense of Tikhonov [40, 41]: An arbitrarily
38 small change of the input data may change the answer. Similarly, even for a symmetric
39 matrix it is not possible to compute narrow error bounds for an individual eigenvector
40 to a double eigenvalue, see (6.1). However, verified bounds are possible for a basis of
41 the invariant subspace.

42 In this note we are interested in fast verification methods for the factors of stan-

*Submitted to the editors DATE.

Funding: This work was partially supported by JSPS KAKENHI Grant Number JP23H03410.

[†]Institute for Reliable Computing, Hamburg University of Technology, Am Schwarzenberg-
Campus 3, Hamburg 21073, Germany, and Faculty of Science and Engineering, Waseda University,
3-4-1 Okubo, Shinjuku-ku, Tokyo 169-8555, Japan (rump@tuhh.de).

[‡]Faculty of Science and Engineering, Waseda University, 3-4-1 Okubo, Shinjuku-ku, Tokyo 169-
8555, Japan (t.ogita@waseda.jp).

43 dard matrix decompositions with emphasis on the complete matrix decomposition.
 44 For example, methods are known [18, 4] to compute error bounds for a single eigen-
 45 pair. For an $n \times n$ matrix those methods might be applied to each individual eigenpair,
 46 however, resulting in totally $\mathcal{O}(n^4)$ operations. In contrast, [25, 24] give methods to
 47 compute error bounds for the complete eigendecomposition in $\mathcal{O}(n^3)$ operations in-
 48 cluding the treatment of clustered and/or multiple eigenvalues.

49 In [34] fast methods are described to compute error bounds for the complete
 50 singular value decomposition of an $m \times n$ matrix with special emphasis on clustered
 51 and/or multiple singular values. Here “fast” means $\mathcal{O}(Pp^2)$ operations with $P :=$
 52 $\max(m, n)$ and $p := \min(m, n)$.

53 However, no verification methods are known for other standard matrix decompo-
 54 sitions. We close this gap by giving fast algorithms for the LU, Cholesky, QR, and
 55 Schur decomposition. In addition to “fast” in terms of $\mathcal{O}(Pp^2)$ operations we aim on
 56 inclusions being accurate for all solution components, i.e., the entrywise relative error
 57 between lower and upper bound should be close to the relative rounding error unit \mathbf{u}
 58 of the floating-point arithmetic in use.

59 Let $\mathbb{K} \in \{\mathbb{R}, \mathbb{C}\}$. We use the notation $M_{m,n}$ for the set of matrices in $\mathbb{K}^{m \times n}$, and
 60 shortly M_n if $m = n$. For $A \in M_{m,n}$ we denote by $A_k \in M_k$ the upper left $k \times k$
 61 principal submatrix of A . We adopt the convention that inverses are assumed to exist
 62 if used. The $n \times n$ identity matrix is denoted by I_n , where the index is omitted if clear
 63 from the context. Moreover, $I_{m,n} \in M_{m,n}$ is the matrix with I_p for $p := \min(m, n)$ in
 64 the upper left corner and zero elsewhere.

65 Any method to compute rigorous error bounds for scalar, vector and matrix op-
 66 erations is suitable for the algorithms to be presented. We use interval arithmetic
 67 [26] because it is simple and intuitive to use, in particular in INTLAB [31], the MAT-
 68 LAB/Octave toolbox for reliable computing. We use the interval notation [15], where
 69 in particular boldface letters indicate interval quantities.

70 Not much knowledge about verification methods and/or interval arithmetic is
 71 necessary to follow this note, basically familiarity with MATLAB notation. The rep-
 72 resentation of intervals like infimum-supremum or midpoint-radius is not important:
 73 throughout this note we only use the *inclusion property*, namely, that interval opera-
 74 tions $\circ \in \{+, -, \cdot, / \}$ are defined such that for compatible interval quantities \mathbf{A}, \mathbf{B}

$$75 \quad (1.1) \quad \forall A \in \mathbf{A} \forall B \in \mathbf{B}: \quad A \circ B \in \mathbf{A} \circ \mathbf{B}$$

76 is satisfied. For details see [26, 32, 28]. For $M \in M_n(\mathbb{K})$ and non-negative $R \in$
 77 $M_n(\mathbb{R})$ the command `midrad(M,R)` is a superset of $\{A \in M_n(\mathbb{K}) : |A - M| \leq R\}$
 78 with entrywise comparison and absolute value. Moreover, $\mathbf{X} = f(\mathbf{A})$ for an interval
 79 quantity \mathbf{A} and the induced function f implies that $f(A) \in \mathbf{X}$ for all $A \in \mathbf{A}$.

80 For an interval \mathbf{X} , the magnitude is defined by $\text{mag}(\mathbf{X}) := \max\{|x| : x \in \mathbf{X}\} \geq 0$.
 81 The definition applies entrywise to vectors and matrices, so that $B = \text{mag}(\mathbf{A})$ satisfies
 82 $|A_{ij}| \leq B_{ij}$ for all i, j , cf. [26]. The result B is a non-negative vector/matrix.

83 Throughout this note we use the new definition [35] of the relative error of an
 84 interval quantity \mathbf{X} , which is basically $\text{diam}(\mathbf{X})/\text{mag}(\mathbf{X})$ for $\text{diam}(\mathbf{X})$ denoting the
 85 diameter of \mathbf{X} . The definition applies to vectors and matrices entrywise.

86 We use some notations in MATLAB-style, in particular

$A^{[l]}$	the strictly lower triangular part of A
$A^{[u]}$	the upper triangular part of A
$\max(A)$	the row vector of columnwise maxima of A
$\text{sum}(A, 2)$	the column vector of rowwise sums of A

88 for square A . For the maximum and sum we use the typewriter font to avoid confusion
 89 with the mathematical terms. The MATLAB notation for $A^{[\ell]}$ and $A^{[u]}$ is `tril(A,-1)`
 90 and `triu(A)`, respectively. We introduce this short notation because we use them
 91 frequently when developing the algorithms for the LU decomposition. The maximum
 92 and sum apply to $A \in M_{m,n}$ as well, and also the triangular parts apply to $A \in M_{m,n}$
 93 in the sense that, for example, $B = \text{tril}(A)$ is the lower triangular part of A . That
 94 means for matrices of any dimension $A = A^{[\ell]} + A^{[u]}$.

95 For $A, B, C \in M_n$ we note entrywise upper bounds for $|ABC|$. Denote by $\mu :=$
 96 `max(|C|)` the row vector of columnwise maxima of $|C|$, and by $\sigma := \text{sum}(|A|,2)$ the
 97 column vector of rowwise sums of $|A|$, i.e., $\mu_\ell = \max_{1 \leq k \leq n} |C_{k\ell}|$ and $\sigma_i = \sum_{j=1}^n |A_{ij}|$.
 98 Then

$$99 \quad (1.2) \quad \begin{aligned} |ABC|_{i\ell} &= \sum_{j=1}^n \sum_{k=1}^n |A_{ij} B_{jk} C_{k\ell}| \leq \sum_{j=1}^n \sum_{k=1}^n |A_{ij} B_{jk} \mu_\ell| \\ &\leq \|B\|_\infty \mu_\ell \sum_{j=1}^n |A_{ij}| \leq \sigma_i \mu_\ell \|B\|_\infty, \end{aligned}$$

100 so that entrywise upper bounds

$$101 \quad (1.3) \quad |ABC| \leq \text{sum}(|A|,2) \text{max}(|C|) \|B\|_\infty \quad \text{and} \quad |AB| \leq \text{sum}(|A|,2) \text{max}(|B|)$$

102 by outer products follow. Note that the computational cost is $\mathcal{O}(n^2)$. With the same
 103 complexity the bounds

$$104 \quad |ABC| \leq \text{sum}(|A|,2) \text{max}(|B|) \cdot |C| \quad \text{and} \quad |ABC| \leq |A| \cdot \text{sum}(|B|,2) \text{max}(|C|)$$

105 follow; however, in our applications $B = (I + F)^{-1}$ is a perturbation of the identity
 106 matrix, so that the entries of B are not available but an upper bound for $\|B\|_\infty$ is.
 107 These estimates are true for real and complex matrices, as well as for any compatible
 108 matrix dimensions of A, B , and C .

109 All computational results use MATLAB [22] and double precision (binary64), i.e.,
 110 some 53 bits in the mantissa with the relative rounding error unit $\mathbf{u} = 2^{-53} \approx 10^{-16}$.
 111 We use directed rounding which is part of the IEEE 754 arithmetic standard [1].
 112 The statement `setround(-1)` implies that from now until the next call of `setround`
 113 the rounding mode is downwards, i.e., towards $-\infty$. As a consequence the result of
 114 floating-point operations is the largest floating-point number being less than or equal
 115 to the true real result. Similarly, `setround(1)` switches the rounding upwards so that
 116 the smallest floating-point number being greater than or equal to the true real result
 117 is computed. These statements are true for vector and matrix operations as well. As
 118 a consequence, the code sequence

```
119   setround(-1), Cinf = A*B;  
120   setround(+1), Csup = A*B;  
121   flpt = isequal(Cinf,Csup)
```

122 produces the result `flpt = true` if, and only if, all entries of the product AB are
 123 exactly representable in floating-point.

124 We aim on producing accurate bounds, i.e., the lower and upper bound often differ
 125 by few bits. Our main tool to achieve this are accurate dot products, either purely
 126 approximate or with error bound. To that end there are many techniques. Some
 127 of the early references are [42, 23, 21]. Later so-called “error-free transformations”
 128 [16, 7] were used to transform a pair (a, b) of floating-point numbers into a new pair
 129 (x, y) such that, e.g., x is the floating-point product ab and y is the error in the sense
 130 $ab = x + y$ and similarly for sum, quotient, and square root. That technique was used
 131 in [27] to introduce “error-free vector transformations” where a vector v of n floating-
 132 point numbers is transformed into a vector w of the same length such that w_n is the

133 floating-point sum of the v_i and $\sum v_i = \sum w_i$. In that paper the term “error-free
 134 transformations” was coined which was the start of a revival of such methods. Using
 135 error-free vector transformations, sums and dot products of arbitrarily large condition
 136 number can be computed with maximal precision [27].

137 The mentioned error-free transformations are based on a relative splitting of the
 138 input data. Yet a completely different method was introduced in [43] where an ab-
 139 solute splitting of vectors was introduced. That method was analyzed in [36] and
 140 is also used for reproducible results [3]. Moreover, this method was used to develop
 141 very efficient algorithm for accurate matrix multiplication [29], with or without error
 142 bounds. In our note we use such algorithms. They work in ordinary double preci-
 143 sion floating-point arithmetic but produce a result “as if” computed with doubled
 144 precision, i.e., some 32 decimal digits, or more. That allows to store the result of a
 145 matrix product in two terms, a higher and a lower order part. We use that technique
 146 occasionally. To that end some double-double arithmetic as in [5] or [19] could be
 147 used as well.

148 To be more precise the function `prodK` computes matrix products in $(1 + \frac{k}{2})$ -fold
 149 precision and rounds the result into working precision. In fact `prodK` is a very versatile
 150 function, but we need only few functionalities. For example, for a given square matrix
 151 A the code

```
152 k = 2;  
153 [L,U,p] = lu(A,'vector');  
154 R = prodK(L,U,-1,A(p,:),k);
```

155 computes the residual $LU - A(p, :)$ in two-fold precision and rounds the result into R
 156 in working precision. The call

```
157 [R,E] = prodK(L,U,-1,A(p,:),k);
```

158 produces the same R but in addition an error matrix E such that

$$159 \quad |LU - A(p, :) - R|_{ij} \leq E_{ij}.$$

160 for all indices i, j . Larger values of k are possible, but not used in this note. When
 161 calculating triple products ABC it may be useful to compute the first product in
 162 two-fold precision but also store it as an unevaluated sum $P_1 + P_2$. For example,

```
163 P = prodK(A,B,k,'OutputTerms',2);  
164 X = prodK(P,C,k);
```

165 stores P as a cell array, and the second call of `prodK` computes $P_1 * C + P_2 * C$ in two-
 166 fold precision and stores the result in X in working precision.

167 The key to our verification methods will be to transform the problem into the
 168 problem for a perturbed identity matrix. In particular in combination with extra-
 169 precise dot products that technique turns out to be effective. The transformation
 170 uses approximate inverses of approximate factors. These are usually full, also for
 171 sparse input matrix. Therefore applying our methods to sparse matrices is prohibitive
 172 because of expected fill-in. For some factorizations such as LU and at least the R -
 173 factor of QR are usually sparse for sparse input. Verified inclusions for these cases
 174 are open problems.

175 We begin with an investigation of the sensitivity of matrix factors. In particular
 176 the fact that, in case of an ill-conditioned input matrix A , well-conditioned factors
 177 tend to be sensitive to perturbations of A seems unknown in numerical analysis. In
 178 the following sections verification methods for the factors of the LU decomposition
 179 of a square and rectangular matrix, Cholesky- and QR decomposition, singular value
 180 decomposition for square and rectangular matrices, eigen- and Schur decomposition

181 are presented, accompanied by numerical results. As an application of the symmetric
 182 eigendecomposition we show how to compute inclusions for the Takagi factors.

183 Throughout the note random matrices $A \in \mathbb{F}^{m \times n}$ with specified condition number
 184 $\text{cond}(A) \approx 10^k$ are generated by

```
185     mn = min(m,n); s = logspace(0,k,mn); S = diag(s(randperm(mn)));
186     if m~n, S(m,n)=0; end; A = orth(randn(m)) * S * orth(randn(n));
```

187 which is for square matrices equivalent to MATLAB's `gallery/randsvd`.

188 We present numerical evidence that mostly our method compute error bounds
 189 with an accuracy close to the relative rounding error unit \mathbf{u} of the floating-point
 190 arithmetic in use. All our algorithms are given and implemented in pure MATLAB
 191 code, therefore suffering from interpretation overhead. Therefore we restrict timing
 192 information to the QR decomposition in Section 5 together with accuracy information
 193 of the built-in (approximate) MATLAB routines; the time ratio of other verification
 194 methods is similar.

195 **2. Sensitivity of factors in a decomposition.** For any of the matrix decom-
 196 positions under investigation we made a general observation which seems to be known
 197 in the literature [13, 12, 8] but not so much in numerical analysis. Some perturbation
 198 bounds for LU, Cholesky, and QR decompositions can be found in [37], see also [11],
 199 however they overestimate the sensitivity of ill-conditioned factors.

200 Let X be a factor of some decomposition of a matrix A . Denote by $A + \Delta A$ a
 201 small perturbation of A such that $\frac{\|\Delta A\|}{\|A\|} \sim \mathbf{u}$ for some matrix norm, and let \tilde{X} be the
 202 corresponding factor of $A + \Delta A$. Then numerical evidence (cf. Tables 1–10) suggests
 203 that often the sensitivity of X satisfies

$$204 \quad (2.1) \quad \text{sensitivity}(X) := \frac{\|\tilde{X} - X\|}{\|X\|} \sim \mathbf{u} \frac{\text{cond}(A)}{\text{cond}(X)},$$

205 where $\text{cond}(B) := \|B\| \cdot \|B^{-1}\|$ for a nonsingular square matrix B .

206 For example, suppose $A = LU$ with L being unit lower triangular. Then the
 207 U -factor of the LU decomposition has usually the same condition number as A . Al-
 208 though the L -factor is usually well conditioned by the widely accepted rule of thumb,
 209 numerical evidence (cf. Tables 1–4) suggests that its sensitivity grows with the con-
 210 dition number of A . That can be seen as follows. Let $A + \Delta A = (L + \Delta L)(U + \Delta U)$,
 211 then to first order

$$212 \quad (2.2) \quad L^{-1} \cdot \Delta A \cdot U^{-1} = L^{-1} \cdot \Delta L + \Delta U \cdot U^{-1}.$$

213 The matrices ΔL and $L^{-1} \cdot \Delta L$ are strictly lower triangular, whereas $\Delta U \cdot U^{-1}$ is
 214 upper triangular. Thus, taking the strictly lower triangular and upper triangular part
 215 from matrices of both sides of (2.2) implies

$$216 \quad \Delta L = L [L^{-1} \cdot \Delta A \cdot U^{-1}]^{[l]} \quad \text{and} \quad \Delta U = [L^{-1} \cdot \Delta A \cdot U^{-1}]^{[u]} U.$$

217 Numerical evidence suggests that the elements of each row of (the upper triangular
 218 part of) U are often of similar magnitude¹, so that $U \approx DX$ for diagonal D with
 219 elements decreasing in magnitude with $\|D^{-1}\| \sim \|A^{-1}\|$ and well-conditioned X with
 220 the upper triangular part of entries close to 1 in magnitude. Hence

$$221 \quad L^{-1} \cdot \Delta A \cdot U^{-1} \approx L^{-1} \cdot \Delta A \cdot X^{-1} D^{-1} =: Y D^{-1}$$

¹That is true due to our practical experience, and it is also satisfied for matrices generated by `randsvd` from MATLAB's matrix gallery in any of the 5 modes. However, for ill-conditioned matrices generated by `sprand` with density 1, i.e., full matrices, it is sometimes not true.

222 for some Y with entries of the size of those of ΔA , i.e., $\|Y\| \sim \|\Delta A\|$. Then

$$223 \quad \Delta L = L [L^{-1} \cdot \Delta A \cdot U^{-1}]^{[\ell]} \approx L [YD^{-1}]^{[\ell]} = LY^{[\ell]}D^{-1}$$

224 and

$$225 \quad \Delta U = [L^{-1} \cdot \Delta A \cdot U^{-1}]^{[u]} U \approx [YD^{-1}]^{[u]} DX = Y^{[u]}X .$$

226 Now $\|L\|$ and $\|X\|$ are small because both are usually well conditioned, so that
227 $\|D^{-1}\| \sim \|A^{-1}\|$ and $\|Y\| \sim \|\Delta A\| \sim \mathbf{u}\|A\|$ imply

$$228 \quad \|\Delta L\| \sim \|\Delta A\| \cdot \|A^{-1}\| \sim \mathbf{u} \cdot \text{cond}(A) \quad \text{and} \quad \|\Delta U\| \sim \|\Delta A\| \sim \mathbf{u}\|A\|$$

229 and explain (2.1) for the LU decomposition.

230 For the QR decomposition (2.1) is mentioned in [20]. Let $A = QR$ and $A + \Delta A =$
231 $(Q + \Delta Q)(R + \Delta R)$, so that to first order

$$232 \quad M := Q^* \cdot \Delta A \cdot R^{-1} = Q^* \cdot \Delta Q + \Delta R \cdot R^{-1} \quad \text{and} \quad [Q^* \cdot \Delta A \cdot R^{-1}]^{[\ell]} = [Q^* \cdot \Delta Q]^{[\ell]} .$$

233 Using $(Q + \Delta Q)^*(Q + \Delta Q) = I$ implies that $C := Q^* \cdot \Delta Q$ is skew-Hermitian, so that
234 $M^{[\ell]} = C^{[\ell]}$ yields

$$235 \quad C = C^{[\ell]} - (C^{[\ell]})^* = M^{[\ell]} - (M^{[\ell]})^* \quad \text{and} \quad \Delta Q = Q [M^{[\ell]} - (M^{[\ell]})^*]$$

236 and explains (2.1) for the Q -factor. The perturbation of the R -factor satisfies

$$237 \quad (2.3) \quad \Delta R = [Q^* \cdot \Delta A]^{[u]} - [M^{[\ell]}R]^{[u]} + [(M^{[\ell]})^* R]^{[u]} .$$

238 The first summand support (2.1), i.e., that R is not very sensitive to perturbations of
239 A , the second and third one need some extra consideration. In a similar way to the
240 LU decomposition, numerical evidence suggests that $R \approx DX$ for diagonal D with
241 elements decreasing in magnitude and well-conditioned X with entries close to 1 in
242 magnitude. Hence

$$243 \quad [M^{[\ell]}R]^{[u]} \approx [[Q^* \cdot \Delta A \cdot X^{-1}D^{-1}]^{[\ell]} DX]^{[u]} = [[Q^* \cdot \Delta A \cdot X^{-1}]^{[\ell]} X]^{[u]}$$

244 which is of the order $\|\Delta A\|$. For the third summand of (2.3) we have

$$245 \quad [(M^{[\ell]})^* R]^{[u]} \approx [([Q^* \cdot \Delta A \cdot X^{-1}D^{-1}]^{[\ell]})^* DX]^{[u]} \approx [(D^{-1})^* YDX]^{[u]}$$

246 for some Y with entries of the size of $\|\Delta A\|$. Since the D_{ii} decrease in magnitude,
247 that supports (2.1) for R as well. For the Cholesky decomposition $A = R^T R$ the
248 ansatz

$$249 \quad R^{-T} \Delta A \cdot R^{-1} = R^{-T} \Delta R^T + \Delta R \cdot R^{-1} \quad \text{implies} \quad \Delta R = [R^{-T} \Delta A \cdot R^{-1}]^{[u]} R$$

250 and explains (2.1) along the same lines as well. The condition number of the Cholesky
251 factor is $\text{cond}(A)^{1/2}$, and numerical evidence suggests indeed that the sensitivity is
252 always of the order $\mathbf{u} \text{cond}(A)^{1/2}$ in accordance with (2.1). Similarly, for the QR
253 decomposition $\text{cond}(Q) = 1$ and $\text{cond}(R) = \text{cond}(A)$, so that (2.1) suggests that Q is
254 sensitive to perturbations of A while R is not. Numerical evidence supporting these
255 statements will be presented in the following sections.

256 As a consequence and from a numerical standpoint of view to our surprise, we
257 may expect that accurate inclusions are more demanding for well conditioned factors.

258 **3. LU decomposition.** If all upper left principal minors $\det(A_k)$ of $A \in M_{m,n}$
 259 are nonzero, then there is a unique LU decomposition of A . That is true for square as
 260 well as for rectangular matrices. If the first $m-1$ minors are nonzero but $\det(A_m) = 0$,
 261 then the decomposition exists but is not unique [11, Theorem 9.1].

262 A verification method for computing inclusions of the L - and U -factor of a matrix
 263 A asserts, with mathematical certainty, that the decomposition exists and is unique.
 264 Thus, a necessary condition is that A has full rank.

265 Given a matrix $A \in M_{m,n}$, the following MATLAB code in Algorithm 3.1 (`getL`)
 computes the L -factor of A , cf. [10, p. 35], [11, (9.2a)].

Algorithm 3.1 Computation of the L -factor

```

function L = getL(A)
    [m,n] = size(A);
    mn = min(m,n);
    L = eye(m,mn);
    for k=1:mn-1
        v = 1:k;
        w = k+1:m;
        Bv = inv(A(v,v)); % last column of inv(A(v,v)) needed
        L(w,k) = A(w,v)*Bv(:,end);
    end
    
```

266 For square $A \in M_n$, this requires to compute the last column of the inverse of A_k
 267 for $1 \leq k \leq n-1$. To that end we see no other way than to compute the inverses
 268 individually at the cost of $\mathcal{O}(k^3)$ operations each, so that totally prohibitive $\mathcal{O}(n^4)$
 269 operations are necessary.

271 Let $A \in \mathbb{K}^{m \times n}$ be given and denote $P := \max(m, n)$ and $p := \min(m, n)$. Our
 272 goal is to compute verified and sharp error bounds for the factors L and U of A with a
 273 total computing time of $\mathcal{O}(Pp^2)$ operations. This will be achieved by preconditioners
 274 X_L, X_U such that $X_L A X_U$ is a perturbed identity matrix I_E .

275 We first show how to compute the LU decomposition of a perturbed identity
 276 matrix, followed by the cases $m = n$, $m > n$, and $m < n$ for the LU decomposition of
 277 a general matrix.

278 **3.1. LU decomposition of a perturbed identity matrix.** Let $A \in M_{m,n}$
 279 with $m \geq n$ be given, denote $E := A - I_{m,n}$ and assume that E_n is convergent. Fix
 280 k with $1 \leq k \leq n$, let i satisfy $k+1 \leq i \leq m$, and denote $B := A_k^{-1} = (I_k + E_k)^{-1}$.
 281 Then, according to Algorithm 3.1 (`getL`), using $(I + E_k)^{-1} = I - (I + E_k)^{-1} E_k$ and
 282 $A_{i\nu} = E_{i\nu}$ for $i > k$ and $\nu \leq k$ yields

$$283 \quad L_{ik} = \sum_{\nu=1}^k A_{i\nu} B_{\nu k} = \sum_{\nu=1}^k E_{i\nu} [I_k - (I_k + E_k)^{-1} E_k]_{\nu k} .$$

284 Hence, denoting the k -th column of I_k by $e^{(k)}$ and using $i > \nu$ gives

$$285 \quad L_{ik} - E_{ik} = - \sum_{\nu=1}^k (E^{[\ell]})_{i\nu} [(I_k + E_k)^{-1} E_k e^{(k)}]_{\nu} .$$

286 Using (1.2) it follows

$$287 \quad |L_{ik} - E_{ik}| \leq \sum_{\nu=1}^k \left| E^{[\ell]} \right|_{i\nu} \|(I_k + E_k)^{-1} E_k e^{(k)}\|_{\infty} \leq \sum_{\nu=1}^k \left| E^{[\ell]} \right|_{i\nu} \frac{\|E_k e^{(k)}\|_{\infty}}{1 - \|E_k\|_{\infty}} .$$

288 The k -th component of the row vector $\max(|E_n^{[u]}|)$ is equal to $\|E_k e^{(k)}\|_\infty$, so that the
 289 strictly lower triangular part of the difference between L and E is bounded above by
 290 the outer product

$$291 \quad (3.1) \quad \left| L^{[\ell]} - E^{[\ell]} \right| \leq \frac{\left(\text{sum}(|E^{[\ell]}|, 2) \cdot \max(|E_n^{[u]}|) \right)^{[\ell]}}{1 - \|E_n\|_\infty} =: \Delta = \Delta^{[\ell]} ,$$

292 and there exists a strictly lower triangular matrix $C = C^{[\ell]}$ with

$$293 \quad (3.2) \quad L = I + E^{[\ell]} + C^{[\ell]} \quad \text{and} \quad |C^{[\ell]}| \leq \Delta^{[\ell]} .$$

294 In other words, the strictly lower triangular part of L is essentially equal to the strictly
 295 lower triangular part of E . The computational cost is $\mathcal{O}(mn)$ operations. For $m < n$
 296 the factor L is square, and along the same lines we deduce

$$297 \quad (3.3) \quad \left| L^{[\ell]} - E_m^{[\ell]} \right| \leq \frac{\left(\text{sum}(|E_m^{[\ell]}|, 2) \cdot \max(|E_m^{[u]}|) \right)^{[\ell]}}{1 - \|E_m\|_\infty} .$$

298 If L is square, i.e., $m \leq n$, an inclusion of L^{-1} can be obtained using verification
 299 methods [32], however, we may proceed directly by using the Neumann expansion

$$300 \quad (I + F)^{-1} = I - F(I + F)^{-1} = I - (I + F)^{-1}F = I - F + F(I + F)^{-1}F .$$

301 Then (1.3) implies

$$302 \quad |(I + F)^{-1} - I + F| \leq \frac{\text{sum}(|F|, 2) \max(|F|)}{1 - \|F\|_\infty}$$

303 provided that $\|F\|_\infty < 1$. Using (3.2), $F := E^{[\ell]} + C^{[\ell]}$ and $G := |E^{[\ell]}| + \Delta^{[\ell]}$ yields

$$304 \quad (3.4) \quad L^{-1} = I - E^{[\ell]} + \delta \quad \text{with} \quad |\delta| \leq \Delta^{[\ell]} + \left[\frac{\text{sum}(G, 2) \max(G)}{1 - \|G\|_\infty} \right]^{[\ell]} .$$

305 Note that $G = |L - I|$ and $\delta = \delta^{[\ell]}$, and error bounds are only needed for the strictly
 306 lower triangular part. The estimate may be improved by using more terms of the
 307 Neumann expansion, however, it seems hardly worth the effort.

308 In order to compute an inclusion of U we may use, regardless whether $m \geq n$ or
 309 $m < n$, the L -factor of the upper left square matrix of $I + E$, an inclusion of which
 310 can be computed as described before. For the case $m \geq n$ we have $(I + E)_n = L_n U$,
 311 and the uniqueness of the LU decomposition, (3.2) and $(I + F)^{-1}x = x - (I + F)^{-1}Fx$
 312 imply that

$$\begin{aligned} U &= (I_n + E_n^{[\ell]} + C_n^{[\ell]})^{-1}(I_n + E_n) \\ &= (I_n + E_n^{[\ell]} + C_n^{[\ell]})^{-1}(I_n + E_n^{[\ell]} + C_n^{[\ell]} + E_n^{[u]} - C_n^{[\ell]}) \\ 313 \quad &= I_n + (I_n + E_n^{[\ell]} + C_n^{[\ell]})^{-1}(E_n^{[u]} - C_n^{[\ell]}) \\ &= I_n + E_n^{[u]} - C_n^{[\ell]} - (I_n + E_n^{[\ell]} + C_n^{[\ell]})^{-1}(E_n^{[\ell]} + C_n^{[\ell]})(E_n^{[u]} - C_n^{[\ell]}) . \end{aligned}$$

314 Note that the rightmost factor $E_n^{[u]} - C_n^{[\ell]}$ is a composition into upper and strictly
 315 lower part. Now $(C^{[\ell]})^{[u]} = O$ because $U = U^{[u]}$ is upper triangular, and similar to

316 the estimate for the L -factor it follows

$$317 \quad (3.5) \quad \left| U - (I_n + E_n^{[u]}) \right| \leq \frac{\left(\text{sum}(|E_n^{[\ell]} + C_n^{[\ell]}|, 2) \cdot \max(|E_n^{[u]} - C_n^{[\ell]}|) \right)^{[u]}}{1 - \|E_n^{[\ell]} + C_n^{[\ell]}\|_\infty}.$$

318 For the case $m < n$ we identify the matrix dimensions by adding subscripts. For
 319 example, E_m denotes the left upper $m \times m$ principal submatrix of E and indicates
 320 the dimension as well. For $A \in M_{m,n}$ we have $A = I_{m,n} + E_{m,n} = L_m U_{m,n}$, so that
 321 $C_m^{[\ell]} = C_m$ with (3.2) for L_m , and for $P \in M_{m,n_1}, Q \in M_{m,n_2}$ with matrix block
 322 notation $[P, Q] \in M_{m,n_1+n_2}$ it follows

$$\begin{aligned} 323 \quad U_{m,n} &= (I_m + E_m^{[\ell]} + C_m^{[\ell]})^{-1} (I_{m,n} + E_{m,n}) \\ 324 \quad &= (I_m + E_m^{[\ell]} + C_m^{[\ell]})^{-1} \left(I_{m,n} + [E_m^{[\ell]} + C_m^{[\ell]}, O_{m,n-m}] + [E_m^{[u]} - C_m^{[\ell]}, E_{m,n-m}] \right) \\ 325 \quad &= I_{mn,n} + (I_m + E_m^{[\ell]} + C_m^{[\ell]})^{-1} [E_m^{[u]} - C_m^{[\ell]}, E_{m,n-m}] \\ 326 \quad &= I_{mn,n} + [E_m^{[u]} - C_m^{[\ell]}, E_{m,n-m}] \\ 327 \quad &\quad + (I_m + E_m^{[\ell]} + C_m^{[\ell]})^{-1} (E_m^{[\ell]} + C_m^{[\ell]}) [E_m^{[u]} - C_m^{[\ell]}, E_{m,n-m}]. \end{aligned}$$

328 Since $U \in M_{m,n}$ is upper triangular we obtain, similar to the previous estimate,

$$329 \quad \left| U - \left(I_{m,n} + E_{m,n}^{[u]} \right) \right| \leq \frac{\left(\text{sum}(|E_m^{[\ell]} + C_m^{[\ell]}|, 2) \cdot \max(B_{m,n}) \right)^{[u]}}{1 - \|E_m^{[\ell]} + C_m^{[\ell]}\|_\infty} =: \Delta$$

330 using $B_{m,n} := [|E_m^{[u]}| + |C_m^{[\ell]}|, |E_{m,n-m}|]$ which is $|E|$ with $E_m^{[\ell]}$ replaced by $|C_m^{[\ell]}|$.
 331 The computational cost is again $\mathcal{O}(mn)$ operations.

332 If $m \leq n$ we may derive an inclusion of U^{-1} as well. We rewrite (3.5) into

$$333 \quad (3.6) \quad U = I + E^{[u]} + C^{[u]} \quad \text{and} \quad |C^{[u]}| \leq \Delta^{[u]} = \Delta$$

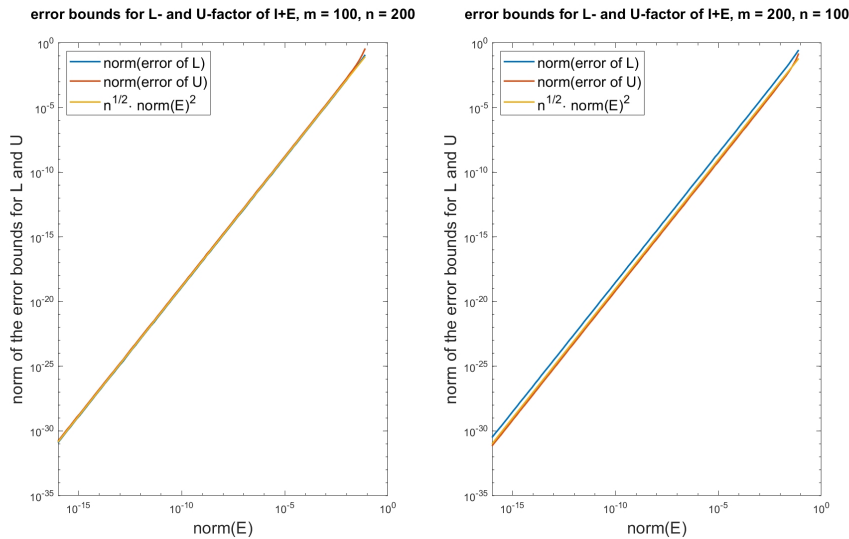
334 analogously to (3.2), and for $G := |E^{[u]}| + \Delta^{[u]}$, similar to (3.4), it follows

$$335 \quad (3.7) \quad U^{-1} = I - E^{[u]} + \delta \quad \text{with} \quad |\delta| \leq \Delta^{[u]} + \left[\frac{\text{sum}(G, 2) \max(G)}{1 - \|G\|_\infty} \right]^{[u]}.$$

336 In Figure 1 the norm of the right hand side of (3.1) and (3.5) together with $\sqrt{n}\|E\|^2$
 337 is displayed for dimensions $(m, n) = (100, 200)$ and $(m, n) = (200, 100)$ for different
 338 norms of E . As can be seen the norm of the error of the bounds for L and U grow with
 339 $\sqrt{n}\|E\|^2$. Hence, for $\|E\| \lesssim 10^{-8}$ the inclusions of the factors are maximally accurate,
 340 with errors of the size of the relative rounding error unit $\mathbf{u} \approx 10^{-16}$. Actually Figure 1
 341 displays results for complex input matrix; the results for real input are similar.

342 The graph of the norm of the error as in (3.4) and (3.7) of L^{-1} and U^{-1} , respectively,
 343 looks exactly like Figure 1, so the error for the inclusions of L^{-1} and U^{-1}
 344 grows with $\sqrt{n}\|E\|^2$ as well.

345 We close this subsection with giving executable MATLAB/INTLAB code in Al-
 346 gorithm 3.2 for the computation of inclusions of the L - and U -factor as well as of
 347 their inverses for a perturbed identity matrix $I + E$. Input is the perturbation E , and
 348 inclusions are stored as perturbations of the identity matrix as well, i.e., $I + E = LU$
 349 implies $L \in I + \text{LE}$, $U \in I + \text{UE}$, $L^{-1} \in I + \text{LinVE}$, and $U^{-1} \in I + \text{UinvE}$.

FIG. 1. Norm of error bound for the L - and U -factor of $I + E$

Algorithm 3.2 Inclusion of L - and U -factor of $I + E$ and their inverses

```

1  function [LE,UE,LinvE,UinvE] = LU_E(E)
2      setround(1)
3      magE = mag(E);
4      DeltaL = tril(sum(tril(magE,-1),2)*max(magE),-1);
5      DeltaL = - ( DeltaL/(norm(magE,inf)-1) );
6      LE = tril(E,-1) + midrad(0,DeltaL);
7      GL = mag(LE);
8      delta = tril(sum(tril(GL,-1),2)*max(GL),-1);
9      delta = DeltaL - delta/(norm(GL,inf)-1);
10     LinvE = -tril(E,-1) + midrad(0,delta);
11     B = triu(magE) + tril(DeltaL,-1);
12     DeltaU = triu(sum(GL,2)*max(B));
13     DeltaU = - ( DeltaU/(norm(GL,inf)-1) );
14     UE = triu(E) + midrad(0,DeltaU);
15     GU = mag(UE);
16     delta = triu(sum(triu(GU),2)*max(GU));
17     delta = DeltaU - delta/(norm(GU,inf)-1);
18     UinvE = -triu(E) + midrad(0,delta);
19     setround(0)

```

350 Throughout the code we use from line 2 rounding upwards, i.e., the computed
 351 floating-point result is always an upper bound of the true result (see Section 1). That
 352 holds true for vector and matrix operations as well. Thus, for example, the sum in
 353 the computation of `DeltaL` in line 4 is an upper bound of the row sums of absolute
 354 values of the strictly lower triangular part of E . The code is simplified in the sense
 355 that in lines 5, 9, 13, and 17 it is assumed that the upper bounds for the norms are

356 strictly less than 1. Then the denominator in line 5 is negative and larger or equal
 357 to $\|E\|_\infty - 1$, so that the negative of the division produces a correct upper bound
 358 **DeltaL**. Similar arguments show the correctness of the code in lines 9, 13, and 17. For
 359 an interval matrix \mathbf{E} , real or complex, the assertions hold true for every $I + \tilde{E} \in I + \mathbf{E}$.

360 For given A , the product $U = L^{-1}A$ can be enclosed by $(I + \mathbf{LE}) \setminus A$ or $A + \mathbf{Lin}v\mathbf{E} \cdot A$,
 361 and similarly $L = AU^{-1}$ can be enclosed² by $A/(I + \mathbf{LU})$ or $A + A \cdot \mathbf{Lin}v\mathbf{U}$. For not
 362 too large $\|E\|$ the latter formulation is advantageous. To that end we compare the
 363 relative error of the former and the latter method. For different values of ε we choose
 364 perturbations $I + E$ with $\|E\| = \varepsilon$ and $A \in M_{100}$ with fixed condition number 10^8 .
 365 The results are shown in Figure 2 for L and U in the left and right graph, respectively.

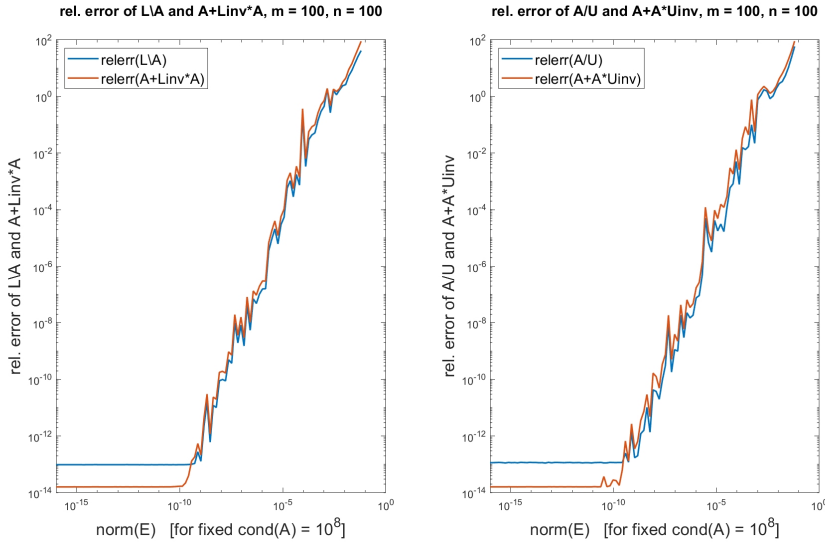


FIG. 2. Relative error of $L \setminus A$ vs. $A + \mathbf{Lin}v \cdot A$ and A/U vs. $A + A \cdot \mathbf{U}inv$

366 As can be seen both the results for $L^{-1}A$ and AU^{-1} improve for perturbations E of
 367 the identity matrix up to 10^{-9} . For larger E there is not much difference.
 368

369 **3.2. LU decomposition of general $A \in M_{m,n}$ with $m = n$.** Let $A \in M_n$ be
 370 given. We first compute approximate factors \tilde{L} and \tilde{U} by an LU decomposition with
 371 partial pivoting and permute the rows of A accordingly.

372 Let $\tilde{X}_L \approx \tilde{L}^{-1}$ and $\tilde{X}_U \approx \tilde{U}^{-1}$ be approximate preconditioners, so that $I_E :=$
 373 $\tilde{X}_L A \tilde{X}_U$ is a perturbed identity matrix. In MATLAB we may compute \tilde{X}_L by `inv(Ls)`
 374 as a left inverse, however, use `eye(n)/Us` to compute \tilde{X}_U as a right inverse, cf. [11].
 375 Then the uniqueness of the LU decomposition and $I_E = L_E U_E$ imply

376
$$L = \tilde{X}_L^{-1} L_E \quad \text{and} \quad U = U_E \tilde{X}_U^{-1} .$$

377 Note that in Algorithm 3.2 (**LU_E**) the offsets \mathbf{LE} and \mathbf{UE} of L_E and U_E to the identity
 378 matrix are computed, respectively. The computational effort is $\mathcal{O}(n^3)$ operations.

²Recall that A/B is the MATLAB notation for AB^{-1} .

TABLE 1
Condition number and sensitivity of the LU-factors for $n = 100$ and different condition numbers

cond(A)	10^2	10^5	10^8	10^{11}	10^{14}
cond(L)	$1.4 \cdot 10^2$	$1.4 \cdot 10^2$	$1.4 \cdot 10^2$	$1.4 \cdot 10^2$	$1.4 \cdot 10^2$
cond(U)	$2.7 \cdot 10^2$	$4.8 \cdot 10^4$	$2.4 \cdot 10^7$	$1.7 \cdot 10^{10}$	$1.3 \cdot 10^{13}$
sensitivity(L)	$5.2 \cdot 10^{-15}$	$7.7 \cdot 10^{-13}$	$3.4 \cdot 10^{-10}$	$2.0 \cdot 10^{-7}$	$1.4 \cdot 10^{-4}$
sensitivity(U)	$3.3 \cdot 10^{-15}$	$5.2 \cdot 10^{-15}$	$6.9 \cdot 10^{-15}$	$8.1 \cdot 10^{-15}$	$9.4 \cdot 10^{-15}$

379 It is well known that the condition number of A moves into the U -factor, i.e.,
 380 the factor L will be well conditioned whereas $\text{cond}(U) \sim \text{cond}(A)$. Thus we might
 381 expect the factor U to be sensitive to perturbations in A whereas L is not so sensitive.
 382 However, the opposite is true, see Table 1 for square matrices of dimension $n = 100$.
 383 As can be seen the condition number of L is small, that of U is of the order of $\text{cond}(A)$.
 384 For the sensitivity displayed in the last two rows we perturb the matrix A into \tilde{A} by
 385 changing each entry of A randomly by 1 bit in the mantissa and display $\|\tilde{L} - L\|/\|L\|$,
 386 and for U correspondingly. As can be seen, both L and U are insensitive for small
 387 condition number, however, for ill-conditioned A a perturbation of the last bit of A
 388 changes L relatively by about 10^{-4} , whereas U changes only about in the last bit.
 389 This is in accordance with our rule of thumb (2.1). The numbers are the median of
 390 1000 samples.

391 The reasoning for this rule of thumb (2.1) in Section 2 relied on the relation of
 392 the magnitude of the elements of U . It was also mentioned in a footnote that this
 393 relation is often not true for ill-conditioned matrices generated randomly by `sprand`
 394 with density 1, i.e., dense matrices. All entries of the factor U of such matrices
 395 are often not far from 1 in magnitude except one or two very small entries on the
 396 diagonal, often not U_{nn} . We measured the sensitivity of L and U for square matrices
 397 of dimension $n = 100$ similar to Table 1 and display the results in Table 2.

TABLE 2
Condition number and sensitivity of the LU-factors for $n = 100$ and matrices generated by `sprand`

cond(A)	10^2	10^5	10^8	10^{11}	10^{14}
cond(L)	$7.0 \cdot 10^1$	$7.2 \cdot 10^1$	$7.3 \cdot 10^1$	$7.4 \cdot 10^1$	$7.6 \cdot 10^1$
cond(U)	$5.4 \cdot 10^2$	$4.1 \cdot 10^5$	$3.9 \cdot 10^8$	$3.8 \cdot 10^{11}$	$4.0 \cdot 10^{14}$
sensitivity(L) median	$2.0 \cdot 10^{-15}$	$1.7 \cdot 10^{-13}$	$2.4 \cdot 10^{-10}$	$3.4 \cdot 10^{-8}$	$2.3 \cdot 10^{-4}$
sensitivity(L) mean	$2.3 \cdot 10^{-15}$	$6.0 \cdot 10^{-13}$	$5.7 \cdot 10^{-10}$	$4.3 \cdot 10^{-7}$	$6.9 \cdot 10^{-4}$
sensitivity(U) median	$1.5 \cdot 10^{-15}$	$3.5 \cdot 10^{-15}$	$5.9 \cdot 10^{-15}$	$7.5 \cdot 10^{-15}$	$7.9 \cdot 10^{-15}$
sensitivity(U) mean	$1.7 \cdot 10^{-15}$	$1.5 \cdot 10^{-13}$	$5.6 \cdot 10^{-11}$	$1.2 \cdot 10^{-7}$	$1.4 \cdot 10^{-4}$

398 As before the input matrices were perturbed entrywise and randomly by one bit, and
 399 the decompositions were performed using the multiple precision package [2] to avoid
 400 distortion of the sensitivity by rounding errors. For the median of the sensitivities of
 401 L and U there is not too much difference to Table 1. The mean and median of the
 402 sensitivity of L are similar³, so that seems to support (2.1). However, the mean of the
 403 sensitivities of U is larger than the median. Thus a few entries of U seem sensitive
 404 to perturbations, but the majority is not. So basically the rule of thumb (2.1) seems
 405 still applicable, but we don't have an explanation for that behavior.

406 The quality of the bounds depend on how close I_E is to the identity matrix, i.e., for

³All medians and means in Table 1 are all similar, so only the medians are displayed.

407 $I + E := I_E$ we want $\|E\|$ to be as small as possible. The median of the relative errors
 408 of all inclusion components of L and of U is displayed in Figure 3 from left to right,
 409 respectively, for condition numbers from 1 to 10^{15} . We first use interval arithmetic for
 410 the computation of $I_E := \tilde{X}_L A \tilde{X}_U$ and for L and U and display the relative errors in
 411 red. As expected, the error grows with the condition number. For condition numbers
 412 close to 10^{15} the inclusions are still accurate to about 8 to 10 decimal figures. The
 413 spikes for very large condition number indicate that the verification failed.

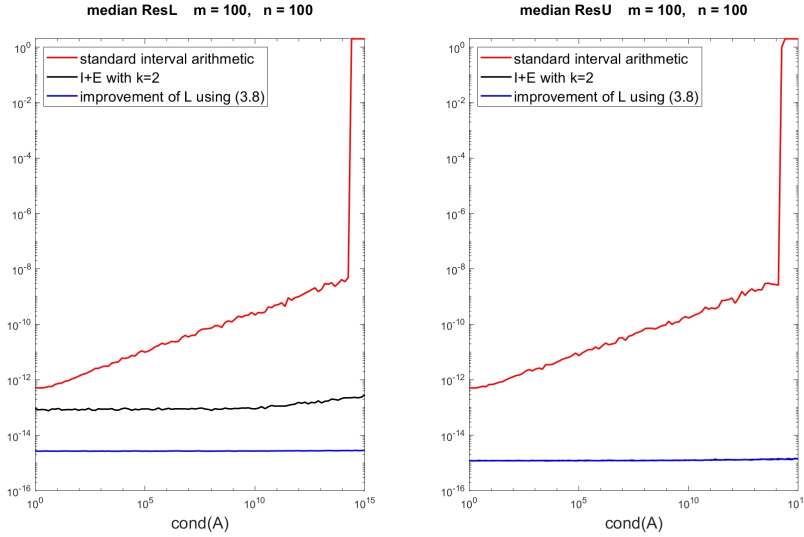


FIG. 3. Norm of error bounds for the L - and U -factor for different condition numbers

414 To achieve more accurate bounds we compute I_E as $\tilde{X}_L(A\tilde{X}_U)$ and use for both
 415 products two-fold precision, equivalent to double-double precision. In the legend of
 416 Figure 3 this is indicated by $k = 2$. The result is displayed in Figure 3 in black, where
 417 in the right picture the black curve is identical to the blue curve to be defined. As can
 418 be seen the accuracy of U is now close to maximal precision equivalent to 16 decimal
 419 places for all condition numbers, and for L the accuracy is a little bit less. That
 420 corresponds to Table 1, i.e., we expect better inclusions for U rather than that of L .
 421 In order to improve the accuracy of the L -factor, we use $L = LU\tilde{X}_U U_E^{-1}$ implying

422 (3.8)
$$L = A\tilde{X}_U U_E^{-1} \quad \text{and} \quad U = U_E \tilde{X}_U^{-1}$$

423 for the square case $m = n$. The product $A\tilde{X}_U$ is computed in doubled precision, the
 424 result is displayed as the blue curve in Figure 3. Now for all condition numbers and
 425 all entries of the factors L and U the bounds are of almost maximal accuracy, except
 426 for $\text{cond}(A) = 10^{15}$ where the verification failed.

427 **3.3. LU decomposition of general $A \in M_{m,n}$ with $m > n$.** Let $A \in M_{m,n}$ be
 428 given with $m > n$. We first compute an approximate LU decomposition with partial
 429 pivoting and permute the rows of A accordingly. Thus we may assume that the upper
 430 square block A_n has an LU decomposition. Following the approach discussed at the

431 beginning of this section

$$432 \quad A = LU = \begin{pmatrix} A_n \\ A_{\bar{n}} \end{pmatrix} = \begin{pmatrix} L_n \\ L_{\bar{n}} \end{pmatrix} U$$

433 so that

$$434 \quad X_L := \begin{pmatrix} L_n^{-1} & O_{n,m-n} \\ -L_{\bar{n}}L_n^{-1} & I_{m-n} \end{pmatrix}, \quad X_U := U^{-1}$$

435 implies

$$436 \quad X_L A X_U = \begin{pmatrix} I_n \\ O_{m-n,n} \end{pmatrix}.$$

437 We compute approximations $\tilde{X}_L \approx X_L$ and $\tilde{X}_U \approx X_U$ using an approximate LU
 438 decomposition $A \approx \tilde{L}\tilde{U}$, so that $I_E := \tilde{X}_L A \tilde{X}_U$ is a perturbed identity matrix. Then
 439 $I_E = L_E U_E$ implies

$$440 \quad \tilde{X}_L := \begin{pmatrix} P & O_{n,m-n} \\ Q & I_{m-n} \end{pmatrix} \Rightarrow L = \begin{pmatrix} P^{-1} & O_{n,m-n} \\ -QP^{-1} & I_{m-n} \end{pmatrix} L_E$$

441 and $U = U_E \tilde{X}_U^{-1}$. The computational effort is $\mathcal{O}(P^2 p)$ operations for $P = \max(m, n)$
 442 and $p = \min(m, n)$.

TABLE 3
 Condition number and sensitivity of the LU-factors for $m = 200$ and $n = 100$

cond(A)	10^2	10^5	10^8	10^{11}	10^{14}
cond(L)	$8.6 \cdot 10^1$	$8.5 \cdot 10^1$	$8.6 \cdot 10^1$	$8.6 \cdot 10^1$	$8.5 \cdot 10^1$
cond(U)	$3.4 \cdot 10^2$	$9.1 \cdot 10^4$	$5.5 \cdot 10^7$	$4.1 \cdot 10^{10}$	$3.5 \cdot 10^{13}$
sensitivity(L_1)	$6.4 \cdot 10^{-15}$	$1.3 \cdot 10^{-12}$	$6.6 \cdot 10^{-10}$	$4.3 \cdot 10^{-7}$	$3.1 \cdot 10^{-4}$
sensitivity(L_2)	$1.3 \cdot 10^{-14}$	$3.6 \cdot 10^{-12}$	$2.3 \cdot 10^{-9}$	$1.8 \cdot 10^{-6}$	$1.5 \cdot 10^{-3}$
sensitivity(U)	$2.7 \cdot 10^{-15}$	$4.5 \cdot 10^{-15}$	$5.9 \cdot 10^{-15}$	$7.1 \cdot 10^{-15}$	$8.5 \cdot 10^{-15}$

443 In Table 3 we show the sensitivity of the upper square block L_n , the remaining
 444 part $L_{\bar{n}}$ of L and of the U -factor, again the median over 1000 samples. Similar
 445 to the square case and as predicted by (2.1), with increasing condition number the
 446 components of the L -factor are getting sensitive to perturbations in A , while those of
 447 U are not. Hence, as for square A , we expect less accurate inclusions of L .

448 The median of the relative errors of all inclusion components of the upper square
 449 part L_n of L , the remaining part $L_{\bar{n}} = L(n+1 : m, :)$ and of U is displayed in Figure
 450 4 from left to right, respectively, for condition numbers from 1 to 10^{15} . As before we
 451 first use interval arithmetic for the computation of $I_E := \tilde{X}_L A \tilde{X}_U$ and to compute
 452 L and U and display the relative errors in red. As expected, the error grows with
 453 the condition number. For condition numbers close to 10^{15} the inclusions are still
 454 accurate to about 8 to 10 decimal figures.

455 To achieve more accurate bounds we compute I_E as $\tilde{X}_L(A\tilde{X}_U)$ and use for both
 456 products two-fold precision, equivalent to double-double precision. The result is dis-
 457 played in Figure 4 in black. As can be seen the accuracy of U is now close to maximal
 458 precision equivalent to 16 decimal places for all condition numbers, for the upper part
 459 of L it improved significantly, and for the lower part of L the accuracy decreases from
 460 condition number 10^9 .

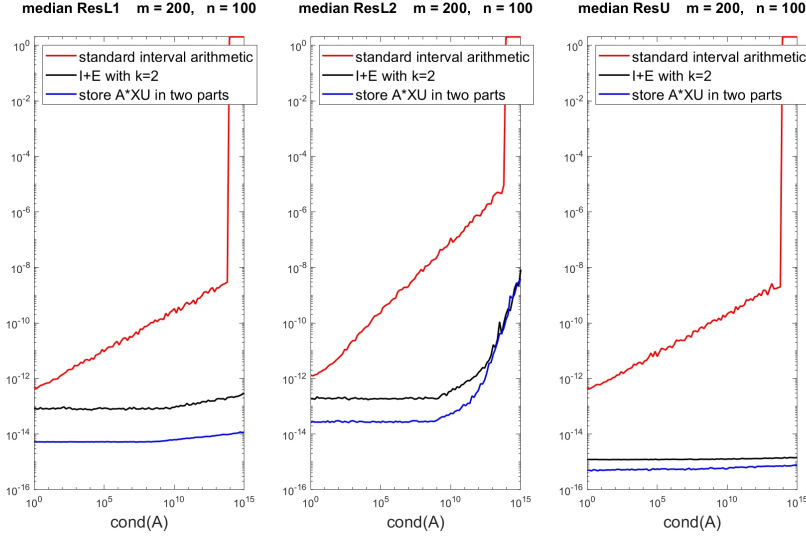


FIG. 4. First methods: Norm of error bounds for the L - and U -factor based on A

461 In order to obtain flat curves in all three pictures, i.e., close to maximal accuracy
 462 for all components of the L - and the U -factor, we compute the product $A\tilde{X}_U$ again in
 463 doubled precision but store it in two parts $C_1 + C_2$, and then compute $\tilde{X}_L C_1 + \tilde{X}_L C_2$
 464 in doubled precision but store it in one matrix I_E . The result displayed in blue in
 465 Figure 4 is better than before, however, there is still a growth of the errors of the
 466 lower part of L from condition number 10^9 .

467 The following alternative approach is faster and better. To that end we use an
 468 approximate LU decomposition $A_n \approx \tilde{L}_n \tilde{U}$ of the upper left square block of A . For
 469 approximations $\tilde{X}_{L_n} \approx \tilde{L}_n^{-1}$ and $\tilde{X}_U \approx \tilde{U}^{-1}$ let $I_E = \tilde{X}_{L_n} A_n \tilde{X}_U = L_E U_E$. Then

470 (3.9)
$$U = U_E \tilde{X}_U^{-1} \quad \text{and} \quad L = A \tilde{X}_U U_E^{-1}$$

471 using $L = LU \tilde{X}_U U_E^{-1}$. Thus, for $m \geq n$ we use the same formula as (3.8) for the
 472 square case. The computational effort is $\mathcal{O}(Pp^2)$ operations.

473 As before display the median of the relative errors of all components of the upper
 474 square part L_n of L , the lower part of L and of U in Figure 5. The red curves are
 475 the results when using interval arithmetic to compute $I_E := \tilde{X}_L A \tilde{X}_U$, L and U and
 476 are similar to those before. Using extra precision to compute $\tilde{X}_L(A\tilde{X}_U)$ is shown in
 477 black and is for both parts of L better than before. Finally, for the blue curve we used
 478 doubled precision to compute $A\tilde{X}_U$ and stored the result in two parts, which are then
 479 multiplied by \tilde{X}_L . That last method yields for all condition numbers and all entries
 480 of the factors L and U bounds of almost maximal accuracy.

481 **3.4. LU decomposition of general $A \in M_{m,n}$ with $m < n$.** Let $A \in M_{m,n}$
 482 with $m < n$ be given. Now the partial pivoting of an approximate LU decomposition
 483 may take only the left square block A_m into account. Therefore, we first compute an
 484 approximate LU decomposition with partial pivoting of A^T and permute the columns
 485 of A accordingly, followed by the computation an approximate LU decomposition of
 486 A with partial pivoting and permute the rows of A accordingly. We may assume that

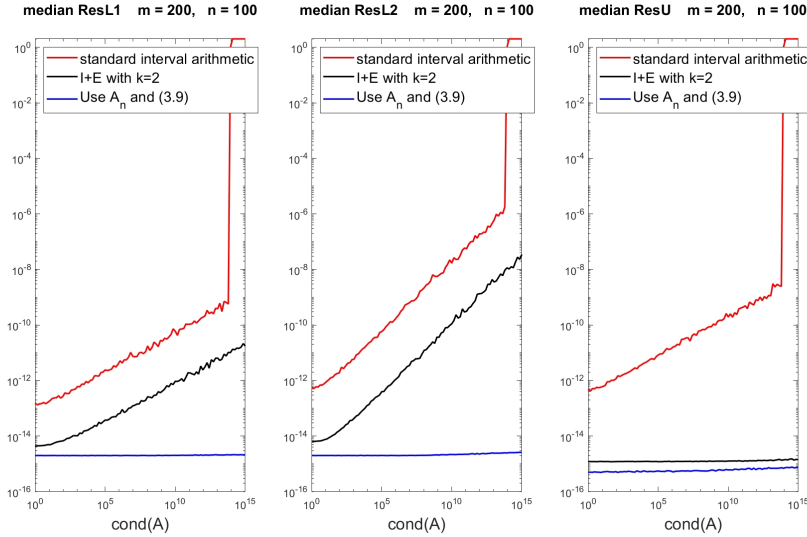


FIG. 5. *Second methods: Norm of error bounds for the L- and U-factor based on A_n*

487 the left square block A_m has an LU decomposition. Then

488
$$A = LU = \begin{pmatrix} A_m & A_{\overline{m}} \end{pmatrix} = L \begin{pmatrix} U_m & U_{\overline{m}} \end{pmatrix}$$

489 and

490
$$X_L := L^{-1}, X_U = \begin{pmatrix} U_m^{-1} & -U_m^{-1}U_{\overline{m}} \\ O_{n-m,m} & I_{n-m} \end{pmatrix}$$

491 implies $X_L A X_U = \begin{pmatrix} I_m & O_{m,n-m} \end{pmatrix}$. We compute approximations $\tilde{X}_L \approx X_L$ and
 492 $\tilde{X}_U \approx X_U$ using an approximate LU decomposition $A \approx \tilde{L}\tilde{U}$, so that again $I_E :=$
 493 $\tilde{X}_L A \tilde{X}_U$ is a perturbed identity matrix. Then $I_E = L_E U_E$ implies

494
$$\tilde{X}_U := \begin{pmatrix} P & Q \\ O_{n-m,m} & I_{n-m} \end{pmatrix} \Rightarrow L = \tilde{X}_L^{-1} L_E$$

495 and

496
$$U = U_E \begin{pmatrix} P^{-1} & -P^{-1}Q \\ O_{n-m,m} & I_{n-m} \end{pmatrix}.$$

497 The computational effort is $\mathcal{O}(P^2 p)$ operations for $P = \max(m, n)$ and $p = \min(m, n)$.

498 The sensitivity of the L -factor, the left square block U_m and the remaining of the
 499 U -factor is displayed in Table 4 and is, as predicted in (2.1), similar to square A or
 500 the case $m > n$. Again we expect it to be more difficult to obtain narrow inclusions
 501 of L .

502 Computational results are shown in Figure 6. The median of relative errors of all
 503 components of the L -factor, the left square part U_m and the remaining of the U -factor
 504 are shown from left to right. The color coding is as in the previous subsections, i.e.,

TABLE 4
Condition number and sensitivity of the LU-factors for $m = 100$ and $n = 200$

$\text{cond}(A)$	10^2	10^5	10^8	10^{11}	10^{14}
$\text{cond}(L)$	$1.4 \cdot 10^2$	$1.4 \cdot 10^2$	$1.4 \cdot 10^2$	$1.4 \cdot 10^2$	$1.4 \cdot 10^2$
$\text{cond}(U)$	$2.7 \cdot 10^2$	$4.8 \cdot 10^4$	$2.4 \cdot 10^7$	$1.7 \cdot 10^{10}$	$1.3 \cdot 10^{13}$
$\text{sensitivity}(L)$	$1.2 \cdot 10^{-14}$	$2.5 \cdot 10^{-12}$	$1.3 \cdot 10^{-9}$	$8.4 \cdot 10^{-7}$	$6.0 \cdot 10^{-4}$
$\text{sensitivity}(U_1)$	$3.8 \cdot 10^{-15}$	$5.4 \cdot 10^{-15}$	$7.0 \cdot 10^{-15}$	$8.3 \cdot 10^{-15}$	$9.6 \cdot 10^{-15}$
$\text{sensitivity}(U_2)$	$9.4 \cdot 10^{-15}$	$1.3 \cdot 10^{-14}$	$1.7 \cdot 10^{-14}$	$1.9 \cdot 10^{-14}$	$2.2 \cdot 10^{-14}$

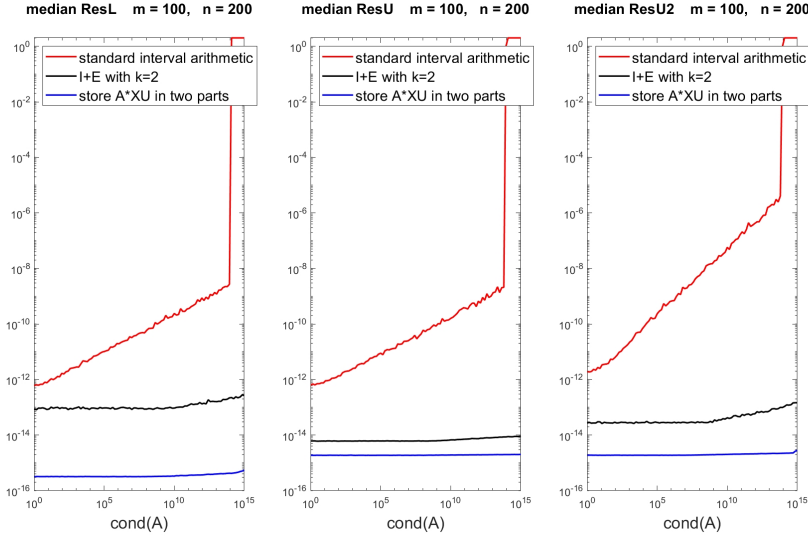


FIG. 6. First methods: Norm of error bounds for the L - and U -factor based on A

505 for the red curve only standard interval arithmetic was used and the expected growth
 506 with the condition number can be seen.

507 For the black curve the two products in $\tilde{X}_L(A\tilde{X}_U)$ are computed in doubled
 508 precision. Supposedly, the results are better than for the case $m \geq n$ because we
 509 performed initially two approximate LU decompositions to identify the permutations
 510 of columns and rows.

511 In order to obtain flat curves close to maximal accuracy for all components of
 512 both the L - and the U -factor, we compute $C_1 + C_2 = A\tilde{X}_U$ in doubled precision and
 513 use two matrices to store the result, and then compute $\tilde{X}_L C_1 + \tilde{X}_L C_2$ in doubled
 514 precision but store the result in one matrix I_E .

515 Now the results are close to maximal accuracy, shown in blue, but the computing
 516 time of $\mathcal{O}(P^2p)$ operations can be improved into $\mathcal{O}(Pp^2)$ operations. Similar as before
 517 we use an approximate LU decomposition $A_m \approx \tilde{L}\tilde{U}_m$ of the left square block of A .
 518 For approximations $\tilde{X}_L \approx \tilde{L}^{-1}$ and $\tilde{X}_{U_m} \approx \tilde{U}_m^{-1}$ let $I_E = \tilde{X}_L A_m \tilde{X}_{U_n} = L_E U_E$. Then

519 (3.10)
$$L = \tilde{X}_L^{-1} L_E \quad \text{and} \quad U = L_E^{-1} \tilde{X}_L A$$

520 using $U = L_E^{-1} \tilde{X}_L L U$ for the latter equality. Now the computational effort is $\mathcal{O}(Pp^2)$
 521 operations.

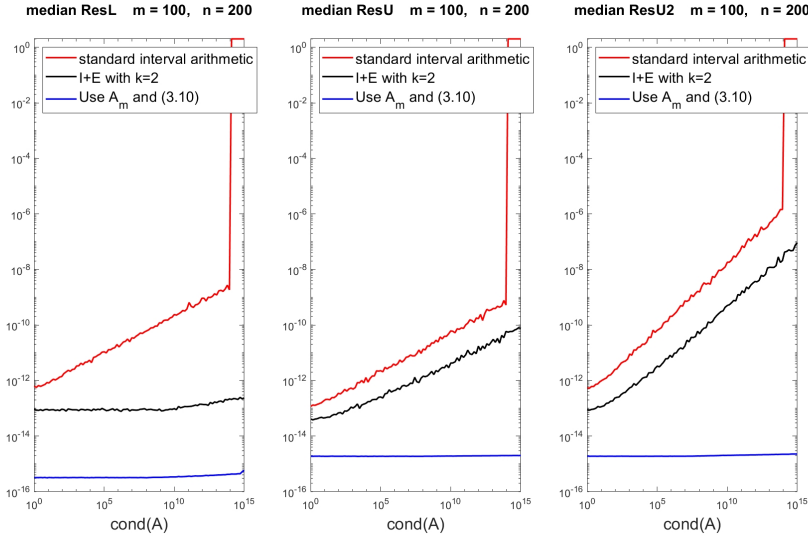


FIG. 7. *Second methods: Norm of error bounds for the L - and U -factor based on A_m*

522 Computational results are shown in Figure 7. Using the same color coding the
 523 main difference is in the black curve, computing the two products $\tilde{X}_L(A\tilde{X}_U)$ in dou-
 524 bled precision but storing either result in one matrix. The results in Figure 6
 525 are better due to the use of the right part $U_{\bar{m}}$ in the computations.

526 The third and best method, shown in blue in Figure 7, is to store $A\tilde{X}_U$ in two
 527 matrices and proceed as before. The resulting curves are flat and close to the relative
 528 rounding error unit 10^{-16} for all condition numbers and all components of the L -
 529 and the U -factor.

530 **4. Cholesky decomposition.** As for the Cholesky decomposition of a sym-
 531 metric positive definite matrix $A \in M_n$ we have all necessary ingredients. We would
 532 like to stress that A being positive definite is not an assumption, because if so, it
 533 would have to be verified *before* starting the computation. In contrast, the property
 534 is verified *a posteriori*, i.e., if successful the matrix has been proved to be positive
 535 definite.

536 To compute bounds for the Cholesky factor, we first use an approximate Cholesky
 537 factor \tilde{G} to precondition A into $I_E := X_G^T A X_G$ for $X_G \approx \tilde{G}^{-1}$. We discussed in Section
 538 3.1 how to obtain verified inclusions for the LU decomposition $I_E = L_E U_E$. That
 539 includes in particular the diagonal D of U_E . The uniqueness of the LU and Cholesky
 540 decomposition implies that $G_E = D^{1/2} L_E^T$ is the Cholesky factor of I_E . Hence, an
 541 inclusion of the Cholesky factor may be computed by

$$542 \quad (4.1) \quad G_E^T G_E = I_E = X_G^T A X_G \quad \Rightarrow \quad G = G_E X_G^{-1} = D^{1/2} L_E^T X_G^{-1} .$$

543 The computational effort is $\mathcal{O}(n^3)$ operations.

544 We first show the median of the sensitivity of the Cholesky factor for 1000 samples
 545 in Table 5. The condition number of G is, of course, the square root of that of A ,
 546 and the sensitivity corresponds to that predicted in (2.1). In some way it seems the
 547 geometric mean between the sensitivity of the L - and the U -factor.

TABLE 5
Condition number and sensitivity of the Cholesky factor for different condition numbers

cond(A)	10 ²	10 ⁵	10 ⁸	10 ¹¹	10 ¹⁴
cond(G)	1.0 · 10 ¹	3.2 · 10 ²	10.0 · 10 ³	3.2 · 10 ⁵	10.0 · 10 ⁶
sensitivity(G)	6.8 · 10 ⁻¹⁶	1.8 · 10 ⁻¹⁴	4.9 · 10 ⁻¹³	1.4 · 10 ⁻¹¹	4.0 · 10 ⁻¹⁰

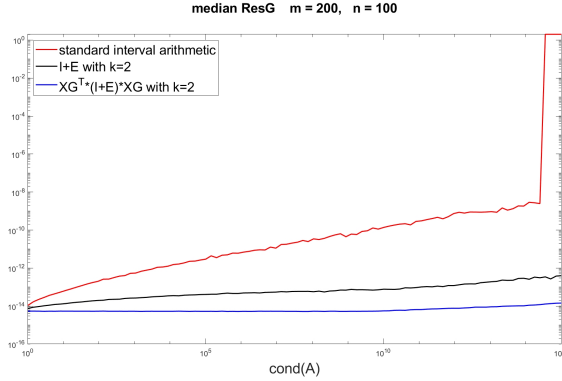


FIG. 8. Norm of error bounds for the Cholesky factor for different condition numbers

548 The computational results for (4.1) are shown in Figure 8. We display the median
 549 of the relative errors of all components of the inclusion of the Cholesky factor for
 550 different condition numbers. The color coding is similar to the previous sections. For
 551 the red curve we use interval arithmetic to compute inclusions of $I_E = X_G^T A X_G$ and
 552 for the inclusion of G according to (4.1). As expected, errors grow with the condition
 553 number but still guaranteeing some 8 correct decimal figures up to condition number
 554 $3 \cdot 10^{14}$ and failure above.

555 The black curve shows the results for computing $X_G^T A X_G$ in doubled precision.
 556 The quality of the inclusion is better and there is no failure.

557 Finally, we compute $C_1 + C_2 = A X_G$ in doubled precision with two results and
 558 $X_G^T C_1 + X_G^T C_2$ again in doubled precision but with one result I_E . Now, shown in
 559 blue, for all condition numbers all components of the Cholesky factor are enclosed
 560 with almost maximal accuracy.

561 **5. QR decomposition.** Assume $A \in M_{m,n}$ with $m \geq n$ to be given with full
 562 rank. Then there is a unique QR decomposition with orthonormal columns and upper
 563 triangular R with non-negative diagonal entries [14, Theorem 2.1.14]. Consider

564
$$A := \begin{pmatrix} 1 & e \\ 1 & e + \varphi e^2 \end{pmatrix}$$

565 with

566
$$Q := \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & -\varphi \\ 1 & \varphi \end{pmatrix} \quad \text{and} \quad R := \frac{1}{\sqrt{2}} \begin{pmatrix} 2 & (2 + \varphi e)e \\ 0 & e^2 \end{pmatrix}.$$

567 Then $A = QR$ for $e, \varphi \in \mathbb{R}$, and for $\varphi = -1$ and $\varphi = 1$ this is the unique QR decom-
 568 position of $A(\varphi)$ for any $e \in \mathbb{R}$. Hence, the computation of the QR decomposition is
 569 ill-posed at $e = 0$ because an arbitrary small perturbation causes a finite change in

570 Q . As a consequence, a verification method is only applicable to matrices with full
571 rank.

572 Let $A \approx \tilde{Q}\tilde{R}$ be an approximate “economy-size” QR decomposition, i.e., $\tilde{Q} \in$
573 $M_{m,n}$ and $\tilde{R} \in M_n$. For $\tilde{X}_R \approx \tilde{R}^{-1}$ we expect $C := A\tilde{X}_R$ to be close to unitary, so
574 that $C^T C$ will be a small perturbation of the identity matrix I_n .

575 In fact, $C^T C \approx \tilde{X}_R^T A^T A \tilde{X}_R =: I + E$ is the same as preconditioning $A^T A$ by
576 the inverse of the approximate Cholesky factor \tilde{R} of $A^T A$. However, the formulation
577 $C^T C \approx (\tilde{X}_R^T A^T) \cdot (A \tilde{X}_R)$ avoids to form the matrix $A^T A$ with squared condition
578 number $\text{cond}(A)^2$.

579 We use the method described in the previous section to compute an inclusion
580 G_E of the Cholesky factor of $I + E$, so that $\tilde{X}_R^T A^T A \tilde{X}_R = I + E = G_E^T G_E$. Hence
581 $R := G_E \tilde{X}_R^{-1}$ is the Cholesky factor of $A^T A$, which in turn is the R -factor of the QR
582 decomposition of A . An inclusion of the economy-size Q -factor is obtained by $Q_1 =$
583 AR^{-1} provided that R is non-singular. A second possibility is to use $Q_1 = A\tilde{X}_R G_E^{-1}$.

584 For the full-size QR -factors note that $Q = (Q_1 \ Q_2)$ where Q_2 is the orthogonal
585 complement of Q_1 and, provided that A has full rank, a basis for the null space of A^* .
586 In [17] several methods are discussed to compute an inclusion of a basis of the null
587 space of a rectangular matrix. For \tilde{Q}_2 denoting an approximation of the orthogonal
588 complement of Q_1 , the solution X of the square linear system

$$589 \quad (5.1) \quad \begin{pmatrix} A^* \\ \alpha \tilde{Q}_2^* \end{pmatrix} X = \begin{pmatrix} O_{n,m-n} \\ \alpha I_{m-n} \end{pmatrix}$$

590 does the job [17], i.e., X spans the orthogonal complement Q_2 of Q_1 . If \mathbf{A} is an
591 interval matrix, then this is true for every $A \in \mathbf{A}$. We choose α within $[\sigma_n(A), \sigma_1(A)]$
592 to ensure that the condition number of the system matrix in (5.1) is equal to that of
593 A .

594 We can expect X to be numerically unitary, but not mathematically. The follow-
595 ing lemma from [34] estimates the distance to an orthonormal basis.

596 LEMMA 5.1. *Let $X, Y \in M_{m,n}$ with $m \geq n$ be given. Define $\alpha := \|I - X^* X\|$ and*
597 *$\delta := \|X - Y\|$. Let \mathcal{V} be an n -dimensional subspace of \mathbb{K}^m that contains all columns*
598 *of Y . Then there exists $Q \in M_{m,n}$ with $Q^* Q = I$ whose columns span \mathcal{V} and*

$$599 \quad \|Q - X\| \leq \alpha + \sqrt{2}\delta.$$

600 The bound is sharp. Note that the bound remains true even if $\alpha \geq 1$, although that
601 may not be very useful. In our practical application α is of the order of the relative
602 rounding error unit $\mathbf{u} \approx 10^{-16}$.

603 The application of Lemma 5.1 is as follows. A very good approximate solution to
604 (5.1) is $X := \tilde{Q}_2$, for which also $\alpha := \|I - X^* X\|$ is close to the relative rounding error
605 unit. Define Y to be the true solution of (5.1). An inclusion of Y is computed by
606 verification methods. In fact, the inclusion will be of the form $\tilde{Q}_2 + \Delta$ for an interval
607 matrix Δ with small norm [32], so that $\delta = \|\Delta\|$. It follows that $\tilde{Q}_2 \pm \delta$ is an inclusion
608 of the orthogonal complement Q_2 to Q_1 . Hence, $(Q_1 \ Q_2)$ is the full Q -factor of the

609 QR decomposition of A , where the full R -factor is $\begin{pmatrix} R \\ O_{m-n,n} \end{pmatrix}$.

610 For $A \in M_{m,n}$ with $m < n$ we compute inclusions of the (full) QR decomposition
611 of the square matrix A_m , so that $A_m = QR_m$ implies $R = Q^* A$. The computational
612 effort for the inclusion is the same as to compute an approximate decomposition.

613 To judge the computational results we first check on the median of the sensitivity
 614 of R and the two parts of the factor Q of $A \in M_{200,100}$ for 1000 samples. The results
 615 are displayed in Table 6. The factor Q is perfectly conditioned, however, sensitive
 616 to perturbations in A . The factor R has the same condition number as A , but is
 617 insensitive to small perturbations in A in accordance with (2.1). The corresponding
 618 data for $m < n$ is completely similar, only Q is sensitive to perturbations in A . Thus
 619 we may expect more problems in the computation of narrow bounds of Q .

TABLE 6
 Sensitivity of the two parts of the factors Q and R for different condition numbers

cond(A)	10^2	10^5	10^8	10^{11}	10^{14}
sensitivity(Q_1)	$1.2 \cdot 10^{-14}$	$6.0 \cdot 10^{-12}$	$4.5 \cdot 10^{-9}$	$3.7 \cdot 10^{-6}$	$3.3 \cdot 10^{-3}$
sensitivity(Q_2)	$1.7 \cdot 10^{-14}$	$8.2 \cdot 10^{-12}$	$6.1 \cdot 10^{-9}$	$5.0 \cdot 10^{-6}$	$4.5 \cdot 10^{-3}$
sensitivity(R)	$4.8 \cdot 10^{-16}$	$5.6 \cdot 10^{-16}$	$6.2 \cdot 10^{-16}$	$7.0 \cdot 10^{-16}$	$7.8 \cdot 10^{-16}$

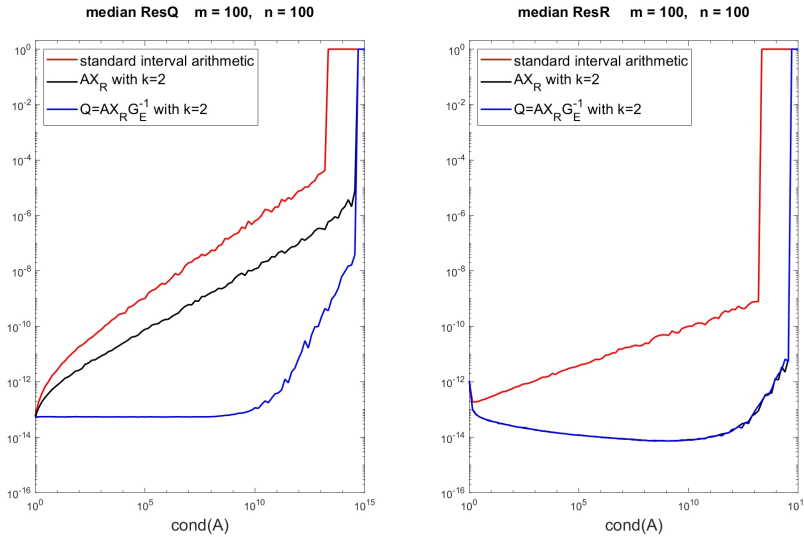


FIG. 9. Norm of error bounds for the Q - and R -factor for different condition numbers for $m = n$

620 Next we show the median of the relative errors of all entries of the inclusions of Q
 621 and R . We start with a square matrix $A \in M_n$. We first compute $C = A\tilde{X}_R$ as well
 622 the inclusions $G_E\tilde{X}_R^{-1}$ and $Q = AR^{-1}$ in standard interval arithmetic. The result for
 623 different condition numbers is the red curve in Figure 9. We observe an increase of the
 624 relative errors proportional to the condition number, and as predicted less accurate
 625 bounds for Q .

626 Secondly, we compute an inclusion $C = A\tilde{X}_R$ with doubled precision with one
 627 output result. The product C^*C is computed in doubled precision as well, otherwise
 628 we use standard interval arithmetic. The result is the black curve in Figure 9. It is
 629 much better than before, in particular the inclusion of R .

630 Finally, we use the second possibility $Q_1 = A\tilde{X}_R G_E^{-1}$ for the inclusion of Q_1 ,
 631 where the first product $Q_1 = A\tilde{X}_R$ is computed in doubled precision. The result is

632 shown as the blue curve in Figure 9, where the black and blue curves are practically
 633 identical for R .

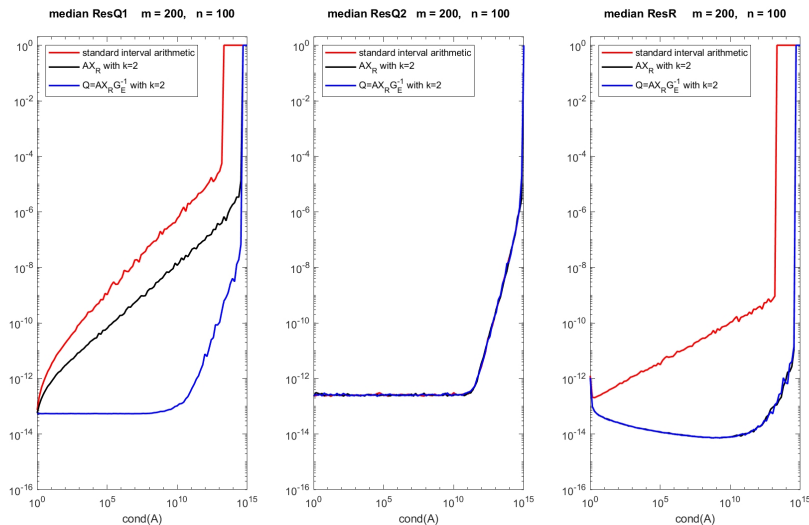


FIG. 10. Norm of error bounds for the Q - and R -factor for different condition numbers for $m > n$

634 The results for $A \in M_{m,n}$ with $m > n$ are shown in Figure 10. They look quite
 635 similar to those in Figure 9 for square A . In particular the quality of the two parts of
 636 Q shows no surprises. That is also true for the case $m < n$, so we omit to show that
 637 graph.

638 Until now we refrained from giving computing times of our inclusion methods,
 639 mainly because those are essentially dominated by the interpretation overhead in
 640 MATLAB. At least for one example, the QR decomposition, we show the ratio of
 641 computing time compared to the built-in MATLAB routine. For most problems we
 642 gave three inclusion methods, as for example in Figure 9. However, the timing is not
 643 too different, therefore we give only the time ratio for our best method compared to `qr`.
 644 As has been said this ratio is biased by the interpretation overhead, and in particular
 645 by the fact that our inclusion methods aim on highly accurate results. Therefore we
 646 also show the median relative error of the MATLAB result.

647 For the following Table 7 we generated real and complex square random matrices
 648 with condition number 10^{10} . From left the dimension and ratio of computing time of
 649 our best inclusion method together with the median relative error of all components
 650 of the floating-point approximation produced by $[Q,R] = \text{qr}(A)$; are displayed, for
 651 real matrices on the left and for complex matrices on the right.

652 As can be seen the verification method (in pure MATLAB code) is significantly
 653 slower than the built-in `qr`, but, according to Figure 9 and Table 7 also more accurate.
 654 Note that the accuracy of the MATLAB approximation is different from the sensitivity
 655 as displayed in Table 6.

656 We finally show the same table for rectangular matrices $A \in M_{m,n}$. We set
 657 $m := 2n$ and display the results in Table 8.

658 The median relative error of Q computed by `qr` is slightly weaker than for square
 659 matrices, more according to Table 6. Otherwise there is not too much difference in

TABLE 7
Ratio of computing times to MATLAB's `qr` for different condition numbers

n	real input			complex input		
	time ratio	Q	R	time ratio	Q	R
100	49.3	$1.2 \cdot 10^{-11}$	$6.8 \cdot 10^{-14}$	47.1	$1.2 \cdot 10^{-11}$	$6.8 \cdot 10^{-14}$
200	33.5	$1.4 \cdot 10^{-11}$	$8.3 \cdot 10^{-14}$	44.2	$1.3 \cdot 10^{-11}$	$7.9 \cdot 10^{-14}$
500	64.6	$6.9 \cdot 10^{-10}$	$4.5 \cdot 10^{-11}$	61.4	$1.2 \cdot 10^{-11}$	$1.1 \cdot 10^{-13}$
1000	58.3	$2.6 \cdot 10^{-9}$	$1.7 \cdot 10^{-8}$	48.2	$1.4 \cdot 10^{-11}$	$1.8 \cdot 10^{-13}$

TABLE 8
Ratio of computing times to MATLAB's `qr` for different condition numbers

n	real input			complex input		
	time ratio	Q	R	time ratio	Q	R
100	50.5	$1.7 \cdot 10^{-8}$	$4.4 \cdot 10^{-14}$	63.3	$2.3 \cdot 10^{-8}$	$5.6 \cdot 10^{-14}$
200	45.0	$1.2 \cdot 10^{-8}$	$5.2 \cdot 10^{-14}$	67.3	$1.9 \cdot 10^{-8}$	$6.6 \cdot 10^{-14}$
500	66.9	$9.6 \cdot 10^{-9}$	$3.5 \cdot 10^{-11}$	52.3	$1.8 \cdot 10^{-8}$	$1.0 \cdot 10^{-13}$
1000	52.3	$8.5 \cdot 10^{-9}$	$3.0 \cdot 10^{-8}$	56.6	$2.1 \cdot 10^{-8}$	$1.7 \cdot 10^{-13}$

660 the ratio of computing times or accuracy.

661 **6. Eigendecomposition.** A verified inclusion of an individual eigenvector to a
 662 multiple eigenvalue matrix is out of the scope of verification methods because the
 663 problem is ill-posed. In case of a non-trivial Jordan block of size k , there may be only
 664 one eigenvector which, after an arbitrarily small perturbation, changes into up to k
 665 individual eigenvectors.

666 Hence, the problem of computing a verified error bound for an individual eigen-
 667 vector is ill-posed, as well as to certify that an eigenvalue is not simple. However, an
 668 inclusion of a cluster and/or multiple eigenvalue becomes well posed if it is separated
 669 from the remaining spectrum. Then computing a basis for the corresponding invariant
 670 subspace is well posed as well.

671 There are approaches to compute error bounds for one cluster of eigenvalues
 672 together with invariant subspace [32, Theorem 13.9], however, the computing time for
 673 the complete eigendecomposition by applying that method to each cluster is $\mathcal{O}(n^4)$
 674 operations.

675 There are papers for computing inclusions of all eigenvalues and -vectors in $\mathcal{O}(n^3)$
 676 operations [25] for the symmetric positive definite case and [24] for general matrices.
 677 However, the given practical implementations face some problems. The algorithms
 678 in [33] for general and [34] for Hermitian matrices for the complete eigendecomposi-
 679 tion also require $\mathcal{O}(n^3)$ operations and are numerically stable. Of course, there are
 680 natural limitations for many large clusters. General matrices in [33] are, similar to
 681 the methods in this paper, transformed into a perturbed identity matrix, whereas the
 682 symmetric and Hermitian case in [34] is treated by generalized perturbation bounds.
 683 Both algorithms handle multiple or clustered eigenvalues as follows.

684 The output of either algorithm is an interval vector \mathbf{L} , an interval matrix \mathbf{X} and a
 685 cell array μ . If a cell element consists of a single element $\{k\}$, then \mathbf{L}_k is an inclusion
 686 of a simple eigenvalue and $\mathbf{X}(:, k)$ an inclusion of a corresponding eigenvector. A
 687 challenge for both algorithms is to identify clusters. To that end a threshold κ can be
 688 specified accepting eigenvalues with distance below κ to be a cluster. For $\kappa = 0$ the
 689 algorithms try, if possible, to produce individual inclusions for all eigenvalues.

690 If a cell element is a set μ_ℓ of indices, then the \mathbf{L}_k for $k \in \mu_\ell$ are identical and
 691 contain exactly $|\mu_\ell|$ eigenvalues, where the set of columns \mathbf{X}_k for $k \in \mu_\ell$ span the
 692 corresponding invariant subspace. For symmetric or Hermitian matrix, Lemma 5.1 is
 693 used to ensure that \mathbf{X} contains a unitary eigenvector basis.

694 There is a difference between the results for symmetric/Hermitian and for general
 695 matrices. In the former case the matrix is diagonalizable so that there exist $L \in \mathbf{L}$ and
 696 $X \in \mathbf{X}$ with $AX = XL$. Thus \mathbf{L} and \mathbf{X} are inclusions of eigenvalues and eigenvectors.
 697 A general matrix A may not be diagonalizable. If a cell element μ_ℓ consists of more
 698 than one index, i.e., $m := |\mu_\ell| > 1$, then the identical elements $\mathbf{A} := \mathbf{L}_k$ for $k \in \mu_\ell$
 699 contain m eigenvalues. That may be an m -fold or m distinct eigenvalues or any
 700 combination. In any case, the set of columns $\{\mathbf{X}_k : k \in \mu_\ell\}$ contains a basis Y of an
 701 invariant subspace of A . That implies existence of a matrix $M \in M_m$ with $AY = YM$,
 702 but there may be no diagonal M with this property.

703 However, the methods in [33] allow to compute a block diagonal interval matrix
 704 \mathbf{D} with the property that there exist $D \in \mathbf{D}$ and $X \in \mathbf{X}$ with $AX = XD$. For the
 705 Schur decomposition discussed in Section 8 it would be important to include upper
 706 triangular \mathbf{T} with the property that there exist $T \in \mathbf{T}$ and $X \in \mathbf{X}$ with $AX = XT$.
 707 However, that is not possible as eigenvectors of multiple eigenvalues need not be
 708 continuous, even for symmetric matrices. It was shown in [30] that the local behavior
 709 of an eigendecomposition of a matrix depending on several parameters may be quite
 710 different from the case of one parameter. The following example is adapted from [38]:

$$711 \quad (6.1) \quad A(e, f) := \begin{pmatrix} 1+f & e \\ e & 1 \end{pmatrix}.$$

712 The two matrices $A_1 := A(e, e)$ and $A_2 := A(e, 2e)$ have eigenvalues $1 + e/2 \pm e\sqrt{5}/2$
 713 and $1 + e \pm e\sqrt{2}$, respectively. So the eigenvalues depend continuously on e at $e = 0$.
 714 However, the corresponding orthogonal eigenvectors do not depend on e and are

$$715 \quad \begin{pmatrix} (1 - \sqrt{5})/2 \\ 1 \end{pmatrix}, \begin{pmatrix} (1 + \sqrt{5})/2 \\ 1 \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} 1 + \sqrt{2} \\ 1 \end{pmatrix}, \begin{pmatrix} 1 - \sqrt{2} \\ 1 \end{pmatrix}$$

716 for A_1 and A_2 , respectively. In other words, the computation of eigenvectors for
 717 multiple eigenvalues, even for symmetric matrices, is an ill-posed problem and outside
 718 the scope of verification methods.

719 The inclusion of the eigendecomposition offers a simple way to compute the matrix
 720 exponential and other matrix functions, however, only for non-defective matrices.

721 For detailed computational tests of inclusions for the eigenproblem see [33] and
 722 [34]. Here we only mention that error bounds of high quality are computed for a
 723 general real or complex, or symmetric or Hermitian matrix. The algorithms are
 724 applicable to interval matrices \mathbf{A} as well, where the inclusions are true for every
 725 $A \in \mathbf{A}$. The total computational effort is $\mathcal{O}(n^3)$ operations.

726 **7. Singular value and polar decomposition.** As for the eigendecomposition,
 727 the computation of singular vectors becomes ill-posed for multiple singular values, see
 728 example (6.1). Hence, as for the eigenproblem, inclusions for the subspaces span-
 729 ning the singular vectors corresponding to a multiple or cluster of singular values is
 730 computed.

731 For square matrices the perturbation bounds for symmetric/Hermitian matrices
 732 can be adapted without too much difficulty. For $A \in M_{m,n}$ with $m > n$ this is still

733 true for the right singular vectors. However, the left singular vectors to the $m - n$
 734 extra zero singular values need some special attention. If 0 or a numerical zero is a
 735 singular value, the singular vectors cannot be distinguished from those of the extra
 736 $m - n$ trivial singular values. They have to be clustered in order to obtain a basis for
 737 a singular subspace.

TABLE 9
Sensitivity of the singular value decomposition for different condition numbers

cond(A)	10^2	10^5	10^8	10^{11}	10^{14}
sensitivity(U)	$9.7 \cdot 10^{-14}$	$1.5 \cdot 10^{-11}$	$5.9 \cdot 10^{-9}$	$3.2 \cdot 10^{-6}$	$2.0 \cdot 10^{-3}$
sensitivity(Σ)	$7.5 \cdot 10^{-16}$	$7.2 \cdot 10^{-16}$	$6.5 \cdot 10^{-16}$	$6.0 \cdot 10^{-16}$	$3.6 \cdot 10^{-16}$
sensitivity(V)	$9.7 \cdot 10^{-14}$	$1.5 \cdot 10^{-11}$	$5.9 \cdot 10^{-9}$	$3.2 \cdot 10^{-6}$	$2.0 \cdot 10^{-3}$

738 As before we verify the rule of thumb (2.1) for the sensitivity of the singular values
 739 and -vectors, the results are displayed in Table 9. In the 1000 test cases we generated
 740 matrices with separated singular values because otherwise the problem to compute
 741 singular vectors becomes ill-posed. Again, the orthogonal/unitary factors become
 742 more and more sensitive for increasing condition number, whereas the singular values
 743 seem insensitive, even for large condition numbers. So extra attention seems necessary
 744 for the singular vectors.

745 To our knowledge [34] is the only paper for computing verified bounds for the
 746 complete singular value decomposition of $A \in M_{m,n}$ in $\mathcal{O}(Pp^2)$ operations for $P :=$
 747 $\max(m, n)$ and $p := \min(m, n)$. Detailed computational results can be found in [34].
 748 The quality of the bounds is often close to maximal accuracy, and even for large
 749 clusters still some 8 decimal figures can be guaranteed.

750 Bounds for the factors of the polar decomposition $A = QP$ with unitary Q and
 751 positive semidefinite P follow by $Q = UV^*$ and $P = V\Sigma V^*$ using the singular value
 752 decomposition $A = U\Sigma V^*$.

753 **8. Schur decomposition.** Let $A = XJX^{-1}$ denote a Jordan decomposition of
 754 A , and let $X = QR$ be the QR decomposition of X . Then

755 (8.1)
$$A = QTQ^*, \quad T := RJR^{-1}$$

756 is a Schur decomposition because J and R are upper triangular. The eigenvalues in
 757 T are sorted according to the diagonal of J .

758 The real Schur decomposition $A = Q'UQ'^T$ for orthogonal Q' and block upper
 759 triangular U becomes ill-posed for double eigenvalues. Consider

760
$$A := \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix}$$

761 with double real eigenvalue 1. An arbitrary small perturbation of A_{12} produces two
 762 simple real or a pair of complex eigenvalues, thus changing the block size of the factor
 763 U of the real Schur decomposition. Hence we restrict our attention to the complex
 764 Schur decomposition $A = QTQ^*$.

765 However, for the symmetric parameterized matrix in (6.1), which is normal,
 766 the Schur decomposition is the eigendecomposition with discontinuous eigenvectors.
 767 Hence, certified bounds for the Schur decomposition are restricted to matrices with
 768 simple eigenvalues – otherwise facing the ill-posed Jordan decomposition.

769 For diagonalizable A the algorithms discussed in Section 6 yield inclusions of an
 770 eigendecomposition $A = XDX^{-1}$. Combining this with the algorithm in Section 5 for
 771 an inclusion of a QR decomposition of X yields inclusions for a Schur decomposition
 772 according to (8.1). However, only inclusions \mathbf{X} and \mathbf{D} of X and D are available,
 773 so verified error bounds for the QR decomposition of \mathbf{X} are to be computed, which
 774 include those of the true X . The factor T is equal to the solution of the linear
 775 system $TR = RD$. We have to replace R and D by their computed inclusions,
 776 introducing an additional source of overestimation. That is also the reason why only
 777 for $\text{cond}(A) \lesssim 10^{14}$ verified inclusions of the Schur factors can be calculated. That
 778 does not apply to the other decompositions, including the eigendecomposition needed
 779 here.

780 The Schur decomposition offers the possibility to compute an inclusion of the
 781 departure from normality of A . To that end, only the inclusion of $T = RDR^{-1}$ is
 782 needed.

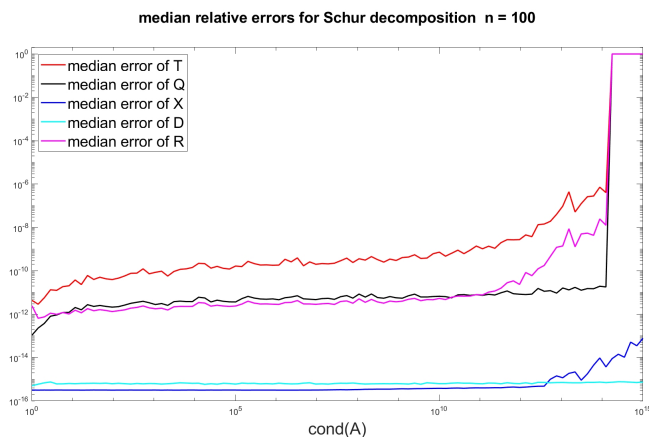


FIG. 11. Error bounds for the Schur decomposition for different condition numbers

783 We show some computational results in Figure 11, where the median relative errors
 784 of all components of T is the red line and those of Q the black line. As has been
 785 mentioned, the accuracy of the results suffers severely from the fact that only an
 786 inclusion \mathbf{X} of X is available.

787 To see that effect we also show the median relative errors of X and D in blue and
 788 cyan, respectively. As can be seen the eigenvalues are enclosed with almost maximal
 789 accuracy, the eigenvectors for condition numbers up to 10^{12} , beyond condition number
 790 10^{12} the quality of the eigenvector inclusions decreases slightly. The errors in R , shown
 791 in magenta, are close to those of Q for condition numbers up to 10^{12} .

792 **9. Takagi decomposition.** We close this note by an application of the inclu-
 793 sion of the factors of a symmetric eigendecomposition. A complex symmetric matrix
 794 $A \in M_n(\mathbb{C})$ with $A^T = A$ allows for a Takagi factorization $A = U\Sigma U^T$, also called
 795 Autonne-Takagi or symmetric singular value decomposition [14], with unitary U and
 796 diagonal Σ with non-negative diagonal elements. The factor Σ comprises of the sin-
 797 gular values of A and is unique if the diagonal elements are in nonincreasing order.
 798 The factor U may be replaced by US for diagonal S with $S^2 = I$.

799 Although less known, the Takagi factorization is used in several applications in
 800 physics and chemistry, including for example the diagonalization of mass matrices of

801 Majorana fermions, quadratic fermionic Hamiltonians, the Bloch-Messiah reduction,
 802 cf. [6, 39] and the literature cited over there.

803 It is well known that the factor U is not continuous for singular A . Consider

804
$$A := \begin{pmatrix} 1 & 0 \\ 0 & e \end{pmatrix}$$

805 with

806
$$U := \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \quad \text{and} \quad \Sigma := \begin{pmatrix} 1 & 0 \\ 0 & e \end{pmatrix} \quad \text{if } e > 0,$$

807 and

808
$$U := \begin{pmatrix} 1 & 0 \\ 0 & \sqrt{-1} \end{pmatrix} \quad \text{and} \quad \Sigma := \begin{pmatrix} 1 & 0 \\ 0 & -e \end{pmatrix} \quad \text{if } e < 0.$$

809 The discontinuity is forced by the non-negativity of Σ . Hence, as for the QR de-
 810 composition, the computation of the Takagi factorization is ill-posed at $e = 0$. As a
 811 consequence, a verification method is only applicable to matrices with full rank.

812 First, we verify the rule of thumb (2.1) for the sensitivity of the Takagi factors,
 813 the results are displayed in Table 10. Here we use $\mathbf{Q} = \text{orth}(\text{randn}(n))$ to generate
 814 a random complex symmetric matrices of size $n = 100$ with $\text{cond}(A) \approx 10^k$ by

815 $\mathbf{D} = \text{diag}(\text{logspace}(0, k, n)); \mathbf{A} = \mathbf{Q}.' * \mathbf{D} * \mathbf{Q}; \mathbf{A} = \mathbf{A} + \mathbf{A}.'$;
 where the last statement symmetrizes the matrix taking care of rounding errors.

TABLE 10
Sensitivity of the Takagi decomposition for different condition numbers

cond(A)	10 ²	10 ⁵	10 ⁸	10 ¹¹	10 ¹⁴
sensitivity(U)	2.3 · 10 ⁻¹³	2.9 · 10 ⁻¹²	9.7 · 10 ⁻¹⁰	5.3 · 10 ⁻⁷	3.3 · 10 ⁻⁴
sensitivity(Σ)	4.4 · 10 ⁻¹⁵	3.3 · 10 ⁻¹⁵	2.4 · 10 ⁻¹⁵	2.2 · 10 ⁻¹⁵	1.8 · 10 ⁻¹⁵

816

817 As anticipated, the factor U is sensitive to perturbations of the matrix A , while Σ is
 818 not. Of course, the insensitivity of Σ follows by well known perturbation results for
 819 singular values.

820 There are several methods known in the literature to approximate the Takagi
 821 factors. For our purposes, the computation of verified bounds, one possibility is the
 822 following [9]. Let $A^T = A$ have full rank and denote the singular value decomposition
 823 by $A = U\Sigma V^*$. Then $D := U^* A \bar{U} \Sigma^{-1}$ is diagonal, and a computation shows that
 824 $UD^{1/2}$ and Σ are the Takagi factors. This is our first method.

825 One drawback is that the computation of D involves two matrix multiplications.
 826 Despite the computational effort this is a source of overestimation because the two
 827 factors are interval matrices \mathbf{U} and \mathbf{V} including the true factors U and V , respectively.
 828 The inclusions \mathbf{U} and \mathbf{V} may be computed by the methods in the previous section.

829 Overestimation can be reduced by using $U^* A \bar{U} \Sigma^{-1} = V^* \bar{U}$ which is again diago-
 830 nal. Now only one multiplication of interval matrices is necessary, and we may expect
 831 better results. That is our second method.

832 We finally transform the problem into a real symmetric eigenproblem, see also
 833 [14, 4.4.P2]. Let nonsingular $A^T = A \in M_n(\mathbb{C})$ be given, and denote $A = E + iF$

834 with $E^T = E, F^T = F$ and $E, F \in M_n(\mathbb{R})$. A direct computation shows that the
 835 eigenvalues of the symmetric matrix

$$836 \quad M := \begin{pmatrix} E & F \\ F & -E \end{pmatrix}$$

837 come in \pm pairs. If $\begin{pmatrix} x \\ y \end{pmatrix}$ is an eigenvector to $\lambda \in \mathbb{R}$, then $\begin{pmatrix} y \\ -x \end{pmatrix}$ is an eigenvector
 838 to $-\lambda$. After suitable renumbering the eigendecomposition of M is

$$839 \quad M \begin{pmatrix} X & Y \\ Y & -X \end{pmatrix} = \begin{pmatrix} X & Y \\ Y & -X \end{pmatrix} \begin{pmatrix} \Sigma & 0 \\ 0 & -\Sigma \end{pmatrix}.$$

840 A direct computation verifies that $U := X + iY$ and $V := X - iY$ are unitary and
 841 $AV = U\Sigma$. Hence the diagonal elements of Σ are the singular values of A . Finally

$$842 \quad A\bar{U} = (E + iF)(X - iY) = AV = U\Sigma$$

843 verifies that U and Σ are the factors of the Takagi decomposition. Inclusions of X
 844 and Y are computed with the methods for symmetric eigendecomposition presented
 845 in Section 6.

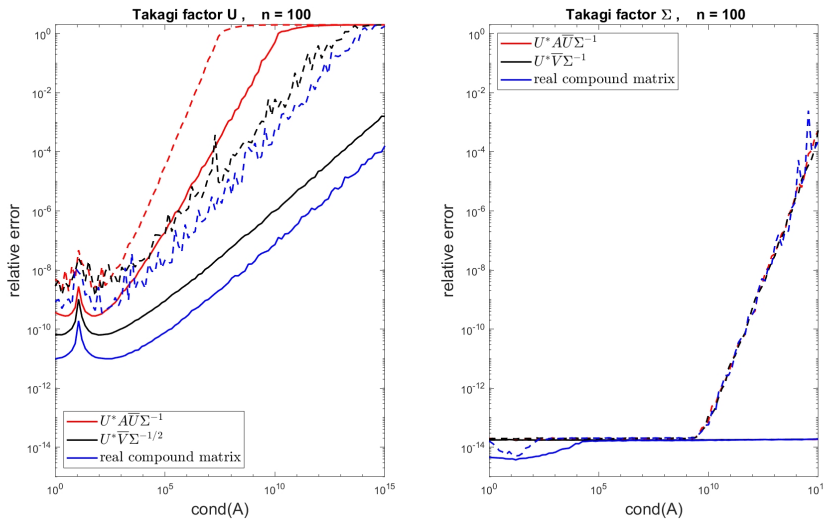


FIG. 12. The Takagi decomposition, solid line median, dashed line maximum relative error

846 We show some computational results in Figure 12. The median relative errors of U in
 847 the left and Σ in the right graph computed by the three methods are the solid lines in
 848 red, black, and blue, respectively. As for the inclusions of U the third method seems
 849 best. In any case, as expected by our rule of thumb (2.1), the relative error increases
 850 with the condition number. The reason for the small peak at $\text{cond}(A) \approx 10$ is not
 851 clear to us, it may be due to the construction of the test matrices. For Σ all three
 852 methods compute bounds of almost maximal accuracy.

853 However, the numbers are slightly misleading because, for example, most elements
 854 of U are enclosed with high accuracy, and only few corresponding to the smallest sin-
 855 gular values are weaker. Thus the median reflects mostly the relative error of the
 856 better inclusions. Therefore, we display for both U and Σ the maximum relative
 857 errors as well, the dashed lines. Figure 12 shows that for the first method and condi-
 858 tion number beyond about 10^9 some inclusions have relative error close to 1, for the
 859 second method beyond 10^{14} , where for the third method even for $\text{cond}(A) \lesssim 10^{15}$ the
 860 inclusions seem to contain some information.

861 For the singular values Σ below condition number 10^9 all bounds are of maximal
 862 accuracy, with some deterioration for larger condition numbers. The results of all
 863 three methods are of similar quality.

864 **Acknowledgment.** The authors wish to express their heartfelt thanks to the
 865 two anonymous referees for their thorough reading and most constructive and helpful
 866 comments.

867

REFERENCES

- 868 [1] IEEE Standard for Floating-point Arithmetic. *IEEE Std 754-2019 (Revision of IEEE 754-*
 869 *2008)*, pages 1–84, 2019.
- 870 [2] Advanpix: Multiprecision Computing Toolbox for MATLAB, 2024. Code and documentation
 871 available at <http://www.advanpix.com/>.
- 872 [3] P. Ahrens, J. Demmel, H. D. Nguyen: Algorithms for efficient reproducible floating-point
 873 summation. *ACM Trans. Math. Software*, 46: 1–49, 2020.
- 874 [4] G. Alefeld, H. Spreuer: Iterative improvement of componentwise error bounds for invariant
 875 subspaces belonging to a double or nearly double eigenvalue. *Computing*, 36: 321–334,
 876 1986.
- 877 [5] D. H. Bailey: A Fortran-90 double-double precision library.
 878 <https://www.davidhbailey.com/dhbsoftware/>.
- 879 [6] G. Cariolaro and G. Pierobon. Bloch-messiah reduction of gaussian unitaries by Takagi factor-
 880 ization. *Phys. Rev. A*, 94: Article 062109, 2016.
- 881 [7] T. J. Dekker: A floating-point technique for extending the available precision. *Numer. Math.*,
 882 18: 224–242, 1971.
- 883 [8] L. Dieci, T. Eirola: On smooth decompositions of matrices. *SIAM J. Matrix Anal. Appl.*,
 884 20(3): 800–819, 1999.
- 885 [9] L. Dieci, A. Papini, A. Pugliese: Takagi factorization of matrices depending on parameters and
 886 locating degeneracies of singular values. *SIAM J. Matrix Anal. Appl.*, 43(3): 1148–1161,
 887 2022.
- 888 [10] F. R. Gantmacher: *The Theory of Matrices*, volume 1, Chelsea, New York, 1959.
- 889 [11] N. J. Higham: *Accuracy and Stability of Numerical Algorithms*, SIAM Publications, Philadel-
 890 phia, 2nd edition, 2002.
- 891 [12] N. J. Higham: private communication, 2023
- 892 [13] F. R. de Hoog, R. S. Anderssen, M. A. Lukas: Differentiation of matrix functionals using
 893 triangular factorization. *Math. Comp.*, 80(275): 1585–1600, 2011.
- 894 [14] R. A. Horn, C. R. Johnson: *Matrix Analysis*, second edition. Cambridge University Press,
 895 2013.
- 896 [15] R. B. Kearfott, M. T. Nakao, A. Neumaier, S. M. Rump, S. P. Shary, P. van Hentenryck:
 897 Standardized notation in interval analysis, *Computational Technologies*, 15(1): 7–13, 2010.
- 898 [16] D. E. Knuth: *The Art of Computer Programming: Seminumerical Algorithms*, volume 2.
 899 Addison Wesley, Reading, Massachusetts, 1969.
- 900 [17] R. Kobayashi, M. Lange, A. Minamihata, S. M. Rump: Verified inclusion of a basis of the null
 901 space. *Reliable Computing*, 27: 26–41, 2020.
- 902 [18] R. Krawczyk: Fehlerabschätzung reeller Eigenwerte und Eigenvektoren von Matrizen. *Com-*
 903 *puting*, 4: 281–293, 1969.
- 904 [19] M. Lange, S. M. Rump: Faithfully rounded floating-point computations. *ACM Trans. Math.*
 905 *Softw.*, 46(3): Article 21, 2020.
- 906 [20] J. de Leeuw: Differentiating the QR decomposition, Preprint, 2023 (DOI: 10.13140/RG.2.2.
 907 19201.12640).

- 908 [21] M. Malcolm: On accurate floating-point summation. *Comm. ACM*, 14(11): 731–736, 1971.
- 909 [22] MATLAB. User’s Guide, Version 2023b, the MathWorks Inc., 2023.
- 910 [23] O. Möller: Quasi double precision in floating-point arithmetic. *BIT Numerical Mathematics*,
- 911 5: 37–50, 1965.
- 912 [24] S. Miyajima: Fast enclosure for all eigenvalues and invariant subspaces in generalized eigenvalue
- 913 problems. *SIAM J. Matrix Anal. Appl.*, 35(3): 1205–1225, 2014.
- 914 [25] S. Miyajima, T. Ogita, S. M. Rump, S. Oishi: Fast verification of all eigenpairs in symmetric
- 915 positive definite generalized eigenvalue problem. *Reliable Computing*, 14: 24–45, 2010.
- 916 [26] A. Neumaier: *Interval methods for systems of equations*. Encyclopedia of Mathematics and its
- 917 Applications. Cambridge University Press, 1990.
- 918 [27] T. Ogita, S. M. Rump, S. Oishi: Accurate sum and dot product. *SIAM J. Sci. Comput.*, 26(6):
- 919 1955–1988, 2005.
- 920 [28] S. Oishi, K. Ichihara, M. Kashiwagi, T. Kimura, X. Liu, H. Masai, Y. Morikura, T. Ogita,
- 921 K. Ozaki, S. M. Rump, K. Sekine, A. Takayasu, N. Yamanaka: *Principle of Verified*
- 922 *Numerical Computations*. Corona Publisher, Tokyo, Japan, 2018. [in Japanese].
- 923 [29] K. Ozaki, T. Ogita, S. Oishi: Error-free transformation of matrix multiplication with a poste-
- 924 riori validation. *Numer. Linear Alg. Appl.*, 23(5): 931–946, 2016.
- 925 [30] F. Rellich: Störungstheorie der Spektralzerlegung. *Math. Ann.*, 113: 79–91, 1994.
- 926 [31] S. M. Rump: INTLAB – INTerval LABORatory. In Tibor Csendes, editor, *Developments in*
- 927 *Reliable Computing*, pages 77–104. Springer Netherlands, Dordrecht, 1999.
- 928 [32] S. M. Rump: Verification methods: Rigorous results using floating-point arithmetic. *Acta*
- 929 *Numerica*, 19: 287–449, 2010.
- 930 [33] S. M. Rump: Verified error bounds for all eigenvalues and eigenvectors of a matrix. *SIAM J.*
- 931 *Matrix Anal. Appl.*, 43(4): 1736–1754, 2022.
- 932 [34] S. M. Rump, M. Lange: Fast computation of error bounds for all eigenpairs of a Hermitian
- 933 and all singular pairs of a rectangular matrix with emphasis on eigen- and singular value
- 934 clusters. *J. Comput. Appl. Math.*, 434: Article 115332, 2023.
- 935 [35] S. M. Rump, T. Ogita: On a quality measure for interval inclusions. *BIT Numerical Mathe-*
- 936 *matics*, 64: Article 22, 2024.
- 937 [36] S. M. Rump, T. Ogita, S. Oishi: Accurate floating-point summation part I: Faithful rounding.
- 938 *SIAM J. Sci. Comput.*, 31(1): 189–224, 2008.
- 939 [37] G. W. Stewart: On the perturbation of LU, Cholesky, and QR factorizations. *SIAM J. Matrix*
- 940 *Anal. Appl.*, 14: 1141–1146, 1993.
- 941 [38] J. Sun: Multiple eigenvalue sensitivity analysis. *Linear Alg. Appl.*, 137/138: 183–211, 1990.
- 942 [39] A. E. Teretenkov. Singular value decomposition for skew-Takagi factorization with quantum
- 943 applications. *Linear and Multilinear Algebra*, 70(22): 7762–7769, 2022.
- 944 [40] A. N. Tikhonov: Regularization of incorrectly posed problems. *Dokl. Akad. Nauk SSSR*, 153(1):
- 945 49–52, 1963.
- 946 [41] A. N. Tikhonov, V. L. Arsenin: *Solutions of Ill-posed Problems*. John Wiley & Sons, New
- 947 York, 1977.
- 948 [42] J. M. Wolfe: Reducing truncation errors by programming. *Comm. ACM*, 7(6): 355–356, 1964.
- 949 [43] G. Zielke, V. Drygalla: Genaue Lösung linearer Gleichungssysteme. *GAMM Mitt. Ges. Angew.*
- 950 *Math. Mech.*, 26: 7–108, 2003.