

Floating-point arithmetic and error analysis (AFAE)

Fast verification methods for linear systems – Part I

Stef Graillat

LIP6/PEQUAN – Sorbonne University

Lecture Master 2 SFPN – MAIN5



Presentation of the problem

- Let $A \in \mathbb{F}^{n \times n}$ be a nonsingular matrix, and $b \in \mathbb{F}^n$ be a vector. We will be interested in a validated solution of a linear system $Ax = b$:
 - computation of an approximate solution $\hat{x} \in \mathbb{F}^n$;
 - computation of a certified bound $\delta \in \mathbb{F}$ such that $\|\hat{x} - x\|_{\infty} \leq \delta$.
- In this lecture, we will present some methods proposed by Oishi, Ogita and Rump in [OR02, OR05]:
 - 1 What is the cost for obtaining such a verification?
 - 2 How evolves the quality of the bound with respect to the cost of the methods?

Classical iterations

Numerous verification algorithms for linear systems are based on iterations with interval arithmetic.

Theorem 1 (Rump 1980, Theorem 10.6 in [Rum10])

Let $A, R \in \mathbb{R}^{n \times n}$, $b \in \mathbb{R}^n$ and $x \in \mathbb{IR}^n$. If $Rb + (I - RA)x \subset \text{int}(x)$, then A and R are nonsingular, and $x := A^{-1}b \in Rb + (I - RA)x$.

This theorem can be used through an iteration of the type

$$x^{(i)} = Rb + (I - RA)x^{(i-1)},$$

by taking $R \approx A^{-1}$ computed with finite precision.

Cost of this method: $6n^3 + O(n^2)$ floating-point operations (flops), i.e. **9 times the cost of solving the system without certification** by Gaussian Elimination (GE).

The different approaches presented here will make it possible to lower the cost of the verification by obtaining a cost less than $6n^3$, allowing a loss in the tightness of the bound.

Assumptions, notations

- Here we assume no underflow occurs. Nevertheless it is possible to deal with them (see [OR02, OR05]).
- u is the unit rounding for rounding to nearest,
- The notation $\text{fl}(\cdot)$ is used to indicate that all the computations within the parenthesis are performed in floating-point arithmetic with rounding to nearest.
- In a similar way, $\text{fl}_{\downarrow}(\cdot)$ et $\text{fl}_{\uparrow}(\cdot)$ means that we use floating-point arithmetic respectively with rounding toward $-\infty$ and with rounding toward ∞ .

Outline

1 Preliminary results

- Triangular linear systems
- Inverse of triangular matrices
- LU factorization and solution of linear systems
- Normwise conditioning of linear systems

2 Verified solution of linear systems

- General principle
- Four methods to upper bound α
- Upper bound for β
- Synthesis

3 Numerical experiments

Outline

1

Preliminary results

- **Triangular linear systems**
- Inverse of triangular matrices
- LU factorization and solution of linear systems
- Normwise conditioning of linear systems

2

Verified solution of linear systems

- General principle
- Four methods to upper bound α
- Upper bound for β
- Synthesis

3

Numerical experiments

Triangular linear systems

Let $T \in \mathbb{F}^{n \times n}$ be a nonsingular triangular matrix and $b \in \mathbb{F}^n$. Without loss of generality, we assume that T is lower triangular. The linear system $Tx = b$ can be solved by forward substitution:

For $i = 1 \dots n$,

$$x_i = \frac{1}{T_{i,i}} \left(b_i - \sum_{j=1}^{i-1} T_{i,j} x_j \right).$$

Cost: n^2 flops.

Algorithm 1 (Substitution)

```
function x = Substitution(T,b)
  for i = 1 : n
    s = b_i
    for j = 1 : i - 1
      s = s - T_{i,j}x_j
    end
    x_i = s/T_{i,i}
  end
```

Assuming the algorithm is performed in floating-point arithmetic: **find the backward error for the computed solution \hat{x} .**

Backward stability

We have the following lemma:

Lemma 1 (Thm 8.5, p. 142 in [Hig02])

If $y = \frac{1}{b_k} (c - \sum_{i=1}^{k-1} a_i b_i)$ is evaluated in floating-point arithmetic with rounding to nearest, then, no matter what the order of evaluation,

$$(1 + \delta_k) b_k \hat{y} = c - \sum_{i=1}^{k-1} (1 + \delta_i) a_i b_k,$$

where $|\delta_i| \leq \gamma_k$, for $1 \leq i \leq k$.

In the substitution algorithm, we evaluate, for $i = 1 \dots n$,

$$x_i = \frac{1}{T_{i,i}} \left(b_i - \sum_{j=1}^{i-1} T_{i,j} x_j \right).$$

Backward stability

So the components $\hat{\mathbf{x}}_i$ of the computed solution satisfy

$$(1 + \delta_{i,i})T_{i,i} \hat{\mathbf{x}}_i = \mathbf{b}_i - \sum_{j=1}^{i-1} (1 + \delta_{i,j})T_{i,j} \hat{\mathbf{x}}_j,$$

with $|\delta_{i,j}| \leq \gamma_i \leq \gamma_n$. We deduce the following result:

Theorem 2 (p. 142, [Hig02])

Let the triangular system $T\mathbf{x} = \mathbf{b}$ where $T \in \mathbb{F}^{n \times n}$ is nonsingular, be solved by substitution with any ordering. Then the computed solution $\hat{\mathbf{x}}$ satisfies

$$(T + \Delta T) \hat{\mathbf{x}} = \mathbf{b}, \quad \text{with} \quad |\Delta T| \leq \gamma_n |T|.$$

Remark: this is true for lower and upper triangular matrix T .

Outline

1 Preliminary results

- Triangular linear systems
- **Inverse of triangular matrices**
- LU factorization and solution of linear systems
- Normwise conditioning of linear systems

2 Verified solution of linear systems

- General principle
- Four methods to upper bound α
- Upper bound for β
- Synthesis

3 Numerical experiments

Inversion of triangular matrices

Let $T \in \mathbb{F}^{n \times n}$ be a nonsingular triangular matrix, and \hat{X} approximate inverse of T . Quantity needed to verify a solution of a linear system: **an a priori bound on $|\hat{X}T - I|$** .

There exist many methods in $\frac{1}{3}n^3 + O(n^2)$ operations to compute an approximate inverse of T [Hig02, chap. 14]. These methods are based

- either on the equation $T\hat{X} = I$,
- or on the equation $\hat{X}T = I$.

Depending on the case, we will be able to bound

- either $|T\hat{X} - I|$, and we talk about “**right residual**”,
- or $|\hat{X}T - I|$, and we talk about “**left residual**”.

Inversion of triangular matrices

Let us denote $X = [X_1, \dots, X_n]$, with $X_j \in \mathbb{R}^n$.

- If \hat{X} is computed as an approximate solution of $TX = I$:
for each triangular system $TX_j = e_j$, we compute \hat{X}_j such
that $(T + \Delta T_j)\hat{X}_j = e_j$ with $|\Delta T_j| \leq \gamma_n |T|$, so
 $|T\hat{X}_j - e_j| \leq \gamma_n |T| |\hat{X}_j|$, and

$$|T\hat{X} - I| \leq \gamma_n |T| |\hat{X}|.$$

We obtain an upper bound for the “right residual”

- If \hat{X} is computed as an approximate solution of $\hat{X}T = I$:
by solving $T^T \hat{X}^T = I$, we have $|T^T \hat{X}^T - I| \leq \gamma_n |T^T| |\hat{X}^T|$, so

$$|\hat{X}T - I| \leq \gamma_n |\hat{X}| |T|.$$

We obtain an upper bound for the “left residual”

Outline

1 Preliminary results

- Triangular linear systems
- Inverse of triangular matrices
- **LU factorization and solution of linear systems**
- Normwise conditioning of linear systems

2 Verified solution of linear systems

- General principle
- Four methods to upper bound α
- Upper bound for β
- Synthesis

3 Numerical experiments

LU factorization

We say that $A \in \mathbb{R}^{n \times n}$ admits a LU factorization of $A = LU$ where

- $L \in \mathbb{R}^{n \times n}$ is a lower triangular matrix with all diagonal entries equal to 1;
- $U \in \mathbb{R}^{n \times n}$ is an upper triangular matrix.

A matrix A admits a unique LU factorization iff $A(1 : k, 1 : k)$ is nonsingular for $k = 1 : n - 1$ ([Hig02, Thm 9.1, p. 161]).

The factorization can be computed by Gaussian elimination (GE) in $\frac{2}{3}n^3$ operations.

We show here how to obtain a **a result on backward analysis for factors \hat{L} and \hat{U} computed by GE in floating-point arithmetic.**

Initialize U to A

for $k = 1 : n$

$L_{k,k} \leftarrow 1$

for $i = k + 1 : n$

$L_{i,k} = U_{i,k}/U_{k,k}; U_{i,k} \leftarrow 0$

for $j = k + 1 : n$

$U_{i,j} \leftarrow U_{i,j} - L_{i,k}U_{k,j}$

LU factorization

Suppose we have already calculated $U(1 : k - 1, :)$ and $L(:, 1 : k - 1)$.

We can determine $U(k, :)$ and $L(:, k)$ by:

$$\begin{aligned} \text{for } j = k : n, \quad U_{k,j} &= A_{k,j} - \sum_{i=1}^{k-1} L_{k,i} U_{i,j}, \\ \text{for } i = k + 1 : n, \quad L_{i,k} &= (A_{i,k} - \sum_{j=1}^{k-1} L_{i,j} U_{j,k}) / U_{k,k}. \end{aligned}$$

All variants of GE calculate L and U according to the two previous equations, even if the order of the operations is different from one variant to another.

In particular, we obtain the Doolittle algorithm:

```
for k = 1 : n
  for j = k : n
     $U_{k,j} = A_{k,j} - \sum_{i=1}^{k-1} L_{k,i} U_{i,j}$ 
  for i = k + 1 : n
     $L_{i,k} = (A_{i,k} - \sum_{j=1}^{k-1} L_{i,j} U_{j,k}) / U_{k,k}$ 
```

LU factorization

For the exact factors L and U , we have:

$$U_{k,j} = A_{k,j} - \sum_{i=1}^{k-1} L_{k,i} U_{i,j}, \quad j \geq k,$$

$$L_{i,k} = (A_{i,k} - \sum_{j=1}^{k-1} L_{i,j} U_{j,k}) / U_{k,k} \quad i > k.$$

If we consider the factors \hat{L} and \hat{U} calculated in floating-point arithmetic,

$$(1 + \delta_k) \hat{U}_{k,j} = A_{k,j} - \sum_{i=1}^{k-1} (1 + \delta_i) \hat{L}_{k,i} \hat{U}_{i,j}, \quad j \geq k,$$

$$(1 + \delta'_k) \hat{L}_{i,k} \hat{U}_{k,k} = A_{i,k} - \sum_{j=1}^{k-1} (1 + \delta'_j) \hat{L}_{i,j} \hat{U}_{j,k}, \quad i > k,$$

with $|\delta_i|, |\delta'_i| \leq \gamma_k$. So

$$|A_{k,j} - \sum_{i=1}^k \hat{L}_{k,i} \hat{U}_{i,j}| \leq \gamma_k \sum_{i=1}^k |\hat{L}_{k,i}| |\hat{U}_{i,j}|, \quad j \geq k,$$

$$|A_{i,k} - \sum_{j=1}^k \hat{L}_{i,j} \hat{U}_{j,k}| \leq \gamma_k \sum_{j=1}^k |\hat{L}_{i,j}| |\hat{U}_{j,k}|, \quad i > k.$$

Recognizing the products $\hat{L}\hat{U}$ and $|\hat{L}||\hat{U}|$, and as $\gamma_k \leq \gamma_n$, we have:

$$|A - \hat{L}\hat{U}| \leq \gamma_n |\hat{L}||\hat{U}|.$$

Backward stability

Theorem 3 (Thm 9.3, p. 164 in [Hig02])

If GE applied to $A \in \mathbb{F}^{n \times n}$ runs to completion then the computed LU factors $\hat{L}, \hat{U} \in \mathbb{F}^{n \times n}$ satisfy

$$\hat{L}\hat{U} = A + \Delta A, \quad \text{with} \quad |\Delta A| \leq \gamma_n |\hat{L}| |\hat{U}|.$$

This result of backward analysis is valid for all variants of GE performing the same operations as the Doolittle algorithm.

We use Gaussian Elimination with Partial Pivoting (GEPP) :

- to be able to treat the case where one meets a null pivot ;
- to improve the numerical stability of the elimination.

This is the case LAPACK function xGETRF for LU factorization and for functions xGESV and xGESVX for solving linear systems. What does the previous result give for GEPP?

LUP factorization

GEPP

Initialize U to A

Initialize P to I

for $k = 1 : n$

Let i such that $|U_{i,k}|$ be the greatest entry in $|U(k : n, k)|$

Swap lines k and i in L , U in P

$L_{k,k} \leftarrow 1$

for $i = k + 1 : n$

$L_{i,k} = U_{i,k}/U_{k,k}; U_{i,k} \leftarrow 0$

for $j = k + 1 : n$

$U_{i,j} \leftarrow U_{i,j} - L_{i,k}U_{k,j}$

We obtain the LU factorization of the matrix PA . In this case, we get:

$$|PA - \widehat{L}\widehat{U}| \leq \gamma_n |\widehat{L}| |\widehat{U}|.$$

Solution of a linear system

Theorem 4

Let $A \in \mathbb{F}^{n \times n}$ be a nonsingular matrix, and $b \in \mathbb{F}^n$. An approximate solution of $Ax = b$ is computed in floating-point arithmetic with precision u as follows:

- let \hat{L} and \hat{U} be the LU factors of PA computed via GEPP;
- let \hat{y} be an approximate solution of $\hat{L}y = Pb$ computed by substitution;
- let \hat{x} be an approximate solution of $\hat{U}x = \hat{y}$ computed by substitution.

We define $\rho(A) := \frac{\|\hat{U}\|_M}{\|A\|_M}$ the **growth factor** associated with matrix A for GEPP. Then

$$(A + \Delta A)\hat{x} = b, \quad \text{with} \quad \|\Delta A\|_\infty \leq n^2 \gamma_{3n} \rho(A) \|A\|_\infty.$$

Remark: We use the same definition of $\rho(A)$ as [Dem97, p. 49]. There exist other definitions: see [Gv96, p.116] and [Hig02,

What about the backward stability of GEPP?

$$(\mathbf{A} + \Delta\mathbf{A})\hat{\mathbf{x}} = \mathbf{b}, \quad \text{with} \quad \|\Delta\mathbf{A}\|_{\infty} \leq n^2 \gamma_{3n} \rho(\mathbf{A}) \|\mathbf{A}\|_{\infty}.$$

We can not guarantee a priori $n^2 \gamma_{3n} \rho(\mathbf{A})$ to be “small” for any matrix \mathbf{A} .

The stability of GEPP depends on $\rho(\mathbf{A})$: the only bound independent of \mathbf{A} is $\rho(\mathbf{A}) \leq 2^{n-1}$, it can be reached for all n (see[Dem97, p. 49]).

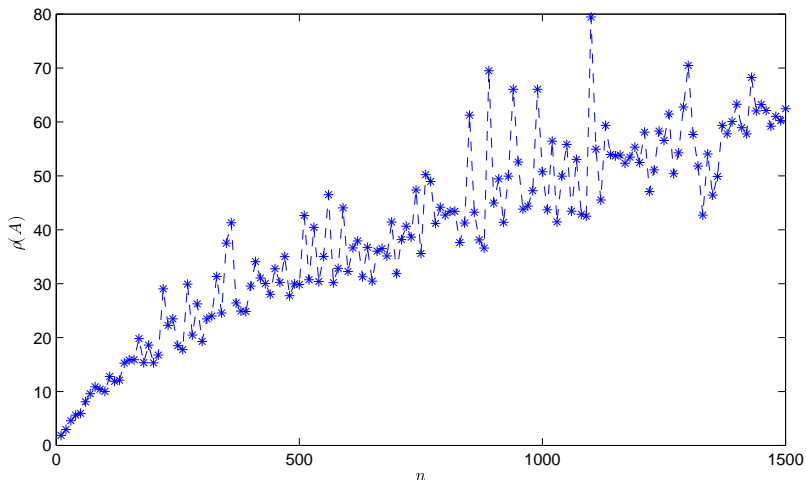
The use of GEPP is generally justified as follows:

- the growth factors encountered in practice are usually “small”,
- the bound of the backward error obtained is often very rough,
- it is possible to improve the backward stability by iterative refinement.

We can sometimes read that GEPP is “inverse-stable in practice”.

What about the backward stability of GEPP?

An example, with matrices whose coefficients are uniformly distributed over $[-1, 1]$:



Outline

1 Preliminary results

- Triangular linear systems
- Inverse of triangular matrices
- LU factorization and solution of linear systems
- **Normwise conditioning of linear systems**

2 Verified solution of linear systems

- General principle
- Four methods to upper bound α
- Upper bound for β
- Synthesis

3 Numerical experiments

Normwise conditioning of linear systems

Let $A \in \mathbb{F}^{n \times n}$, $b \in \mathbb{F}^n$, and $\hat{x} \in \mathbb{F}^n$ be an approximate solution of $Ax = b$ computed with GEPP. Then \hat{x} satisfies a backward analysis relation of the form:

$$(A + \Delta A)\hat{x} = b, \quad \text{with} \quad \|\Delta A\|_{\infty} \leq \varepsilon \|A\|_{\infty}.$$

In the lecture, we will use the condition numbers defined by:

$$\lim_{\varepsilon \rightarrow 0} \sup_{\|\Delta A\|_{\infty} \leq \varepsilon \|A\|_{\infty}} \left\{ \frac{\|\hat{x} - x\|_{\infty}}{\|x\|_{\infty}}, (A + \Delta A)\hat{x} = b \right\}$$

This is the definition of the relative conditioning of $F_b : X \mapsto X^{-1}b$ in A , using $\|\cdot\|_{\infty}$ for input and output. Following Rice's theorem (see [Hig02, p. 121]), one can show that

$$\kappa_{\infty}(A) := \|A^{-1}\|_{\infty} \|A\|_{\infty}.$$

Normwise conditioning of linear systems

We assume the following result:

[Hig02, p. 121]

Let $A, A + \Delta A \in \mathbb{R}^{n \times n}$ be nonsingular matrices, and $b \in \mathbb{R}^n$. Let us denote $x = A^{-1}b$ and $\hat{x} = (A + \Delta A)^{-1}b$. If $\|\Delta A\|_\infty \leq \varepsilon \|A\|_\infty$ and $\varepsilon \kappa(A) < 1$, then

$$\frac{\|x - \hat{x}\|_\infty}{\|x\|_\infty} \leq \frac{\varepsilon \kappa(A)}{1 - \varepsilon \kappa(A)} = \varepsilon \kappa(A) + O(\varepsilon^2).$$

If $\hat{x} \in \mathbb{F}^n$ is an approximate solution computed by GEPP of a linear system $Ax = b$ with floating-point numbers, we must be able to observe that the relative error

$$\frac{\|x - \hat{x}\|_\infty}{\|x\|_\infty}$$

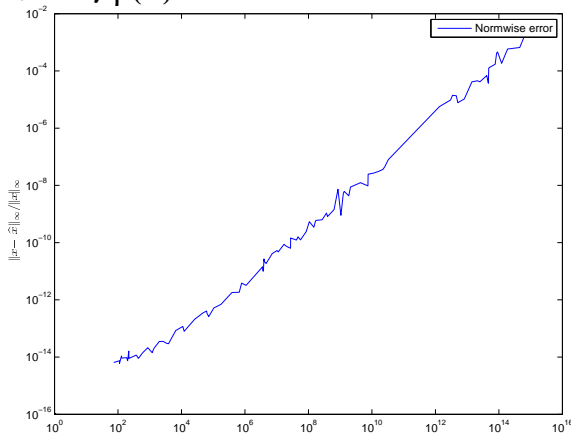
is in the worst case of the order of $\kappa(A)$.

CNormwise conditioning of linear systems

Generator of ill-conditioned matrices is available in MATLAB:

```
gallery('randsvd', 100, 10^(15*rand));
```

In this simulation, $\rho(A) \leq 27$.



Outline

1 Preliminary results

- Triangular linear systems
- Inverse of triangular matrices
- LU factorization and solution of linear systems
- Normwise conditioning of linear systems

2 Verified solution of linear systems

- General principle
- Four methods to upper bound α
- Upper bound for β
- Synthesis

3 Numerical experiments

Outline

1

Preliminary results

- Triangular linear systems
- Inverse of triangular matrices
- LU factorization and solution of linear systems
- Normwise conditioning of linear systems

2

Verified solution of linear systems

- **General principle**
- Four methods to upper bound α
- Upper bound for β
- Synthesis

3

Numerical experiments

Matrix norms

The norm on $\mathbb{R}^{n \times n}$ subordinated to the p -norm $\|\cdot\|_p$ on \mathbb{R}^n is defined by:

$$\|A\|_p := \max_{x \neq 0} \frac{\|Ax\|_p}{\|x\|_p} = \max_{\|x\|_p=1} \|Ax\|_p.$$

By definition, for all $A \in \mathbb{R}^{n \times n}$ and for all $x \in \mathbb{R}^n$,

$$\|Ax\|_p \leq \|A\|_p \|x\|_p.$$

Moreover, for all $A, B \in \mathbb{R}^{n \times n}$, $\|AB\|_p \leq \|A\|_p \|B\|_p$.

Indeed, it exists $x^* \in \mathbb{R}^n$ with $\|x^*\|_p = 1$ such that

$$\|AB\|_p = \|ABx^*\|_p, \text{ so}$$

$$\|AB\|_p \leq \|A\|_p \|Bx^*\|_p \leq \|A\|_p \|B\|_p \|x^*\|_p = \|A\|_p \|B\|_p.$$

Let us recall that $\lambda \in \mathbb{C}$ is an eigenvalue of A , then

$$|\lambda| \leq \|A\|_p.$$

Indeed, if $Ax = \lambda x$ for $x \in \mathbb{R}^n - \{0\}$, then

$$\|A\|_p \|x\|_p \geq \|Ax\|_p = \|\lambda x\|_p = |\lambda| \|x\|_p.$$

Matrix norms

In the cas of the ∞ -norm:

- the sub-multiplicative inequality can be attained:

$$\underbrace{\begin{bmatrix} 2 & 2 \\ 2 & 2 \end{bmatrix}}_{\|\cdot\|_\infty=4} = \underbrace{\begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}}_{\|\cdot\|_\infty=2} \times \underbrace{\begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}}_{\|\cdot\|_\infty=2}$$

- but it can also be “arbitrarily pessimistic”:

$$\underbrace{\begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}}_{\|\cdot\|_\infty=0} = \underbrace{\begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}}_{\|\cdot\|_\infty=2} \times \underbrace{\begin{bmatrix} 1 & 1 \\ -1 & -1 \end{bmatrix}}_{\|\cdot\|_\infty=2}$$

Same behavior as the triangular inequality ...

Inversion and norm

Lemma 2

Let $A, R \in \mathbb{R}^{n \times n}$. If $\|RA - I\|_p < 1$ then R and A are nonsingular.

Lemma 3 (Lemma 2.3.3, p. 58 in [Gv96])

Let $F \in \mathbb{R}^{n \times n}$. If $\|F\|_\infty < 1$ then $I - F$ is nonsingular and

$$\|(I - F)^{-1}\|_\infty \leq \frac{1}{1 - \|F\|_\infty}.$$

Proof: we prove that $I - F$ is nonsingular as in the previous lemma. Let $G = (I - F)^{-1}$. We have $(I - F)G = I$, so

$$1 = \|I\|_\infty = \|G - FG\|_\infty \geq \|G\|_\infty - \|FG\|_\infty \geq \|G\|_\infty(1 - \|F\|_\infty).$$

As $\|F\|_\infty < 1$, we obtain the result. \square

Inversion and residual

Interpretation of the sufficient condition for the nonsingularity of A :

- Assume that $R = (A + \Delta A)^{-1}$, with $|\Delta A| \leq u|A|$.
Then, $R = A^{-1}(I + \Delta A A^{-1})^{-1} = A^{-1} - A^{-1}\Delta A A^{-1} + O(u^2)$. So
 $RA - I = -A^{-1}\Delta A + O(u^2)$ and

$$\|RA - I\|_{\infty} \leq u\kappa_{\infty}(A) + O(u^2).$$

- If R is computed with LAPACK function `xGETRI`, which is the case in the `inv` function of MATLAB. In this case (see [Hig02, p. 268]),

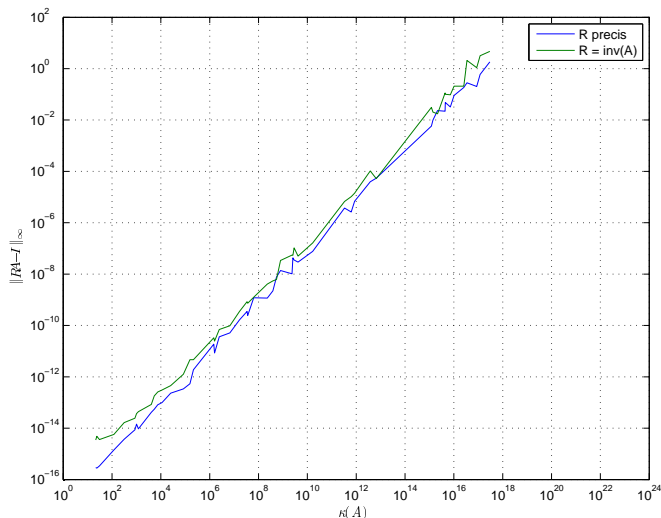
$$|RA - I| \leq p(n)u|R\|L\|U|,$$

where $p(n)$ is a polynomial in n . As a consequence,

$$\|RA - I\|_{\infty} \leq p(n)u\|R\|_{\infty}\|L\|U\|_{\infty}.$$

Conditioning and residual

Example, with 50 matrices 20×20 generated by « randsvd » :



Upper bound

Theorem 5

Let $A, R \in \mathbb{R}^{n \times n}$, and $b, \hat{x} \in \mathbb{R}^n$. If $\|RA - I\|_\infty < 1$, then A is nonsingular, and

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

Proof: It is sufficient to prove the upper bound.

We have $A^{-1} = (RA)^{-1}R = (I - (I - RA))^{-1}R$, so

$$\|A^{-1}b - \hat{x}\|_\infty = \|A^{-1}(b - A\hat{x})\|_\infty \leq \|(I - (I - RA))^{-1}\|_\infty \|R(b - A\hat{x})\|_\infty$$

This gives the desired inequality. \square

Remark: we also have

$$\|A^{-1}b - \hat{x}\|_\infty \leq \frac{\|R\|_\infty \|A\hat{x} - b\|_\infty}{1 - \|RA - I\|_\infty},$$

but this upper bound uses a “sub-multiplicative inequality”.

Verified upper bound

Let $\hat{\mathbf{x}}$ be an approximate solution of $A\mathbf{x} = \mathbf{b}$ and R be an approximate inverse of A . If we can check that $\|R\mathbf{A} - \mathbf{I}\|_\infty < 1$, then

$$\|A^{-1}\mathbf{b} - \hat{\mathbf{x}}\|_\infty \leq \frac{\|R(A\hat{\mathbf{x}} - \mathbf{b})\|_\infty}{1 - \|R\mathbf{A} - \mathbf{I}\|_\infty}.$$

The algorithm will therefore have the general form:

function $[\hat{\mathbf{x}}, \bar{\delta}] = \text{certifLSV}(A, \mathbf{b})$

$\hat{\mathbf{x}} = \text{fl}(A^{-1}\mathbf{b})$ $\{ O(n^3) \}$

$R = \text{fl}(A^{-1})$ $\{ O(n^3) \}$

Compute an upper bound $\bar{\alpha}$ of $\|R\mathbf{A} - \mathbf{I}\|_\infty$ $\{ O(n^3) \}$

if $\bar{\alpha} \geq 1$ then error('Certification failed')

Compute an upper bound $\bar{\beta}$ of $\|R(A\hat{\mathbf{x}} - \mathbf{b})\|_\infty$ $\{ O(n^2) \}$

$\bar{\delta} = \text{fl}_\uparrow(\bar{\alpha} / \text{fl}_\downarrow(1 - \bar{\beta}))$

In the sequel, we denote $\alpha := \|R\mathbf{A} - \mathbf{I}\|_\infty$ and $\beta := \|R(A\hat{\mathbf{x}} - \mathbf{b})\|_\infty$.

Outline

1

Preliminary results

- Triangular linear systems
- Inverse of triangular matrices
- LU factorization and solution of linear systems
- Normwise conditioning of linear systems

2

Verified solution of linear systems

- General principle
- **Four methods to upper bound α**
- Upper bound for β
- Synthesis

3

Numerical experiments

Upper bound for $\alpha = \|RA - I\|_\infty$: method 1

We simply use the techniques of fast implementation of interval arithmetic for operations on the matrices (see previous lecture):

Cost: $4n^3 + O(n^2)$

function $\bar{\alpha}_1 = \text{bndAlpha1}(A, R)$

$$\underline{G} = \text{fl}_\downarrow(RA - I) \quad \{ 2n^3 + O(n^2) \}$$

$$\overline{G} = \text{fl}_\uparrow(RA - I) \quad \{ 2n^3 + O(n^2) \}$$

$$T = \max(|\underline{G}|, |\overline{G}|) \quad \{ O(n^2) \}$$

$$\bar{\alpha}_1 = \max(\text{fl}_\uparrow(Te)) \quad \{ O(n^2) \}$$

We have $\underline{G} \leq RA - I \leq \overline{G}$, so $|RA - I| \leq T$, $|RA - I|e \leq \text{fl}_\uparrow(Te)$ and as $\|RA - I\|_\infty = \max(|RA - I|e)$, we can deduce that $\alpha \leq \max(\text{fl}_\uparrow(Te))$.

Upper bound for $\alpha = \|RA - I\|_\infty$: method 2

A solution to reduce the cost of calculating an upper bound of α is to use an a priori **error analysis to upper bound $|RA - I|$** .
With rounding to nearest,

$$|(RA - I) - \text{fl}(RA - I)| \leq \gamma_{n+1}(|R||A| + I).$$

We deduce that: $|RA - I|e \leq |\text{fl}(RA - I)|e + \gamma_{n+1}(|R||A| + I)e$.

So, by denoting $\bar{\gamma}_k := \text{fl}_\uparrow(ku/(1 - ku))$, we obtain the following algorithm:

Cost: $2n^3 + O(n^2)$

function $\bar{\alpha}_2 = \text{bndAlpha2}(A, R)$

if $(n + 1)u \geq 1$ then error("Verification failed")

$H = \text{fl}(|RA - I|)$ $\{ 2n^3 + O(n^2) \}$

$\bar{h} = \text{fl}_\uparrow(H e + \bar{\gamma}_{n+1}(|R|(|A|e) + e))$ $\{ O(n^2) \}$

$\bar{\alpha}_2 = \max(\bar{h})$

Upper bound for $\alpha = \|\mathbf{RA} - \mathbf{I}\|_\infty$: method 3

If $\mathbf{PA} = \mathbf{LU}$ is an LU factorization of \mathbf{A} , we can replace \mathbf{A}^{-1} by $\mathbf{U}^{-1}\mathbf{L}^{-1}\mathbf{P} \dots$

Let $\mathbf{X}_L, \mathbf{X}_U \in \mathbb{R}^{n \times n}$ be approximate inverse of \mathbf{L} and \mathbf{U} .

We denote $\mathbf{R} := \mathbf{X}_U \mathbf{X}_L \mathbf{P}$, and we want to bound

$$\alpha = \|\mathbf{X}_U \mathbf{X}_L \mathbf{PA} - \mathbf{I}\|_\infty.$$

We can use some a priori bounds. Let us make appear $|\mathbf{X}_U \mathbf{U} - \mathbf{I}|$:

$$|\mathbf{X}_U \mathbf{X}_L \mathbf{PA} - \mathbf{I}| = |\mathbf{X}_U (\mathbf{X}_L \mathbf{PA} - \mathbf{U}) + (\mathbf{X}_U \mathbf{U} - \mathbf{I})| \leq |\mathbf{X}_U| |\mathbf{X}_L \mathbf{PA} - \mathbf{U}| + |\mathbf{X}_U \mathbf{U} - \mathbf{I}|.$$

We assume that \mathbf{X}_U is computed by solving $\mathbf{XU} = \mathbf{I}$:

$$|\mathbf{X}_U \mathbf{U} - \mathbf{I}| \leq \gamma_n |\mathbf{X}_U| |\mathbf{U}|.$$

For $|\mathbf{X}_L \mathbf{PA} - \mathbf{U}|$, we still use an a priori error analysis,

$$|\mathbf{X}_L \mathbf{PA} - \mathbf{U}| \leq |\text{fl}(\mathbf{X}_L \mathbf{PA} - \mathbf{U})| + \gamma_{n+1} (|\mathbf{X}_L| |\mathbf{PA}| + |\mathbf{U}|),$$

Upper bound for $\alpha = \|RA - I\|_\infty$: methode 3

We also obtain

$$|X_U X_L PA - I| \leq |X_U| (|fl(X_L PA - U)| + \gamma_{n+1} (|X_L| |PA| + 2|U|)),$$

and $\|X_U X_L PA - I\|_\infty \leq \max(\bar{h})$, with

$$\bar{h} = fl_\uparrow (|X_U| (|fl(X_L PA - U)|e + \bar{\gamma}_{n+1} (|X_L| (|PA|e) + 2|U|e)))$$

Cost: $n^3 + O(n^2)$

function $\bar{\alpha}_3 = \text{bndAlpha3}(A, P, L, U, X_L, X_U)$

if $(n+1)u \geq 1$ then error('Verification failed')

$H = fl(|X_L PA|)$ $\{ n^3 + O(n^2) \}$

$\bar{h} = fl_\uparrow (|X_U| (He + \bar{\gamma}_{n+1} (|X_L| (|PA|e) + 2|U|e)))$ $\{ O(n^2) \}$

$\bar{\alpha}_3 = \max(\bar{h})$

Upper bound for $\alpha = \|RA - I\|_\infty$: method 4

We can use the way we computed the factors L and U.
Assuming it was with GEPP,

$$|LU - PA| \leq \gamma_n |L| |U|.$$

Let us make appear $|LU - PA|$ in the bound for $|X_U X_L PA - I|$:

$$\begin{aligned} |X_U X_L PA - I| &\leq |X_U X_L (PA - LU) + X_U X_L LU - I| \\ &\leq |X_U X_L (PA - LU) + X_U (X_L L - I) U + (X_U U - I)| \\ &\leq |X_U| |X_L| |PA - LU| + |X_U| |X_L L - I| |U| + |X_U U - I|. \end{aligned}$$

We assume that X_L and X_U were computed "by solving $XT = I$ ":

$$\begin{aligned} |X_U X_L PA - I| &\leq \gamma_n |X_U| |X_L| |L| |U| + \gamma_n |X_U| |X_L| |L| |U| + \gamma_n |X_U| |U| \\ &\leq 2\gamma_n |X_U| |X_L| |L| |U| + \gamma_n |X_U| |U|. \end{aligned}$$

Upper bound for $\alpha = \|RA - I\|_\infty$: method 4

We obtain $\|X_U X_L P A - I\|_\infty \leq \max(\bar{h})$, with

$$\bar{h} = \text{fl}_\uparrow(2 \bar{\gamma}_n |X_U| |X_L| |L| |U| e + \gamma_n |X_U| |U| e).$$

Cost: $O(n^2)$

```
function  $\bar{\alpha}_4 = \text{bndAlpha4}(A, P, L, U, X_L, X_U)$   
  if  $\text{nu} \geq 1$  then error('Verification failed')  
   $\bar{h} = \text{fl}_\uparrow(2 \bar{\gamma}_n |X_U| |X_L| |L| |U| e + \gamma_n |X_U| |U| e)$       {  $O(n^2)$  }  
   $\bar{\alpha}_4 = \max(\bar{h})$ 
```

Outline

1

Preliminary results

- Triangular linear systems
- Inverse of triangular matrices
- LU factorization and solution of linear systems
- Normwise conditioning of linear systems

2

Verified solution of linear systems

- General principle
- Four methods to upper bound α
- **Upper bound for β**
- Synthesis

3

Numerical experiments

Upper bound for β

Let us recall that $\beta := \|R(A\hat{x} - b)\|_\infty$. Two remarks:

- An upper bound for β can be computed with $O(n^2)$ flops.
- The global cost will be in $O(n^3)$. So we can use more accurate algorithms to compute the residual $A\hat{x} - b$: for example `CompDot2`.

Upper bound for β

When an approximate inverse R of A is computed explicitly:

Compute $\bar{\beta}_1$ such that $\beta \leq \bar{\beta}_1$ in $O(n^2)$ flops.

```
function  $\bar{\beta}_1 = \text{bndBeta1}(A, b, \hat{x}, R)$   
   $\langle m_{\text{res}}, r_{\text{res}} \rangle = \text{resLinSys2}(A, b, \hat{x})$   
   $\underline{g} = \text{fl}_{\downarrow}(Rm_{\text{res}}); \bar{g} = \text{fl}_{\uparrow}(Rm_{\text{res}})$   
   $\underline{h} = \text{fl}_{\uparrow}(\max(|\bar{g}|, |\underline{g}|) + |R|r_{\text{res}})$   
   $\bar{\beta}_1 = \max(\underline{h})$ 
```

Proof: Denote $A\hat{x} - b = m_{\text{res}} + \rho_{\text{res}}$, with $|\rho_{\text{res}}| \leq r_{\text{res}}$. Then,
 $R(A\hat{x} - b) = Rm_{\text{res}} + \rho$, with $|\rho| \leq |R|r_{\text{res}}$, and
 $|R(A\hat{x} - b)| \leq |Rm_{\text{res}}| + |R|r_{\text{res}}$.

Upper bound for β

If R is only known in the form of $X_U X_L P$:

Compute $\bar{\beta}_2$ such that $\beta \leq \bar{\beta}_2$ in $O(n^2)$ flops.

```
function  $\bar{\beta}_2 = \text{bndBeta2}(A, b, \hat{x}, X_L, X_U)$ 
```

```
   $\langle m_{\text{res}}, r_{\text{res}} \rangle = \text{resLinSys2}(A, b, \hat{x})$ 
```

```
   $\bar{g} = \sup(|X_U(X_L(P \langle m_{\text{res}}, r_{\text{res}} \rangle))|)$       { interval arithmetic }
```

```
   $\bar{\beta}_2 = \max(\bar{g})$ 
```

Outline

1

Preliminary results

- Triangular linear systems
- Inverse of triangular matrices
- LU factorization and solution of linear systems
- Normwise conditioning of linear systems

2

Verified solution of linear systems

- General principle
- Four methods to upper bound α
- Upper bound for β
- **Synthesis**

3

Numerical experiments

Synthesis

```
function  $[\hat{x}, \bar{\delta}] = \text{certifLSV}\{\mathbf{1}, \mathbf{2}\}(A, b)$   
   $[L, U, P] = \text{fl}(\text{xGETRF}(A))$   $\{ \mathbf{2/3} \cdot n^3 + O(n^2) \}$   
   $\hat{x} = \text{fl}(\text{xGETRS}(P, L, U, b))$   $\{ O(n^2) \}$   
   $R = \text{fl}(\text{xGETRI}(P, L, U))$   $\{ \mathbf{4/3} \cdot n^3 + O(n^2) \}$   
   $\bar{\alpha} = \text{bndAlpha}\{\mathbf{1}, \mathbf{2}\}(A, R)$   $\{ \mathbf{4} \cdot n^3 + O(n^2) / \mathbf{2} \cdot n^3 + O(n^2) \}$   
  if  $\bar{\alpha} \geq 1$  then error('Certification failed')  
   $\bar{\beta} = \text{bndBeta1}(A, b, \hat{x}, R)$   $\{ O(n^2) \}$   
   $\bar{\delta} = \text{fl}_{\uparrow}(\bar{\alpha} / \text{fl}_{\downarrow}(1 - \bar{\beta}))$ 
```

Algorithm	Global cost	Certif. cost	Ref.
certifLSV1	$\mathbf{6} \cdot n^3 + O(n^2)$	$\mathbf{16/3} \cdot n^3 + O(n^2)$	[ORO05]
certifLSV2	$\mathbf{4} \cdot n^3 + O(n^2)$	$\mathbf{10/3} \cdot n^3 + O(n^2)$	[ORO05]

Synthesis

```

function  $[\hat{x}, \bar{\delta}] = \text{certifLSV}\{\mathbf{3}, \mathbf{4}\}(A, b)$ 
   $[L, U, P] = \text{fl}(\text{xGETRF}(A))$   $\{ \mathbf{2/3} \cdot n^3 + O(n^2) \}$ 
   $\hat{x} = \text{fl}(\text{xGETRS}(P, L, U, b))$   $\{ O(n^2) \}$ 
   $X_L = \text{fl}(\text{xTRTRI}(L))$   $\{ \mathbf{1/3} \cdot n^3 + O(n^2) \}$ 
   $X_U = \text{fl}(\text{xTRTRI}(U))$   $\{ \mathbf{1/3} \cdot n^3 + O(n^2) \}$ 
   $\bar{\alpha} = \text{bndAlpha}\{\mathbf{3}, \mathbf{4}\}(A, P, L, U, X_L, X_U)$   $\{ \mathbf{n^3} + O(n^2) / O(n^2) \}$ 
  if  $\bar{\alpha} \geq 1$  then error('Certification failed')
   $\bar{\beta} = \text{bndBeta2}(A, b, \hat{x}, X_L, X_U)$   $\{ O(n^2) \}$ 
   $\bar{\delta} = \text{fl}_{\uparrow}(\bar{\alpha} / \text{fl}_{\downarrow}(1 - \bar{\beta}))$ 

```

Algorithm	Global cost	Certif. cost	Ref.
certifLSV3	$\mathbf{7/3} \cdot n^3 + O(n^2)$	$\mathbf{5/3} \cdot n^3 + O(n^2)$	cf. [OR02]
certifLSV4	$\mathbf{4/3} \cdot n^3 + O(n^2)$	$\mathbf{2/3} \cdot n^3 + O(n^2)$	[OR02]

Algorithm	Global cost	Certif. cost	Ref.
certifLSV1	$6 \cdot n^3 + O(n^2)$	$16/3 \cdot n^3 + O(n^2)$	[ORO05]
certifLSV2	$4 \cdot n^3 + O(n^2)$	$10/3 \cdot n^3 + O(n^2)$	[ORO05]
certifLSV3	$7/3 \cdot n^3 + O(n^2)$	$5/3 \cdot n^3 + O(n^2)$	Cf. [OR02]
certifLSV4	$4/3 \cdot n^3 + O(n^2)$	$2/3 \cdot n^3 + O(n^2)$	[OR02]

What about the numerical quality of those algorithms?

Outline

1

Preliminary results

- Triangular linear systems
- Inverse of triangular matrices
- LU factorization and solution of linear systems
- Normwise conditioning of linear systems

2

Verified solution of linear systems

- General principle
- Four methods to upper bound α
- Upper bound for β
- Synthesis

3

Numerical experiments

Numerical experiments

- All the experiments were performed with MATLAB, and INTLAB toolbox¹ for:
 - dealing with rounding modes;
 - using interval arithmetic.
- We use IEEE-754 double precision arithmetic.
- The ill-conditioned linear systems are generated as follows:

```
A = gallery('randsvd', n, 10^(k*rand));  
b = A*ones(n,1);
```

We obtain some condition numbers between 1 and 10^k .
The exact condition number is then recomputed later.

1. <http://www.ti3.tu-harburg.de/rump/intlab/>

Numerical experiments

Implementation of function `bndAlpha1`:

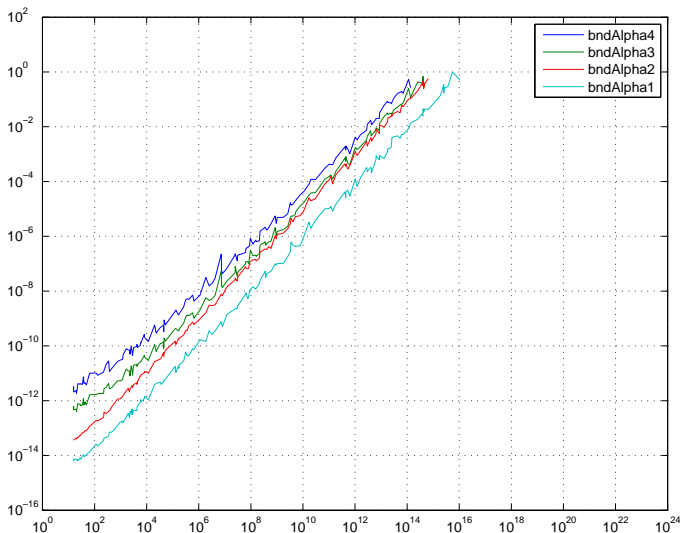
```
function alpha = bndAlpha1(A, R, n)
    setround(-1);          % rounding toward -Inf
    infG = R*A-eye(n);
    setround(1);           % rounding toward +Inf
    supG = R*A-eye(n);
    T = max(abs(infG), abs(supG));
    alpha = max(T*ones(n,1));
```

Implementation of function `bndBeta2`:

```
function beta = bndBeta2(P, A, b, x, XL, XU)
    [mres, rres] = resLinSys2(b, A, x);
    r = midrad(mres, rres);          % midpoint/radius
    g = sup(abs(XU*(XL*(P*r)))); % interval arithmetic
    beta = max(g);
```

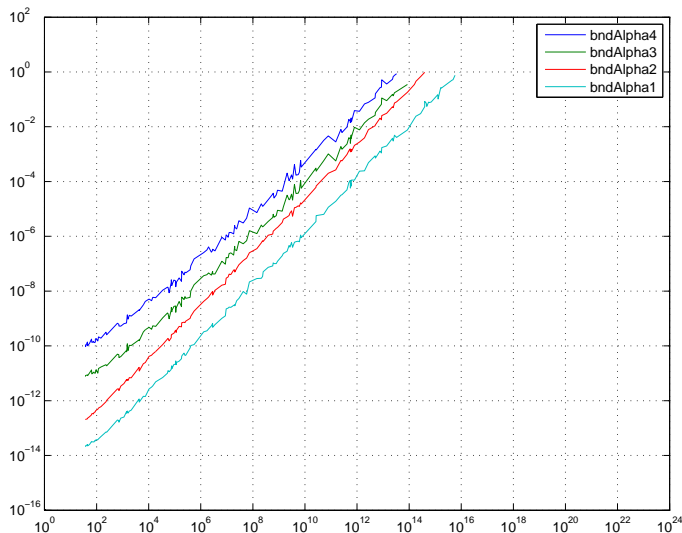
Numerical experiments

System of 20 equations, bound for $\alpha = \|RA - I\|_\infty$:



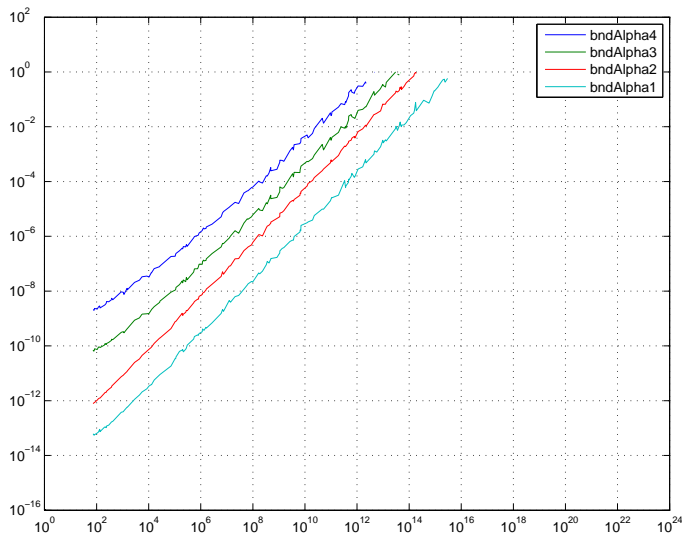
Numerical experiments

System of 50 equations, bound for $\alpha = \|RA - I\|_\infty$:



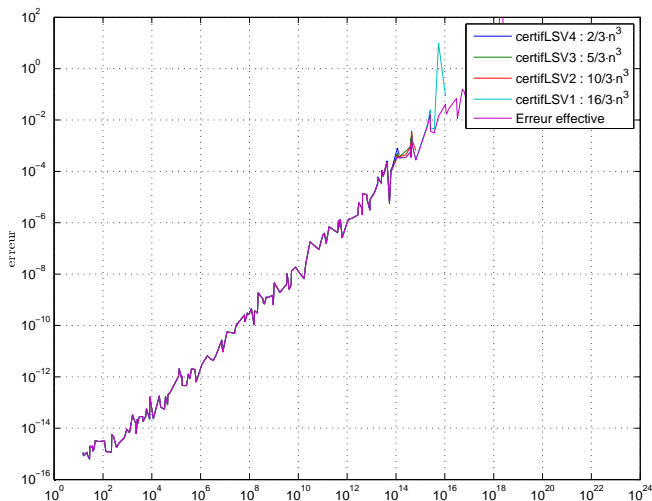
Numerical experiments

System of 100 equations, bound for $\alpha = \|RA - I\|_\infty$:



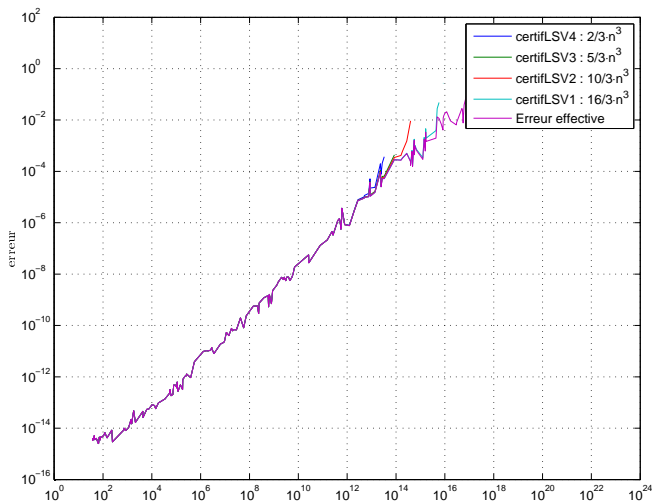
Numerical experiments

System of 20 equations, bound δ for $\|\hat{\mathbf{x}} - \mathbf{x}\|_\infty$, we report $\delta/\|\mathbf{x}\|_\infty$:



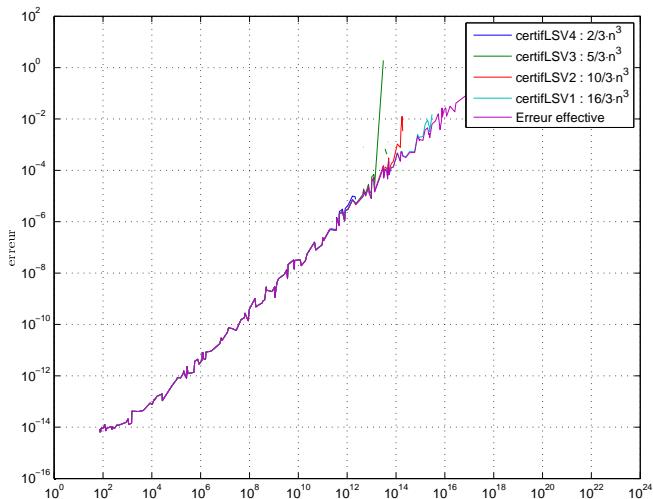
Numerical experiments

System of 50 equations, bound δ for $\|\hat{\mathbf{x}} - \mathbf{x}\|_\infty$, we report $\delta/\|\mathbf{x}\|_\infty$:



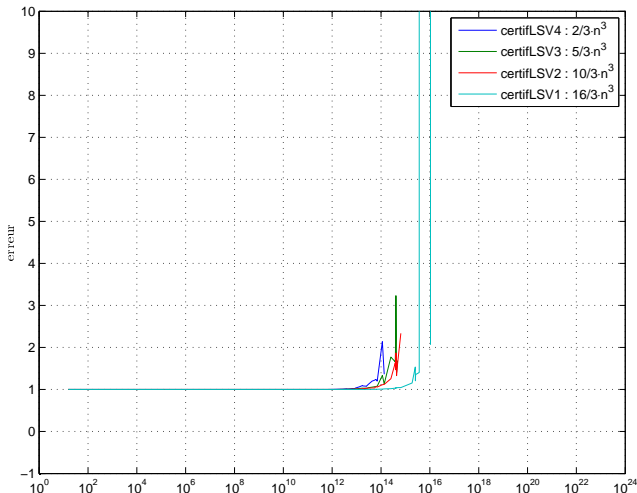
Numerical experiments

System of 100 equations, bound δ for $\|\hat{\mathbf{x}} - \mathbf{x}\|_\infty$, we report $\delta/\|\mathbf{x}\|_\infty$:



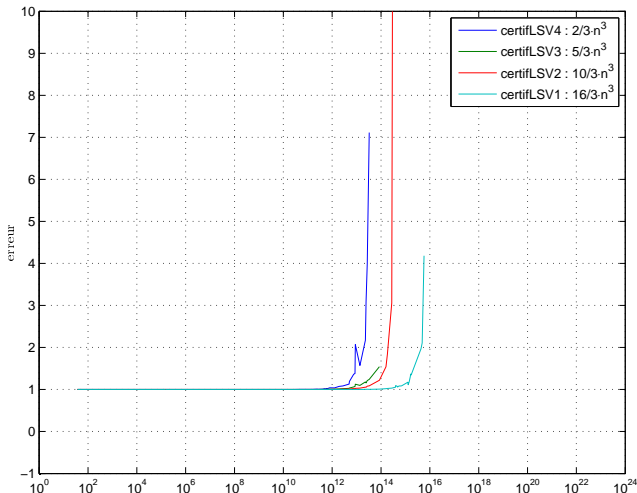
Numerical experiments

System of 20 equations, bound δ for $\|\hat{\mathbf{x}} - \mathbf{x}\|_\infty$, we report $\delta / \|\hat{\mathbf{x}} - \mathbf{x}\|_\infty$



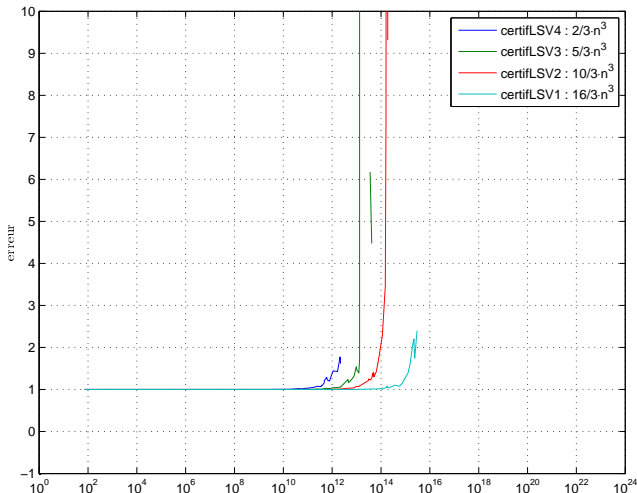
Numerical experiments

System of 50 equations, bound δ for $\|\hat{\mathbf{x}} - \mathbf{x}\|_\infty$, we report $\delta / \|\hat{\mathbf{x}} - \mathbf{x}\|_\infty$



Numerical experiments





System of 100 equations, bound δ for $\|\hat{\mathbf{x}} - \mathbf{x}\|_\infty$, we report $\delta/\|\hat{\mathbf{x}} - \mathbf{x}\|_\infty$





Summary

- Upper bound for $\alpha = \|RA - I\|_\infty$:
 - The more expensive the method is, the more one is able to verify the solution for large condition numbers.
 - When n increases, the condition number below which we can verify a solution decreases
- Upper bound for the error $\delta = \|x - \hat{x}\|_\infty$:
 - If $\alpha < 1$, the computed upper bound δ is often good, independently of the algorithm used: this is due to the fact that the residual $A\hat{x} - b$ is computed with twice the working precision.
 - On the other hand, the upper bound δ becomes more pessimistic for “large condition numbers”: the quality of R as approximate inverse gets worse.
- It is possible to bound α with heuristic in $O(n^2)$ flops: the idea is to use such estimates to select at runtime the fastest method for verification.

References I

-  J. W. Demmel, Applied numerical linear algebra, SIAM, Philadelphia, PA, 1997.
-  G. H. Golub and C. van Loan, Matrix computations, 3 ed., The Johns Hopkins University Press, London, 1996.
-  Nicholas J. Higham, Accuracy and stability of numerical algorithms, second ed., Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2002.
-  Shin'Ichi Oishi and Siegfried M. Rump, Fast verification of solutions of matrix equations, Numerische Mathematik **90** (2002), no. 4, 755–773.

References II

-  Takeshi Ogita, Siegfried M. Rump, and Shin'ichi Oishi, Verified solution of linear systems without directed rounding, Technical Report No. 2005-04, Advanced Research Institute for Science and Engineering, Waseda University, 2005.
-  Siegfried M. Rump, Verification methods: Rigorous results using floating-point arithmetic, Acta Numerica **19** (2010), 287-449.