# Iterative Refinement and Verified Numerical Linear Algebra

## OGITA, Takeshi

Division of Mathematical Sciences

Tokyo Woman's Christian University

Japan

SIAM LA21, Virtual Conference

(Originally scheduled in New Orleans, Louisiana, U.S.)

May 19 (CDT) [May 20 (JST)], 2021

# Outline

Question · How do we know the quality of computed error bounds for solutions in numerical linear algebra?

$\Longrightarrow$ We consider error of error for solutions of linear systems.

Purpose · We answer the question using the following two tools:

- Verified numerical computations
  $\Longrightarrow$ error bounds of computed solutions.

- Iterative refinement
  $\Longrightarrow$ error reduction of computed solutions.

# (Usual) verified computation

Notation: For $x = (x_1, \ldots, x_n)^T \in \mathbb{R}^n$, $|x| = (|x_1|, \ldots, |x_n|)^T$.

Given an approximate solution $\widetilde{x}$ of $Ax = b$, the usual verified computation gives an upper bound of the error or its norm:

$$|\widetilde{x} - A^{-1}b| \leq \epsilon \in \mathbb{R}^n \quad \text{or} \quad \|\widetilde{x} - A^{-1}b\|_\infty \leq \max_{1 \leq i \leq n} \epsilon_i = \varepsilon \in \mathbb{R}$$

$\Longrightarrow$ At least, $\widetilde{x}_i$ has correct digits (accuracy) corresponding to $\epsilon_i$.

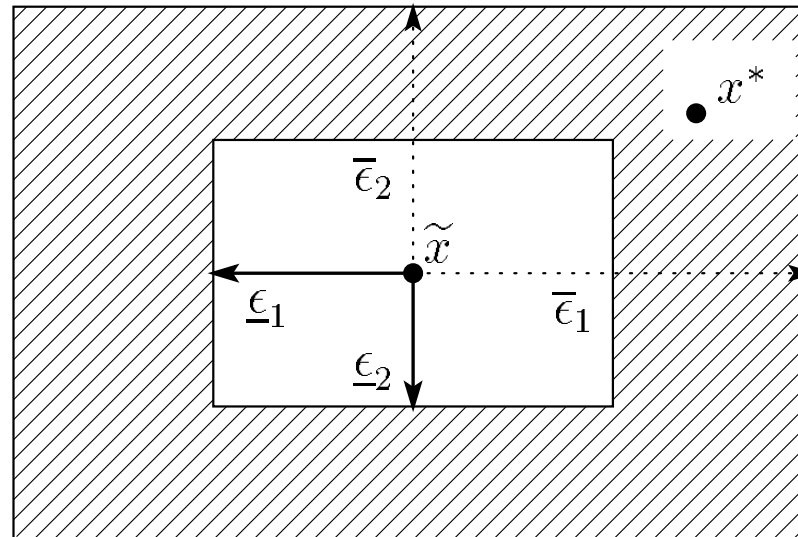$\Longrightarrow$ However, $\epsilon_i$ may be overestimated (too pessimistic).

$\Longrightarrow$ The quality of the verification is still not known!

# Quality of the verification

How (and whether) can we know it?

# Compute both lower and upper error bounds

If both $\underline{\epsilon}$ and $\overline{\epsilon}$ s.t. $\underline{\epsilon} \leq |\widetilde{x} - A^{-1}b| \leq \overline{\epsilon}$ and $\overline{\epsilon} \approx \underline{\epsilon}$ are obtained, then the quality of the verification (evaluation) can be confirmed!



 : possible domain where the exact solution $x^*$ exists

Figure 1: Inner and outer enclosure of the exact solution (two-dimensional case)

Question: Is it possible to obtain such $\underline{\epsilon}$ and $\bar{\epsilon}$ without much computational cost?

Answer: Yes. It is not so difficult! Let's see how to do it.

# Nonsingularity of $A$ and upper bound of $\|A^{-1}\|$

It needs some effort in terms of computational cost. For example,

- Let $R$ be an approximate inverse of $A$. If $\|I - RA\| < 1$, then $A$ is proved to be nonsingular and

$$\|A^{-1}\| \leq \frac{\|R\|}{1 - \|I - RA\|}.$$

- computing a lower bound $\underline{\sigma}$ of the smallest singular value of $A \implies$ If $\underline{\sigma} > 0$, then $\|A^{-1}\|_2 \leq 1/\underline{\sigma}$.

# Basic theorem for tight error bounds

**Theorem 1. [O. et al., 2003]** *Let $A$ be a real $n \times n$ matrix and $b$ be a real $n$-vector. Let $\widetilde{x}$ be an approximate solution of $Ax = b$ and $r := b - A\widetilde{x}$. Let $\widetilde{y}$ be an approximate solution of $Ay = r$. If $A$ is nonsingular, then it holds for $p \in \{1, 2, \infty\}$ that*

$$|A^{-1}b - \widetilde{x}| \le |\widetilde{y}| + \|A^{-1}\|_p \|r - A\widetilde{y}\|_p e, \tag{1}$$

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

# Tight enclosure of the solution

For an arbitrary $y \in \mathbb{R}^n$, we have

$$A^{-1}b - \widetilde{x} = A^{-1}b - (\widetilde{x} + y) + y.$$

It follows that

$$|y| - \epsilon_y \leq |A^{-1}b - \widetilde{x}| \leq |y| + \epsilon_y \quad \text{with} \quad \epsilon_y := |A^{-1}b - (\widetilde{x} + y)|.$$

Using this and Theorem 1, we have the following proposition.

**Proposition 1.** *Let $A, b, \widetilde{x}$ and $r$ be as in Theorem 1. Let $\widetilde{y}$ be an approximate solution of $Ay = r$. Assume that $A$ is nonsingular and $\rho$ satisfies $\|A^{-1}\|_p \leq \rho$ for any $p \in \{1, 2, \infty\}$. Then*

$$\max(|\widetilde{y}| - \epsilon, \mathbf{o}) \leq |A^{-1}b - \widetilde{x}| \leq |\widetilde{y}| + \epsilon, \qquad (2)$$

*where $\epsilon := \rho\|r - A\widetilde{y}\|_p e$ and $\mathbf{o} = (0, \ldots, 0)^T \in \mathbb{R}^n$.*

$\implies$ If $|\widetilde{y}_i| \gg \epsilon_i$, the error bounds are very tight!

$\implies$ Such $|\widetilde{y}|$ can be obtained by the iterative refinement method.

# Iterative refinement and staggered correction

To obtain a tight enclosure of an approximate solution $\widetilde{x}$ of a linear system $Ax = b$, we introduce a so-called "staggered correction".

$\mathbb{F}$: a set of floaing-point numbers

Using iterative refinements, we can obtain $\widetilde{x} + y$ with arbitrarily higher precision: For $R \approx A^{-1}$

$$y^{(\ell+1)} = R * (b - A(\widetilde{x} + y^{(\ell)})),$$

where $y^{(\ell)} = \sum_{k=1}^{M} y_k^{(\ell)}$ with $y_k^{(\ell)} \in \mathbb{F}^n$. $\implies$ The correction term $y$ can be expressed by the sum of floating-point vectors.

This makes only sense for calculating the residual $b - A(\widetilde{x} + y^{(\ell)})$ accurately.

$\implies$ Fortunately, we have accurate dot product algorithms.

[1] O., Rump, Oishi: *Accurate sum and dot product*, SISC, 26:6 (2005).

[2] Rump, O., Oishi: *Accurate floating-point summation: Part I/Part II*, SISC, 31:1/2 (2008).

On the other hand, to obtain tight error bounds, we need to compute

$$\epsilon_i = \rho \| r - A\widetilde{y} \|_p = \rho \| b - A(\widetilde{x} + \widetilde{y}) \|_p.$$

This is compatible with the iterative refinements!

# Behavior of iterative refinement

Assume that an approximate inverse $R \in \mathbb{F}^{n \times n}$ of $A$ is computed by a backward stable algorithm, e.g. LU factorization with partial pivoting. Then, the following is known as a rule of thumb: For $\mu := \mathrm{cond}_\infty(A) < \mathbf{u}^{-1}$ and $G := I - RA$,

$$\alpha := \|G\|_\infty = \mathcal{O}(n\mathbf{u})\mu. \tag{3}$$

Let $\widetilde{x} = Rb$ and $e := (1, \ldots, 1)^T$. Since

$$|A^{-1}b - \widetilde{x}| = |A^{-1}b - Rb| = |(I - RA)A^{-1}b| \leq |G||A^{-1}b|,$$

it holds that

$$|A^{-1}b - \widetilde{x}| \leq \|A^{-1}b\|_\infty |G|e. \tag{4}$$

After an iterative refinement by using $y^{(1)} = R(b - A\widetilde{x})$, it follows that

$$
\begin{aligned}
|A^{-1}b - (\widetilde{x} + y^{(1)})| &= |A^{-1}b - \widetilde{x} - R(b - A\widetilde{x})| \\
&= |(I - RA)(A^{-1}b - \widetilde{x})| \\
&\leq |G||A^{-1}b - \widetilde{x}|. \quad\quad (5)
\end{aligned}
$$

Inserting (4) into (5) yields

$$
|A^{-1}b - (\widetilde{x} + y^{(1)})| \leq \|A^{-1}b\|_{\infty}|G|^2 e.
$$

For $k \geq 2$, it can inductively be proved for $y^{(k)} = y^{(k-1)} + R(b - A(\widetilde{x} + y^{(k-1)}))$ that

$$
|A^{-1}b - (\widetilde{x} + y^{(k)})| \leq \|A^{-1}b\|_{\infty}|G|^{k+1} e
$$

and
$$|A^{-1}b - (\widetilde{x} + y^{(k)})| \leq \alpha^{k+1}\|A^{-1}b\|_\infty e. \qquad (6)$$
Therefore, if $\alpha < 1$, then the iterative refinement converges with the factor $\alpha = \mathcal{O}(n\mathbf{u})\mu$ for each iteration.

In practice, due to the rounding error, we have $\widetilde{x}^{(k)} = \mathrm{fl}\left(\widetilde{x} + y^{(k)}\right)$ with $\widetilde{x}^{(0)} = \widetilde{x}$ and

$$|A^{-1}b - \widetilde{x}^{(k)}| \leq \mathbf{u}|A^{-1}b| + \mathcal{O}(\alpha^{k+1})\|A^{-1}b\|_\infty e. \qquad (7)$$

This is a *componentwise* error bound and explains the behavior of the iterative refinement.

# Example

Let us consider the case where $A \in \mathbb{F}^{5 \times 5}$ with $\mathrm{cond}_\infty(A) \approx 10^{10}$ and the *exact* solution $A^{-1}b = (1, 10^3, 10^6, 10^9, 134217728)^T$.

All computations are done in double precision arithmetic on Matlab, so that $\mathbf{u} = 2^{-53} \approx 10^{-16}$.

An approximate inverse $R$: computed by a Matlab function `inv` $\implies$
$\alpha = \|I - RA\|_\infty \approx 10^{-6}$   $(\mathbf{u} \cdot \mathrm{cond}_\infty(A) \approx 10^{-6})$

$\widetilde{x}^{(0)} = \mathrm{fl}\,(Rb)$

Let's see the result of iterative refiments.

## Table 1: History of iterative refinement for $k = 0, 1, 2$

| $i$ | $\widetilde{x}^{(0)}$ | $\widetilde{x}^{(1)}$ |
|---|---|---|
| 1 | $-1.711885408 \cdot 10^2$ | $\underline{0.99993}5694692795$ |
| 2 | $\underline{1.02}1301738 \cdot 10^3$ | $\underline{1.00000001}1437893 \cdot 10^3$ |
| 3 | $\underline{1.000}55792 \cdot 10^6$ | $\underline{0.9999999997}9213 \cdot 10^6$ |
| 4 | $\underline{1.000000}83 \cdot 10^9$ | $\underline{0.99999999999980} \cdot 10^9$ |

| $i$ | $\widetilde{x}^{(2)}$ |
|---|---|
| 1 | $\underline{1.0000000000}2729$ |
| 2 | $\underline{1.00000000000000} \cdot 10^3$ |
| 3 | $\underline{1.00000000000000} \cdot 10^6$ |
| 4 | $\underline{1.00000000000000} \cdot 10^9$ |

## Table 2: True absolute errors $\epsilon$ and tight error bounds

| $i$ | $\epsilon^{(0)}$ | $\epsilon^{(1)}$ | $\epsilon^{(2)}$ |
|---|---|---|---|
| 1 | $1.7218854 \cdot 10^2$ | $6.430530 \cdot 10^{-5}$ | $2.728484 \cdot 10^{-12}$ |
| 2 | $2.1301738 \cdot 10^1$ | $1.143789 \cdot 10^{-5}$ | $3.410605 \cdot 10^{-13}$ |
| 3 | $5.5792398 \cdot 10^1$ | $2.078711 \cdot 10^{-5}$ | $0$ |
| 4 | $8.3648966 \cdot 10^1$ | $2.026557 \cdot 10^{-5}$ | $0$ |

| $i$ | $[\underline{\epsilon}_i^{(0)}, \overline{\epsilon}_i^{(0)}]$ | $[\underline{\epsilon}_i^{(1)}, \overline{\epsilon}_i^{(1)}]$ | $[\underline{\epsilon}_i^{(2)}, \overline{\epsilon}_i^{(2)}]$ |
|---|---|---|---|
| 1 | $1.72188^{55}_{53} \cdot 10^2$ | $6.4305^{34}_{28} \cdot 10^{-5}$ | $2.728^6_4 \cdot 10^{-12}$ |
| 2 | $2.13017^{39}_{37} \cdot 10^1$ | $1.1437^{92}_{86} \cdot 10^{-5}$ | $3.410^9_4 \cdot 10^{-13}$ |
| 3 | $5.57923^{99}_{97} \cdot 10^1$ | $2.0787^{15}_{09} \cdot 10^{-5}$ | $[0, 2.05] \cdot 10^{-17}$ |
| 4 | $8.36489^{67}_{65} \cdot 10^1$ | $2.0265^{61}_{55} \cdot 10^{-5}$ | $[0, 2.05] \cdot 10^{-17}$ |

# Thanks!