

Unit 6: Numerical Stability

Any real number can be written as: $x = u \times \beta^e$

where: β is the base (an integer > 1)

$u \in (-1, 1)$ is the mantissa

e is the exponent

For our discussion, $x = u \times \beta^e$ such that:

- $\beta = 2$ (computer use binary system)
- $u = \pm \delta_0 \delta_1 \dots \delta_{t-1}$ (u has only + binary digits, $\delta_i \in \{0, 1\}$)
- $\delta_0 = 0$ iff $u=0$ (the mantissa is normalized)
- $-1 \leq e \leq 0$

Example: Let $\beta = 2$, $t = 3$, $u = 0.101$, $e = -1$

$$\begin{aligned} x &= u \times \beta^e \\ &= .101 \times 2^{-1} \\ &= 0.125 \end{aligned}$$

Example of overflow

Since $\|x\|_2 = \sqrt{\sum_{i=0}^{n-1} x_i^2} \Rightarrow \|x\|_2 \leq \sqrt{n} \max_{i=0}^{n-1} |x_i|$

If more than one x_i is close to overflowing,

the result will overflow since \sqrt{t} is overflowed.

The solution is to find k such that:

$$\|x_k\| = \max_{i=0}^{n-1} |x_i|$$

Giấy dó trắng tự nhiên
không hiat mat

And then instead compute

$$\|x\|_2 = |x_k| \sqrt{\sum_{i=0}^{n-1} \left(\frac{x_i}{x_k} \right)^2}$$

Remark E.2.t.3.

Any real number stored in our computer is stored as nearby floating point number (element in F) either through rounding or truncation.

Error in storing real number as floating point number

The error is bounded by:

$$|\tilde{x}| \leq \epsilon_{\text{mach}} |x|$$

where ϵ_{mach} is known as the machine epsilon or unit roundoff. When single precision floating point numbers are used $\epsilon_{\text{mach}} \approx 10^{-8}$ (eight decimal digits of accuracy).

When double precision floating point numbers are used $\epsilon_{\text{mach}} \approx 10^{-16}$

If \tilde{x} is computed by truncating:

$$|\tilde{x}| \leq 2^{-(t-k)} |x|$$

If \tilde{x} is computed by rounding:

$$|\tilde{x}| \leq 2^{-k} |x|$$

Giấy dó trắng tự nhiên
không hiat mat

fahasa

Machine epsilon (unit roundoff)

ϵ_{mach} is defined as the smallest positive floating point number x such that the floating point number that represents $1 + x$ is greater than 1.

Python code to understand this:

```
epsilon = 2**-52 # machine epsilon for double precision
print(1 + epsilon > 1) # True
print(1 + (epsilon/2) > 1) # False
```

Ex 6.2.2.1

Given system with $b=2$, a with t digits, truncation when storing

Number as floating point number in this system

$$0.\underbrace{10 \dots 0}_{t \text{ digits}} \times 2^1$$

What ϵ_{mach} of this system?

$$\epsilon_{\text{mach}} = 2^{-(t-1)}$$

Standard Computational Model (SCM)

$$f_1(x \oplus \psi) = (x \oplus \psi)(1 + \epsilon), | \epsilon | \leq \epsilon_{\text{mach}}$$

and $\oplus \in \{+, -, \times, /, \}$

Remark 6.2.3.2

SCM can be interpreted as: These operations are performed exactly and it is only in storing the result that a roundoff error occurs.

Example

Consider $k = \frac{4}{3}$. Our floating point system with $b=2$, $t=4$. Our real number $\frac{4}{3} = 1.333\ldots$ is stored as $.1010 \times 2^1$, if $t=4$ and truncation is deployed. This equals 1.25 in decimal, the relative error was 0.0625

$$k = f_1(\frac{4}{3})$$

$$= 1.25$$

$$= 1.333\ldots + (-0.0833\ldots)$$

$$= 1.333\ldots \times (1 + \frac{-0.0833\ldots}{1.333\ldots})$$

$$= \frac{4}{3} \times (1 + (-0.0625))$$

$$= k(1 + \epsilon) \text{ where } | \epsilon | = 0.0625 \leq \epsilon_{\text{mach}}$$

$$\frac{4}{3} \times \frac{1}{8} = \frac{1}{6}$$

Alternative Computational Model (ACM)

$$f(x \text{ op } \epsilon) = \frac{x \text{ op } \epsilon}{1 + \epsilon}, |\epsilon| \leq \epsilon_{\text{mach}}, \text{op} \in \{+, -, *, /\}$$

Notice that the Taylor series expansion at $\frac{1}{1+\epsilon}$ is given by:

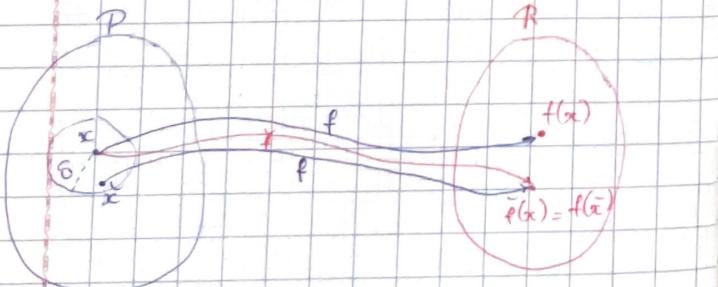
$$\frac{1}{1+\epsilon} = \underbrace{1}_{(1-\epsilon)} + \underbrace{\epsilon^2}_{\epsilon^2 - \epsilon^3} + \dots$$

ACM is useful when analyzing algorithms that involve division.

Number instability

We wish to evaluate the mapping $f: D \rightarrow R$.

Let $\tilde{f}: D \rightarrow R$ is the computer implementation of f



Here, \tilde{f} represent the function that use floating point arithmetic, thus incurring errors. The fact that

for a nearby value, \tilde{x} , the computed value equals the exact function applied to the slightly perturbed input x

$$\tilde{f}(\tilde{x}) = f(x)$$

means that the error in computation can be attributed to the small change in the input. It is said that

\tilde{f} is a (numerically) stable implementation of f

Backward stable implementation

Given mapping $f: D \rightarrow R$, where $D \subset \mathbb{R}^n$, $R \subset \mathbb{R}^m$

Let \tilde{f} be the computer implementation of f .

\tilde{f} is backward stable (numerically stable)

implementation of f if for all $x \in D$, exists \tilde{x} "close" to x such that:

$$\tilde{f}(x) = f(\tilde{x})$$

Forward stable: $\tilde{f}(x) \approx f(x)$

Condition vs Stability

- Conditioning is a property of the problem. A problem is well-conditioned if a small change in the input is guaranteed to only result in a small change in output.

Stability is a property of an implementation. A numerically stable implementation when executed with an input always yields an output that can be attributed to slightly changed input.

Now,

- Well-conditioned problem + numerically stable implementation = answer that is close to the actual answer
- Else, you may not get a good answer.

Definition 6.2.6.1 Absolute value of vector & matrix

Given $x \in \mathbb{R}^n$ and $A \in \mathbb{R}^{m \times n}$

$$|x| = \begin{pmatrix} |x_0| \\ |x_1| \\ \vdots \\ |x_{n-1}| \end{pmatrix} \quad |A| = \begin{pmatrix} |a_{0,0}| & \dots & |a_{0,n-1}| \\ |a_{1,0}| & \dots & |a_{1,n-1}| \\ \vdots & & \vdots \\ |a_{m-1,0}| & \dots & |a_{m-1,n-1}| \end{pmatrix}$$

Definition 6.2.6.2

Let $\Delta \in \{\leq, \geq, =, \geq, \geq\}$ and $x, y \in \mathbb{R}^n$, then:

- $|x| \Delta |y|$ iff $|x_i| \Delta |y_i|$
- $|A| \Delta |B|$ iff $|a_{ij}| \Delta |b_{ij}|$

Giấy dó trắng tự nhiên
không hại mắt

fahasa

thm 6.2.6.1

$$\text{Let } A \in \mathbb{R}^{m \times k}, B \in \mathbb{R}^{k \times n}$$

$\langle \text{ALWAYS} \rangle$ $|AB| \leq |A| |B|$

Theorem 6.2.6.3

$$\text{Let } A, B \in \mathbb{R}^{m \times m}. \text{ If } |A| \leq |B| \text{ then } \|A\|_F \leq \|B\|_F$$

$\|A\|_F \leq \|B\|_F$, and $\|A\|_\infty \leq \|B\|_\infty$

Backward error analysis of dot product: general case

Lemma 6.3.2.1

Let $\epsilon_i \in \mathbb{R}$ ($0 \leq i \leq n-1$)

$\epsilon_n \in \epsilon_{\text{mach}}$

$|\epsilon_i| \leq \epsilon_{\text{mach}}$. Then $\exists \theta_n \in \mathbb{R}$ such that:

$$\prod_{i=0}^{n-1} (1 + \epsilon_i)^{\pm 1} = 1 + \theta_n$$

with $|\theta_n| \leq \left(\frac{n}{(1 - n\epsilon_{\text{mach}})} \right) \epsilon_{\text{mach}}$

Remark 6.3.2.2

The quantity θ_n is the number of error that has been accumulated. Example: $\epsilon_0^0 + \epsilon_1^1 + \epsilon_2^2 + \dots + \epsilon_{n-1}^{n-1}$

Giấy dó trắng tự nhiên
không hại mắt

fahasa

The bound of $\hat{\theta}_n$

$$|\hat{\theta}_n| \leq \frac{n \epsilon_{\text{mach}}}{1 - n \epsilon_{\text{mach}}} = \gamma_n$$

For $n \geq 1$ and $n \epsilon_{\text{mach}} < 1$

With these notation, dot product simplified to:

$$\begin{aligned} \tilde{k} &= x_0 \psi_0 (1 + \theta_0) + x_1 \psi_1 (1 + \theta_1) + x_2 \psi_2 (1 + \theta_2) \\ &\quad \dots + x_{n-1} \psi_{n-1} (1 + \theta_{n-1}) \\ &= \begin{pmatrix} x_0 \\ \vdots \\ x_{n-1} \end{pmatrix}^T \begin{pmatrix} (1 + \theta_0) & & & \\ & (1 + \theta_1) & & \\ & & \ddots & \\ & & & (1 + \theta_{n-1}) \end{pmatrix} \begin{pmatrix} \psi_0 \\ \vdots \\ \psi_{n-1} \end{pmatrix} \\ &= \begin{pmatrix} x_0 \\ \vdots \\ x_{n-1} \end{pmatrix}^T \left(I + \begin{pmatrix} \theta_0 & & & \\ & \ddots & & \\ & & \theta_{n-1} & \\ & & & \theta_n \end{pmatrix} \right) \begin{pmatrix} \psi_0 \\ \vdots \\ \psi_{n-1} \end{pmatrix} \\ &\approx \boxed{x^T (I + \Sigma^{(n)}) y} \quad \text{Theorem 6.3.3.1} \end{aligned}$$

where $|\theta_j| \leq \gamma_j$, $j = 2, \dots, n$

Lemma 6.2.3.5

If $n, b \geq 1 \Rightarrow \gamma_n \leq \gamma_{n+b}$

and $\gamma_n + \gamma_b + \gamma_n \gamma_b \leq \gamma_{n+b}$

This lemma is used when we want to bound $|\tilde{k}|$

Giấy dó trắng tự nhiên
không hai mặt

such that $1 + \epsilon = (1 + \epsilon_1)(1 + \epsilon_2)$

$$= 1 + (\epsilon_1 + \epsilon_2 + \epsilon_1 \epsilon_2)$$

knowing that $|\epsilon_1| \leq \gamma_n$, $|\epsilon_2| \leq \gamma_b$

Theorem 6.3.3.1

Let $k \geq 0$, assume that $|\epsilon_1|, |\epsilon_2| \leq \epsilon_{\text{mach}}$,
with $\epsilon_1 = 0$ if $k = 0$

$$\left(\begin{array}{cc} I + \Sigma^{(k)} & 0 \\ 0 & (1 + \epsilon) \end{array} \right) (1 + \epsilon_2) = (I + \Sigma^{(k+1)})$$

Corollary 6.3.3.2

Under the assumption of Theorem 6.3.3.1, the following relations hold:

$$R-1B \quad \tilde{k} = (x + \delta x)^T y, \text{ where } |\delta x| \leq \gamma_n |x|$$

$$R-2B \quad \tilde{k} = x^T (y + \delta y), \text{ where } |\delta y| \leq \gamma_n |y|$$

$$R-1F \quad \tilde{k} = x^T y + \delta k, \text{ where } |\delta k| \leq \gamma_n |x|^T |y|$$

Matrix-vector multiplication

→ next page

Giấy dó trắng tự nhiên
không hai mặt

Theorem 6.3.1

Let $A \in \mathbb{R}^{m \times n}$, $x \in \mathbb{R}^n$, $y \in \mathbb{R}^m$, consider $y = Ax$

$$R-1B \quad \tilde{y} = (A + \Delta A)x, \text{ where } |\Delta A| \leq \gamma_n |A|$$

$$R-1F \quad \tilde{y} = Ax + \Delta y, \text{ where } |\Delta y| \leq \gamma_n |A| \|x\|$$

Note: $\tilde{y} = A(x + \Delta x)$ is not possible in most cases since Δx is different for every row \tilde{a}_i^T of A

Matrix-matrix multiplication Theorem 6.3.5.1

Let $C \in \mathbb{R}^{m \times n}$, $A \in \mathbb{R}^{m \times k}$ and $B \in \mathbb{R}^{k \times n}$ and consider $C = AB$, then $\exists \Delta C$ such that

$$\tilde{C} = AB + \Delta C, \text{ where } |\Delta C| \leq \gamma_k |A| |B|$$

Numerical stability of triangular solve

Lemma 6.4.1.2

Let $n \geq 1$, $\lambda, v \in \mathbb{R}$ and $x, y \in \mathbb{R}^n$. Assume $\lambda \neq 0$ and consider the computation

$$z = (\alpha - x^T y) / \lambda$$

$$\text{then: } (\lambda + \delta \lambda) z = \alpha - (x - \delta x)^T y$$

$$\text{where } |\delta x| \leq \gamma_n |x|, |\delta \lambda| \leq \gamma_2 |L|$$

Theorem 6.4.1.3 Triangular solve

Let $L \in \mathbb{R}^{n \times n}$ be nonsingular lower triangular matrix

\tilde{x} is the computed result when executing $Lx = y$

Then \exists matrix ΔL such that

$$(L + \Delta L) \tilde{x} = y, \text{ where } |\Delta L| \leq \max(\gamma_2, \gamma_{n+1}) \|L\|$$

We can relax the bound of the above theorem about:

$$|\Delta L| \leq \gamma_n \|L\|$$

Note that γ_n does not equals $\max(\gamma_2, \gamma_{n+1})$, especially for cases when $n=1$

Theorem 6.4.2.1 Backward error of Crout variant for LU factorization

Let $A \in \mathbb{R}^{n \times n}$ and let the LU factorization of A be computed via the Crout variant, yielding approx factors \tilde{L} and \tilde{U} . Then:

$$(A + \Delta A) = \tilde{L} \tilde{U} \text{ with } |\Delta A| \leq \gamma_n \|\tilde{L}\| \|\tilde{U}\|$$

Numerical stability of linear solve via LU factorization

From what we have known:

$$\cdot (A + \Delta A) - L\bar{U}, \|AA\| \leq \gamma_n \|L\| \|U\|$$

$$\cdot (L + \Delta L) \bar{x} = y, \|LL\| \leq \gamma_n \|L\|$$

$$\cdot (U + \Delta U) \bar{x} = z, \|UU\| \leq \gamma_n \|U\|$$

$$(A + \Delta A) \bar{x} = y, \|AA\| \leq (3\gamma_n + \gamma_n^2) \|L\| \|U\|$$

Note: ΔA from solving Ax = y is different from ΔA from LU factorization

Thm 6.4.3.1

Consider $Ax = y$ and $(A + \Delta A)(x + \delta x) = y$

$$\frac{\|\delta x\|}{\|x\|} \leq k(A) \frac{\|\Delta A\|}{\|A\|}$$

Theorem 6.4.3.2

Let A be nonsingular, $\|\cdot\|$ be a subordinate norm, and

$$\frac{\|\Delta A\|}{\|A\|} < \frac{1}{k(A)}$$

Then $A + \Delta A$ is nonsingular

How partial pivoting is important for numerical stability

When $A = LU$ without pivoting, we might encounter multiplier that is "small" $\rightarrow \|L\|$ and $\|U\|$ very large $\rightarrow \|AA\| \leq (3\gamma_n + \gamma_n^2) \|L\| \|U\|$ very very large

When $PA = LU$ with pivoting, we pivot the largest multiplier on top and divide by it $\rightarrow \|L\|$ and $\|U\|$ have elements $\leq 1 \rightarrow \|AA\| \leq (3\gamma_n + \gamma_n^2) \|L\| \|U\|$ is more reasonable

$$P(A + \Delta A) = L\bar{U}, \text{ where } P(\Delta A) \leq \gamma_n \|L\| \|U\|$$

Thm 6.4.5.2

$$A = \begin{pmatrix} 1 & 0 & \dots & 0 & 1 \\ -1 & 1 & & & \vdots \\ -1 & \vdots & \ddots & \ddots & \vdots \\ -1 & \vdots & \vdots & \ddots & 1 \\ -1 & -1 & \dots & -1 & 1 \end{pmatrix} \quad \text{LU-factorize } A, \text{ we have:}$$
$$U = \begin{pmatrix} 1 & \dots & 1 \\ 0 & 1 & 2 \\ \vdots & \ddots & 2^{n-2} \\ 0 & 0 & 2^{n-1} \end{pmatrix}$$

Therefore, in theory even with partial pivoting, U result can get very large.

However, in practice, this does not happen and LU factorization (especially with pivoting) is stable.