

Round-off error analysis

Input error $x \rightarrow fl(x)$

$$fl(x) = \pm 0.d_1d_2 \dots d_n \beta^b \xrightarrow{\text{Chopping}} \text{round-off}$$

$$x = \pm 0.d_1d_2 \dots d_n d_{n+1} \dots \beta^b \xrightarrow{\text{rounding}}$$

Round-off $|x - fl(x)| \leq \frac{1}{2} \beta^{b-n}$ $x \neq 0, d_1 \neq 0$

$$x \neq 0, \left| \frac{x - fl(x)}{x} \right| \leq \frac{\frac{1}{2} \beta^{b-n}}{0.1 \cdot \beta^b} = \boxed{\frac{1}{2} \beta^{-n+1}} = \epsilon$$

ϵ is called the precision, epsilon, accuracy.
 $fl(x)$ is the normalized floating number

$$\beta = 10, \quad \epsilon = \frac{1}{2} 10^{-n+1} \quad n = 23 \quad \epsilon$$

How to represent computers in terms of real number system.

Thm: (Wilkinson) If $x \in \mathbb{R}$, $\underline{fl}(x) = \frac{x(1+\delta_x)}{1+\delta_x}$
 is a perturbation of true value of x

$$x - \underline{fl}(x) = -x\delta$$

$$|x - \underline{fl}(x)| = |x||\delta|$$

relative error $\frac{x - \underline{fl}(x)}{x} = \frac{-x\delta}{x} = -\delta$ $x = 3.14$

$$|\delta| \leq \frac{1}{2} \beta^{-n+1} = \epsilon$$

$$\frac{|x - \underline{fl}(x)|}{|x|} \leq \epsilon = \frac{1}{2} \beta^{-n+1}$$

Conclusion: When we input a number into a computer system, the relative error is bounded by the precision

Operation error analysis $\pm \%$ / logic oper.

Let x and y are numbers in a computer system, and, if, or, no error

$$fl(x \circ y) = (x \circ y) (1 + \delta_1) \quad \text{Theorem.}$$

$$|\delta_1| \leq \epsilon. \quad \circ, +, -, \times, \div, \sqrt{\quad}$$

$$x = 0.d_1d_2\dots d_n \cdot \beta^{b_1}$$

$$y = 0.\bar{d}_1\bar{d}_2\dots\bar{d}_n \beta^{b_2}$$

$$b_1 \neq b_2, d_1 \neq 0$$

$$0.d_1d_2\dots d_n 10^{+k} \bar{d}_1 \neq 0$$

$$0.\bar{d}_1\bar{d}_2\dots\bar{d}_n \cdot 10^p$$

x, y are not in the computer system, $x \rightarrow 0.00\bar{d}_1\dots\bar{d}_{n+2} \cdot 10^2$

$$x \rightarrow fl(x) \quad y \rightarrow fl(y)$$

$$fl(x \cdot y) = fl(x(1+\delta_1) \cdot y(1+\delta_2)), \quad |\delta_i| \leq \epsilon$$

$$= xy (1+\delta_1)(1+\delta_2)(1+\delta_3)$$

relative error

input x

input y

multiplication of

$fl(x) \cdot fl(y)$

$$= xy (1 + \underbrace{\delta_1 + \delta_2 + \delta_3}_{\leq \epsilon} + \underbrace{\delta_1\delta_2 + \delta_1\delta_3 + \delta_2\delta_3}_{\epsilon^2} + \underbrace{\delta_1\delta_2\delta_3}_{\epsilon^3})$$

$$= xy (1 + \delta_4)$$

$$|\delta_4| \leq 3\epsilon$$

$$\epsilon = 10^{-8}$$

$$\epsilon = 10^{-16}$$

$$\epsilon^2 = 10^{-16}$$

Subtraction, $x \in \mathbb{R}, y \in \mathbb{R}$.

$$fl(x-y) = fl \left(\underset{\text{input } x}{x(1+\delta_1)} - \underset{\text{input } y}{y(1+\delta_2)} \right)$$

$$= (x(1+\delta_1) - y(1+\delta_2))(1+\delta_3) \quad \delta_3 \text{ from the subtraction}$$

$$= (x-y)(1+?) \quad |?| \leq$$

$$= (x + x\delta_1, -y - y\delta_2)(1+\delta_3)$$

$$= (x-y + x\delta_1 - y\delta_2)(1+\delta_3) + (x-y)(1+\delta_3)$$

$$= (x-y) \left(1 + \frac{x\delta_1 - y\delta_2}{x-y} + \delta_3 \right) + O(\epsilon^2)$$

$$\delta_4 = \frac{x\delta_1 - y\delta_2}{x-y} + \delta_3, \quad |\delta_4| \leq ? \epsilon$$

If x and y have different signs $x = -y$

$$|\delta_4| \leq \frac{|x|\delta_1 + |y|\delta_2}{|x| + |y|} + \epsilon \leq \frac{|x|\epsilon + |y|\epsilon}{|x| + |y|} = 3\epsilon$$

If $x-y$ is very small, the relative error can be large, very close, we may have large relative error, Cancellation of significant digits. Catastrophic

Review: machine precision $\frac{1}{2} \epsilon^{-n+1}$
 32 bits, $n=52$.

Ex 2 $f(x) = \frac{xy}{z} = \frac{xy}{z} (1+\delta)$ $|\delta| \leq 5\epsilon$

$$= \frac{x(1+\delta_1)y(1+\delta_2)(1+\delta_3)}{z(1+\delta_4)(1+\delta_5)}$$

$$= \frac{xy(1+\tilde{\delta}_1)}{z(1+\tilde{\delta}_2)}$$

$$\approx \frac{xy}{z} (1+\tilde{\delta}_1)(1-\tilde{\delta}_2)$$

$$\approx \frac{xy}{z} (1+\delta)$$

$\frac{1}{1+\tilde{\delta}_2} = 1 - \tilde{\delta}_2 + \tilde{\delta}_2^2 - \tilde{\delta}_2^3 + \dots + (-1)^n \tilde{\delta}_2^n$

$\Rightarrow \frac{xy}{z} \frac{1}{z + \epsilon \sin(\theta)}$ $\frac{\partial \phi}{|\partial \phi|}$

Aug 29-9:34 AM

How to minimize round-off errors.
 $ax^2+bx+c=0$ $x = \frac{-b \pm \sqrt{b^2-4ac}}{2a}$

Algorithm:
 A1: $x_1 = \frac{-b + \sqrt{b^2-4ac}}{2a}$ $x_2 = \frac{-b - \sqrt{b^2-4ac}}{2a}$
 A2: $x_1 = \frac{-b + \sqrt{b^2-4ac}}{2a}$ $x_2 = \frac{c}{ax_1}$ $b < 0$
 A3: $x_1 = \frac{-b - \sqrt{b^2-4ac}}{2a}$ $x_2 = \frac{c}{ax_1}$ $b > 0$ ✓

Vieta's formula $x_1+x_2 = -\frac{b}{a}$, $x_1x_2 = \frac{c}{a}$

Let $a=1$, $c=\epsilon$, $b=2$ small number $(1-\epsilon)^x$
 $x_1 = \frac{-2 - \sqrt{4-4\epsilon}}{2} = -1 - \sqrt{1-\epsilon} = -1 - 1 + \frac{\epsilon}{2} = -2 + \frac{\epsilon}{2}$
 $x_2 = \frac{-2 + \sqrt{4-4\epsilon}}{2} = -1 + 1 - \frac{\epsilon}{2} = -\frac{\epsilon}{2}$

$\frac{x_1 - x_2}{x_1 - x_2}$

Aug 29-9:54 AM

How to avoid/minimize $(a+b)(a-b) = a^2 - b^2$
 $-b + \sqrt{b^2-4ac}$

if $b > 0$, $\left(\frac{-b + \sqrt{b^2-4ac}}{-b - \sqrt{b^2-4ac}} \right) \left(-b - \sqrt{b^2-4ac} \right)$

$$= \frac{(-b) - (b^2-4ac)}{-b - \sqrt{b^2-4ac}} = \frac{4ac}{-b - \sqrt{b^2-4ac}}$$
 ✓

Ex: $1 - \cos x$

$$= \begin{cases} 1 - \cos x & \text{if } |x| > \delta \\ 1 - (1 - \frac{x^2}{2} + \frac{x^4}{24}) & |x| \leq \delta \end{cases}$$

$$= \begin{cases} 1 - \cos x \\ \frac{x^2}{2} \end{cases}$$

$u'(x) = \frac{u(x+h) - u(x)}{h}$
 $+ O(h)$
 method error $O(h)$
 method error $O(h)$
 The best h
 $\delta = \sqrt{\epsilon}$

Aug 29-10:10 AM

Some basic programming skills.

$\sum_{i=1}^n a_i$, $S=0$ ✓ initialization.
 for $i=1:n$
 $S = S + a(i)$
 end

$\prod_{i=1}^n a_i$
 $p=1$
 for $i=1:n$
 $p = p * a(i)$
 end

At Matrix-vector multiplication,
 $A \in \mathbb{R}^{m \times n}$, $x \in \mathbb{R}^{n \times 1}$ $Ax \in \mathbb{R}^m$ $y = Ax$

for $i=1:m$
 $y(i) = 0$
 for $j=1:n$
 $y(i) = y(i) + a(i,j) * x(j)$
 end
 end

$\begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$

Aug 29-10:18 AM

$p(x) = a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0 x^0$

$p=0$
 for $i=0:n$
 $p = p + a(i) * x^i(i)$
 end

Complexity ?

Aug 29-10:25 AM

Algorithm efficiency: Accuracy, fast, (FFT)
 Simple, storage.

Research topic: How to obtain triple accuracy based on the current computer system.

Evaluate a polynomial $p_n(x) = \sum_{i=0}^n a_i x^i$

$$p_3(x) = a_3 x^3 + a_2 x^2 + a_1 x + a_0$$

$$= x(a_3 x^2 + a_2 x + a_1) + a_0$$

$$= x(x(a_3 x + a_2) + a_1) + a_0$$

Nexted algorithm

A pseudo-code

```

P = a(n)
for i = n-1, -1, 0
    P = P * x + a(i);
end
    
```

Complexity: $+ - \cdot n$
 $X/\pm n$

The best for sequential machines.

Aug 31-9:33 AM

Vector and matrix norms, a review.

$x \rightarrow f(x) \quad f(x) = x(x+s) \quad |s| \leq \epsilon$
 $\vec{x} \rightarrow f(\vec{x}) \quad \frac{1}{\beta} \beta^{-n+1}$

Definition of a vector norm.

Given a \mathbb{R}^n , $\mathbb{R}^n = \{\vec{x}, \vec{y} = \begin{bmatrix} x \\ x \\ \vdots \\ x \end{bmatrix}, x \in \mathbb{R}\}$

$x \in \mathbb{R}, |x|, \vec{x} \in \mathbb{R}^3, x = \begin{bmatrix} 1 \\ -1 \\ 2 \end{bmatrix}, \|\vec{x}\|_2 = \sqrt{1^2 + (-1)^2 + 2^2} = \sqrt{6}$

A norm is a special function $f(\vec{x})$ that satisfies

(i) $f(\vec{x}) \geq 0, f(\vec{x}) = 0$ if and only if $\vec{x} = \vec{0}$.

(ii) $f(\alpha \vec{x}) = |\alpha| f(\vec{x})$ iff

(iii) $f(\vec{x} + \vec{y}) \leq f(\vec{x}) + f(\vec{y})$

Then $f(\vec{x})$ is a norm in \mathbb{R}^n , can be written as $f(x) = \|x\|$

Aug 31-9:50 AM

Ex: 1. $f(x) = \max_{1 \leq i \leq n} \{ |x_i| \}$

Check: 1. $f(\vec{x}) \geq 0$. If $f(\vec{x}) = 0$, $\max_{1 \leq i \leq n} \{ |x_i| \} = 0$

2. $f(\alpha \vec{x}) = \max_{1 \leq i \leq n} \{ |\alpha x_i| \} = |\alpha| \max_{1 \leq i \leq n} \{ |x_i| \} = |\alpha| f(\vec{x})$

3. $f(\vec{x} + \vec{y}) = \max_{1 \leq i \leq n} \{ |x_i + y_i| \} = \max_{1 \leq i \leq n} \{ |x_i| + |y_i| \}$

$$\leq \max_{1 \leq i \leq n} \{ |x_i| \} + \max_{1 \leq i \leq n} \{ |y_i| \}$$

$$\leq f(\vec{x}) + f(\vec{y})$$

Aug 31-10:01 AM

Are the following vector norms?

1. $f(x) = \max_{1 \leq i \leq n} \{ |x_i| \} + s, f(\vec{0}) \neq 0$

2. $f(x) = \frac{\max_{1 \leq i \leq n} \{ |x_i| \}}{x_i^2}$, No $f(\vec{x}) \neq |k| f(\vec{x})$

3. $f(x) = (\sum_{i=1}^n x_i^2)^k$, Yes.

(i) $f(\vec{x}) \geq 0, f(\vec{x}) = 0$ iff $\vec{x} = \vec{0}$

(ii) $f(\alpha \vec{x}) = (\sum_{i=1}^n (\alpha x_i)^2)^k = (\sum_{i=1}^n \alpha^2 x_i^2)^k = |\alpha|^{2k} (\sum_{i=1}^n x_i^2)^k = |\alpha|^{2k} f(\vec{x})$

(iii) $f(\vec{x} + \vec{y}) \leq f(\vec{x}) + f(\vec{y})$

$$\sqrt{\sum_{i=1}^n (x_i + y_i)^2} \leq \sqrt{\sum_{i=1}^n x_i^2} + \sqrt{\sum_{i=1}^n y_i^2}$$

$$\sum_{i=1}^n (x_i + y_i)^2 \leq \sum_{i=1}^n x_i^2 + \sum_{i=1}^n y_i^2 + 2 \sqrt{\sum_{i=1}^n x_i^2 \sum_{i=1}^n y_i^2}$$

$$\sum_{i=1}^n (x_i^2 + 2x_i y_i + y_i^2) \leq \sum_{i=1}^n x_i^2 + \sum_{i=1}^n y_i^2 + 2 \sqrt{\sum_{i=1}^n x_i^2 \sum_{i=1}^n y_i^2}$$

$$\sum_{i=1}^n 2x_i y_i \leq 2 \sqrt{\sum_{i=1}^n x_i^2 \sum_{i=1}^n y_i^2}$$

Aug 31-10:06 AM

Cauchy-Schwartz inequality $\checkmark \checkmark$

Introduce $f(\lambda) = \sum_{i=1}^n (x_i - \lambda y_i)^2 \geq 0$

$$= \sum_{i=1}^n (x_i^2 - 2\lambda x_i y_i + \lambda^2 y_i^2)$$

$$= \sum_{i=1}^n x_i^2 - 2\lambda \sum_{i=1}^n x_i y_i + \lambda^2 \sum_{i=1}^n y_i^2$$

$$= a + b\lambda + c\lambda^2$$

$b^2 - 4ac \leq 0$

$$\left(-2 \sum_{i=1}^n x_i y_i \right)^2 \leq 4 \sum_{i=1}^n x_i^2 \sum_{i=1}^n y_i^2$$

$$\left| \sum_{i=1}^n x_i y_i \right| \leq \sqrt{\sum_{i=1}^n x_i^2} \cdot \sqrt{\sum_{i=1}^n y_i^2}$$

$\sum_{i=1}^n y_i^2 \leq f(\vec{x}) = \sqrt{\sum_{i=1}^n x_i^2}$ is a vector norm

$\|\vec{x}\|_2$ Euclidean Norm

Aug 31-10:17 AM