Elements of Floating-point Arithmetic

Sanzheng Qiao

Department of Computing and Software McMaster University

September, 2011

◆□▶ ◆□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Outline

Floating-point Numbers

- Representations
- IEEE Floating-point Standards
- Underflow and Overflow
- Correctly Rounded Operations

2 Sources of Errors

- Rounding Error
- Truncation Error
- Discretization Error
- Stability of an Algorithm
 - 4 Sensitiviy of a Problem

5 Fallacies

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

Outline

Floating-point Numbers

- Representations
- IEEE Floating-point Standards
- Underflow and Overflow
- Correctly Rounded Operations
- 2 Sources of Errors
 - Rounding Error
 - Truncation Error
 - Discretization Error
- 3 Stability of an Algorithm
- 4 Sensitiviy of a Problem
- 5 Fallacies

Two ways of representing floating-point

On paper we write a floating-point number in the format:

 $\pm d_1.d_2\cdots d_t \times \beta^e$

$$0 < d_1 < \beta, 0 \le d_i < \beta \ (i > 1)$$

- t: precision
- β: base (or radix), almost universally 2, other commonly used bases are 10 and 16

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

e: exponent, integer

Two ways of representing floating-point (cont.)

- Examples: 1.0×10^{-1}
- t = 2 (the last zero counts), $\beta = 10$, e = -1

Two ways of representing floating-point (cont.)

- Examples:
- 1.0×10^{-1}
- t = 2 (the last zero counts), $\beta = 10$, e = -1
- $1.234 imes 10^2$
- $t = 4, \beta = 10, e = 2$

Two ways of representing floating-point (cont.)

- Examples:
- 1.0×10^{-1}
- t = 2 (the last zero counts), $\beta = 10$, e = -1
- $1.234 imes 10^2$
- $t = 4, \, \beta = 10, \, e = 2$
- 1.10011×2^{-4}
- $t = 6, \beta = 2$ (binary), e = -4

Two ways of representing floating-point (cont.)

Examples: 1.0×10^{-1} t = 2 (the last zero counts), $\beta = 10$, e = -1 1.234×10^{2} t = 4, $\beta = 10$, e = 2 1.10011×2^{-4} t = 6, $\beta = 2$ (binary), e = -4

The precision *t*, the base β , and the range of the exponent *e* determine a floating-point number system.

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

◆□▶ ◆□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

In memory, a floating-point number is stored in three consecutive fields:

sign (1 bit) exponent (depends on the range)

fraction (depends on the precision)

In memory, a floating-point number is stored in three consecutive fields:

sign (1 bit) exponent (depends on the range) fraction (depends on the precision)

In order for a memory representation to be useful, there must be a standard.

◆□▶ ◆□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

In memory, a floating-point number is stored in three consecutive fields:

```
sign (1 bit)
exponent (depends on the range)
fraction (depends on the precision)
```

In order for a memory representation to be useful, there must be a standard.

◆□▶ ◆□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

IEEE floating-point standards: single precision and double precision.

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

Characteristics

A floating-point number system is characterized by four (integer) parameters:

- base β (also called radix)
- precision t
- exponent range $e_{min} \le e \le e_{max}$

Machine precision

A real number representing the accuracy.

Machine precision

Denoted by ϵ_M , defined as the distance between 1.0 and the next larger floating-point number, which is $0.0...01 \times \beta^0$.

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

Machine precision

A real number representing the accuracy.

Machine precision

Denoted by ϵ_M , defined as the distance between 1.0 and the next larger floating-point number, which is $0.0...01 \times \beta^0$.

Thus, $\epsilon_M = \beta^{1-t}$.

Equivalently, the distance between two consecutive floating-point numbers between 1.0 and β . (The floating-point numbers between 1.0 and β are evenly spaced, 1.0...000, 1.0...001, 1.0...010, ..., 1.1...111.)

How would you compute the underlying machine precision?

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

How would you compute the underlying machine precision?

◆□▶ ◆□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

```
The smallest \epsilon such that 1.0 + \epsilon > 1.0.
```

How would you compute the underlying machine precision?

◆□▶ ◆□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

```
The smallest \epsilon such that 1.0 + \epsilon > 1.0.
```

For $\beta = 2$:

```
eps = 1.0;
while (1.0 + eps > 1.0)
        eps = eps/2;
end
2*eps,
```

How would you compute the underlying machine precision?

```
The smallest \epsilon such that 1.0 + \epsilon > 1.0.
```

For $\beta = 2$:

```
eps = 1.0;
while (1.0 + eps > 1.0)
        eps = eps/2;
end
2*eps,
```

```
Examples. (\beta = 2)
When t = 24, \epsilon_M = 2^{-23} \approx 1.2 \times 10^{-7}
When t = 53, \epsilon_M = 2^{-52} \approx 2.2 \times 10^{-16}
```

Approximations of real numbers

Since floating-point numbers are discrete, a real number, for example, $\sqrt{2}$, may not be representable in floating-point. Thus real numbers are approximated by floating-point numbers.

We denote

 $\mathrm{fl}(\mathbf{x}) \approx \mathbf{x}.$

as a floating-point approximation of a real number x.

Approximations of real numbers (cont.)

Example

The floating-point number 1.10011001100110011001101 \times 2⁻⁴ can be used to approximate 1.0 \times 10⁻¹. The best single precision approximation of decimal 0.1.

 1.0×10^{-1} is not representable in binary.

Approximations of real numbers (cont.)

Example

The floating-point number 1.10011001100110011001101 \times 2⁻⁴ can be used to approximate 1.0 \times 10⁻¹. The best single precision approximation of decimal 0.1.

 1.0×10^{-1} is not representable in binary.

When approximating, some kind of rounding is involved.

Error measurements: ulp and *u*

If the nearest rounding is applied and $fl(x) = d_1.d_2...d_t \times \beta^e$, then the absolute error is bounded by

$$|\mathrm{fl}(\mathbf{x}) - \mathbf{x}| \leq rac{1}{2} eta^{1-t} eta^{\mathsf{e}},$$

half of the unit in the last place (ulp);

Error measurements: ulp and u

If the nearest rounding is applied and $fl(x) = d_1.d_2...d_t \times \beta^e$, then the absolute error is bounded by

$$|\mathrm{fl}(\mathbf{x}) - \mathbf{x}| \leq \frac{1}{2} \beta^{1-t} \beta^{\mathsf{e}},$$

half of the unit in the last place (ulp);

the relative error is bounded by

$$rac{|\mathrm{fl}(x)-x|}{|\mathrm{fl}(x)|} \leq rac{1}{2}eta^{1-t}, ext{ since } |\mathrm{fl}(x)| \geq 1.0 imes eta^{\mathrm{e}},$$

called the *unit of roundoff* denoted by *u*.

Unit of roundoff u

When $\beta = 2$, $u = 2^{-t}$.



Unit of roundoff u

When $\beta = 2$, $u = 2^{-t}$. How would you compute *u*?



Unit of roundoff u

When $\beta = 2$, $u = 2^{-t}$. How would you compute *u*?

The largest number such that 1.0 + u = 1.0. Also, when $\beta = 2$, the distance between two consecutive floating-point numbers between 1/2 and 1.0 $(1.0...0 \times 2^{-1}, ..., 1.1...1 \times 2^{-1}, 1.0.)$ $1.0 + 2^{-t} = 1.0$ (Why?)

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

Unit of roundoff u

When $\beta = 2$, $u = 2^{-t}$. How would you compute *u*?

The largest number such that 1.0 + u = 1.0. Also, when $\beta = 2$, the distance between two consecutive floating-point numbers between 1/2 and 1.0 $(1.0...0 \times 2^{-1}, ..., 1.1...1 \times 2^{-1}, 1.0.)$ $1.0 + 2^{-t} = 1.0$ (Why?)

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

Four parameters

Base $\beta = 2$.

	single	double
precision t	24	53
e _{min}	-126	-1022
e _{max}	127	1023

Formats:

	single	double
Exponent width	8 bits	11 bits
Format width in bits	32 bits	64 bits

 $x \neq y \Rightarrow 1/x \neq 1/y$?

How many single precision floating-point numbers in [1,2)?

 $x \neq y \Rightarrow 1/x \neq 1/y$?

How many single precision floating-point numbers in [1,2)?

$1.00...00 \to 1.11...11$

2²³, evenly spaced.

 $x \neq y \Rightarrow 1/x \neq 1/y$?

How many single precision floating-point numbers in [1,2)?

 $1.00...00 \to 1.11...11$

2²³, evenly spaced.

How many single precision floating-point numbers in (1/2, 1]?

$$1.00...01 \times 2^{-1} \to 1.00...00$$

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

2²³, evenly spaced.

Floating-point Numbers Sources of Errors Stability of an Algorithm

Sensitivity of a Problem Fallacies Summary

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

$x \neq y \Rightarrow 1/x \neq 1/y$? (cont.)

How many single precision floating-point numbers in [3/2, 2)? $(1/2)\times 2^{23}$

$$x \neq y \Rightarrow 1/x \neq 1/y$$
? (cont.)

How many single precision floating-point numbers in [3/2,2)? (1/2) $\times 2^{23}$ How many single precision floating-point numbers in (1/2,2/3]?

 $(1/3) \times 2^{23}$.

$$x \neq y \Rightarrow 1/x \neq 1/y$$
? (cont.)

How many single precision floating-point numbers in [3/2, 2)? $(1/2)\times 2^{23}$

How many single precision floating-point numbers in (1/2, 2/3]?

 $(1/3) \times 2^{23}$.

Since $(1/2) \times 2^{23} > (1/3) \times 2^{23}$, there exist $x \neq y \in [3/2, 2)$ such that $1/x = 1/y \in (1/2, 2/3]$.

◆□▶ ◆□▶ ◆□▶ ◆□▶ ◆□▶ ◆□

Hidden bit and biased representation

Since the base is 2 (binary), the integer bit is always 1. This bit is not stored and called *hidden bit*.

Hidden bit and biased representation

Since the base is 2 (binary), the integer bit is always 1. This bit is not stored and called *hidden bit*.

The exponent is stored using the biased representation. In single precision, the bias is 127. In double precision, the bias is 1023.
Hidden bit and biased representation

Since the base is 2 (binary), the integer bit is always 1. This bit is not stored and called *hidden bit*.

The exponent is stored using the biased representation. In single precision, the bias is 127. In double precision, the bias is 1023.

Example Single precision 1.1001100110011001101 \times 2⁻⁴ is stored as

0 01111011 10011001100110011001101

Special quantities

The special quantities are encoded with exponents of either $e_{max} + 1$ or $e_{min} - 1$. In single precision, 11111111 in the exponent field encodes $e_{max} + 1$ and 00000000 in the exponent field encodes $e_{min} - 1$.

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

Special quantities

The special quantities are encoded with exponents of either $e_{max} + 1$ or $e_{min} - 1$. In single precision, 11111111 in the exponent field encodes $e_{max} + 1$ and 00000000 in the exponent field encodes $e_{min} - 1$.

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

Signed zeros: ± 0

Binary representation:



When testing for equal, +0 = -0, so the simple test if (x == 0) is predictable whether x is +0 or -0.

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

- The relation 1/(1/x) = x holds when $x = \pm \infty$.
- $\log(+0) = -\infty$ and $\log(-0) = NaN$; sign(+0) = 1 and sign(-0) = -1.

Signed zeros

If
$$z = -1$$
, $\sqrt{1/z} = i$, but $1/\sqrt{z} = -i$.
 $\sqrt{1/z} \neq 1/\sqrt{z}!$



Signed zeros

If
$$z = -1$$
, $\sqrt{1/z} = i$, but $1/\sqrt{z} = -i$.
 $\sqrt{1/z} \neq 1/\sqrt{z}!$

Why? Square root is multivalued, can't make it continuous in the entire complex plane. However, it is continuous for $z = \cos \theta + i \sin \theta$, $-\pi \le \theta \le \pi$, if a brabch cut consisting of all negative real numbers is excluded from the consideration. With signed zeros, for the numbers with negative real part, -x + i(+0), x > 0, has a square root of $i\sqrt{x}$; -x + i(-0) has a square root of $-i\sqrt{x}$.

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

Signed zeros

If
$$z = -1$$
, $\sqrt{1/z} = i$, but $1/\sqrt{z} = -i$.
 $\sqrt{1/z} \neq 1/\sqrt{z}!$

Why? Square root is multivalued, can't make it continuous in the entire complex plane. However, it is continuous for $z = \cos \theta + i \sin \theta$, $-\pi \le \theta \le \pi$, if a brabch cut consisting of all negative real numbers is excluded from the consideration. With signed zeros, for the numbers with negative real part, -x + i(+0), x > 0, has a square root of $i\sqrt{x}$; -x + i(-0) has a square root of $-i\sqrt{x}$.

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

$$z = -1 = -1 + i(+0), 1/z = -1 + i(-0)$$
, then
 $\sqrt{1/z} = -i = 1/\sqrt{z}$

Signed zeros

If
$$z = -1$$
, $\sqrt{1/z} = i$, but $1/\sqrt{z} = -i$.
 $\sqrt{1/z} \neq 1/\sqrt{z}!$

Why? Square root is multivalued, can't make it continuous in the entire complex plane. However, it is continuous for $z = \cos \theta + i \sin \theta$, $-\pi \le \theta \le \pi$, if a brabch cut consisting of all negative real numbers is excluded from the consideration. With signed zeros, for the numbers with negative real part, -x + i(+0), x > 0, has a square root of $i\sqrt{x}$; -x + i(-0) has a square root of $-i\sqrt{x}$.

$$z = -1 = -1 + i(+0), 1/z = -1 + i(-0)$$
, then
 $\sqrt{1/z} = -i = 1/\sqrt{z}$

However, +0 = -0, and $1/(+0) \neq 1/(-0)$. (Shortcoming)

Infinities

Infinities

Infinities: $\pm\infty$

Binary Representation:

x 11111111 000000000000000000000000

- Provide a way to continue when exponent gets too large, $x^2 = \infty$, when x^2 overflows.
- When $c \neq 0$, $c/0 = \pm \infty$.
- Avoid special case checking, 1/(x + 1/x), a better formula for $x/(x^2 + 1)$, with infinities, there is no need for checking the special case x = 0.

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

NaN

NaNs (not a number)

Binary representation:

X 11111111 nonzero fraction



NaN

NaNs (not a number)

Binary representation:

X 11111111 nonzero fraction

Provide a way to continue in situations like

Operation	NaN Produced By
+	$\infty + (-\infty)$
*	$0 * \infty$
/	0/0, ∞/∞
REM	$x \text{ REM 0}, \infty \text{ REM } y$
sqrt	sqrt(x) when $x < 0$

Example for NaN

The function zero(f) returns a zero of a given quadratic polynomial f.

lf

$$f = x^2 + x + 1,$$

 $d = 1 - 4 < 0$, thus $\sqrt{d} = NaN$ and
 $rac{-b \pm \sqrt{d}}{2a} = NaN,$

no zeros.

Denormalized numbers

Denormalized Numbers

Binary representation:

X 00000000 nonzero fraction

- * ロ * * @ * * 注 * 注 * の < @

Denormalized numbers

Denormalized Numbers

Binary representation:

X 00000000 nonzero fraction

When $e = e_{\min} - 1$ and the bits in the fraction are $b_2, b_3, ..., b_t$, the number being represented is $0.b_2b_3...b_t \times 2^{e+1}$ (no hidden bit)

Denormalized numbers

Denormalized Numbers

Binary representation:

X 00000000 nonzero fraction

When $e = e_{\min} - 1$ and the bits in the fraction are $b_2, b_3, ..., b_t$, the number being represented is $0.b_2b_3...b_t \times 2^{e+1}$ (no hidden bit)

- Guarantee the relation: $x = y \iff x y = 0$
- Allow gradual underflow. Without denormals, the spacing abruptly changes from $\beta^{-t+1}\beta^{e_{\min}}$ to $\beta^{e_{\min}}$, which is a factor of β^{t-1} .

Complex division

$$\frac{a+ib}{c+id} = \frac{ac+bd}{c^2+d^2} + i\frac{bc-ad}{c^2+d^2}.$$

▲□▶▲圖▶▲≣▶▲≣▶ ≣ のQ@

Complex division

$$\frac{a+ib}{c+id}=\frac{ac+bd}{c^2+d^2}+i\frac{bc-ad}{c^2+d^2}.$$

Underflows when *a*, *b*, *c*, and *d* are small.

▲ロ▶▲母▶▲目▶▲目▶ 目 のへで

Smith's formula

$$\frac{a+b(d/c)}{c+d(d/c)} + i\frac{b-a(d/c)}{c+d(d/c)} \quad \text{if } |d| < |c|$$
$$\frac{b+a(c/d)}{d+c(c/d)} + i\frac{-a+b(c/d)}{d+c(c/d)} \quad \text{if } |d| \ge |c|$$

|▲□▶ ▲圖▶ ▲≣▶ ▲≣▶ = 通 = のへで

Example for denormalized numbers

Smith's formula

$$\frac{a+b(d/c)}{c+d(d/c)} + i\frac{b-a(d/c)}{c+d(d/c)} \quad \text{if } |d| < |c|$$
$$\frac{b+a(c/d)}{d+c(c/d)} + i\frac{-a+b(c/d)}{d+c(c/d)} \quad \text{if } |d| \ge |c|$$

For $a = 2\beta^{e_{\min}}$, $b = \beta^{e_{\min}}$, $c = 4\beta^{e_{\min}}$, and $d = 2\beta^{e_{\min}}$, the result is 0.5 with denormals $(a + b(d/c) = 2.5\beta^{e_{\min}})$ or 0.4 without denormals $(a + b(d/c) = 2\beta^{e_{\min}})$.

Smith's formula

$$rac{a+b(d/c)}{c+d(d/c)}+irac{b-a(d/c)}{c+d(d/c)}$$
 if $|d|<|c|$

$$rac{b+a(c/d)}{d+c(c/d)}+irac{-a+b(c/d)}{d+c(c/d)}$$
 if $|d|\geq |c|$

For $a = 2\beta^{e_{\min}}$, $b = \beta^{e_{\min}}$, $c = 4\beta^{e_{\min}}$, and $d = 2\beta^{e_{\min}}$, the result is 0.5 with denormals $(a + b(d/c) = 2.5\beta^{e_{\min}})$ or 0.4 without denormals $(a + b(d/c) = 2\beta^{e_{\min}})$.

It is typical for denormalized numbers to guarantee error bounds for arguments all the way down to $\beta^{e_{\min}}$.

IEEE floating-point representations

Exponent	Fraction	Represents
$e = e_{min} - 1$	<i>f</i> = 0	±0
$e = e_{min} - 1$	f eq 0	$0.f imes 2^{e_{\min}}$
$\mathbf{e}_{min} \leq \mathbf{e} \leq \mathbf{e}_{max}$		$1.f imes 2^{e}$
$e = e_{max} + 1$	<i>f</i> = 0	$\pm\infty$
$e = e_{max} + 1$	f eq 0	NaN

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Examples (IEEE single precision)

Underflow

An arithmetic operation produces a number with an exponent that is too small to be represented in the system.

Example. In single precision,

$$a = 3.0 \times 10^{-30},$$

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

a * *a* underflows. By default, it is set to zero.



An arithmetic operation produces a number with an exponent that is too large to be represented in the system.

Example. In single precision,

$$a=3.0\times10^{30},$$

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

a * a overflows. In IEEE standard, the default result is ∞ .

Avoiding unnecessary underflow and overflow

Sometimes, underflow and overflow can be avoided by using a technique called scaling.

Avoiding unnecessary underflow and overflow

Sometimes, underflow and overflow can be avoided by using a technique called scaling. Given $x = (a, b)^T$, $a = 1.0 \times 10^{30}$, b = 1.0, compute $c = ||x||_2 = \sqrt{a^2 + b^2}$.

Avoiding unnecessary underflow and overflow

Sometimes, underflow and overflow can be avoided by using a technique called scaling.

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

Given
$$x = (a, b)^{T}$$
, $a = 1.0 \times 10^{30}$, $b = 1.0$, compute
 $c = ||x||_{2} = \sqrt{a^{2} + b^{2}}$.
scaling: $s = \max\{|a|, |b|\} = 1.0 \times 10^{30}$
 $a \leftarrow a/s (1.0)$,
 $b \leftarrow b/s (1.0 \times 10^{-30})$
 $t = \sqrt{a * a + b * b} (1.0)$
 $c \leftarrow t * s (1.0 \times 10^{30})$

Example: Computing 2-norm of a vector

Compute

$$\sqrt{x_1^2 + x_2^2 + \ldots + x_n^2}$$

Example: Computing 2-norm of a vector

Compute

$$\sqrt{x_1^2 + x_2^2 + ... + x_n^2}$$

Summary

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

Efficient and robust:

Avoid multiple loops:

searching for the largest; Scaling; Summing.

Example: Computing 2-norm of a vector

Compute

$$\sqrt{x_1^2 + x_2^2 + \ldots + x_n^2}$$

Summary

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

Efficient and robust:

Avoid multiple loops:

searching for the largest; Scaling; Summing.

Result: One single loop

Technique: Dynamic scaling

Summary

◆□ > ◆□ > ◆三 > ◆三 > ・三 のへで

Example: Computing 2-norm of a vector

```
scale = 0.0i
ssq = 1.0;
for i=1 to n
   if (x(i) != 0.0)
      if (scale<abs(x(i))
         tmp = scale/x(i);
         ssq = 1.0 + ssq*tmp*tmp;
         scale = abs(x(i));
      else
         tmp = x(i)/scale;
         ssq = ssq + tmp*tmp;
      end
   end
end
nrm2 = scale*sqrt(ssq);
```

◆□▶ ◆□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Correctly rounded operations

Correctly rounded means that result must be the same as if it were computed exactly and then rounded, usually to the nearest floating-point number. For example, if \oplus denotes the floating-point addition, then given two floating-point numbers *a* and *b*,

$$a \oplus b = \mathrm{fl}(a+b).$$

Correctly rounded operations

Correctly rounded means that result must be the same as if it were computed exactly and then rounded, usually to the nearest floating-point number. For example, if \oplus denotes the floating-point addition, then given two floating-point numbers *a* and *b*,

$$a \oplus b = \mathrm{fl}(a+b).$$

Example $\beta = 10, t = 4$ $a = 1.234 \times 10^{0}$ and $b = 5.678 \times 10^{-3}$

Correctly rounded operations

Correctly rounded means that result must be the same as if it were computed exactly and then rounded, usually to the nearest floating-point number. For example, if \oplus denotes the floating-point addition, then given two floating-point numbers *a* and *b*,

$$a \oplus b = \mathrm{fl}(a+b).$$

Example $\beta = 10, t = 4$ $a = 1.234 \times 10^{0}$ and $b = 5.678 \times 10^{-3}$ Exact: a + b = 1.239678

Correctly rounded operations

Correctly rounded means that result must be the same as if it were computed exactly and then rounded, usually to the nearest floating-point number. For example, if \oplus denotes the floating-point addition, then given two floating-point numbers *a* and *b*,

$$a \oplus b = \mathrm{fl}(a+b).$$

Example $\beta = 10, t = 4$ $a = 1.234 \times 10^{0}$ and $b = 5.678 \times 10^{-3}$ Exact: a + b = 1.239678Floating-point: fl $(a + b) = 1.240 \times 10^{0}$
◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

Correctly rounded operations

IEEE standards require the following operations are correctly rounded:

- arithmetic operations +, -, *, and /
- square root and remainder
- conversions of formats (binary, decimal)

・ ロ ト ・ 『 ト ・ 日 ト ・ 日 ト

Outline

- Floating-point Numbers
 - Representations
 - IEEE Floating-point Standards
 - Underflow and Overflow
 - Correctly Rounded Operations

2 Sources of Errors

- Rounding Error
- Truncation Error
- Discretization Error
- 3 Stability of an Algorithm
- 4 Sensitiviy of a Problem

5 Fallacies

Rounding error

Due to finite precision arithmetic, a computed result must be rounded to fit storage format.

Example $\beta = 10, p = 4 (u = 0.5 \times 10^{-3})$ $a = 1.234 \times 10^{0}, b = 5.678 \times 10^{-3}$ $x = a + b = 1.239678 \times 10^{0}$ (exact) $\hat{x} = fl(a + b) = 1.240 \times 10^{0}$

the result was rounded to the nearest computer number.

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

Rounding error

Due to finite precision arithmetic, a computed result must be rounded to fit storage format.

Example $\beta = 10, p = 4 (u = 0.5 \times 10^{-3})$ $a = 1.234 \times 10^{0}, b = 5.678 \times 10^{-3}$ $x = a + b = 1.239678 \times 10^{0}$ (exact) $\hat{x} = fl(a + b) = 1.240 \times 10^{0}$ the result was rounded to the nearest computer number.

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

Rounding error: $fl(a + b) = (a + b)(1 + \epsilon), |\epsilon| \le u$.

Rounding error

Due to finite precision arithmetic, a computed result must be rounded to fit storage format.

Example $\beta = 10, p = 4 (u = 0.5 \times 10^{-3})$ $a = 1.234 \times 10^{0}, b = 5.678 \times 10^{-3}$ $x = a + b = 1.239678 \times 10^{0}$ (exact) $\hat{x} = fl(a + b) = 1.240 \times 10^{0}$ the result was rounded to the nearest computer number. Rounding error: $fl(a + b) = (a + b)(1 + \epsilon), |\epsilon| \le u$. $1.240 = 1.239678(1 + 2.59... \times 10^{-4}), |2.59... \times 10^{-4}| < u$

Effect of rounding errors

Top:
$$y = (x - 1)^6$$

Bottom: $y = x^6 - 6x^5 + 15x^4 - 20x^3 + 15x^2 - 6x + 7$



Two ways of evaluating the polynomial $(x - 1)^6$

・ロト・日本・日本・日本・日本・日本

Summary

Real to floating-point

double x = 0.1;

What is the value of x stored?



Real to floating-point

double x = 0.1;

What is the value of x stored?

 $1.0\times 10^{-1} = 1.100110011001100110011...\times 2^{-4}$

Real to floating-point

double x = 0.1;

What is the value of x stored?

 $1.0\times 10^{-1} = 1.100110011001100110011...\times 2^{-4}$

Decimal 0.1 cannot be exactly represented in binary. It must be rounded to

 $\begin{array}{rl} 1.10011001100...110011010 \times 2^{-4} \\ > & 1.10011001100...11001100110011... \end{array}$

slightly larger than 0.1.

Real to floating-point

double x, y, h; x = 1/2; h = 0.1; for i=1 to 5 x = x + h; end y = 1.0 - x;

y > 0 or y < 0 or y = 0?

Real to floating-point

double x, y, h; x = 1/2; h = 0.1; for i=1 to 5 x = x + h; end y = 1.0 - x;

$$y > 0$$
 or $y < 0$ or $y = 0$?
Answer: $y \approx 1.1 \times 10^{-16} > 0$

◆□▶ ◆□▶ ◆ □▶ ◆ □▶ → □ ● ● ● ●

Real to floating-point (cont.)

Why?

$$\begin{array}{rcl} 0.5 & = & 1.0000000...00 \times 2^{-1} \\ h & = & 0.00110011...11010 \times 2^{-1} \end{array}$$

▲□▶▲□▶▲□▶▲□▶ □ のへで

Floating-point Numbers Sources of Errors Stability of an Algorithm Sensitiviy of a Problem

Fallacies Summary

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Real to floating-point (cont.)

Why?

$$\begin{array}{rcl} 0.5 & = & 1.0000000...00 \times 2^{-1} \\ h & = & 0.00110011...11010 \times 2^{-1} \end{array}$$

Rounding errors in floating-point addition.

Integer to floating-point

Fallacy

Java converts an integer into its mathematically equivalent floating-point number.

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

```
long k = 1801439850948199; \\
long d = k - (long)((double) k);
```

Note $1801439850948199 = 2^{54} + 1$

d = 0?

Integer to floating-point

Fallacy

Java converts an integer into its mathematically equivalent floating-point number.

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

```
long k = 1801439850948199; \\
long d = k - (long)((double) k);
```

Note $1801439850948199 = 2^{54} + 1$

d = 0?

No, d = 1!

Integer to floating-point

Why?

|▲□▶ ▲圖▶ ▲≣▶ ▲≣▶ = 通 = のへで

Integer to floating-point

Why?

$k = 1.00...0001 \times 2^{54}$ (double) $k = 1.00...00 \times 2^{54}$

▲□▶▲□▶▲□▶▲□▶ □ ● ● ●

Truncation error

When an infinite series is approximated by a finite sum, truncation error is introduced.

Example. If we use

$$1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots + \frac{x^n}{n!}$$

to approximate

$$e^{x} = 1 + x + \frac{x^{2}}{2!} + \frac{x^{3}}{3!} + \dots + \frac{x^{n}}{n!} + \dots,$$

then the truncation error is

$$\frac{x^{n+1}}{(n+1)!} + \frac{x^{n+2}}{(n+2)!} + \cdots$$

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

Discretization error

When a continuous problem is approximated by a discrete one, discretization error is introduced. Example. From the expansion

$$f(x+h) = f(x) + hf'(x) + \frac{h^2}{2!}f''(\xi),$$

for some $\xi \in [x, x + h]$, we can use the following approximation:

$$y_h(x) = rac{f(x+h)-f(x)}{h} \approx f'(x).$$

The discretization error is $E_{dis} = |f''(\xi)|h/2$.

Example

Let
$$f(x) = e^x$$
, compute $y_h(1)$.



Example

Let $f(x) = e^x$, compute $y_h(1)$. The discretization error is

$$\mathcal{E}_{ ext{dis}} = rac{h}{2} |f''(\xi)| \leq rac{h}{2} \mathrm{e}^{1+h} pprox rac{h}{2} \mathrm{e}^{-1}$$
 for small h .

Example

Let $f(x) = e^x$, compute $y_h(1)$. The discretization error is

$$E_{
m dis}=rac{h}{2}|f''(\xi)|\leq rac{h}{2}{
m e}^{1+h}pproxrac{h}{2}{
m e}~~{
m for~small}~~h.$$

The computed $y_h(1)$:

$$\widehat{y}_h(1) = \frac{(e^{(1+h)(1+\epsilon_1)}(1+\epsilon_2) - e(1+\epsilon_3))(1+\epsilon_4)}{h}(1+\epsilon_5),$$

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

 $|\epsilon_i| \leq u.$

Example

Let $f(x) = e^x$, compute $y_h(1)$. The discretization error is

$$E_{
m dis}=rac{h}{2}|f''(\xi)|\leq rac{h}{2}{
m e}^{1+h}pproxrac{h}{2}{
m e}~~{
m for~small}~~h.$$

The computed $y_h(1)$:

$$\widehat{y}_h(1) = \frac{(e^{(1+h)(1+\epsilon_1)}(1+\epsilon_2) - e(1+\epsilon_3))(1+\epsilon_4)}{h}(1+\epsilon_5),$$

 $|\epsilon_i| \leq u$. The rounding error is

$$egin{aligned} E_{ ext{round}} &= \widehat{y}_h(1) - y_h(1) pprox rac{7u}{h} e. \end{aligned}$$

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

Example (cont.)

The total error:

$$E_{\text{total}} = E_{\text{dis}} + E_{\text{round}} \approx \left(\frac{h}{2} + \frac{7u}{h}\right) e.$$

$$\int_{\frac{1}{2}} \frac{1}{\sqrt{1-\frac{1}{2}}} \frac{1}$$

Example (cont.)

The total error:

$$E_{\text{total}} = E_{\text{dis}} + E_{\text{round}} \approx \left(\frac{h}{2} + \frac{7u}{h}\right) e.$$

$$\int_{a}^{a} \int_{a}^{a} \int_{$$

The optimal *h*: $h_{\rm opt} = \sqrt{12u} \approx \sqrt{u}$.

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

Outline

- Floating-point Numbers
 - Representations
 - IEEE Floating-point Standards
 - Underflow and Overflow
 - Correctly Rounded Operations
- 2 Sources of Errors
 - Rounding Error
 - Truncation Error
 - Discretization Error
- Stability of an Algorithm
 - Sensitiviy of a Problem
 - 5 Fallacies

Backward errors

Recall that

$$a \oplus b = \mathrm{fl}(a+b) = (a+b)(1+\eta), \quad |\eta| \le u$$

Backward errors

Recall that

$$\boldsymbol{a} \oplus \boldsymbol{b} = \mathrm{fl}(\boldsymbol{a} + \boldsymbol{b}) = (\boldsymbol{a} + \boldsymbol{b})(1 + \eta), \quad |\eta| \leq u$$

In other words,

$$m{a} \oplus m{b} = ilde{m{a}} + ilde{m{b}}$$

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

where $\tilde{a} = a(1 + \eta)$ and $\tilde{b} = b(1 + \eta)$, for $|\eta| \le u$, are slightly different from *a* and *b* respectively.

Backward errors

Recall that

$$\boldsymbol{a} \oplus \boldsymbol{b} = \mathrm{fl}(\boldsymbol{a} + \boldsymbol{b}) = (\boldsymbol{a} + \boldsymbol{b})(1 + \eta), \quad |\eta| \leq u$$

In other words,

$$\pmb{a} \oplus \pmb{b} = ilde{\pmb{a}} + ilde{\pmb{b}}$$

where $\tilde{a} = a(1 + \eta)$ and $\tilde{b} = b(1 + \eta)$, for $|\eta| \le u$, are slightly different from *a* and *b* respectively.

The computed sum (result) is the exact sum of slightly different *a* and *b* (inputs).

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

Example

$$\begin{split} \beta &= 10, \, p = 4 \, (u = 0.5 \times 10^{-3}) \\ a &= 1.234 \times 10^{0}, \, b = 5.678 \times 10^{-3} \\ a \oplus b &= 1.240 \times 10^{0}, \, a + b = 1.239678 \\ 1.240 &= 1.239678(1 + 2.59... \times 10^{-4}), \, |2.59... \times 10^{-4}| < u \\ 1.240 &= a(1 + 2.59... \times 10^{-4}) + b(1 + 2.59... \times 10^{-4}) \end{split}$$

(日)

æ

Example

$$\begin{split} \beta &= 10, \, p = 4 \, (u = 0.5 \times 10^{-3}) \\ a &= 1.234 \times 10^{0}, \, b = 5.678 \times 10^{-3} \\ a \oplus b &= 1.240 \times 10^{0}, \, a + b = 1.239678 \\ 1.240 &= 1.239678(1 + 2.59... \times 10^{-4}), \, |2.59... \times 10^{-4}| < u \\ 1.240 &= a(1 + 2.59... \times 10^{-4}) + b(1 + 2.59... \times 10^{-4}) \end{split}$$

The computed sum (result) is the exact sum of slightly different *a* and *b* (inputs).

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

Backward errors (cont.)

A general example

$$\mathbf{s}_n = \mathbf{x}_1 \oplus \mathbf{x}_2 \oplus \cdots \oplus \mathbf{x}_n$$

The computed result $(x_1 \oplus \cdots \oplus x_n)$ is the exact result of the problem with slightly perturbed data. $(x_1(1 + \eta_1), ..., x_n(1 + \eta_n))$. Backward errors:

$$egin{aligned} &|\eta_1| \leq 1.06(n-1)u \ &|\eta_i| \leq 1.06(n-i+1)u, \, i=2,3,...,n \end{aligned}$$

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

Backward errors (cont.)

A general example

$$\mathbf{s}_n = \mathbf{x}_1 \oplus \mathbf{x}_2 \oplus \cdots \oplus \mathbf{x}_n$$

The computed result $(x_1 \oplus \cdots \oplus x_n)$ is the exact result of the problem with slightly perturbed data. $(x_1(1 + \eta_1), ..., x_n(1 + \eta_n))$. Backward errors:

$$|\eta_1| \le 1.06(n-1)u$$

 $|\eta_i| \le 1.06(n-i+1)u, i = 2, 3, ..., n$

If the backward errors are small, then we say that the algorithm is backward stable.

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

Outline

- Floating-point Numbers
 - Representations
 - IEEE Floating-point Standards
 - Underflow and Overflow
 - Correctly Rounded Operations
- 2 Sources of Errors
 - Rounding Error
 - Truncation Error
 - Discretization Error
- 3 Stability of an Algorithm
 - Sensitiviy of a Problem
- 5 Fallacies

Introduction

Example: a + ba = 1.23, b = 0.45, s = a + b = 1.68

Introduction

Example: a + b a = 1.23, b = 0.45, s = a + b = 1.68Slightly perturbed $\hat{a} = a(1 + 0.01), \hat{b} = b(1 + 0.001), \hat{s} = \hat{a} + \hat{b} = 1.69275$

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで
Example: a + b a = 1.23, b = 0.45, s = a + b = 1.68Slightly perturbed $\hat{a} = a(1 + 0.01), \hat{b} = b(1 + 0.001), \hat{s} = \hat{a} + \hat{b} = 1.69275$

Relative perturbations in data (a and b) are at most 0.01.

Example: a + ba = 1.23, b = 0.45, s = a + b = 1.68Slightly perturbed

 $\hat{a} = a(1+0.01), \hat{b} = b(1+0.001), \hat{s} = \hat{a} + \hat{b} = 1.69275$

Relative perturbations in data (*a* and *b*) are at most 0.01. Causing a relative change in the result $|\hat{s} - s|/|s| \approx 0.0076$,

Example: a + ba = 1.23, b = 0.45, s = a + b = 1.68Slightly perturbed

 $\hat{a} = a(1+0.01), \hat{b} = b(1+0.001), \hat{s} = \hat{a} + \hat{b} = 1.69275$

Relative perturbations in data (*a* and *b*) are at most 0.01. Causing a relative change in the result $|\hat{s} - s|/|s| \approx 0.0076$, which is about the same as the perturbation 0.01

Example: a + b a = 1.23, b = 0.45, s = a + b = 1.68Slightly perturbed $\hat{a} = a(1 + 0.01), \hat{b} = b(1 + 0.001), \hat{s} = \hat{a} + \hat{b} = 1.69275$

Relative perturbations in data (*a* and *b*) are at most 0.01. Causing a relative change in the result $|\hat{s} - s|/|s| \approx 0.0076$, which is about the same as the perturbation 0.01

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

The result is insensitive to the perturbation in data.

Fallacies Summary

Introduction

a = 1.23, b = -1.21, s = a + b = 0.02

Floating-point Numbers Sources of Errors Stability of an Algorithm Sensitivity of a Problem Fallacies Summary

Introduction

$$a = 1.23, b = -1.21, s = a + b = 0.02$$

Slightly perturbed
 $\widehat{a} = a(1 + 0.01), \widehat{b} = b(1 + 0.001), \widehat{s} = \widehat{a} + \widehat{b} = 0.03109$

◆□▶ ◆□▶ ◆巨▶ ◆巨▶

∃ 990

Floating-point Numbers Sources of Errors Stability of an Algorithm Sensitivity of a Problem Fallacies Summary

Introduction

$$a = 1.23, b = -1.21, s = a + b = 0.02$$

Slightly perturbed
 $\hat{a} = a(1 + 0.01), \hat{b} = b(1 + 0.001), \hat{s} = \hat{a} + \hat{b} = 0.03109$

Relative perturbations in data (a and b) are at most 0.01.

$$a = 1.23, b = -1.21, s = a + b = 0.02$$

Slightly perturbed $\hat{a} = a(1 + 0.01), \ \hat{b} = b(1 + 0.001), \ \hat{s} = \hat{a} + \hat{b} = 0.03109$

Relative perturbations in data (*a* and *b*) are at most 0.01. Causing a relative change in the result $|\hat{s} - s|/|s| \approx 0.5545$,

$$a = 1.23, b = -1.21, s = a + b = 0.02$$

Slightly perturbed $\widehat{a} = a(1+0.01), \ \widehat{b} = b(1+0.001), \ \widehat{s} = \widehat{a} + \widehat{b} = 0.03109$

Relative perturbations in data (*a* and *b*) are at most 0.01. Causing a relative change in the result $|\hat{s} - s|/|s| \approx 0.5545$, which is more than 55 times as the perturbation 0.01

$$a = 1.23, b = -1.21, s = a + b = 0.02$$

Slightly perturbed $\hat{a} = a(1 + 0.01), \ \hat{b} = b(1 + 0.001), \ \hat{s} = \hat{a} + \hat{b} = 0.03109$

Relative perturbations in data (*a* and *b*) are at most 0.01. Causing a relative change in the result $|\hat{s} - s|/|s| \approx 0.5545$, which is more than 55 times as the perturbation 0.01

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

The result is sensitive to the perturbation in the data.

Perturbation analysis

Example: *a* + *b*

$$\begin{array}{l} \displaystyle \frac{|\boldsymbol{a}(\boldsymbol{1}+\delta_{\boldsymbol{a}})+\boldsymbol{b}(\boldsymbol{1}+\delta_{\boldsymbol{b}})-(\boldsymbol{a}+\boldsymbol{b})|}{|\boldsymbol{a}+\boldsymbol{b}|} \\ \leq \quad \displaystyle \frac{|\boldsymbol{a}|+|\boldsymbol{b}|}{|\boldsymbol{a}+\boldsymbol{b}|}\,\delta, \quad \delta=\max(\delta_{\boldsymbol{a}},\delta_{\boldsymbol{b}}). \end{array}$$

Perturbation analysis

Example: a + b

$$\frac{|\boldsymbol{a}(1+\delta_{\boldsymbol{a}})+\boldsymbol{b}(1+\delta_{\boldsymbol{b}})-(\boldsymbol{a}+\boldsymbol{b})|}{|\boldsymbol{a}+\boldsymbol{b}|} \\ \leq \frac{|\boldsymbol{a}|+|\boldsymbol{b}|}{|\boldsymbol{a}+\boldsymbol{b}|} \delta, \quad \delta = \max(\delta_{\boldsymbol{a}},\delta_{\boldsymbol{b}}).$$

Condition number: (|a| + |b|)/|a + b|, magnification of the relative error.

$$\frac{\text{relative error in result}}{\text{relative error in data}} \le \text{cond}$$

Condition number is a measurement (an upper bound) of the sensitivity of the problem to changes in data.





Two methods for calculating z(x + y):

 $z \otimes x \oplus z \otimes y$ and $z \otimes (x \oplus y)$



Example

Two methods for calculating z(x + y):

```
z \otimes x \oplus z \otimes y and z \otimes (x \oplus y)
```

$$\beta = 10, t = 4$$

 $x = 1.002, y = -0.9958, z = 3.456$
Exact $z(x + y) = 2.14272 \times 10^{-2}$
 $z \otimes (x \oplus y) = \text{fl}(3.456 * 6.200 \times 10^{-3})$
 $= 2.143 \times 10^{-2}$
error: 2.8×10^{-6}

Example

Two methods for calculating z(x + y):

```
z \otimes x \oplus z \otimes y and z \otimes (x \oplus y)
```

$$\beta = 10, t = 4$$

 $x = 1.002, y = -0.9958, z = 3.456$
Exact $z(x + y) = 2.14272 \times 10^{-2}$
 $z \otimes (x \oplus y) = fl(3.456 * 6.200 \times 10^{-3})$
 $= 2.143 \times 10^{-2}$
error: 2.8×10^{-6}
 $(z \otimes x) \oplus (z \otimes y)) = fl(3.463 - 3.441)$
 $= 2.200 \times 10^{-2}$
error: 5.7×10^{-4}
More than 200 times!

Backward error analyses

 $z \otimes x \oplus z \otimes y$

$$= (zx(1+\epsilon_1)+zy(1+\epsilon_2))(1+\epsilon_3)$$

$$= z(1+\epsilon_3)(x(1+\epsilon_1)+y(1+\epsilon_2)), \ |\epsilon_i| \leq u$$

▲ロト▲御▶▲臣⊁▲臣⊁ ―臣 - のへで

Backward error analyses

 $z \otimes x \oplus z \otimes y$

$$= (\mathbf{z}\mathbf{x}(1+\epsilon_1)+\mathbf{z}\mathbf{y}(1+\epsilon_2))(1+\epsilon_3)$$

$$= z(1+\epsilon_3)(x(1+\epsilon_1)+y(1+\epsilon_2)), \ |\epsilon_i| \leq u$$

$$z \otimes (x \oplus y)$$

= $z((x + y)(1 + \epsilon_1))(1 + \epsilon_3)$
= $z(1 + \epsilon_3)(x(1 + \epsilon_1) + y(1 + \epsilon_1)), |\epsilon_i| \le u$

◆ロ〉 ◆御〉 ◆臣〉 ◆臣〉 「臣」 のへで

Backward error analyses

 $z \otimes x \oplus z \otimes y$

- $= (zx(1+\epsilon_1)+zy(1+\epsilon_2))(1+\epsilon_3)$
- $= z(1+\epsilon_3)(x(1+\epsilon_1)+y(1+\epsilon_2)), \ |\epsilon_i| \leq u$

$$z \otimes (x \oplus y)$$

= $z((x + y)(1 + \epsilon_1))(1 + \epsilon_3)$
= $z(1 + \epsilon_3)(x(1 + \epsilon_1) + y(1 + \epsilon_1)), |\epsilon_i| \le u$

◆□▶ ◆□▶ ◆□▶ ◆□▶ ▲□ ◆ ○ ◆

Both methods are backward stable.

Perturbation analysis

$$z(1 + \delta_z)(x(1 + \delta_x) + y(1 + \delta_y))$$

$$\approx zx(1 + \delta_z + \delta_x) + zy(1 + \delta_z + \delta_y)$$

$$= z(x + y) + zx(\delta_z + \delta_x) + zy(\delta_z + \delta_y)$$

$$= z(x + y)(1 + (\delta_z + \delta_x) + (\delta_y - \delta_x)/(x/y + 1))$$

$$\frac{|z(1+\delta_z)(x(1+\delta_x)+y(1+\delta_y))-z(x+y)|}{|z(x+y)|} \\ \leq \left(2+\frac{2}{|\frac{x}{y}+1|}\right)\delta, \quad \delta = \max(|\delta_x|, |\delta_y|, |\delta_z|$$

▲ロト ▲御 ト ▲ 臣 ト ▲ 臣 ト ○ 臣 - のへで

Perturbation analysis

$$z(1 + \delta_z)(x(1 + \delta_x) + y(1 + \delta_y))$$

$$\approx zx(1 + \delta_z + \delta_x) + zy(1 + \delta_z + \delta_y)$$

$$= z(x + y) + zx(\delta_z + \delta_x) + zy(\delta_z + \delta_y)$$

$$= z(x + y)(1 + (\delta_z + \delta_x) + (\delta_y - \delta_x)/(x/y + 1))$$

$$egin{aligned} & rac{|z(1+\delta_{z})(x(1+\delta_{x})+y(1+\delta_{y}))-z(x+y)|}{|z(x+y)|} \ & \leq & \left(2+rac{2}{|rac{x}{y}+1|}
ight)\delta, \quad \delta=\max(|\delta_{x}|,|\delta_{y}|,|\delta_{z}|) \end{aligned}$$

The condition number can be large if $y \approx -x_a$ and $\delta_x \neq \delta_y$.

Floating-point Numbers Sources of Errors Stability of an Algorithm Sensitivity of a Problem Fallacies Summary

Example (cont.)

Forward error analysis

$$z \otimes x \oplus z \otimes y$$

$$= z(1 + \epsilon_3)(x(1 + \epsilon_1) + y(1 + \epsilon_2))$$

$$\approx z(x + y)(1 + (\epsilon_3 + \epsilon_1) + (\epsilon_2 - \epsilon_1)/(x/y + 1)), \quad |\epsilon_i| \le u$$

$$\frac{|(z \otimes x \oplus z \otimes y) - z(x + y)|}{|z(x + y)|} \le \left(2 + \frac{2}{|\frac{x}{y} + 1|}\right)u$$

▲□▶▲□▶▲□▶▲□▶ □ のへで

Floating-point Numbers Sources of Errors Stability of an Algorithm **Sensitiviy of a Problem** Fallacies Summary

Example (cont.)

Forward error analysis (cont.)

$$egin{aligned} & z \otimes (x \oplus y) pprox z(x+y)(1+\epsilon_1+\epsilon_3), \quad |\epsilon_i| \leq u \ & rac{|z \otimes (x \oplus y) - z(x+y)|}{|z(x+y)|} \leq 2u \end{aligned}$$

Floating-point Numbers Sources of Errors Stability of an Algorithm Sensitivity of a Problem Fallacies Summary



forward error \leq cond \cdot backward error

- If we can prove the algorithm is stable, in other words, the backward errors are small, say, no larger than the measurement errors in data, then we know that large forward errors are due to the ill-conditioning of the problem.
- If we know the problem is well-conditioned, then large forward errors must be caused by unstable algorithm.
- Condition number is an upper bound. It is possible that a well-designed stable algorithm can produce good results even the problem is ill-conditioned.

Example revisited

$$\beta = 10, t = 4$$

 $x = 1.002, y = -0.9958, z = 3.456$
Exact $z(x + y) = 2.14272 \times 10^{-2}$
 $z \otimes (x \oplus y) = fl(3.456 * 6.200 \times 10^{-3})$
 $= 2.143 \times 10^{-2}$
error: 2.8×10^{-6}
 $(z \otimes x) \oplus (z \otimes y)) = fl(3.463 - 3.441)$
 $= 2.200 \times 10^{-2}$
error: 5.7×10^{-4}
More than 200 times!
Why?

Example revisited

$$(z \otimes x) \oplus (z \otimes y)) = fl(3.463 - 3.441)$$

= 2.200 × 10⁻²
error: 5.7 × 10⁻⁴

Cancellation in subtracting two computed (contaminated) numbers. (Catastrophic)

◆□▶ ◆□▶ ◆三▶ ◆三▶ → 三 のへぐ

Example revisited

$$(z \otimes x) \oplus (z \otimes y)) = fl(3.463 - 3.441)$$

= 2.200 × 10⁻²
error: 5.7 × 10⁻⁴

Cancellation in subtracting two computed (contaminated) numbers. (Catastrophic)

$$z \otimes (x \oplus y) = \text{fl}(3.456 * 6.200 \times 10^{-3})$$

= 2.143 × 10⁻²
error: 2.8 × 10⁻⁶

Cancellation in subtracting two original (not contaminated) numbers. (Benign)

Example revisited

$$(z \otimes x) \oplus (z \otimes y)) = fl(3.463 - 3.441)$$

= 2.200 × 10⁻²
error: 5.7 × 10⁻⁴

Cancellation in subtracting two computed (contaminated) numbers. (Catastrophic)

$$z \otimes (x \oplus y) = \text{fl}(3.456 * 6.200 \times 10^{-3})$$

= 2.143 × 10⁻²
error: 2.8 × 10⁻⁶

Cancellation in subtracting two original (not contaminated) numbers. (Benign)

Catastrophic cancellation v.s. benign cancellation.

A classic example of avoiding cancellation

Solving quadratic equation

$$ax^2 + bx + c = 0$$

Text book formula:

$$x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

Computational method:

$$x_1 = \frac{2c}{-b - \operatorname{sign}(b)\sqrt{b^2 - 4ac}}, \quad x_2 = \frac{c}{ax_1}$$

▲□▶▲圖▶▲圖▶▲圖▶ 圖 のQ@

Floating-point Numbers Sources of Errors Stability of an Algorithm Sensitivity of a Problem Fallacies Summary

Question

Suppose $\beta = 10$ and t = 8 (single precision), solve

$$ax^2 + bx + c = 0,$$

where

$$a = 1$$
, $b = -10^5$, and $c = 1$,

using the both methods.

Floating-point Numbers Sources of Errors Stability of an Algorithm Sensitivity of a Problem Fallacies Summary

Outline

- Floating-point Numbers
 - Representations
 - IEEE Floating-point Standards
 - Underflow and Overflow
 - Correctly Rounded Operations
- 2 Sources of Errors
 - Rounding Error
 - Truncation Error
 - Discretization Error
- 3 Stability of an Algorithm
- 4 Sensitiviy of a Problem



Fallacies

- Cancellation in the subtraction of two nearly equal numbers is always bad.
- The final computed answer from an algorithm cannot be more accurate than any of the intermediate quantities, that is, errors cannot cancel.
- Arithmetic much more precise than the data it operates upon is needless and wasteful.
- Classical formulas taught in school and found in handbooks and software must have passed the Test of Time, not merely withstood it.

Summary

- A computer number system is determined by four parameters: Base, precision, *e*_{min}, and *e*_{max}
- IEEE floating-point standards, single precision and double precision. Special quantities: Denormals, ±∞, NaN, ±0, and their binary representations.
- Error measurements: Absolute and relative errors, unit of roundoff, unit in the last place (ulp)
- Sources of errors: Rounding error (computational error), truncation error (mathematical error), discretization error (mathematical error). Total error (combination of rounding error and mathematical errors)
- Issues in floating-point computation: Overflow, underflow, cancellations (benign and catastrophic)
- Error analysis: Forward and backward errors, sensitivity of a problem and stability of an algorithm