

Numerical Error Analysis

→ Analytic representation of numbers

1) Integers : $0, \pm 1, \pm 2, \pm 3, \dots$. Let x be an integer.

$$\rightarrow x = (\alpha_0 + \alpha_1 b + \alpha_2 b^2 + \dots + \alpha_n b^n) s \quad \text{where}$$

- a) b natural number > 1 (2, 3, 4, ...). \leftrightarrow base of representation
- b) For all k : $0 \leq \alpha_k < b$, and natural \leftrightarrow digits
- c) $s = +1$ or $s = -1$ \leftrightarrow sign.

example

$$666 = 6 + 6 \cdot 10 + 6 \cdot 100 \quad (\text{decimal})$$

$$666 = 0 + 1 \cdot 2 + 0 \cdot 4 + 1 \cdot 8 + 1 \cdot 16 + 0 \cdot 32 + 0 \cdot 64 + 1 \cdot 128 + 0 \cdot 256 + 1 \cdot 512 = (1010011010)_2 \quad (\text{binary})$$

$$666 = 10 + 9 \cdot 16 + 2 \cdot 16^2 = (29A)_{16} \quad (\text{hexadecimal}).$$

↳ use A, B, C, D, E, F = 10, 11, 12, ..., 15.

► note the relation between binary and hexadecimal:

$$\begin{array}{cccc} (1010011010)_2 & & & \\ \downarrow \downarrow \downarrow \downarrow & & & \\ (29A)_{16} & & & \end{array}$$

$$\begin{array}{l} (1010)_2 = 10 = (A)_{16} \\ (1001)_2 = 9 = \\ (10)_2 = 2 = \end{array}$$

Every 4 digits of the binary representation correspond to 1 digit of hexadecimal. This can be used in binary \leftrightarrow hexadecimal conversion. ◀

► integer division : $\boxed{\frac{a}{b} = q + \frac{r}{b}}$ q = quotient
r = remainder $< b$

$q = a / b$ \leftarrow integer division] convention in
 $r = a \% b$ \leftarrow remainder operation. C and C++ ◀

Converting to decimal is easy. Converting from decimal as follows

(2)

Algorithm : Let $x > 0$ be a given integer
 $n = \max \{ k \in \mathbb{N} \mid b^k < x \}$
 $x_n = x$
 for $k = n, n-1, \dots, 1, 0$
 || $a_k = x_k / b^k$
 || $x_{k-1} = x_k \% b^k$, if $k > 0$

2) Floating point numbers : 3.45, π , $\sqrt{2}$

$$x = s b^e \left[\frac{m_1}{b} + \frac{m_2}{b^2} + \frac{m_3}{b^3} + \dots \right]$$

- a) b natural number > 1 (2, 3, 4, ...) \leftrightarrow base
- b) e integer (0, $\pm 1, \pm 2, \dots$) \leftrightarrow exponent
- c) $s = 1$ or $s = -1$ \leftrightarrow sign
- d) For all k : $m_k \in \mathbb{N}$ and $0 \leq m_k < b$ \leftrightarrow mantissa

To convert to base b :

Algorithm :
 $e = \min \{ k \in \mathbb{Z} \mid b^k > x \}$
 $x_0 = x / b^e$
 for $k = 0, 1, 2, \dots$
 || $m_k = \lfloor x_k b \rfloor$
 || $x_{k+1} = b x_k - \lfloor b x_k \rfloor$

- Note that if $x_j = x_k$ with $j < k$ then, the sequence (m_k) will repeat after the k . ◀
- If $x_k = 0$ for some k , then the algorithm may terminate. Then we say that x is terminating with respect to b . ◀

(3)

→ Numerical representation of numbers.

- 1) Short integers (16-bit) $\rightarrow x = s(a_0 + a_1 2 + \dots + a_{14} 2^{14})$
 - 2) Integers (32 bit) $\rightarrow x = s(a_0 + a_1 2 + \dots + a_{30} 2^{30})$
 - 3) Long integers (64 bit) $\rightarrow x = s(a_0 + a_1 2 + \dots + a_{62} 2^{62})$
 - 4) Unsigned short integer $\rightarrow x = a_0 + a_1 2 + \dots + a_{15} 2^{15}$
 - 5) Unsigned integer $\rightarrow x = a_0 + a_1 2 + \dots + a_{31} 2^{31}$
 - 6) Unsigned long integer $\rightarrow x = a_0 + a_1 2 + \dots + a_{63} 2^{63}$
- with $s = 1, -1$ and $a_k = 0, 1$

► Integer arithmetic is accurate except for overflow and underflow errors. ◀

► In general: k -bit integer $\rightarrow x = s(a_0 + a_1 2 + \dots + a_{k-2} 2^{k-2})$
 k -bit unsigned integer $\rightarrow x = a_0 + a_1 2 + \dots + a_{k-1} 2^{k-1}$

It is also possible to implicitly modify the range by assuming that x is shifted by a given amount. ◀

▼ Numerical floating point representation of numbers.

There are many variants. We discuss the IEEE standard.
 In general:

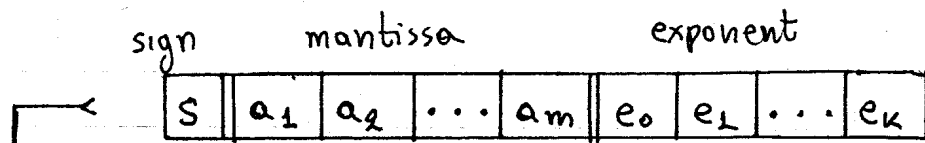
$$x = s b^e \left[\frac{a_1}{b} + \frac{a_2}{b^2} + \dots + \frac{a_m}{b^m} \right]$$

$$e = \mu + e_0 + e_1 \cdot b + \dots + e_k \cdot b^k$$

- | | |
|------------------------|---|
| where, $b =$ base | $k =$ exponent bitlength. |
| $m =$ mantissa length | $\mu =$ smallest exponent |
| $a_k =$ mantissa bits. | $M = \mu + (b-1) + (b-1)b + \dots + (b-1)b^k$ |
| $e_k =$ exponent bits | $=$ biggest exponent |
| $e =$ exponent | $s =$ sign bit. |

The actual information that needs to be stored is:

(4)



$2 \cdot 2^m \cdot 2^{k+1}$ possible numbers can be represented when $b=2$.

→ Rounding function and floating point operations.

Let \mathcal{R} be the set of all the numbers that can be generated by a floating point representation.

A floating point system consists of

- ① A rounding function $rd: \mathbb{R} \rightarrow \mathcal{R}$ that maps every real number to a representable number.
- ② The floating point operations that map a pair of representable numbers to another representable number.

$$x \oplus y \equiv rd(x+y) = \overline{x+y}$$

$$x \ominus y \equiv rd(x-y) = \overline{x-y}$$

$$x \odot y \equiv rd(xy) = \overline{xy}$$

$$x \oslash y \equiv rd(x/y) = \overline{x/y}$$

↖ book notation

The computer should perform these operations in a manner that is always consistent with the rounding function, for all $x, y \in \mathcal{R}$. (pairs of representable numbers).

→ Machine epsilon

Definition: The machine epsilon eps of a floating point system is the largest number such that:

$$1 \oplus eps = 1$$

This is NOT the smallest number that can be represented in the floating point system. It instead measures the ability of the floating point system to handle numbers with very big differences in the order of magnitude.

(5)

The machine epsilon provides an estimate of the relative error that is introduced during a fundamental floating point operation:

$$(x \oplus y) = (x+y)(1+\epsilon_1) \quad , \quad |\epsilon_1| < \text{eps}$$

$$(x \odot y) = (xy)(1+\epsilon_2) \quad , \quad |\epsilon_2| < \text{eps}$$

$$(x \ominus y) = (x-y)(1+\epsilon_3) \quad , \quad |\epsilon_3| < \text{eps}$$

$$(x \oslash y) = (x/y)(1+\epsilon_4) \quad , \quad |\epsilon_4| < \text{eps}$$

Note that ϵ_1 , etc are dependent on the specific $x, y \in \mathbb{F}$, nevertheless are always bound by eps.

→ A note on error propagation

The complete theory of error propagation is discussed in J. Stoer, R. Bulirsch, "Introduction to Numerical Analysis" Floating point operations are fundamentally different from the corresponding analytic ones. For example; if:

$$a = 0.23371258 e - 04$$

$$b = 0.33678429 e + 02$$

$$c = -0.33677811 e + 02$$

then

$$a \oplus (b \oplus c) = 0.64137126 e - 03 \quad \text{accurate}$$

$$(a \oplus b) \oplus c = 0.64100000 e - 03 \quad \text{very inaccurate.}$$

The order with which operations are performed can make a very big difference. There is a formal theory, due to Bauer, with which you can prove whether a computation scheme is or isn't ^{accurate} ~~argument~~. We will not cover this theory in this course. We will only state a prediction:

► Let $y = \varphi(\vec{x})$ be a function of many variables. The inherent error is defined as:

$$\Delta^{(0)} y = [|\nabla \varphi(\vec{x}) \cdot \vec{x}| + |y|] \text{eps}$$

If a computation scheme is "accurate", it should exhibit an

⑥

error that has the same order of magnitude as $\Delta^{\circ}y$ \leftarrow
Therefore, if you observe an error $\gg \Delta^{\circ}y$, then you have a bug.

► Problems with steep gradient $\nabla\varphi(\vec{x})$ have a large inherent error, and it is impossible to eliminate such error.

We say that such problems are ill-conditioned. \leftarrow

The following rules of thumb help.

1) Add large sequences of numbers in ascending order.

2) Avoid $x-y$ if $x-y \ll x$ or $x-y \ll y$.

3) Avoid x/y if $y \ll 1$

4) Avoid $x+y$ if $y \ll x$ or $x \ll y$. (e.g. $1+\text{eps} = 1$).

5) Avoid $\sin x, \cos x, \tan x$ for $x \gg 2\pi$

6) Use Horner scheme to evaluate polynomials.

7) Be paranoid - do formal analysis and proof if your computation is responsible for smart missiles and stuff like that, otherwise you might piss off China.

→ Rounding functions.

Rounding is implicit in the implementation of floating point arithmetic; there should exist a rounding function such that the floating point operations are consistent with it. How that is actually done is messy business. Computers now do it in hardware not software.

Note that although computers grok binary numbers, they may choose to round in a different system, like octal or hexadecimal, instead of binary. I.e. they carry out operations in binary such that rounding in another system is satisfied.

There is two basic types of rounding.

1) Chopping: Simply drop the digits that can not be resolved.
May cause rounding error up to eps .

(7)

2) Symmetric chopping: Drop the digits but round up or down, whichever is nearest.

e.g. $rd(3.14159) = 3.1416$

$rd(3.492) = 3.49$

In case of a tie, pick the even number (i.e. choose "randomly" to minimize statistical bias)

e.g. $rd(3.135) = 3.14$

$rd(3.125) = 3.12$

Other conventions have been to round towards 0, $+\infty$, or $-\infty$.

→ Repeating numbers.

Fact: $0 < |a| < 1 \Rightarrow 1 + a + a^2 + \dots = \frac{1}{1-a}$ (geometric series theorem).

Analytically $1/5 = 0.2$ seems a trivial computation. However, in binary arithmetic $1/5$ is a repeating number so it is rounded, and upon conversion to decimal, you might get 0.1999999 or 0.2000001 . To see that $1/5$ is repeating:

$$\begin{aligned} \frac{1}{5} &= \frac{1}{4-1} = \frac{1}{4} \frac{1}{1-1/4} = \frac{1}{4} \left[1 + \frac{1}{4} + \frac{1}{4^2} + \dots \right] = \\ &= \sum_{k=1}^{+\infty} \frac{1}{4^k} = \sum_{k=1}^{+\infty} \frac{1}{2^{2k}} = 0.01010101\dots = 0.\overline{01} \end{aligned}$$

So repetitiveness depends on the representation system. In general, the following theorem can be proven:

Theorem: Let $0 < b_1 < b_2$ integers such that $b_1 | b_2$ (b_1 divides b_2). Then the following are true:

a) If x repeats in $b_2 \Rightarrow x$ repeats in b_1

b) ~~There~~ There exists an x that does not repeat in b_2 but repeats in b_1 .