

Announcements :

- * Assignment 1 is due Monday 9 pm
- * MarkUs link is posted on the course website.

Last Class

- * Conditioning / sensitivity
- * stability
- * Error : Absolute, Relative, Forward, Backward.

Today

- Floating Point Numbers.
- ↳ Representing real numbers
- ↳ Floating point arithmetic

Floating-Point Numbers

Problem: We need to store real numbers (continuous) using hardware that is discrete.

Ideally, each real number will ~~not~~ take up the same amount of storage space.

Solution : Floating-Point Numbers!

↳ which are like scientific notation.

surface area of Earth: 510072000 km^2

in scientific notation: $5.10072 \times 10^8 \text{ km}^2$

with fewer precision: $5.101 \times 10^8 \text{ km}^2$

Planck constant : $6.626 \times 10^{-34} \text{ m}^2 \text{ kg/s}$

always one digit before the decimal. precision. base exponent denotes the position of decimal

A Floating Point Number System is characterized by ③
4 integers:

β - base (radix)

eg// base 10 for scientific notation

p - precision

base 2 on most computers.

$[L, U]$ - exponent range.

In this system, a real number x is represented as:

$$x = \pm \left(d_0 + \frac{d_1}{\beta} + \frac{d_2}{\beta^2} + \dots + \frac{d_{p-1}}{\beta^{p-1}} \right) \beta^E$$

where $0 \leq d_i \leq \beta-1$ $d_i \in \mathbb{N}$

$L \leq E \leq U$.

0.5101 × 10⁹

eg// surface area of earth $5.101 \times 10^8 \text{ km}^2$.

$$(5 + \frac{1}{10} + \frac{0}{10^2} + \frac{1}{10^3}) 10^8$$

The digits $d_0 d_1 \dots d_{p-1}$ is called the mantissa.

$d_1 \dots d_{p-1}$: : : fraction.

E : : : exponent.

eg// IEEE Single Precision System:

$$\beta = 2.$$

bits we need

$$p = 24$$

mantissa - 24. 23

$$L = -126$$

exponent - 8 There are 254 different numbers between

$$U = 127$$

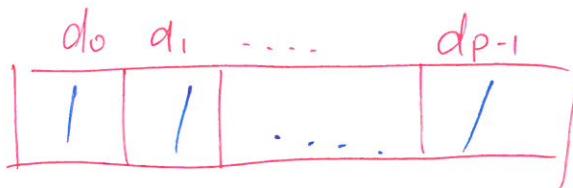
sign - 1 1
33 32 (inclusive).

$$127 - (-126) + 1 = \underline{\underline{254}}$$

Q: What is the largest number we can store using IEEE SP.

2^{128} — too large, $\# 128 > U$

$$2^{127} - 1 + \sum \text{mantissa.} - ?$$



E
127

二

~~2128~~ 2128-24
2128 2128-24

$$\left(1 + \frac{1}{2} + \frac{1}{2^2} + \dots + \frac{1}{2^{23}}\right) \times 2^{127}$$

Normalization:

A Floating-Point System is normalized if the leading digit is $\neq 0$ unless the number represented is 0.

- There is a unique representation for each non-zero value.
 5.101×10^8 ✓
 0.5101×10^9 X
 - If $\beta=2$, we don't have to store the first digit.

Q : listbox

How many floating point numbers are in a floating point system $F(B, p, L, U)$ normalized.

$$2(\beta-1) \underbrace{\beta^{P-1}}_{\text{mantissa}} (\underbrace{U-L+1}_{\text{exponents}}) + \underbrace{1}_{\text{zero}}$$

not all real numbers are exactly representable in our floating-point system.

Numbers that are exactly representable are called machine numbers.

Rounding: $f_1(x) \approx$ approximate $x \in \mathbb{R}$ by a "nearby" machine number
 \hookrightarrow rounding error is introduced here.

Rounding Rule:

1. Chop. truncate the base β expansion of x after $(p-1)$ digits \approx "Round toward zero"
2. Round to nearest.
 $\#$ find $f_1(x)$ nearest to x .

e.g. // $x = 5.10072 \times 10^8 \text{ km}^2$

want to represent x in $F(\beta=10, p=4, L=-5000, U=5000)$

Chop: 5.100×10^8

Nearest: 5.101×10^8

Q: What is the smallest normalized floating-point number

in $F(\beta, p, L, U)$:

$$1 \times \beta^L = \underline{\beta^L} \leftarrow \text{underflow level (UFL)}$$

... largest:

$$\frac{(1-\beta^{-p})\beta^{U+1}}{\underline{}} \leftarrow \text{overflow level (OFL)}$$

~~Defn~~ The machine precision characterizes the accuracy of the floating point system.

Can be defined in a few ways:

- Max. relative error in representing a nonzero real number \underline{x} with $\underline{f(x)}$
- Smallest ϵ s.t. $f(1+\epsilon) > 1$.

If we use rounding by chopping — β^{1-p}
 nearest — $\frac{1}{2}\beta^{1-p}$

$$5.10\underline{09999} \rightarrow 5.100 \quad (\text{in } F/\beta=10, p=4)$$

we lose ~~all~~ of this.
 $\hookrightarrow 10^{-3}$

When we plot floating point (machine) numbers, there is a "gap" around 0. due to normalization.

To remove this gap, we'll relax normalization and allow leading zeros ($a_0 = 0$) when $E = L$

These new floating point numbers are called Subnormal or denormalized F.P. numbers.

This augment system exhibit gradual underflow.

Floating Point Arithmetic.

→ Floating Point Arithmetic is Inexact

Addition/Subtraction: shift the mantissa to match exponents before adding/subtracting

$$x = 1.924 \times 10^2$$

$$y = 6.357 \times 10^{-1}$$

in F($\beta=10$, $p=4$)

$$\begin{array}{r} x+y : \quad 1.924 \times 10^2 \\ \quad \quad \quad 6.357 \times 10^{-1} \\ \hline \boxed{1.930357} \times 10^2 \end{array}$$

↑ we will lose this information
due to rounding

Multiplication: product of two f.p. numbers with precision p will have $2p$ digits

Division: quotient can have many more digits
 $1 \div 3$ has infinite # digits in base 10.

Possible Issues

1. Loss of accuracy.
2. Overflow - exponent > U
3. Underflow - exponent < L.
(silently set to zero)

e.g. $\sum_{n=1}^{\infty} \frac{1}{n}$ should diverge, but converges in F.P. arithmetic.

1. overflow
2. it eventually underflows
3. it becomes insignificant compared to the ϵ so far

example of 3: in $F(\beta=10, p=4)$

$$\left\{ \begin{array}{l} x = 1.924 \times 10^2 \\ y = 6.357 \times 10^{-4} \end{array} \right.$$

$$x+y = x.$$

extension of \oplus :

You can imagine ϵ where $1+\epsilon = 1$

$$\text{so } (1+\epsilon) + \epsilon = 1.$$

$$\text{but } \cancel{\textcircled{+}} \quad 1 + (\epsilon + \epsilon) > 1.$$

$$\Rightarrow (1+\epsilon) + \epsilon \neq 1 + (\epsilon + \epsilon)$$

\Rightarrow Floating point Addition is not associative

$$(a+b)+c \neq a+(b+c)$$

so the order that you add numbers matter.

Normal laws of arithmetic don't hold.

Ideally: a floating point operation flop (F.P. addition)
and its arbitrary precision op. (real addition)

should have:

$$\underline{f_1(x) \text{ flop. } f_1(y) = f_1(x \text{ op } y)}$$

Cancellation: subtracting two p-digit numbers with similar sign & magnitude yield result with much fewer than p digits:

$$\begin{array}{r}
 1.92403 \times 10^2 \\
 - 1.92275 \times 10^2 \\
 \hline
 1.28000 \times 10^{-1}
 \end{array} \quad (p=6) \quad (\beta=10)$$

3 significant digits!

eg// ϵ where $0 < \epsilon < \epsilon_{\text{mach}}$ \leftarrow machine precision $\begin{cases} \beta^{1-p} \text{ chop} \\ \frac{1}{2}\beta^{1-p} \text{ near} \end{cases}$

$$(1+\epsilon) - (1-\epsilon) = 1 - 1 = 0.$$

↳ slightly different issue than cancellation.
but still shows order of operations matter

We don't want to compute small quantities as a difference of 2 large quantities.

eg// Standard deviation; need to compute.

$$\sum_{i=1}^N (x_i - \bar{x})^2 = \sum_{i=1}^N x_i^2 - N\bar{x}^2.$$

\uparrow \uparrow
 data point Mean
 $\underbrace{\qquad\qquad\qquad}_{\text{difference is small}} \quad \underbrace{\qquad\qquad\qquad}_{\text{large #}}$
 compared to the 2 values

Catastrophic cancellation

9.

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

For $x < 0$ this gives disastrous results because of cancellation.

For example, e^{-40} is small, but if we try to compute

$$e^{-40} = 1 + (-40) + \frac{(-40)^2}{2!} + \frac{(-40)^3}{3!} + \dots$$

Each term is large, and adjacent terms have opposite signs, leading to large errors.