# CSC 338 Lecture 2

## Floating-Point Numbers

Advantages: - FP numbers is a way to store <u>continuous</u> quantities using <u>discrete hardware</u>

- Each FP number takes same amount of storage

Disadvantages: - FP arithmetic is inexact
  - ↳ overflow/underflow
  - ↳ loss of accuracy
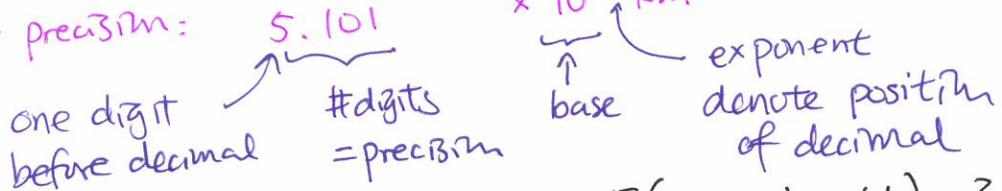  - ↳ cancellation / catastrophic cancellation.

Our treatment of FP numbers will be more general than in CSC 207

<u>Intuition</u>: floating-point numbers are like scientific notation.

S.A. of earth.    510 072 000 km²

scientific notation:    $5.10972 \times 10^8$ km².
                          ← in base 10 ↗

fewer precision:    $5.101 \times 10^8$ km²

one digit                    #digits         base         exponent
before decimal              = precision                  denote position
                                                          of decimal

<u>Def</u>: A Floating-Point Number System $F(\beta, p, L, u)$ is characterized by integers    $\beta$ — <u>base</u>
                            $p$ — <u>precision</u>
                    $[L, u]$ — <u>exponent range</u>.

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

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

where    $d_i \in \{0, 1, \dots, \beta-1\}$

        $E \in \{L, L+1, \dots, u-1, u\}$

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

$d_1 \ldots d_{p-1}$ is called the fraction.

$E$ is called the exponent

eg// In the system $F(\beta=10, \ p=3, \ L=-100, \ U=100)$ we can represent the S.A. of earth as.

$$+ \left(5 + \frac{1}{10} + \frac{0}{10^2}\right) \cdot 10^8 \qquad \Longleftrightarrow \qquad \underline{5.10} \times \underline{10^8}$$

$\qquad \qquad d_0 \quad d_1 \quad d_2 \qquad E. \qquad \qquad \qquad \qquad \leftarrow$ base 10 #s.

Q what about $\quad 0.51 \times 10^9 \qquad$ or $\quad +\left(0 + \frac{5}{10} + \frac{1}{10^2}\right) \cdot 10^9$ ?

Def. A F.P. system is normalized if we enforce the additional rule that the leading digit $d_0 \neq 0$ unless the number represented is 0.

$\Rightarrow$ There is a unique normalized repr for each nonzero value.

$\Rightarrow$ If $\beta = 2$, then $d_0 = 1$ for all nonzero values, so we don't need to store $d_0$.

eg//. IEEE SP System $F(\beta = 2, \ p = 24, \ L = -126, \ U = 127)$.

what are the smallest and largest positive number that we can represent in this system? ($\underleftarrow{\text{normalized .F.P.}}$)

Smallest : set $d_0 = 1$, $d_i = 0$ for $i \neq 0$, $E = -126$.

$$\Rightarrow \quad +\left(1 + \frac{0}{2} + \ldots + \frac{0}{2^{23}}\right) \times 2^{-126} = \cdot 2^{-126}.$$

largest : Set $d_i = 1$, $E = 127$

$$\Rightarrow \quad +\left(1 + \frac{1}{2} + \ldots + \frac{1}{2^{23}}\right) \times 2^{127} \qquad \text{positive}$$

Def. The underflow level (UFL) is the smallest normalized FP number in $F(\beta, p, L, U)$. Can show $\quad UFL = \beta^L$

The overflow level (OFL) is the largest, normalized positive FP number in $F(\beta, p, L, U)$. Can show $OFL = (1 - \beta^{-p}) \beta^{U+1}$

eg// How many FP numbers are there in a normalized FP System $F(\beta, p, L, U)$?

$$\underbrace{2}_{sign} \times \underbrace{(\beta-1)}_{\substack{leading \\ digit}} \underbrace{\beta^{p-1}}_{\substack{rest \, of \\ mantissa}} \underbrace{(U-L+1)}_{exponent} + \underbrace{1}_{zero}$$

<u>Def</u> : Numbers that are exactly representable in a ~~normalized~~ F.P. system are called <u>machine numbers</u>.

when we plot machine numbers on the number line, there is a "gap" around zero due to normalization.

<u>Def</u>. In a <u>subnormal</u> or <u>denormalized</u> FP system, we relax normalization and allow leading zeros $(d_0 = 0)$ when $E = L$. These new FP numbers are called <u>subnormal</u> or <u>denormalized</u>. This augmented system exhibit <u>gradual underflow</u>.

(Break : 4:05)

<u>Rounding</u>

We represent $x \in \mathbb{R}$ in a F.P. system by approximately $x$ with a "nearby" machine number $fl(x)$ via <u>rounding</u>.

Two method of rounding.
1. Chopping "Round toward zero" ~ truncate the base $\beta$ expansion of $x$ after $(p-1)$ digits
2. Round to nearest. Choose $fl(x)$ to minimize $|fl(x) - x|$

eg// want to represent $x = 5.10072 \times 10^8$ in $F(\beta=10, p=4, L=-100, U=100)$

chop – $5.100 \times 10^8$

Nearest – $5.101 \times 10^8$

The machine precision $\epsilon_{mach}$ characterizes the accuracy of a F.P. system. Can be defined in a few ways.

Def$_1$   $\epsilon_{mach}$ is the maximum relative error in representing a nonzero $x \in \mathbb{R}$ with $fl(x)$

Def$_2$   $\epsilon_{mach}$ is the smallest $\epsilon$ with $fl(1+\epsilon) > 1$

Q: what is the machine precision for $F(\beta, p, L, U)$?

A:   If we round by chopping.   $\epsilon_{mach} = \beta^{1-p} = (\beta^{-p+1})$

to nearest.   $\epsilon_{mach} = \frac{1}{2}\beta^{1-p}$.

## Floating-Point Arithmetics

Floating-Point Arithmetic is inexact.

multiplication. Product of two FP numbers with precision $p$ will have $2p$ digits $\Rightarrow$ need to round.

eg//. In $F(\beta=10, p=2, L=-10, U=10)$

$x = 5.9 \times 10^2$
$y = 6.1 \times 10^1$

$$
\begin{array}{r}
59 \quad \times 10^2 \\
6.1 \quad \times 10^1 \\
\hline
59 \\
354 \\
\hline
3599 \quad \times 10^4 \quad \Rightarrow \quad 3.6 \times 10^4.
\end{array}
$$

$2p = 4$ digits.

DIVISIONS   Also inexact.   Long division.
Quotient can have many more digits, & potentially infinite.

eg// $1 \div 3$ has infinite # digits in $\beta = 10$.

# Addition (-subtraction).

shift the mantissa to match the exponent, then
add/subtract along the column.

eg// In $F(\beta=10, p=4, L=-70, U=10)$

$x = 1.924 \times 10^2$

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

$x+y$ :

$$\begin{array}{r} 1.924 \quad \times 10^2 \\ 6.357 \times 10^{-1} \\ \hline \boxed{1.930}\,357. \quad \times 10^2 \\ = \end{array}$$

lose this info due
to rounding.

$x+y \implies 1.930 \times 10^2$

eg// In $F(\beta=5, p=3, L=-10, U=10)$

base 5

$x = 2.44 \times 5^0$

$y = 3.33 \times 5^1$

$$\begin{array}{r} \overset{1}{.}2.44 \quad \times 5^0 \\ 3.33 \quad \times 5^1 \\ \hline 4.12.4 \quad \times 5^1 \end{array}$$

Round to
nearest $\longrightarrow$ $4.13 \quad \times 5^1$

Ideally a floating-point operation   flop.   (eg// F.P. add.iTm)
        vs.  its  arbitrary precision      op    (eg// real addition)

should have     $fl(x)$ flop $fl(x) = fl(x \text{ op } y)$

In most computers, this is true!
Relative error of each computation is bounded by $\epsilon_{mach}$

But there are still issues to be aware of
    eg// nomal law of arithmetic don't hfd.

$\Rightarrow$ Floating-Point Addition is not associative!

eg// Choose $\epsilon = 0.6\, \epsilon_{mach}$.      So   $1 +_{fl} \epsilon = 1$   (because $\epsilon < \epsilon_{mach}$)

So   $(1 +_{fl} \epsilon) +_{fl} \epsilon = 1$

$1 +_{fl} (\epsilon +_{fl} \epsilon) \gtreqless 1$

eg// Choose $\epsilon = 0.6\ \epsilon$mach.

$$(1 +_{fl} \epsilon) -_{fl} (1 -_{fl} \epsilon) = |-_{fl}.| = 0.$$

eg//

FP arithmetic can cause an <u>overflow</u> when the result has $E > U$, an <u>underflow</u> when the result has $E < L$

eg// $\sum_{i=1}^{\infty} \frac{1}{n}$. should diverge, but converges in FP arithmetic.

Potential Reason :

✗ Overflow✗ because $\Sigma$ gets too large.

✗ underflow✗ because $\frac{1}{n}$ gets too small

✓ because $\frac{1}{n}$ becomes insignificant compare to the $\Sigma$ so far

## Cancellation

Subtracting two p-digit numbers with similar sign and magnitude. yield result with much <u>fewer</u> than $p$ digits

eg//. $F(\beta = 10,\ p = 6,\ L = -10,\ U = 10)$

$$
\begin{array}{r}
1.92403 \times 10^2 \\
- 1.92275 \times 10^2 \\
\hline
1.28000 \times 10^{-1}
\end{array}
$$

only 3 significant digits (started w/ 6)

eg// Computing standard deviation

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

data point   mean       large      large

difference much smaller

# Catastrophic Cancellation

In some cases, cancellation can be so bad that the solution has <u>no</u> correct significant digits!

eg// $e^x = 1 + x + \dfrac{x^2}{2!} + \dfrac{x^3}{3!} + \dots$

For $x$ that is large and negative, computing $e^x$ using its Taylor series give disasterous results.

For $x = 40$, $e^{-40}$ is small but with

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

each term is large, and adjacent terms have opposite signs $\Rightarrow$ cancellation.

moral of today's lecture:
- FP usually good enough
- avoid subtracting two large values to get a small result.