

Welcome to Ma221! Lecture 3 Fall 2024

## Floating Point + Error Analysis

FP: how to represent real numbers

Long ago (<1985) many computers did arithmetic differently  $\Rightarrow$  hard to write portable code  $\Rightarrow$  IEEE 754 FP Standard, led by Prof Kahan

First standard 1985, then 2008, 2019, 2029...

Prizes: Turing Award 1989  
IEEE Milestone 2023

Scientific Notation  $\pm d.d\dots d \cdot \text{radix}^e$

Usually radix = 2 (or 10 for finance)

Store: sign bit ( $\pm$ )

exponent ( $e$  integer)

mantissa: ( $d_1 d_2 \dots d_p$ )

$p = \#$  digits in mantissa

9 decimal digits  $\rightarrow 9 \cdot 4 = 36$  bits  $\rightarrow \left(\frac{10}{16}\right)^4 = .015$  memory usage

9 decimal digits  $\rightarrow 3$  groups of 3 digits

$\rightarrow 3$  groups of 10 bits

$\rightarrow 30$  bits  $\rightarrow \left(\frac{1000}{1024}\right)^3 = .93$  memory usage

Both  $p$  and #bits in  $e$  are limited  
to fit 64 bits, 32 bits, 16 bits, 128 bits

Historically, hardware support for 32, 64 bits  
128 bit usually in software

Now 16 bit popular for ML

Lots of accelerators for matmul using  
16 bits

Some use `bfloat16` ("brain float")

eg Google TPUs (Tensor processing units)

Error analysis will only depend on knowing  
 $p$  and #bits in  $e$  for each format

For now, ignore "exceptions" i.e.  
numbers too big or too small to  
represent using  $e$

Normalization: use

$3.100e0$  or  $.0031e3$  ?

assume leading digit is non zero

Normalization  $\Rightarrow$  unique representation

and in binary  $\rightarrow$  leading digit = 1

$\rightarrow$  don't need to store it

$\rightarrow$  free extra bit!

"hidden bit"

Def:  $\text{rnd}(x) =$  nearest FP number to  $x$

(need to break ties: default rule:

"round to nearest even"  $\rightarrow$  round to nearest number ending in zero

good idea: unbiased, half time round up, half time down, more accurate for long sums

eg google Vancouver Stock Exchange

Def: Relative Representation Error

$$\text{RRE}(x) = \frac{|x - \text{rnd}(x)|}{|\text{rnd}(x)|}$$

Def: Maximum RRE =  $\max_{x \neq 0} \text{RRE}(x)$

aka machine epsilon, macheps,  $\epsilon$

(Matlab: their eps =  $2 \cdot \epsilon$ )

Max RRE = half distance from 1 to next larger FP number =  $1 + 2^{(1-p)}$   
 $= 2^{-p}$  in binary

Round off model (no over/underflow)

$$(*) \text{ fl}(a \text{ op } b) = \text{rnd}(a \text{ op } b)$$

= true result rounded to nearest even

$$= (a \text{ op } b)(1 + \delta), |\delta| \leq \epsilon$$

$$op \in \{+, -, *, /\}$$

(\*) a) so true for complex arithmetic  
with larger  $\varepsilon$  (see Q1.12 for details)

( C standard has 30+ lines of code  
for complex multiply )

Existing IEEE binary Formats

single (S) = 32 bits, double (D) = 64 bits

quad (Q) = 128 bits, half (H) = 16 bits

$$S: 32 \text{ bits} = 1 \text{ sign}$$

$$+ 8 \text{ exponent}$$

$$+ 23 \text{ mantissa}$$

$$\Rightarrow p = 1 + 23 = 24 \Rightarrow \varepsilon = 2^{-24} \sim 6 \cdot 10^{-8}$$

$$OV = \text{overflow threshold} \sim 2^{128} \sim 10^{38}$$

$$UN = \text{underflow} \quad " \quad \sim 2^{-126} \sim 10^{-38}$$

$$D: 64 = 1 + 11 + 52 \Rightarrow p = 53$$

$$\Rightarrow \varepsilon = 2^{-53} \sim 10^{-16}$$

$$OV \sim 2^{1024} \sim 10^{308}$$

$$UN \sim 2^{-1022} \sim 10^{-308}$$

$$Q: 128 = 1 + 15 + 112 \Rightarrow \varepsilon \sim 10^{-34}$$

$$OV \sim 10^{4932} \sim \frac{1}{UN}$$

$$H: 16 = 1 + 5 + 10 \Rightarrow \varepsilon \sim 5 \cdot 10^{-4}$$

$$OV \sim 10^4 \sim \frac{1}{UN}$$

$$\text{Bfloat: } 16 = 1 + \underline{8} + 7 \sim \varepsilon \sim 4 \cdot 10^{-3}$$

same as S

Higher precision packages for arbitrary precision in software, see web page (ARPREC, GMP, ...)

Error analysis for Horner:  $p(x) = \sum_{i=0}^d a_i \cdot x^i$

$$p = a_d$$

$$\text{for } i = d-1 : -1 : 0$$

$$p = x \cdot p + a_i$$

Label intermediate terms

$$p_d = a_d$$

$$\text{for } i = d-1 : -1 : 0$$

$$p_i = x \cdot p_{i+1} + a_i$$

Introduce roundoff

$$p_d = a_d$$

$$\text{for } i = d-1 : -1 : 0$$

$$p_i = [x \cdot p_{i+1} (1 + \delta_i) + a_i] (1 + \delta'_i)$$

$$|\delta_i| \leq \epsilon \quad |\delta'_i| \leq \epsilon$$

Simplify

$$p_0 = \sum_{i=0}^{d-1} a_i \cdot x^i \left[ (1 + \delta'_i) \prod_{j=0}^{i-1} (1 + \delta_j) (1 + \delta'_j) \right]$$

$$+ a_d \cdot x^d \left[ \prod_{j=0}^{d-1} (1 + \delta_j) (1 + \delta'_j) \right]$$

$$= \sum_{i=0}^{d-1} \left( \left[ \text{product of } 2i+1 \text{ terms like } (1 + \delta) \right] a_i \right) x^i$$

$$+ \left( \sum \text{product of } 2d \text{ terms like } (1+\delta) a_d \right) x^d$$

$$= \sum_{i=0}^d a'_i x^i \quad a'_i = a_i \left[ \text{at most } 2d \text{ terms product of like } (1+\delta) \right]$$

In words: Horner's rule is backward stable:  
returns exact value of a polynomial  
with slightly different coefficients

Simplify to get error bound:

$$\prod_{i=1}^n (1+\delta_i) \leq \prod_{i=1}^n (1+\varepsilon) = (1+\varepsilon)^n$$

$$= 1 + n\varepsilon + O(\varepsilon^2)$$

usually ignore  $O(\varepsilon^2)$

$$\leq 1 + \frac{n\varepsilon}{1-n\varepsilon} \quad \text{if } n\varepsilon < 1$$

... proof left to you!

$$\prod_{i=1}^n (1+\delta_i) \geq \prod_{i=1}^n (1-\varepsilon) = (-n\varepsilon + O(\varepsilon^2))$$

$$\geq 1 - \frac{n\varepsilon}{1-n\varepsilon}, \quad n\varepsilon < 1$$

$$\Rightarrow \left| \prod_{i=1}^n (1+\delta_i) - 1 \right| \leq n\varepsilon \quad (+ O(\varepsilon^2))$$

Apply to Horner:  $|\text{computed } p_0 - p(x)|$

$$\leq \sum_{i=0}^{d-1} (2i+1)\varepsilon |a_i x^i| + 2d\varepsilon |a_d \cdot x^d|$$

$$\text{relative error} = \frac{|\text{computed } p_0 - p(x)|}{|p(x)|}$$

$$\leq \frac{\sum_{i=0}^d |a_i x^i|}{|p(x)|} \cdot 2d\varepsilon$$

condition number
backward error

How many correct digits can I guarantee?

$k$  correct digits  $\Leftrightarrow$  relative error bound  $\leq 10^{-k}$

$$\Leftrightarrow -\log_{10}(\text{relative error}) \geq k$$

Modify Horner to get error bound

$$p = a_d, \text{ ebnrd} = |a_d|$$

for  $i = d-1 : -1 : 0$

$$p = x \cdot p + a_i, \text{ ebnrd} = |x| \cdot \text{ebnrd} + |a_i|$$

$$\text{ebnd} = \text{ebnrd} \cdot 2 \cdot d \cdot \varepsilon \dots \text{absolute bound}$$

## Details of FP:

- Important to understand to write reliable software
- lots of recent HW developments
- Analyzing code reliability hard
  - recent work on automatic tools to detect errors
- see posted notes

# ① Exception Handling - IEEE 754 rules:

Underflow: Tiny / Big = 0 or  
"subnormal": special numbers  
with smallest exponent, leading  
zeros in mantissa (not normalized)

$$0.00010.2^{\text{minexp}}$$

no hidden bit

Use case: if  $(x \neq y) \quad z = \frac{1}{x-y}$

no division by zero only if  
you have subnormal numbers

Overflow or divide by 0

$$\pm 1/0 = \pm \text{Inf} = \text{"Infinity"}$$

special bit pattern reserved

Natural rules:  $\text{Big} + \text{Big} = \text{Inf}$

$$3 - \text{Inf} = -\text{Inf}$$

$$7 / \text{Inf} = 0$$

Invalid  $0/0 = \text{"NaN"} = \text{Not a number}$

Rules:  $\text{Inf} - \text{Inf} = \text{NaN}$

$$\sqrt{-1} = \text{NaN}$$

$$3 + \text{NaN} = \text{NaN}$$



Flags available to check if  
Inf or NaN appeared