

Welcome to Ma221! Lec 4, Aug 30

Floating Point + Error Analysis

Last time: evaluated $(x-2)^{13}$ using
Horner's rule: $O(10^{-8})$ errors $\gg |p(x)|$
near root

Floating point: how to represent
real numbers

Long ago (<1985) most computers did
arithmetic differently \Rightarrow hard to
write portable code \Rightarrow IEEE 754
Standard Committee formed to create
a standard, led by Prof. Kahan

First Standard 1985, then 2008, then 2019,
forming committee for 2028 standard

Scientific Notation:

$\pm d.d\ldots d \cdot \text{radix}^e$

Usually radix = 2 (or 10 for finance)

store sign bit \pm

exponent (e)

mantissa d. d...d

$p = \# \text{ digits in mantissa}$

both p and # bits in e are limited to
16, 32, 64, 128 bits

Historically, only hardware for 32, 64

Lately 16 popular, for ML

lots of accelerators under construction

Also bfloat16 from Google

more exponential bits than IEEE16
smaller p

Latest committee: 8 bit arithmetic
 e_4p_4 and e_5p_3

For starters, ignore limits on e,
so assume no overflow, no underflow

Normalization: use

$3.100 \cdot 10^0$ not $0.0031 \cdot 10^3$

\Rightarrow unique representation

+ hardware simpler

+ in binary, leading d = 1

\Rightarrow "free bit", don't store it

(decimal: store 3 digits in 10 bits
not 1 digit in 4 bits)

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

(Default IEEE rounding! round to "nearest even", i.e last bit is zero if there is a tie
i.e $\frac{1}{2}$ time roundup, $\frac{1}{2}$ time down
 \Rightarrow unbiased)

Def: Relative Representation Error (RRE)

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

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

aka machine epsilon, macheps, ε
(note: matlab's $\text{eps} = 2\varepsilon$)

Max RRE = half distance from 1 to next larger floating point number
 $= 1 + (\text{radix})^{1-p}$

$$\text{MaxRRE} = 0.5 \cdot \text{radix}^{1-p} = 2^{-p} \text{ in binary}$$

Round off model (no over/underflow)

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

= true result rounded to
nearest even

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

op can be +, -, *, /

?54 also includes $\sqrt{}$, - . . .

(*) also true for complex arithmetic
with larger ε

Existing IEEE Binary Formats

half(H), single(S), double(D), quad(Q)

S: 32 bits = 1 (sign)

+8 (exponent bits)

+23 (mantissa bits)

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

-126 $\leq e \leq 127 \Rightarrow$ OV = overflow thresh.

$$\approx 2^{128} \approx 10^{38}$$

\Rightarrow UN = underflow thresh

$$= 2^{-126} \approx 10^{-38}$$

$$D: 64 = 1 + 11 + 52 \Rightarrow P = 53 \Rightarrow \varepsilon = 2^{-53} \approx 10^{-16}$$

$$-1022 \leq e \leq 1023 \quad OV \approx 2^{1024} \approx 10^{308}$$

$$UN = 2^{-1022} \approx 10^{-308}$$

$$Q: 128 = 1 + 15 + 112$$

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

$$OV \sim 10^4, VN \sim 10^{-4}$$

$$Bfloat16: 16 = 1 + \underbrace{8+7}_{\text{Same as single S}} \Rightarrow \varepsilon \sim 4 \cdot 10^{-3}$$

$$\text{Use } (*) \quad fl(a \oplus b) = (a \oplus b)(1 + \delta) \quad |\delta| \leq \varepsilon$$

for error analysis

$$\text{Horner's rule for } p(x) = \sum_{i=0}^d a_i \cdot x^i$$

$$p = a_d$$

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

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

Label intermediate terms

$$p_d = a_d$$

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

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

Introduce roundoff

$$p_d = a_d$$

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

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

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

Simplify:

$$P_0 = \sum_{i=0}^{d-1} a_i x^i \left[\left(1 + \delta'_i\right) \prod_{j=0}^{i-1} (1 + \delta'_j)(1 + \delta'_j) \right] \\ + a_d \cdot x^d \left[\prod_{j=0}^{d-1} (1 + \delta'_j)(1 + \delta'_j) \right]$$

How many factors $(1 + \delta)^2$?

$$= \sum_{i=0}^{d-1} a'_i x^i \quad \text{where } \leq 2d \\ a'_i = a_i \cdot (\text{factors } (+\delta))$$

Horner is backwards stable:

exact value or slightly
different polynomial

Simplify:

$$\prod_{i=1}^n (1 + \delta_i) \leq \prod_{i=1}^n (1 + \varepsilon) = (1 + \varepsilon)^n \\ = 1 + n\varepsilon + O(\varepsilon^2) \\ \leq 1 + \frac{n\varepsilon}{1 - n\varepsilon} \quad \text{if } n\varepsilon < 1$$

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

$$|\text{computed } p_d - p(x)| \leq \sum_{c=0}^{d-1} (2\epsilon + 1) |a_c x^c| + 2d\epsilon |a_d x^d|$$

$$\begin{aligned} \text{relerr} &= \frac{|\text{computed } p_d - p(x)|}{|p(x)|} \\ &\leq \frac{\sum_{c=0}^{d-1} |a_c x^c|}{|p(x)|} \cdot 2d\epsilon \\ &\quad \underbrace{|p(x)|}_{\substack{\text{condition} \\ \text{number}}} \quad \underbrace{2d\epsilon}_{\text{backward error}} \end{aligned}$$

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

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

Computer absolute bounds

$$p = ad, e_{bnd} = |ad|$$

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

$$p = x \cdot p + a_i, e_{bnd} = |x| \cdot e_{bnd} + |a_i|$$

$$e_{bnd} = e_{bnd} \cdot 2d\epsilon$$