

CS 450: Numerical Analysis

Chapter 1 – Scientific Computing

Lecture 2

Floating Point

Edgar Solomonik

Department of Computer Science
University of Illinois at Urbana-Champaign

Review

- ▶ Last lecture introduced the notions of *roundoff* and *truncation* error.
 - ▶ roundoff error concerns floating point error due to finite precision
 - ▶ truncation error concerns error incurred due to algorithmic approximation, e.g. the representation of a function by a finite Taylor series

$$f(x+h) \approx g_x(h) = \sum_{i=0}^k \frac{f^{(i)}(x)}{i!} h^i$$

truncation error is $|f(x+h) - g_x(h)| \leq \left| \sum_{i=k+1}^{\infty} \frac{f^{(i)}(x)}{i!} h^i \right| = O(h^{k+1})$

- ▶ To study the propagation of roundoff error in arithmetic we can use the notion of conditioning. as $h \rightarrow 0$

$$\kappa_x(f) = \lim_{h \rightarrow 0} \left| \frac{(f(x+h) - f(x)) / f(x)}{h/x} \right| = \left| \frac{f'(x)x}{f(x)} \right|$$

condition number of f at x

Floating Point Numbers

► Scientific Notation

$$2.13798 \cdot 10^{17}$$

exponent

significant digits (significand)

error limited to variations in the least significant digit

► Significand (Mantissa) and Exponent Given x with s leading bits x_0, \dots, x_{s-1}

$$x_0 \cdot x_1 \cdot x_2 \cdots x_{s-1} \cdot 2^{\underbrace{e_0 \dots e_k}_{k \text{ exponent bits}}} \Rightarrow \text{range } [2^{-2^k}, 2^{2^k}]$$

significand

UFL

representable normalized numbers

Rounding Error

- ▶ Maximum Relative Representation Error (Machine Epsilon)

$$\epsilon_{\text{mach}} = 2^{1-s} \quad s \text{ is \# digits in the significand}$$

$$\forall x \in [2^{-2^k}, 2^{2^k}]$$

$$\frac{|f(x) - x|}{|x|} \leq \epsilon_{\text{mach}}$$

Rounding Error in Operations (I)

► Addition and Subtraction

(subtraction is addition with sign bit flipped)

Catastrophic cancellation:

$$\begin{array}{r} \underline{1.37284104} \times 10^{-6} \\ - \underline{1.37283002} \times 10^{-6} \\ \hline \end{array} \left. \vphantom{\begin{array}{r} \underline{1.37284104} \times 10^{-6} \\ - \underline{1.37283002} \times 10^{-6} \end{array}} \right\} 9 \text{ digits of accuracy}$$

$$\rightarrow \underline{1.10200000} \times 10^{-11} \left. \vphantom{\underline{1.10200000} \times 10^{-11}} \right\} 4 \text{ digits of accuracy}$$

absolute error stayed the same,
but magnitude of value decreased significantly

→ large relative error

rel. err. $\leq 10^{-8}$
rel. err. $\leq 10^{-3}$

Rounding Error in Operations (II)

► Multiplication and Division

$$\begin{array}{r} 2.1314 \times 10^3 \\ \times 7.8912 \times 10^{-6} \\ \hline \end{array}$$

$$\underbrace{x_1 \cdot x_2 \cdot x_3 \cdot x_4 \cdot x_5}_{\text{5 digits}} \times 10^{-3}$$

$$\frac{|f|(f(x) \cdot f(y)) - x \cdot y|}{|x \cdot y|} = \text{rel. err}$$

$$\text{rel. err} \leq \frac{|(1 + \epsilon_{\text{mach}}) \left((1 + \epsilon_{\text{mach}}) x \cdot (1 + \epsilon_{\text{mach}}) y \right) - x y|}{|x y|}$$

$$= \frac{|(1 + \epsilon_{\text{mach}})^3 x y - x y|}{|x y|} \approx \frac{|(1 + 3\epsilon_{\text{mach}}) x y - x y|}{|x y|} \leq 3\epsilon_{\text{mach}}$$

Exceptional and Subnormal Numbers

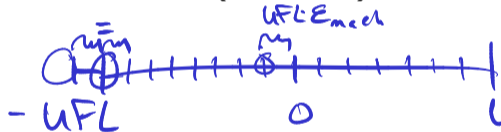
- ▶ **Exceptional Numbers** = not normalized

$$1/0 = \text{INF} \quad | \quad 0$$

$$0/0 = \infty - \infty = \text{NaN}$$

norm. $\rightarrow 1.x_1 \dots x_{s-1} \times 2^s = \# \text{ digits in significant}$

- ▶ **Subnormal (Denormal) Number Range**



$$0.x_1 x_2 x_3 \dots x_{s-1} \times 2^{-2^k}$$

smallest subnormal is $0.0 \dots 01 \cdot 2^{-2^k}$

UFL \leftarrow underflow limit for normalized numbers 2^{-2^k}

- ▶ Gradual Underflow: Avoiding underflow in addition

underflow means result is not a normalized machine number