# CS 450: Numerical Anlaysis

## Chapter 1 – Scientific Computing
### Lecture 2
### Floating Point
### Vector Norms

Edgar Solomonik

Department of Computer Science
University of Illinois at Urbana-Champaign

January 20, 2018

# Floating Point Numbers

- **Scientific Notation**

$$3.124 \times 10^8 \qquad [3.123, \quad 3.125] \times 10^8$$

$$2.103 \times 10^{-3} \qquad [2.102, \quad 2.104] \times 10^8$$

rel error is $10^{-3}$

- **Significand (Mantissa) and Exponent**

$$1.01011 \times 2^4 \qquad \text{normalized}$$

significand

# Rounding Error

▶ **Maximum Relative Representation Error (Machine Epsilon)**

$$\text{\# digits in significand} = \boxed{12}$$

$$\text{exponent range} \sim [-1023, 1024]$$

$1.\underbrace{01101101}_{\phantom{x}}\!\!m$

$\underbrace{\phantom{xxxx}}_{2^{-7}}$
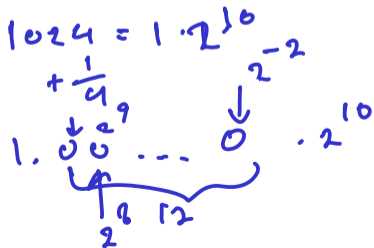
$$UFL = 1 \cdot 2^{-1023}$$

$$OFL = 1 \cdot 2^{1024}$$

$$\text{Machine epsilon} = 2^{-12} = \varepsilon_{mac}$$

smallest number we can
add to 1, so that $fl(1 + \varepsilon_{mac}) \neq 1$

▶ **Rounding Error Analysis**

$$1024 = 1 \cdot 2^{10}$$
$$+ \frac{1}{4} \qquad 2^{-2}$$

$1.\underset{\underset{2^8 \ 12}{\uparrow}}{0 0} \cdots 0 \cdot 2^{10}$

$1.000 \ldots 0$
$+ 0.000 \ldots 1 \quad - \varepsilon_{mac}$
$1.000 \ldots 0 \cdot 2^{-12}$

# Rounding Error in Operations

▶ **Addition and Subtraction**

$$1.0110 \longrightarrow \text{approximate to relative precision}$$
$$-1.0100 \qquad\qquad\qquad\qquad\qquad 2^{-4}$$

$$= 0.0010$$
$$= 1.0000 \cdot 2^{-3}$$

▶ **Multiplication and Division**

$$z = x \cdot y = 1.1010\underset{\sim}{\phantom{0}}$$

$$1.0111 \qquad\qquad 1.1100$$

Overflow leads to a result of $\pm\infty$

$$0/0 = \infty + (-\infty) = NaN$$

# Subnormal Numbers

- **Subnormal (Denormal) Number Range**



smallest subnormal

$$0.00001 \cdot 2^{-\beta}$$

smallest normalized

$$1.0000 \cdot 2^{-\beta} = \varepsilon_{nee} \cdot 2^{-\beta}$$

$$0.00100 \cdot 2^{-\beta}$$

$$0.01001 \cdot 2^{-\beta}$$

- **Gradual Underflow: Avoiding underflow in addition**

$$x - y = 0 \quad \text{if and only if} \quad x = y$$