

# CS 450: Numerical Analysis<sup>1</sup>

## Introduction to Scientific Computing

University of Illinois at Urbana-Champaign

---

<sup>1</sup>*These slides have been drafted by Edgar Solomonik as lecture templates and supplementary material for the book “Scientific Computing: An Introductory Survey” by Michael T. Heath ([slides](#)).*

# Scientific Computing Applications and Context

- ▶ **Mathematical modelling for computational science** *Typical scientific computing problems are numerical solutions to PDEs*
  - ▶ *Newtonian dynamics: simulating particle systems in time*
  - ▶ *Fluid and air flow models for engineering*
  - ▶ *PDE-constrained numerical optimization: finding optimal configurations (used in engineering of control systems)*
  - ▶ *Quantum chemistry (electronic structure calculations): many-electron Schrödinger equation*
- ▶ **Linear algebra and computation**
  - ▶ *Linear algebra and numerical optimization are building blocks for machine learning methods and data analysis*
  - ▶ *Computer architecture, compilers, and parallel computing use numerical algorithms (matrix multiplication, Gaussian elimination) as benchmarks*

## Example: Mechanics<sup>2</sup>

- ▶ Newton's laws provide incomplete particle-centric picture
- ▶ Physical systems can be described in terms of *degrees of freedom* (DoFs)
  - ▶ A piston moving up and down requires \_\_\_\_\_ DoFs
  - ▶ 1-particle system requires \_\_\_\_\_ DoFs
  - ▶ 2-particle system requires \_\_\_\_\_ DoFs
  - ▶ 2-particles at a fixed distance require \_\_\_\_\_ DoFs
- ▶  $N$ -particle system *configuration* described by  $3N$  DoFs

---

<sup>2</sup>*Variational Principles of Mechanics*, Cornelius Lanczos, Dover Books on Physics, 1949.

## **Course Structure**

- ▶ **Complex numerical problems are generally reduced to simpler problems**
  
  
  
  
  
  
  
  
  
  
- ▶ **The course topics will follow this hierarchical structure**

# Numerical Analysis

- ▶ **Numerical Problems involving Continuous Phenomena:**

- ▶ **Error Analysis:**

## Sources of Error

- ▶ **Representation of Numbers:**

- ▶ **Propagated Data Error:**

- ▶ **Computational Error =  $\hat{f}(x) - f(x) = \text{Truncation Error} + \text{Rounding Error}$**

# Error Analysis

▶ **Forward Error:**

▶ **Backward Error:**

# Visualization of Forward and Backward Error

$$A\hat{x} = b$$

$$f(A, b) = x$$

$$\hat{f}(A, b) = \hat{x}$$

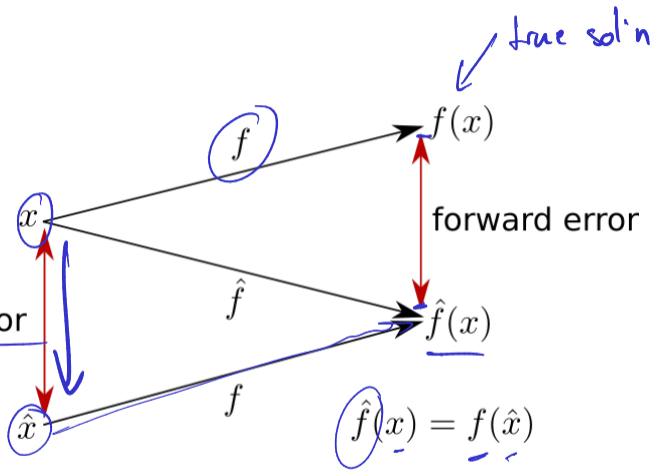
forward error is  $\hat{x} - x$

backward error  
mag... hold

backward error

$$A\hat{x} - b = A\hat{x} - Ax = \delta b$$

$$f(A, b + \delta b) = \hat{x}$$



$$\hat{f}(x) = f(\hat{x})$$



# Conditioning

► Absolute Condition Number:  $\geq 0$

evaluate  $f$  at  $x$ ,

$$k(f, x) = |f'(x)|$$

evaluate  $f$  at  $y$

$$k(f, D) = \max_{x \in D} k(f, x)$$

for some  $x \in D$

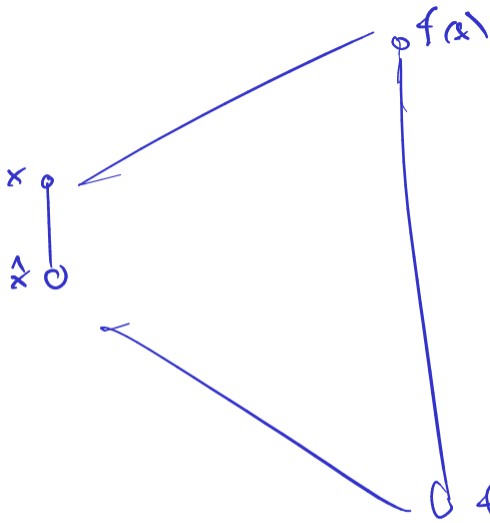
$$\lim_{\Delta x \rightarrow 0} \frac{|f(x+\Delta x) - f(x)|}{\Delta x} \rightarrow f'(x)$$

► (Relative) Condition Number:

relative distance  $\left\{ \begin{array}{l} \uparrow \\ x \\ \uparrow \\ x \\ \uparrow \\ x \end{array} \right\} \rightarrow x - x = \Delta$

$\left\{ \begin{array}{l} \uparrow \\ f(x) \\ \uparrow \\ f(x) \end{array} \right\}$  relative distance

$$\lim_{\Delta x \rightarrow 0} \frac{|f(x+\Delta x) - f(x)| / |f(x)|}{|\Delta x| / |x|}$$



amplification factor  
→ bounded by

$$K(f, x)$$

$$|\text{forward error}| \leq K(f, x) |\text{backward error}|$$

# Posedness and Conditioning

- ▶ What is the condition number of an ill-posed problem?

$$\kappa(f, x) = \infty$$

$f$   
↓  
soln exists is unique and  
•  $f$  changes continuously with  $x$

# Stability and Accuracy

► **Accuracy:** forward error is small

$$\hat{f}(x) \approx f(x) \text{ for all } x \in D$$

↑ algorithm in exact arithmetic (no round-off error)      ↑ input domain

► **Stability:** sensibility of the algorithm to round-off error

$$\tilde{f}(x) \approx \hat{f}(x)$$

↑  $\neq \hat{f}(f(x))$

$\tilde{f}$  is the algorithm in finite precision

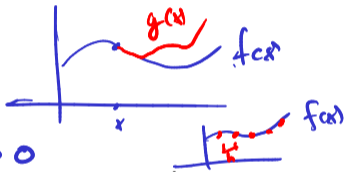
# Error and Conditioning

$$\tilde{f} \approx \hat{f} \quad \hat{f} \approx f$$

- ▶ Two major sources of error: 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

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

$$f(x+h) - q(h) = \sum_{i=k+1}^{\infty} \frac{f^{(i)}(x)}{i!} h^i = O(h^{k+1}) \quad \text{as } h \rightarrow 0$$



- ▶ To study the propagation of roundoff error in arithmetic we can use the notion of conditioning.

$$f(f_1(x))$$

↑

x rounded to  
the nearest floating-point number

# Floating Point Numbers

Demo: Picking apart a floating point number

Demo: Density of Floating Point Numbers

## ► Scientific Notation

$$\underbrace{2.103}_{\text{significant}} \times 10^{\underbrace{7}_{\text{exponent}}}$$

have uniformly low relative error in representation (rounding error)

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

# bits in the significant and the range of exponents

$$1.\underbrace{1100101}_{\text{significant}} \times 2^{\underline{17}}$$

normalized floating point number  
store only bits right of the decimal place

$$[-L, L]$$
$$\frac{-L}{2} \quad \frac{L}{2}$$

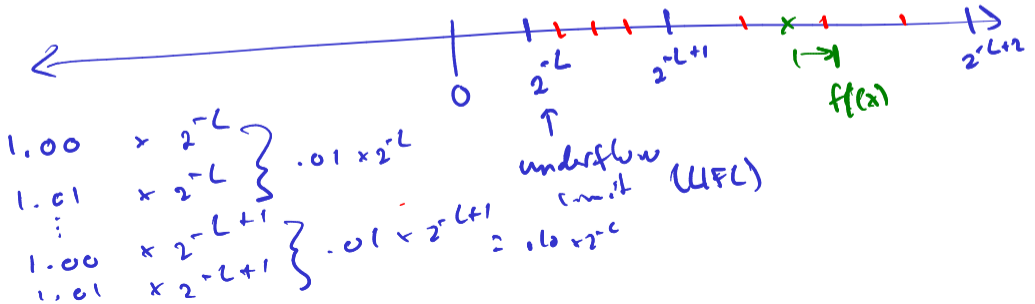
f-p system

1-bit for sign

$$\epsilon_{\text{mach}} = 2^{-b} = .01 \times 2^0 = .01$$

for  $x \in [2^{-L}, 2^L]$ ,  $\left| \frac{f(x) - x}{x} \right| \leq \epsilon_{\text{mach}}$

b-bits for the significand and (b+1 bits of accuracy)  
log<sub>2</sub>(L)-bits of exponent, to represent numbers in  $[2^{-L}, 2^L]$



# Rounding Error

*Demo: Floating point and the Harmonic Series*

*Demo: Floating Point and the Series for the Exponential Function*

## ► Maximum Relative Representation Error (Machine Epsilon)

$$\epsilon_{\text{mach}} = \frac{\text{argmin}_x |x(1+x)|}{x} \neq 1$$

---

$$\begin{array}{r} 1.000 \\ \Downarrow \\ 1.0001 \\ 2^{-6} \end{array} + \begin{array}{r} .0001 \\ 1.0001 \end{array}$$



# Rounding Error in Operations (I)

Demo: Catastrophic Cancellation

Activity: Cancellation in Standard Deviation Computation

## ► Addition and Subtraction

subtraction  $\rightarrow$  addition with a negation of an operand

$$x - y = x + (-y)$$

$$\underbrace{3.124}_{4 \text{ digits of accuracy}} + \underbrace{(-3.102)}_{2 \text{ digits of accuracy}} = \underbrace{.022}_{2 \text{ digits of accuracy}} \Rightarrow \frac{2.2 \times 10^{-2}}{2.00 \times 10^2}$$

catastrophic cancellation

$$\frac{|(x+y) - (f(x) + f(y))|}{|x+y|} \leq \frac{\epsilon(|x| + |y|)}{|x+y|}$$

addition for arbitrary  $x, y$  is ill-posed

## Rounding Error in Operations (II)

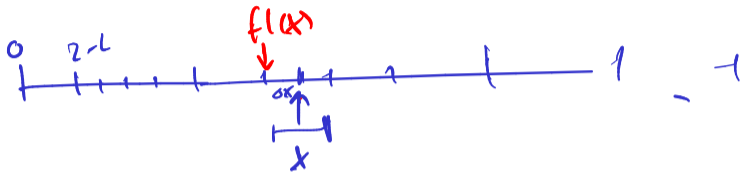
► Multiplication and Division

multiplication by the reciprocal

$$x/y = x \cdot \underline{\underline{1/y}}$$

$$\begin{aligned} & 3.12 \times 10^2 \\ & \cancel{1.23} \times 10^{-3} \\ & = 3. \times 10^{-1} \end{aligned}$$

$$\frac{|xy - f(f(x)f(y))|}{|xy|} = \frac{|xy - (1+\epsilon)^3 xy|}{|xy|} \approx 3\epsilon$$



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

$$\left| \frac{(x+y) - (f(x) + f(y))}{|x+y|} \right| \leq$$

$$\frac{|\delta x| + |\delta y|}{|x+y|} \leq \frac{\epsilon|x| + \epsilon|y|}{|x+y|}$$

# Exceptional and Subnormal Numbers

- ▶ **Exceptional Numbers**
- ▶ **Subnormal (Denormal) Number Range**
- ▶ **Gradual Underflow: Avoiding underflow in addition**

# Floating Point Number Line

