

CS 357 – Numerical Methods 1

Review

Numerical methods

Method = Math + Complexity + Accuracy

Real life
problems



Mathematical
formulation

+

Numerical Experiments



Simulation of physical reality

MODEL, NUMBERS, ERRORS

Random Numbers

Making Models using Monte Carlo simulation

Errors, accuracy, convergence

Truncation (Taylor Series)

Rounding (Floating Point)

NORMS
CONDITIONING

LINEAR SYSTEMS

EIGENVALUE (Markov Chains)

SVD (low-rank app.)

INTERPOLATION (fit data "exactly" to a fn)

LINEAR OPERATORS

MODEL SET OF DATA POINTS

LINEAR LEAST-SQUARES (SVD) ("best" data fit to a trend)

ITERATIVE METHODS

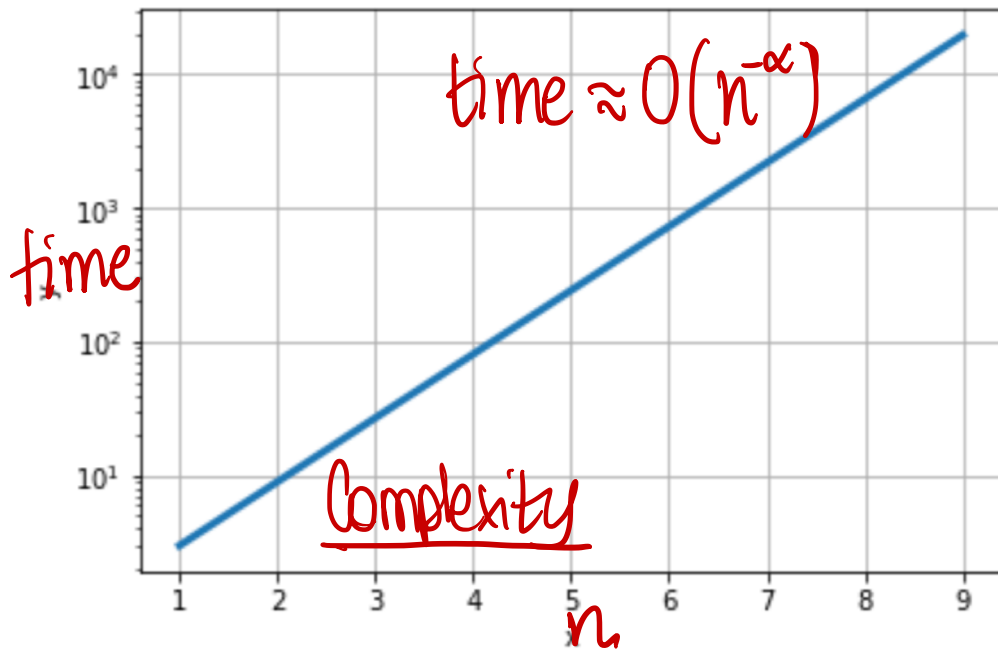
NONLINEAR SYSTEM OF EQUATIONS

- BISECTION) 1D
- SECANT) 1D
- NEWTON) 1D
- BROYDEN) ND
- NEWTON) ND

OPTIMIZATION

- GOLDEN SECTION) 1D
- NEWTON) 1D

- STEEPEST DESCENT/ND
NONLINEAR LEAST-SQUARES

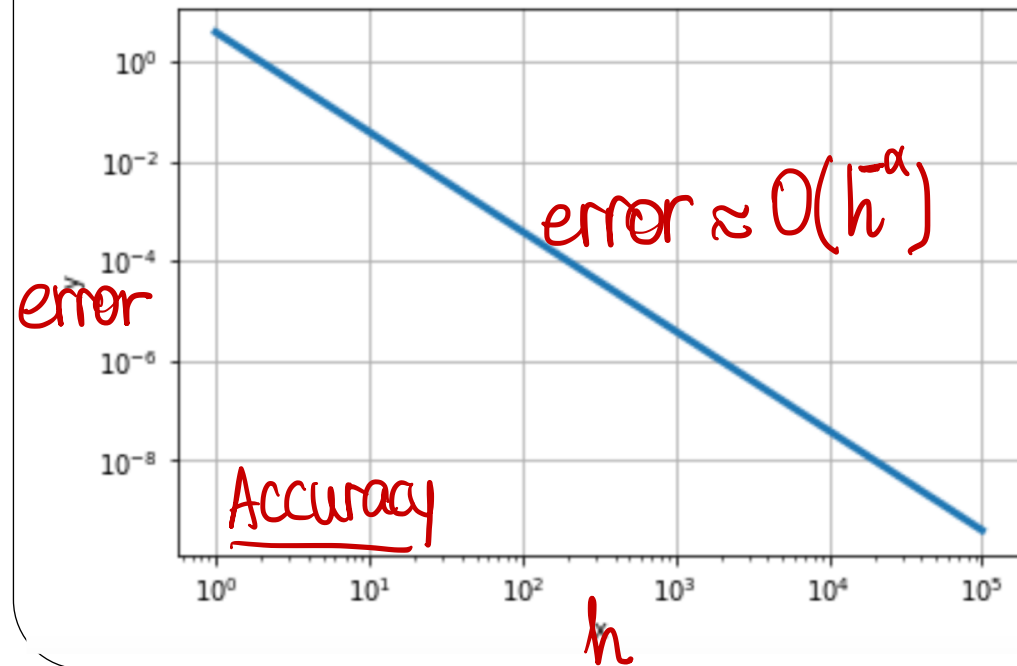


→ Exponential function

$$y = e^{\alpha x}$$

$$\log y = \log(e^{\alpha x})$$

$$\bar{y} = \alpha x$$



→ Power function

$$y = x^{\alpha}$$

$$\log y = \log(x^{\alpha}) = \alpha \log(x)$$

$$\bar{y} = \alpha \bar{x}$$

Let's talk about plots...

- Power functions: $y = a x^b$

$$\log y = \log(a x^b) = \log(a) + \log(x^b) = \log(a) + b \log(x)$$

$$\bar{y} = \bar{a} + b \bar{x}$$

- Exponential functions: $y = a b^x$

$$\log y = \log(a b^x) = \log(a) + \log(b^x) = \log(a) + x \log(b)$$

$$\bar{y} = \bar{a} + \bar{b} x$$

- Log functions: $y = a \log(b x)$

$$y = a \log(b) + a \log(x)$$

$$y = \bar{b} + a \bar{x}$$

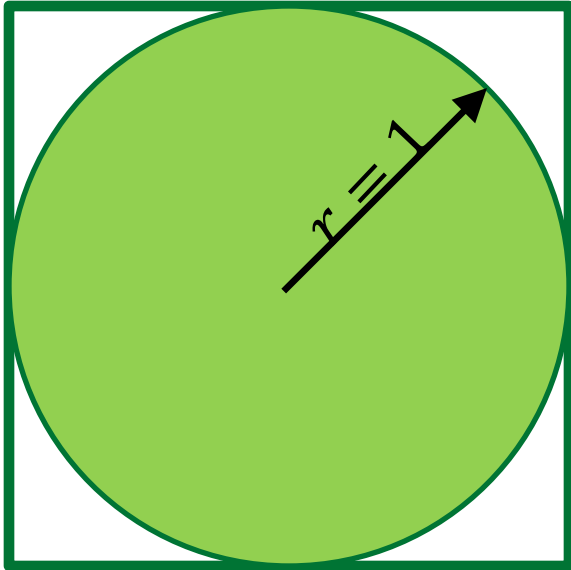
Random numbers

- Pseudo-random numbers
 - Numbers and sequences appear random, but they are in fact reproducible
 - Good for algorithm development and debugging
- Desired properties for a good random number generator
 - Random pattern
 - Long period
 - Efficiency
 - Repeatability
 - Portability

Monte Carlo Methods

- Algorithms that compute approximations of desired quantities based on repeated randomized sampling
- Typically used to model:
 - Nondeterministic problems
 - Complicated deterministic problems, specially in high dimensions
- Asymptotic behavior is $O\left(\frac{1}{\sqrt{n}}\right)$ when $n \rightarrow \infty$, where n is the number of samples
 $\sim O(n^{-0.5})$

Using Monte Carlo to approximate integrals



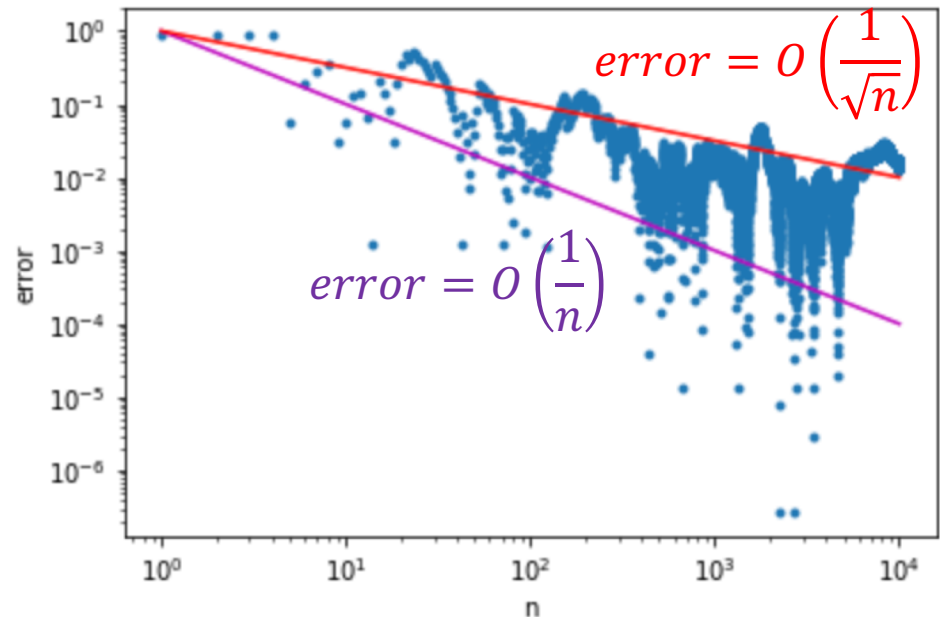
We sample points uniformly inside the domain $D = [x_0, x_1] \times [y_0, y_1]$

N_c : number of points inside circle

N_s : number of points inside square = total number of sample points = n

$$A_c = A_s \left(\frac{N_c}{n} \right)$$

- CONS: Slow convergence rate
- PROS: Efficiency does not degrade when increasing the dimension of the problem



Error in Numerical Methods

Absolute error: $|x - \bar{x}|$

Relative error: $\frac{|x - \bar{x}|}{|x|}$

Accurate to n significant digits means that you can trust a total of n digits. *Accurate digits* is a measure of relative error.

Relative error: $error = \frac{|x_{exact} - x_{approx}|}{|x_{exact}|} \leq 10^{-n+1}$

n is the number of accurate significant digits

3 sig. figures \longrightarrow error $\leq 10^{-3+1}$ \longrightarrow error $\leq 10^{-2}$ (1%)

Rates of convergence or growth

- $error = O(n^{-\alpha})$ or $complexity = O(n^{\alpha})$
Power function – straight line in a log-log plot
Algebraic convergence/growth
- $error = O(e^{-\alpha n})$ or $complexity = O(e^{\alpha n})$
Exponential function – straight line in a linear-log plot
Exponential convergence/growth

Taylor Series

The Taylor Series approximation about point x_0 is given by:

$$f(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2 + \frac{f'''(0)}{3!}(x - x_0)^3 + \dots$$

$$f(x) = \sum_{i=0}^{\infty} \frac{f^{(i)}(x_0)}{i!} (x - x_0)^i$$

The Taylor Series approximation of degree n , $\hat{f}(x)$, is given by:

$$\hat{f}(x) = \sum_{i=0}^n \frac{f^{(i)}(x_0)}{i!} (x - x_0)^i$$

→ truncated function
 $h = (x - x_0)$

And the Taylor error of degree n due the truncation is: $O(h^{n+1})$

$$R = f(x) - \hat{f}(x) = \sum_{i=n+1}^{\infty} \frac{f^{(i)}(x_0) h^i}{i!} \leq \frac{f^{(n+1)}(\xi)}{(n+1)!} h^{(n+1)}$$

Taylor series for e^x about $x_0=0$ is

$$e^x \approx 1 + \frac{x}{1!} + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots$$

1) Approximate e^x at $x=0.2$ using Taylor expansion of degree 2

$$1 + \frac{0.2}{1!} + \frac{(0.2)^2}{2!} = 1.22$$

2) Approximate $\frac{d}{dx}(e^x)$ at $x=0.1$ using Taylor expansion of degree 3

$$1 + \frac{2x}{2!} + \frac{3x^2}{3!} = 1 + 0.1 + \frac{(0.1)^2}{2} = 1.105$$

3) If we want to use the first 4 terms of the expansion, what is the order of the error of the Taylor approximation?

$h = x - x_0$ error $\approx O(h^{n+1})$

(A) 2

(B) 3

(C) 4

(D) 5

(E) 6

Normalized Floating-point numbers

Normalized floating point numbers are expressed as

$$x = \pm q \times 2^m = \pm 1.b_1b_2b_3 \dots b_n \times 2^m = \pm 1.f \times 2^m$$

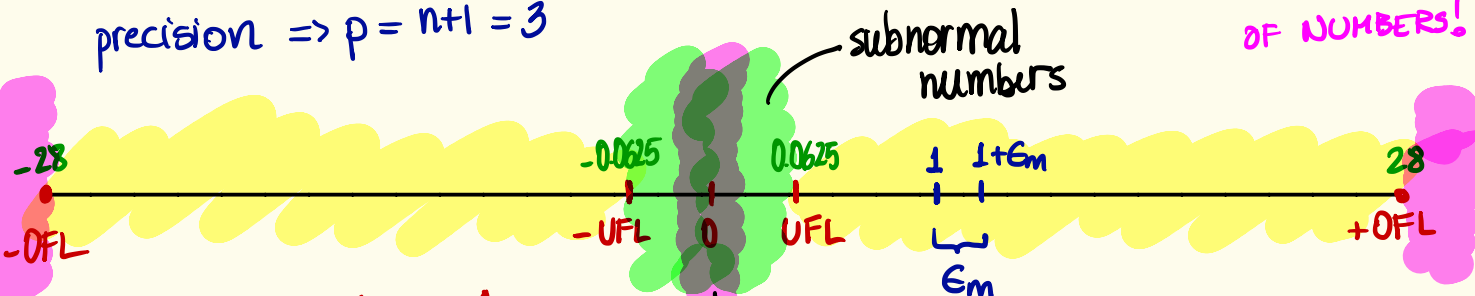
where f is the fractional part of the significand, m is the exponent and $b_i \in \{0,1\}$.

- Exponent range: $[L, U]$
- Precision: $p = n + 1$
- Smallest positive normalized FP number: $\text{UFL} = 2^L$
- Largest positive normalized FP number: $\text{OFL} = 2^{U+1}(1 - 2^{-p})$
- Machine epsilon (ϵ_m): is defined as the distance (gap) between 1 and the next largest floating point number.
- Noticeable gap around zero, present in any floating system, due to normalization
 - ✓ The smallest possible significand is $1.00 \dots 0$
 - ✓ The smallest possible exponent is L
- Relax the requirement of normalization, and allow the leading digit to be zero, only when the exponent is at its minimum ($m = L$)
- Computations with subnormal numbers are often slow.

$$x = \pm 1.b_1 b_2 \times 2^m \quad m \in [-4, 4] \quad b_i \in \{0, 1\}$$

precision $\Rightarrow p = n+1 = 3$

ONLY FINITE REPRESENTATION OF NUMBERS!



$$UFL = (1.00)_2 \times 2^{-4} = 2^{-4} = 0.0625$$

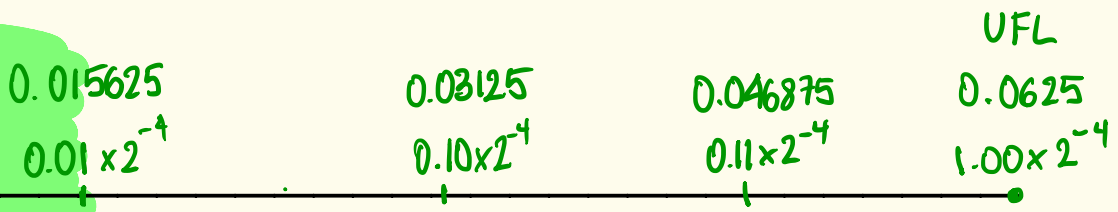
$$OFL = (1.11)_2 \times 2^4 = 28$$

Machine epsilon: gap between 1 and next largest fp number

$$\epsilon_m = 1.00 \times 2^0 - 1.01 \times 2^0 = 0.01 \times 2^0 = 0.25$$

Subnormal numbers \rightarrow let the leading digit $\neq 1$ be zero if $m = L$

reduce the gap around zero



What is machine epsilon double precision?

$$\pm 1. \underbrace{b_1 b_2 b_3 \dots b_{52}}_{52 \text{ bits in the fraction}} \times 2^m \quad \begin{array}{l} \text{1 bit} \\ \text{11 bits} \end{array}$$

$$1.000 \dots 0 \times 2^0 = 1$$

$$1.000 \dots 1 \times 2^0 = \text{next fp number after 1}$$

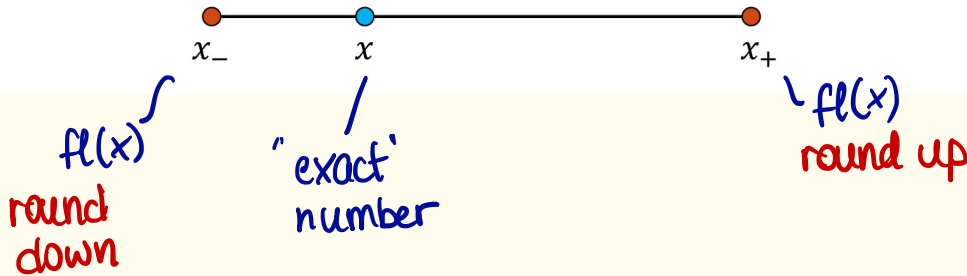
$$\epsilon_m = 0.000 \dots \underset{\substack{\uparrow \\ 2^{-52}}}{1} \times 2^0 \quad \Rightarrow \quad \boxed{\epsilon_m = 2^{-52}}$$

Rounding

The process of replacing x by a nearby machine number \tilde{x} is called rounding, and the error involved is called **roundoff error**.

$$x = \pm 1.b_1b_2b_3 \dots b_n \dots \times 2^m$$

round to nearest



Relative error due to rounding:

$$e_r = \frac{|fl(x) - x|}{|x|} \leq \epsilon_m$$

Absolute error:

$$e_a = |fl(x) - x| \leq \epsilon_m |x|$$

single precision

$$e_r \leq 2^{-23} \approx 10^{-7}$$

(7 decimal accurate digits)

double precision

$$e_r \leq 2^{-52} \approx 10^{-16}$$

(16 decimal accurate digits)

Loss of Significance

Results from floating point arithmetic used by computers, and the numbers of significant digits is substantially reduced.

Consider the decimal representation for π using 20 significant figures

$$a = 3.1415926535897932385;$$

Consider the decimal representation for π using 10 significant figures

$$b = 3.141592654;$$

Now we perform the following calculations:

$$a - 3.14159265300000000000 \text{ (* calculation gives 9 significant digits *)}$$

$$= 5.89793239 \times 10^{-10}$$

$$b - 3.141592653 \text{ (* calculation gives 1 significant digit *)}$$

$$= 1. \times 10^{-9}$$

$$f(x) = \sqrt{x^2 + 1} - 1 = \frac{x^2}{\sqrt{x^2 + 1} + 1}$$

Using five-decimal-digit arithmetic:

$$f(10^{-3}) = \sqrt{(10^{-3})^2 + 1} - 1 = 0$$

$$f(10^{-3}) = \frac{(10^{-3})^2}{\sqrt{(10^{-3})^2 + 1} + 1} = \frac{10^{-6}}{2}$$

Norms

What's a norm?

- A generalization of 'absolute value' to vectors.
- $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}_0^+$, returns a 'magnitude' of the input vector
- In symbols: Often written $\|\mathbf{x}\|$.

Define **norm**.

A function $\|\mathbf{x}\| : \mathbb{R}^n \rightarrow \mathbb{R}_0^+$ is called a norm if and only if

1. $\|\mathbf{x}\| > 0 \Leftrightarrow \mathbf{x} \neq \mathbf{0}$.
2. $\|\gamma\mathbf{x}\| = |\gamma| \|\mathbf{x}\|$ for all scalars γ .
3. Obeys triangle inequality $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$

Properties of Matrix Norms

Matrix norms inherit the vector norm properties:

1. $\|A\| > 0 \Leftrightarrow A \neq \mathbf{0}$.
2. $\|\gamma A\| = |\gamma| \|A\|$ for all scalars γ .
3. Obeys triangle inequality $\|A + B\| \leq \|A\| + \|B\|$

But also some more properties that stem from our definition:

1. $\|A\mathbf{x}\| \leq \|A\| \|\mathbf{x}\|$
2. $\|AB\| \leq \|A\| \|B\|$ (easy consequence)

Both of these are called **submultiplicativity** of the matrix norm.

Vectors \underline{v} :
($n \times 1$)

$$\|\underline{v}\|_p = \left(|v_1|^p + |v_2|^p + \dots + |v_n|^p \right)^{1/p} \quad p \geq 1$$

$$\|\underline{v}\|_1 = |v_1| + |v_2| + \dots + |v_n|$$

$$\|\underline{v}\|_2 = \sqrt{v_1^2 + v_2^2 + \dots + v_n^2}$$

$$\|\underline{v}\|_\infty = \left(|v_1|^\infty + |v_2|^\infty + \dots + |v_n|^\infty \right)^{1/\infty} = \max_i |v_i|$$

Note that we take the absolute value of the components!

Matrices A :

$$\|A\| = \max_{\|x\|=1} \|Ax\| \quad \text{or} \quad \|A\| = \max_{\|x\|=1} \frac{\|Ax\|}{\|x\|}$$

$$A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$$

Row sums: 3, 7
Column sums: 4, 6

$$\|A\|_1 = \max \text{ absolute column sum}$$

$$\|A\|_\infty = \max \text{ absolute row sum}$$

$$\|A\|_2 = \max \text{ singular value}$$

$$\|A\|_1 = 6 \quad \|A\|_\infty = 7$$

$$B = \begin{bmatrix} -1 & 2 \\ 3 & -4 \end{bmatrix}$$

$$\|A\|_1 = \|B\|_1$$


$$\|A\|_\infty = \|B\|_\infty$$

Linear System of Equations


- Solve $\underline{A} \underline{x} = \underline{b}$

- Factorize \underline{A} : LU factorization $\Rightarrow \underline{A} = \underline{L} \underline{U}$ ($O(n^3)$)

$$\underline{L} \underline{U} \underline{x} = \underline{b} \Rightarrow \underline{L} \underline{y} = \underline{b}$$

 $[\underline{y}] = [\underline{b}]$
Forward substitution
 $O(n^2)$

$$\underline{U} \underline{x} = \underline{y}$$

 $[\underline{x}] = [\underline{y}]$
backward substitution
 $O(n^2)$

- Solve $\underline{A} \underline{x}_i = \underline{b}_i$
for different RHS ($i=1, \dots, k$) \Rightarrow cost $\cong O(n^3 + 2kn^2)$ and NOT $O(kn^3)$

- Partial Pivoting: avoids division by zero

Find largest entry in the column (absolute value) and swap it with top row

Partial Pivoting

2 points

Consider LU factorization on the following matrix:

$$A = \begin{bmatrix} 5 & 3 & 2 & 9 \\ -10 & 2 & 1 & 5 \\ 0 & 3 & 3 & 3 \\ 1 & 2 & 1 & 2 \end{bmatrix}$$

If partial pivoting is employed during elimination of the first column, what matrix value should be chosen as the pivot?

Answer*

Let's assume that when solving the system of equations $\mathbf{K} \mathbf{U} = \mathbf{F}$, we observe the following:

- When \mathbf{K} has dimensions (100,100), computing the LU factorization takes about 1 second and each solve (forward + backward substitution) takes about 0.01 seconds.

Estimate the total time it will take to find the solution \mathbf{U} corresponding to 10 different right-hand sides \mathbf{F} when the matrix has dimensions (1000,1000)?

LU factorization : $(100)^3 \text{ — } 1 \text{ sec}$ $\rightarrow x = \left(\frac{1000}{100}\right)^3 = 10^3 \text{ seconds}$
 $(1000)^3 \text{ — } x$

One solve : $(100)^2 \text{ — } 0.01 \text{ sec}$ $\rightarrow y = \left(\frac{1000}{100}\right)^2 0.01 = 10^2 (0.01) = 1 \text{ sec}$
 $(1000)^2 \text{ — } y$

total time = $10^3 \text{ seconds} + 10(1 \text{ sec}) \approx 10^3 \text{ seconds}$

Condition number

Small condition numbers mean not a lot of error amplification. Small condition numbers are good!

$$\frac{\|\Delta \mathbf{x}\|}{\|\mathbf{x}\|} \leq \|\mathbf{A}^{-1}\| \|\mathbf{A}\| \frac{\|\Delta \mathbf{b}\|}{\|\mathbf{b}\|} = \text{cond}(\mathbf{A}) \frac{\|\Delta \mathbf{b}\|}{\|\mathbf{b}\|} \qquad \frac{\|\Delta \mathbf{x}\|}{\|\mathbf{x}\|} \leq \|\mathbf{A}^{-1}\| \|\mathbf{A}\| \frac{\|\mathbf{E}\|}{\|\mathbf{A}\|}$$

When solving linear system of equations, the residual is $\mathbf{r} = \mathbf{b} - \mathbf{A} \hat{\mathbf{x}}$, where $\hat{\mathbf{x}} = (\mathbf{x} + \Delta \mathbf{x})$ is the solution of the perturbed problem. For well-conditioned matrices, small relative residual implies small relative error.

$$\frac{\|\Delta \mathbf{x}\|}{\|\mathbf{x}\|} \leq \text{cond}(\mathbf{A}) \frac{\|\mathbf{r}\|}{\|\mathbf{A}\| \|\mathbf{x}\|}$$

Gaussian elimination with partial pivoting (where c is small) yields **small relative residual regardless of conditioning of the system**

$$\frac{\|\mathbf{r}\|}{\|\mathbf{A}\| \|\hat{\mathbf{x}}\|} \leq c \epsilon_m$$

Rule of thumb: If the entries in \mathbf{A} and \mathbf{b} are accurate to S decimal digits, then a condition number 10^W reduces W digits in the accuracy of the computed solution (i.e., solution will have $(S - W)$ accurate digits).

When solving a system of linear equations via LU with partial pivoting, which of the following is guaranteed to be small?

A) Relative residual: $\frac{\|r\|}{\|A\|\|x\|}$

B) Relative error: $\frac{\|\Delta x\|}{\|x\|}$

C) Neither one of them

D) Both of them

Changing the Right-Hand Side (RHS)

1 point

You performed an experiment using an expensive piece of scientific equipment, then used the resulting \mathbf{b} to solve the linear system $A\mathbf{x} = \mathbf{b}$. Later, you realized that due to a calibration mistake, there may be some error associated with your calculated \mathbf{b} .

Your lab is out of funding and can't afford to run another experiment. Worried about the accuracy of your computed result (say $\hat{\mathbf{x}}$), you calculated:

$$\frac{\|\Delta\mathbf{b}\|}{\|\mathbf{b}\|} = 10^{-3} \quad \text{where} \quad \Delta\mathbf{b} = \hat{\mathbf{b}} - \mathbf{b}.$$

The condition number of your matrix (A), i.e. $\text{cond}(A)$, is 100.

Calculate an upper bound on the relative error e in the true solution given by:

$$e = \frac{\|\Delta\mathbf{x}\|}{\|\mathbf{x}\|} \quad \text{where} \quad \Delta\mathbf{x} = \hat{\mathbf{x}} - \mathbf{x}$$

$$\frac{\|\Delta\mathbf{x}\|}{\|\mathbf{x}\|} \leq \text{cond}(A) \frac{\|\Delta\mathbf{b}\|}{\|\mathbf{b}\|} = 100 (10^{-3}) = 0.1$$

$$\boxed{\frac{\|\Delta\mathbf{x}\|}{\|\mathbf{x}\|} \leq 0.1}$$

Clicker question

Matrix Conditioning: Accurate digits

1 point

Let's say we want to solve the following linear system:

$$Ax = b$$

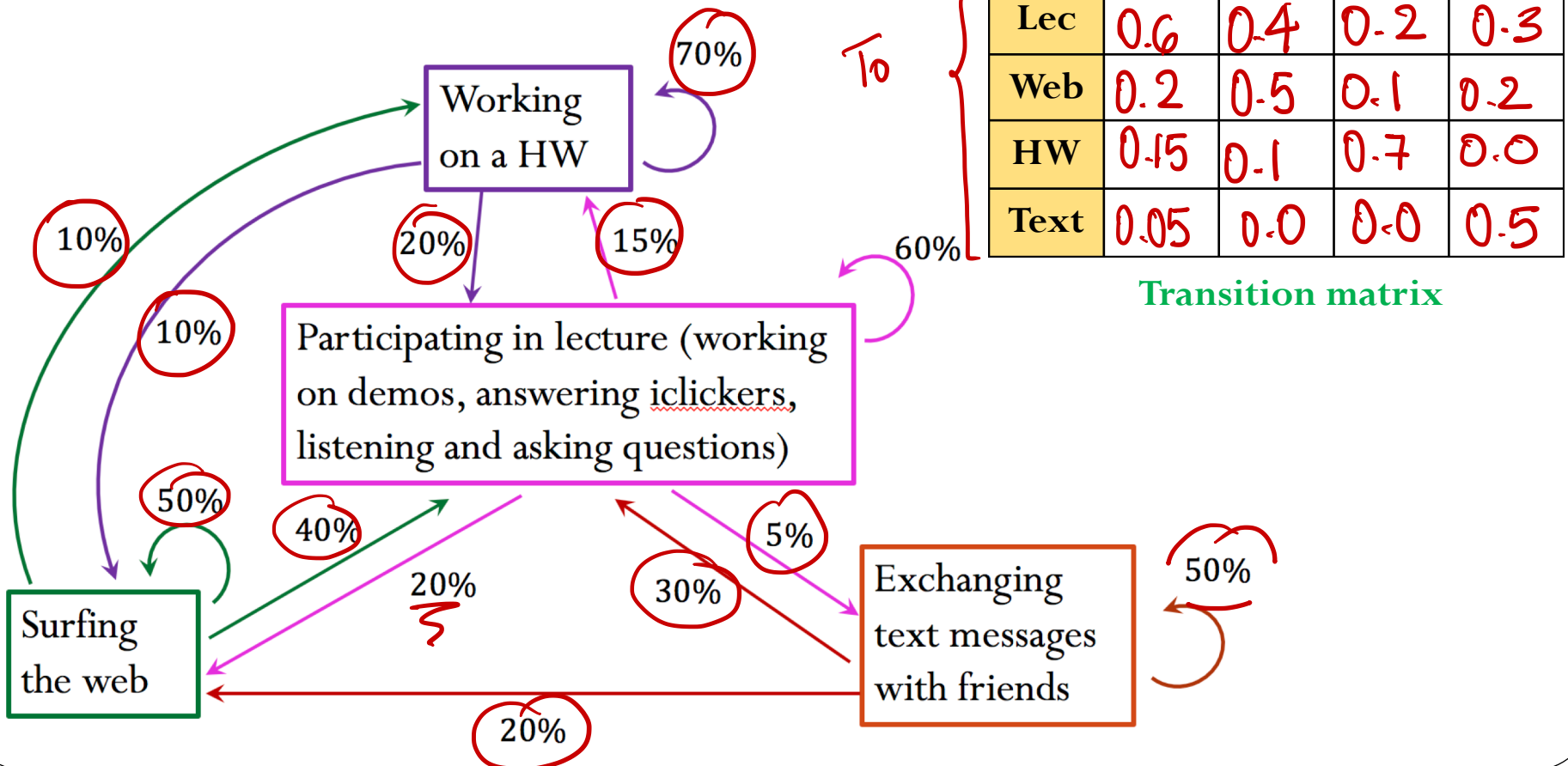
Assuming you are working with IEEE double precision floating point numbers, how many digits of accuracy will your answer have if $\kappa(A) = 1000$?

- A) 3
- B) 10
- C) 13
- D) 16
- E) 32

IEEE double precision \rightarrow input has 16 accurate digits
accuracy is decreased by w , where 10^w is the condition number
output $x \rightarrow (16 - w)$ accurate digits = 13

Markov chain: only the most recent state matters to determine the probability of the next state. $x_n = A x_{n-1}$

Transition (Markov or stochastic) matrix: used to describe the transitions of a Markov chain. Each of its entries is a non-negative real number representing a probability.



Sparse Storage

$$A = \begin{bmatrix} 0 & 0 & 1.3 \\ -1.5 & 0.2 & 0 \\ 5 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0.3 & 3 \end{bmatrix} \begin{matrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \end{matrix}$$

COO Format: row = [0 1 1 2 4 4] nnz integers
col = [2 0 1 0 1 2] nnz integers
data = [1.3 -1.5 0.2 5 0.3 3] nnz doubles

CSR Format: rowptr = [0 1 3 4 4 6] (nrow+1) integers
col = [2 0 1 0 1 2] (nnz) integers
data = [1.3 -1.5 0.2 5 0.3 3] (nnz) doubles

indices: $\underbrace{0}_{\text{row 0}} \quad \underbrace{1 \quad 2}_{\text{row 1}} \quad \underbrace{3}_{\text{row 2}} \quad \underbrace{4}_{\text{row 3}} \quad \underbrace{4 \quad 5}_{\text{row 4}}$

Note: A red arrow points from the text "zero" row to the value 6 in the rowptr array.

MODEL, NUMBERS, ERRORS

- Random Numbers
- Making Models using Monte Carlo simulation
- Errors, accuracy, convergence

- Truncation (Taylor Series)
- Rounding (Floating Point)

LINEAR OPERATORS

- NORMS
- CONDITIONING
- LINEAR SYSTEMS
- EIGENVALUE (Markov Chains)
- SVD (low-rank app.)

MODEL SET OF DATA POINTS

INTERPOLATION (fit data "exactly" to a fn)

LINEAR LEAST-SQUARES (SVD) ("best" data fit to a trend)

ITERATIVE METHODS

NONLINEAR SYSTEM OF EQUATIONS

- BISECTION) 1D
- SECANT) 1D
- NEWTON) 1D
- BROYDEN) ND
- NEWTON) ND

OPTIMIZATION

- GOLDEN SECTION) 1D
- NEWTON) 1D
- STEEPEST DESCENT/ND
- NONLINEAR LEAST-SQUARES

Eigenvalue Problems

• \underline{A} is $n \times n$ matrix, $\underline{x} \neq 0$ is eigenvector of \underline{A} , λ is eigenvalue of \underline{A}

$$\underline{A}\underline{x} = \lambda\underline{x}$$

• \underline{A} is diagonalizable if it has n linearly independent eigenvectors.

$$\underline{A} = \underline{U}\underline{D}\underline{U}^{-1}$$

$$\underline{U} = \begin{bmatrix} | & | & & | \\ \underline{u}_1 & \underline{u}_2 & \dots & \underline{u}_n \\ | & | & & | \end{bmatrix}$$

$$\underline{D} = \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \dots & \\ & & & \lambda_n \end{bmatrix}$$

• Power iteration: \rightarrow finds the largest eigenvalue (λ_1)

$$\underline{\underline{x}}_{k+1} = \underline{\underline{A}} \underline{\underline{x}}_k$$

usually normalize

$$\underline{\underline{x}}_k = \frac{\underline{\underline{A}} \underline{\underline{x}}_{k-1}}{\|\underline{\underline{A}} \underline{\underline{x}}_{k-1}\|}$$

$$\underline{\underline{x}}_k = \lambda_1^k \left[\alpha_1 \underline{\underline{u}}_1 + \underbrace{\alpha_2 \left(\frac{\lambda_2}{\lambda_1} \right)^k}_{e_k} \underline{\underline{u}}_2 + \dots + \alpha_n \left(\frac{\lambda_n}{\lambda_1} \right)^k \underline{\underline{u}}_n \right]$$

$(|\lambda_1| > |\lambda_2| > |\lambda_3| > \dots > |\lambda_n|)$

$\lim_{k \rightarrow \infty} \frac{\underline{\underline{x}}_k}{\lambda_1^k} = \alpha_1 \underline{\underline{u}}_1 \rightarrow \underline{\underline{x}}_k$ converges to a multiple of the first eigenvector $\underline{\underline{u}}_1$!

Once algorithm converges to eigenvector $\underline{\underline{x}}$ then the eigenvalue is evaluated using Rayleigh coefficient:

$$\lambda = \frac{\underline{\underline{x}}^T \underline{\underline{A}} \underline{\underline{x}}}{\underline{\underline{x}}^T \underline{\underline{x}}}$$

- Convergence:

$$\frac{e_{k+1}}{e_k} \approx \frac{|\lambda_2|}{|\lambda_1|}$$

\rightarrow linear convergence

Note that...

Convergence is faster when

$\frac{|\lambda_2|}{|\lambda_1|}$ is small or

$\frac{|\lambda_1|}{|\lambda_2|}$ is large

- Cost per iteration: $O(n^2) \rightarrow$ matrix-vect multiply

- Inverse Power Iteration \rightarrow finds the smallest eigenvalue (λ_n)

$$\underline{\tilde{x}}_{k+1} = \underline{A}^{-1} \underline{\tilde{x}}_k \quad \rightarrow \quad x_0 = \text{initial guess}$$

$$\text{solve } \underline{A} \underline{\tilde{x}}_{k+1} = \underline{\tilde{x}}_k \quad \text{for } \underline{\tilde{x}}_{k+1}$$

$$\text{Factorize } \underline{PA} = \underline{LU} \quad \rightarrow \quad \underline{L} \underline{U} \underline{\tilde{x}}_{k+1} = \underline{P} \underline{\tilde{x}}_k$$

done
only once!

$$\text{solve } \underline{L} \underline{y} = \underline{P} \underline{\tilde{x}}_k$$

$$\text{solve } \underline{U} \underline{\tilde{x}}_{k+1} = \underline{y}$$

} • Cost per iteration
 $O(n^2)$

- Convergence:

$$\frac{e_{k+1}}{e_k} \approx \frac{|\lambda_n|}{|\lambda_{n-1}|}$$

\rightarrow linear convergence

convergence is faster when
this ratio is small

• Shifted-Inverse Power Iteration → finds eigenvalue close to σ

$$\underline{\tilde{x}}_{k+1} = (\underline{A} - \underline{\sigma I})^{-1} \underline{\tilde{x}}_k \longrightarrow x_0 = \text{initial guess}$$
$$\text{solve } (\underline{A} - \underline{\sigma I}) \underline{x}_{k+1} = \underline{\tilde{x}}_k \quad \underline{B} = \underline{A} - \underline{\sigma I}$$

Factorize PB = LU → LU $\underline{x}_{k+1} = \underline{P} \underline{x}_k$

done only once!

$$\left. \begin{array}{l} \text{solve } \underline{L} \underline{y} = \underline{P} \underline{x}_k \\ \text{solve } \underline{U} \underline{x}_{k+1} = \underline{y} \end{array} \right\} \cdot \text{Cost per iteration } O(n^2)$$

λ_{c1} → eigenvalue closest to σ

λ_{c2} → next eigenvalue closest to σ

$$\text{Convergence} : \frac{e_{k+1}}{e_k} \approx \underbrace{\frac{|\lambda_{c1} - \sigma|}{|\lambda_{c2} - \sigma|}}_{\text{Convergence is faster when this rate is small}} \longrightarrow \text{linear convergence}$$

• Rayleigh Quotient Iteration → finds eigenvalue close to σ

$$\underline{\tilde{x}}_{k+1} = (\underline{A} - \underline{\sigma}_k \underline{I})^{-1} \underline{\tilde{x}}_k \longrightarrow x_0 = \text{initial guess}$$

$$\underline{B}_k = \underline{A} - \underline{\sigma}_k \underline{I}$$

$$\text{Factorize } \underline{P} \underline{B}_k = \underline{L} \underline{U} \quad \left. \vphantom{\text{Factorize}} \right\} O(n^3)$$

$$\text{solve } \underline{L} \underline{y} = \underline{P} \underline{x}_k \quad \left. \vphantom{\text{solve}} \right\} O(n^2)$$

$$\text{solve } \underline{U} \underline{x}_{k+1} = \underline{y}$$

$$\text{update } \underline{\sigma}_k = \frac{\underline{x}_{k+1}^T \underline{A} \underline{x}_{k+1}}{\underline{x}_{k+1}^T \underline{x}_{k+1}}$$

cost per
iteration:
 $O(n^3)$

At least quadratic convergence

Compare cost of all methods.

Singular Value Decomposition

$$A = U \Sigma V^T$$

$m \times n$ $m \times m$ $m \times n$ $n \times n$

→ complexity: $O(n^3)$

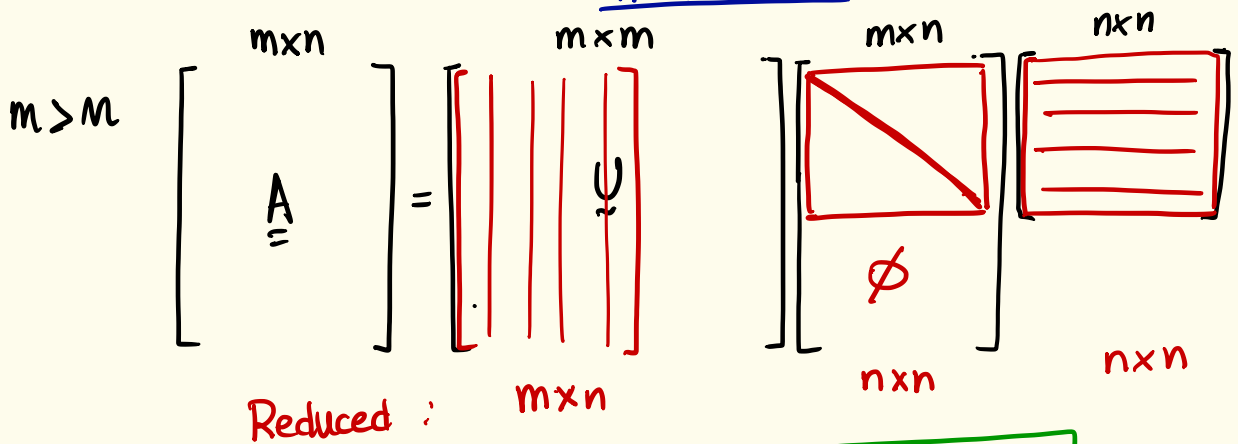
$$A = \begin{bmatrix} | & | & & | \\ u_1 & u_2 & \dots & u_m \\ | & | & & | \end{bmatrix} \begin{bmatrix} \sigma_1 & \sigma_2 & \dots & \sigma_n \\ & & & 0 \end{bmatrix} \begin{bmatrix} - & \sigma_1^T & - \\ & \vdots & \\ - & \sigma_m^T & - \end{bmatrix}$$

$m > n$
 $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n$

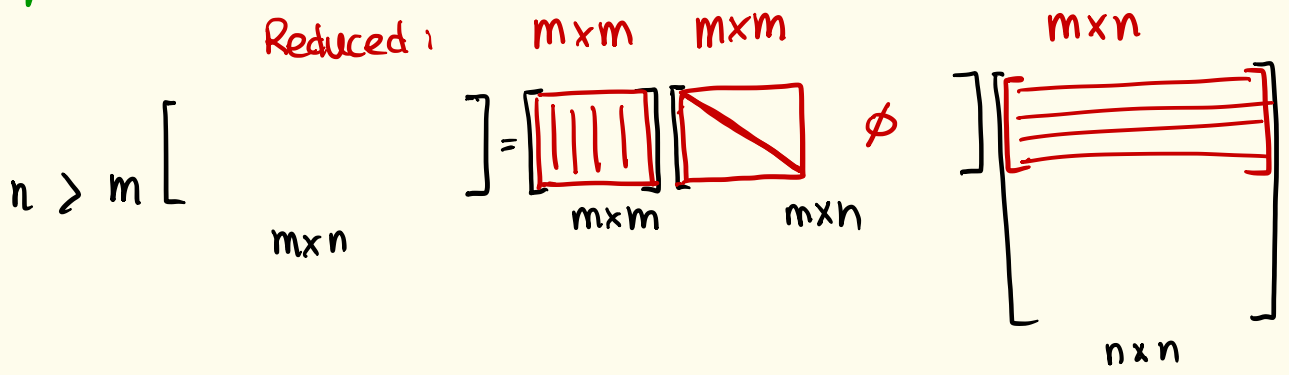
- Columns of V are eigenvectors of $A^T A$ (right singular vectors)
- Columns of U are eigenvectors of AA^T (left singular vectors)
- σ_i^2 are eigenvalues of $A^T A$ (singular values)

Note: $A^T A \rightarrow$ positive definite \rightarrow positive eigenvalues (singular values are always positive!)

"Reduced" SVD



$$k = \min(m, n) \rightarrow (m \times k) \quad (k \times k) \quad (k \times n)$$



Low-rank approximations

$$\underline{A} = \sigma_1 \underline{u}_1 \underline{v}_1^T + \sigma_2 \underline{u}_2 \underline{v}_2^T + \dots + \sigma_n \underline{u}_n \underline{v}_n^T$$

$m \times n$ $m \times 1$ $1 \times n$

($m > n$)
 n singular values

Best rank- k approximation for $m \times n$ matrix \underline{A} is solution of:

$$\min_{\underline{A}_k} \|\underline{A} - \underline{A}_k\|_p$$

s.t. $\text{rank}(\underline{A}_k) \leq k$

when using $p=2$, the best rank- k approximation is:

$$\underline{A}_k = \sigma_1 \underline{u}_1 \underline{v}_1^T + \sigma_2 \underline{u}_2 \underline{v}_2^T + \dots + \sigma_k \underline{u}_k \underline{v}_k^T$$

$$A = \begin{bmatrix} \boxed{-0.817} & \boxed{-0.576} \\ \boxed{-0.576} & \boxed{0.817} \end{bmatrix} \begin{bmatrix} \boxed{5.465} \\ \boxed{0.366} \end{bmatrix} \begin{bmatrix} \boxed{-0.404} & \boxed{-0.914} \\ \boxed{-0.914} & \boxed{0.404} \end{bmatrix}$$

\downarrow \downarrow \downarrow \downarrow \downarrow
 u_1 u_2 σ_1 σ_2 v_1^T v_2^T
 left singular vectors singular values right singular values

① What is $\|A\|_2 = \max \sigma_i = 5.465$

② What is $\text{cond}_2(A) = \|A\|_2 \|A^{-1}\| = \frac{\sigma_{\max}}{\sigma_{\min}} = \frac{5.465}{0.366} = 14.9$

③ Best rank-1 approximation

$$\underline{\underline{A_1}} = \sigma_1 \underline{\underline{u_1}} \underline{\underline{v_1^T}} = 5.465 \begin{bmatrix} -0.817 \\ -0.576 \end{bmatrix} \begin{bmatrix} -0.404 & -0.914 \end{bmatrix}$$

$$\text{error} \approx \|\underline{\underline{A}} - \underline{\underline{A_1}}\|_2 = \sigma_2$$

Linear least-squares

Given m data points : $(t_1, y_1), (t_2, y_2), \dots, (t_m, y_m)$

Find the n coefficients ($n < m$) of the function $f(t)$

that best fits the data

Eg. $f(t) = x_0 + x_1 t + x_2 t^2$

$$\begin{bmatrix} 1 & t_1 & t_1^2 \\ 1 & t_2 & t_2^2 \\ 1 & t_3 & t_3^2 \\ \vdots & \vdots & \vdots \\ 1 & t_m & t_m^2 \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_m \end{bmatrix}$$

$$\Rightarrow \underline{A} \underline{x} \approx \underline{b}$$

more equations than unknowns!

OVERDETERMINED SYSTEM

$$\min_{\underline{x}} \|\underline{A}\underline{x} - \underline{b}\|_2^2$$

solution is unique if $\text{rank}(A) = n$

① $\underline{A}^T \underline{A} \underline{x} = \underline{A}^T \underline{b}$ (Normal equations) \rightarrow Cost: $O(n^3)$ \rightarrow conditioning issues

② Use SVD: $\underline{x} = \underline{V} \underline{\Sigma}^+ \underline{U}^T \underline{b}$ where $\underline{\Sigma}^+ = \begin{cases} 1/\sigma_i, & \sigma_i \neq 0 \\ 0, & \sigma_i = 0 \end{cases}$ \rightarrow Cost: $O(n^3)$ (but more expensive)

$$\textcircled{1} \underline{\underline{A}}^T \underline{\underline{A}} \underline{\underline{x}} = \underline{\underline{A}}^T \underline{\underline{b}} \quad (\text{Normal equations})$$

$$A_{m \times n} \rightarrow A^T A_{n \times n} \quad m > n$$

→ Cost of factorization: $O(n^3)$

→ Cost of solving: $O(n^2)$

→ matrix can be very ill-conditioned

→ solution exists if $\text{rank}(A) = n$

$$\textcircled{2} \text{ Use SVD: } \underline{\underline{x}} = \underline{\underline{V}} \underline{\underline{\Sigma}}^+ \underline{\underline{U}}^T \underline{\underline{b}} \quad \text{where } \underline{\underline{\Sigma}}^+ = \begin{cases} 1/\sigma_i, & \sigma_i \neq 0 \\ 0, & \sigma_i = 0 \end{cases}$$

→ Cost of factorization: $O(n^3)$

→ Cost of solving: $O(mn)$

$$\underline{\underline{z}} = \underline{\underline{U}}^T \underline{\underline{b}} \rightarrow O(mn)$$

$(m \times n)^T \times m \times 1$

$$\underline{\underline{y}} = \underline{\underline{\Sigma}}^T \underline{\underline{z}} \rightarrow O(n)$$

$$\underline{\underline{v}} \underline{\underline{y}} \rightarrow O(n^2)$$

→ solution always exist

• unique if $\text{rank}(A) = n \rightarrow \underline{\underline{x}} = \underline{\underline{V}} \underline{\underline{\Sigma}}^{-1} \underline{\underline{U}}^T \underline{\underline{b}}$

• $\text{rank}(A) < n \rightarrow \underline{\underline{x}} = \underline{\underline{V}} \underline{\underline{\Sigma}}^+ \underline{\underline{U}}^T \underline{\underline{b}}$

Interpolation with monomials

$$p_{n-1}(t) = \sum_{j=0}^{n-1} x_j t^j \quad \rightarrow \text{polynomial of degree } (n-1)$$

in matrix form, this yields the Vandermonde Matrix

If the points t_i are equally spaced on an interval of length h , and the "true function" we are trying to interpolate is sufficiently smooth, then the error of the interpolant of degree $= (n-1)$ is:

$$\text{error} = O(h^n) = O(h^{\text{degree}+1})$$

Nonlinear Equation (1D)

Solve $f(x) = 0$ for $f: \mathbb{R} \rightarrow \mathbb{R}$ (Root Finding)

- 1) Bisection method: $\left\{ \begin{array}{l} \text{no need for derivatives} \\ \text{one function evaluation per iteration} \\ \text{linear convergence} \\ h_{k+1} = 0.5 h_k \end{array} \right.$
- 2) Newton's method: $\left\{ \begin{array}{l} x_{k+1} = x_k + f(x_k)/f'(x_k) \\ \text{need to evaluate one function value and one} \\ \text{derivative for each iteration} \\ \text{quadratic convergence} \end{array} \right.$
- 3) Secant method: $\left\{ \begin{array}{l} \text{need to initial guesses } (x_0, x_1) \\ \hat{f}'(x_k) = (f(x_k) - f(x_{k-1})) / (x_k - x_{k-1}) \rightarrow \text{derivative approx.} \\ x_{k+1} = x_k - f(x_k) / \hat{f}'(x_k) \\ \text{no need for derivatives} \\ \text{one function evaluation per iteration} \\ \text{superlinear convergence} \end{array} \right.$

System Nonlinear Equations

Solve $\underline{f}(\underline{x}) = \underline{0}$ for $f: \mathbb{R}^n \rightarrow \mathbb{R}^n$

1) Newton's Method

\underline{x}_0 : initial guess
solve $\underline{J}(\underline{x}_k) \underline{s}_k = -f(\underline{x}_k)$
 $\underline{x}_{k+1} = \underline{x}_k + \underline{s}_k$
typically quadratic convergence
local convergence
main cost: computing $\underline{J}(\underline{x}_k)$ and solve

$$[\underline{J}(\underline{x}_k)]_{ij} = \frac{\partial f_i}{\partial x_j}$$

$$\underline{J}(\underline{x}_k) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \dots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_n}{\partial x_1} & \frac{\partial f_n}{\partial x_2} & \dots & \frac{\partial f_n}{\partial x_n} \end{bmatrix}$$

2) Secant method (eg. Broyden)

→ use a approximation for the Jacobian that satisfies

$$\hat{\underline{J}}(\underline{x}_{k+1} - \underline{x}_k) = \underline{f}(\underline{x}_{k+1}) - \underline{f}(\underline{x}_k)$$

3) Finite Difference Method

Approximate partial derivatives as: $(\delta_j)_k = \begin{cases} 0, & k \neq j \\ 1, & k = j \end{cases}$

$$\frac{\partial f_i}{\partial x_j} \approx \frac{f_i(\underline{x} + h \delta_j) - f_i(\underline{x})}{h}$$

$$\delta_j = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \leftarrow j^{\text{th}} \text{ component}$$

Optimization 1D

$$\begin{cases} \min_x f(x) \longrightarrow \text{objective (cost) function} \\ \text{s.t. } g(x) = 0 \longrightarrow \text{equality constraint} \\ h(x) \leq 0 \longrightarrow \text{inequality constraint} \end{cases}$$

solution x^* is the minimizer.

Unconstrained optimization

First order necessary condition: $f'(x) = 0 \rightarrow$ if $f'(x^*) = 0 \rightarrow x^*$ is critical point

Second order sufficient condition: if $f''(x^*) > 0 \rightarrow x^*$ is a minimizer

1) Golden Section Search

- guaranteed to converge to global minimum for unimodal functions
- no need for derivatives
- only one function evaluation per iteration
- $h_{k+1} = 0.618 h_k$
- linear convergence

2) Newton's Method

- $x_{k+1} = x_k - f'(x_k)/f''(x_k)$
- typically quadratic convergence
- local convergence
- need to evaluate $f'(x_k)$ and $f''(x_k)$ for each iteration
- may fail to converge (or converge to a max or saddle, since it does not check sign of $f''(x)$)

Optimization ND

Unconstrained Optimization: $\min_{\underline{x}} f(\underline{x})$ $f: \mathbb{R}^n \rightarrow \mathbb{R}$

First order necessary condition: $\nabla f(\underline{x}) = \underline{0} \rightarrow$ if $\nabla f(\underline{x}^*) = \underline{0} \rightarrow \underline{x}^*$ is critical point

Second order sufficient condition: $H_f(\underline{x}) = \nabla^2 f(\underline{x})$

if $H_f(\underline{x}^*)$ positive definite \rightarrow minimum

$H_f(\underline{x}^*)$ negative definite \rightarrow maximum

$H_f(\underline{x}^*)$ indefinite \rightarrow saddle

1) Steepest Descent

\underline{x}_0 : initial guess

$$\underline{s}_k = -\nabla f(\underline{x}_k)$$

$$\alpha_k^* = \underset{\alpha}{\operatorname{argmin}} (f(\underline{x}_k + \alpha \underline{s}_k))$$

$$\underline{x}_{k+1} = \underline{x}_k + \alpha_k^* \underline{s}_k$$

- linear convergence

2) Newton's Method

\underline{x}_0 : initial guess

$$\text{solve } H_f(\underline{x}_k) \underline{s}_k = -\nabla f(\underline{x}_k)$$

$$\underline{x}_{k+1} = \underline{x}_k + \underline{s}_k$$

- typical quadratic convergence
- need Hessian 😞
- local convergence
- high computational cost