

Exam. Fri Nov. 10 16:00 - 17:30 (extendable)
open book / electronic.



5 Questions. (normal questions) \rightarrow 20 pts./

3 bonus \rightarrow 10 pts.

eracery - hand

QR comp.

CS450: Numerical Analysis

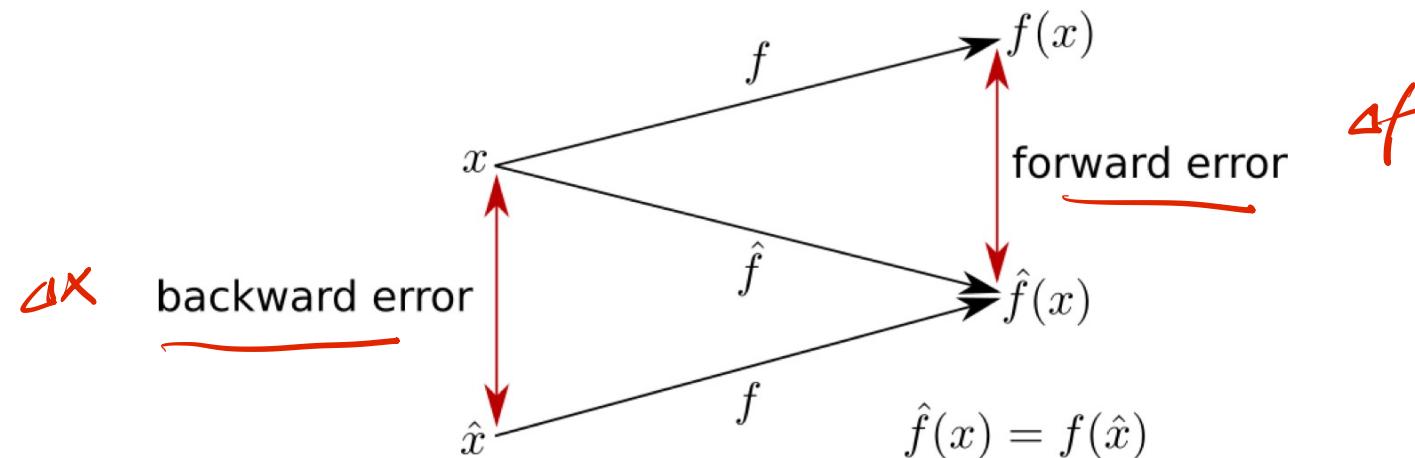
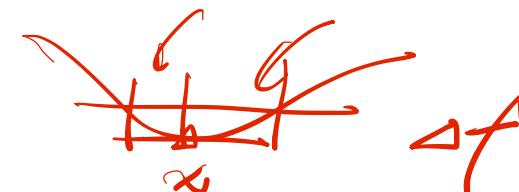
Howard Hao Yang

Assistant Professor, ZJU-UIUC Institute

08/11/2023

Approximation and Errors

- Forward and backward error
 - Suppose we want to compute $f(x)$, where $f: R \rightarrow R$, but obtain an approximate value $\hat{f}(x) = f(\hat{x})$
 - Forward error: $\Delta f = \hat{f}(x) - f(x)$ ↗ unique?
 - Backward error: $\Delta x = \hat{x} - x$ ↗ unique?



Approximation and Errors

- Forward and backward error
 - What is forward error?
 - The computational error of an algorithm/machine
 - Essentially, this is what we really want for an algorithm, but usually hard to obtain...
 - What is backward error?
 - Backward error enables us to measure computational error with respect to data propagation error



Conditioning

- Absolute condition number
 - The absolute condition number is a **property of the problem**, measuring its sensitivity to perturbations in input
 - How much a small change in the input leads to changes in the output
 - Formally, defined by the ratio of absolute errors at output and input

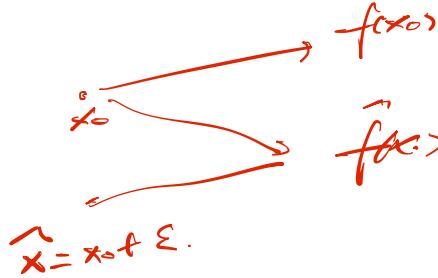
$$\kappa_{abs}(f) = \lim_{\delta \rightarrow 0} \max_{\|\Delta x\| < \delta} \frac{\| \Delta f \|}{\| \Delta x \|}$$

→ f f'

- Linking the forward and backward error

$$\text{forward_error} = \kappa(f) \times \text{backward_error}$$

Condition Number



- Working example

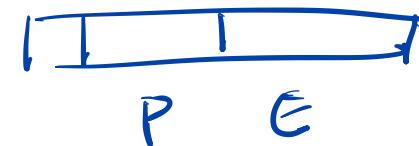
- Suppose we want to evaluate $f(x_0) = e^{\sqrt{x_0}}$ at $x_0 = 4$, but computation results in $\hat{f}(x_0) = f(x_0 + \epsilon)$ where $\epsilon = 10^{-4}$. Give a bound to the forward error.

$$\begin{aligned} |\Delta f| &= |f(x_0 + \epsilon) - f(x_0)| \leq K(f) \cdot |x_0 - \hat{x}| \\ &= f'(x_0) \cdot \epsilon \quad x_0 = 4 \\ &= \frac{1}{2} \cdot \frac{1}{\sqrt{x_0}} \cdot e^{\sqrt{x_0}} \cdot \epsilon \\ &= \frac{1}{2} \cdot \frac{1}{2} \cdot e^2 \times 10^{-4} \quad e \approx 3 \\ &\leq \underline{\underline{2 \cdot 10^{-4}}} \end{aligned}$$

Floating Point Numbers (cont'd)

- Floating point number system is characterized by **four integers**

- Base/radix: β
- Precision: p
- Exponent range: $[L, U]$



- Number x is represented as

$$x = \pm \left(d_0 + \frac{d_1}{\beta} + \frac{d_2}{\beta^2} + \cdots + \frac{d_{p-1}}{\beta^{p-1}} \right) \times \beta^E$$

where $0 \leq d_i \leq \beta - 1$, $i = 0, \dots, p - 1$, and $L \leq E \leq U$

- **Sign**, **exponent**, and **mantissa** are stored in separate fixed-width fields

Rounding

- Not all real numbers can be represented exactly by floating-point systems
- The numbers **representable** are called **machine numbers**
- The **numbers not representable** are approximated by “nearby” floating-point numbers
- Two commonly used rounding rules
 - Chop: truncate base- β expansion of x after $(p - 1)$ -th precision digit; also called round toward zero
 - Round to nearest: $\text{fl}(x)$ is nearest floating-point number to x
- Round to nearest is most accurate and is default in IEEE systems

Machine Precision



- Definition 1: The machine precision or machine epsilon is the smallest number ϵ such that $\text{fl}(1 + \epsilon) > 1$, formally,

$$\epsilon_{\text{mach}} = \operatorname{argmin}_{\epsilon \geq 0} \{ \text{fl}(1 + \epsilon) \neq 1 \}$$

example: $\beta=10$, $p=3$, $E=2$ $0.01 \sim 9.99 \times 10^{-2}$

$$\begin{aligned} 1.00 &\rightarrow 1.00 \times 10^0 \\ 0.009 &\rightarrow 0.09 \times 10^{-1} \end{aligned} \quad \left\{ \rightarrow \begin{aligned} 1.00 + 0.009 &= 1.00 \times 10^0 + 0.09 \times 10^{-1} \\ &= 1.00 \times 10^0 + \cancel{0.009} \times 10^{-1} \\ &= 1.00 \times 10^0 \end{aligned} \right.$$

Floating-Point Arithmetic

- Rule: When adding two numbers, the exponents align before mantissa adds

$$2.414 \times 10^3 + 1.2 \times 10^2 = 2.414 \times 10^3 + 0.12 \times 10^3 = 2.534 \times 10^3$$

- Floating-point addition and multiplication are commutative but not associative

- If ϵ is slightly smaller than the machine precision ϵ_{mach} , then $(1 + \epsilon) + \epsilon = 1$, but $1 + (\epsilon + \epsilon) > 1$

- Q: How to add up $1 + \left(\frac{1}{2} \epsilon_{\text{mach}} + \frac{1}{4} \epsilon_{\text{mach}} + \dots + \frac{1}{2^n} \epsilon_{\text{mach}}\right)$

$$\approx \frac{1}{2} \cdot \frac{1}{1 - \frac{1}{2}} \cdot \epsilon_{\text{mach}}$$

order matters

$$1 + \left(\frac{1}{2} \epsilon_{\text{mach}} + \frac{1}{4} \epsilon_{\text{mach}} + \dots + \frac{1}{2^n} \epsilon_{\text{mach}} \right)$$

$n \rightarrow \infty$

$$\begin{aligned} & \frac{1}{2} \epsilon_{\text{mach}} + \frac{1}{4} \epsilon_{\text{mach}} \\ &= \frac{3}{4} \epsilon_{\text{mach}} \end{aligned}$$

Vector Norms

$$b = \frac{1}{2}a.$$



ZJU-UIUC INSTITUTE
Zhejiang University/University of Illinois at Urbana-Champaign Institute

- What is a (vector) norm?
- A metric to measure the “length” of a vector, or “distance” between two vectors
 - An operator that $\|\cdot\|: R^d \rightarrow R_+$, and satisfies
 - $\|x\| \geq 0$ and $\|x\| = 0$ if and only if $x = \mathbf{0}$
 - $\|a \cdot x\| = |a| \cdot \|x\|$ for any $a \in R$
 - $\|x + y\| \leq \|x\| + \|y\|$

Vector Norms (cont'd)

- Typical norms for vectors
 - The L-p norms

- L-p norm: $\|\beta\|_p = \left(\sum_{j=1}^d \beta_j^p \right)^{1/p}, p \geq 1$

- $p = 1, \|\beta\|_1 = \sum_j |\beta_j|$ ✓

- $p = 2, \|\beta\|_2 = \sqrt{\sum_j |\beta_j|^2}$ ✓

- $p = \infty, \|\beta\|_\infty = \max_j |\beta_j|$

- The L-0 norm: Counts the number of non-zero entries, e.g., if $\beta = (10, 0, 2, 0.01, 0, 1)^T$, then

$$\|\beta\|_0 = 4$$

$$\begin{aligned}
 \|\beta\|_p &= \left(\beta_1^p + \beta_2^p + \cdots + \beta_d^p \right)^{1/p} \\
 &\leq \left(d \cdot (\max_j |\beta_j|)^p \right)^{1/p} \\
 &\leq \underbrace{d^{1/p}}_{p \rightarrow \infty} \cdot \max_j |\beta_j|
 \end{aligned}$$

Matrix Norms

- Given matrix A (square)

- We want an operator that satisfies

① $\|A\| \geq 0$ and $\|A\| = 0$ if and only if $A = 0$

② $\|a \cdot A\| = |a| \cdot \|A\|$ for any $a \in R$

③ $\|A + B\| \leq \|A\| + \|B\|$

④ $\|AB\| \leq \|A\| \cdot \|B\|$ (a property in matrix norms)

if $A = B = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$,

$$\|A\|_{\max} = \|B\|_{\max} \geq 1$$

$$AB = \begin{bmatrix} 2 & 2 \\ 2 & 2 \end{bmatrix} \quad \|AB\|_{\max} = 2$$

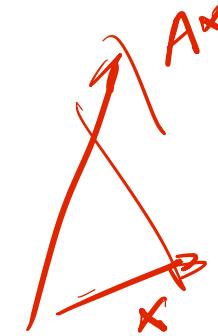
$$\|A\|_F = \sqrt{\sum_{i=1}^n \|a_i\|^2}$$

$$\|A\|_{\max} = \max_{ij} |a_{ij}|$$

Matrix Norms (cont'd)

- Given a vector norm $\|\cdot\|$
 - The matrix norm induced by this vector norm is given as

$$\|A\| = \max_{x \neq 0} \frac{\|Ax\|}{\|x\|} = \max_{\|x\|=1} \|Ax\|$$



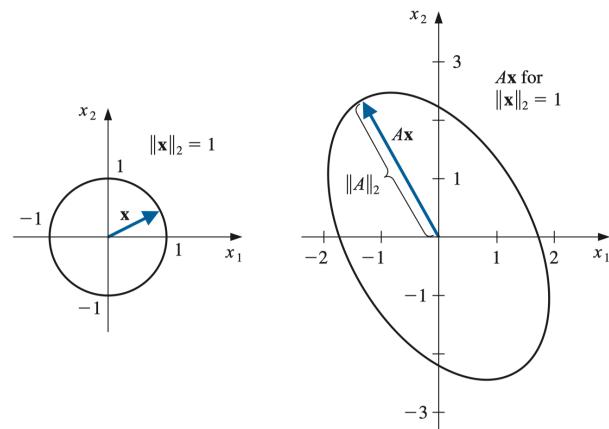
- Different vector norms can induce different matrix norms for the same matrix
- Examples
 - When the norm is L-1 norm, $\|A\|_1 = \max_j \sum_{i=1}^m |a_{ij}|$: the maximum absolute **column sum** of the matrix
 - When the norm is L- ∞ norm, $\|A\|_\infty = \max_i \sum_{j=1}^n |a_{ij}|$: the maximum absolute **row sum** of the matrix

Matrix Norms (cont'd)

- When the norm is the Euclidean norm $\|\cdot\|_2$
 - The matrix norm induced by this vector norm is given as

$$\|A\|_2 = \max_{x \neq 0} \frac{\|Ax\|_2}{\|x\|_2} = \max_{\|x\|_2=1} \|Ax\|_2 \sim \lambda_{\max}(A)$$

- This is called the spectral norm of matrix A



(Matrix) Norms

- Working example

- If A is a square matrix, show that $\|A\|_1 = \|A^T\|_\infty$ and $\|A\|_2^2 \leq \|A\|_1 \|A\|_\infty$?

$$\|A\|_2 = \max_{\|x\|_2=1} \|Ax\|_2 \quad \longleftrightarrow \quad \max_{\|x\|_2=1} x^T A^T A x \rightarrow L(x, \lambda)$$

$$\frac{\partial L(x, \lambda)}{\partial x} \geq 0 \Rightarrow A^T A x - \lambda x \geq 0 \Rightarrow A^T A x = \lambda x \rightarrow \lambda = \lambda_{\max}(A^T A) = \|A\|_2^2$$

$$x^* = \arg \max_{\|x\|_2=1} \|Ax\|_2$$

$$\underbrace{A^T A x^*}_{\|A\|_2^2 \cdot x^*} = \|A\|_2^2 \cdot x^*$$

$$\|A\|_2^2 \cdot \|x^*\|_1$$

$$\begin{aligned} \|\|A\|_2^2 \cdot x^*\|_1 &= \|A^T A x^*\|_1 \\ &\leq \|A^T\|_1 \cdot \|A x^*\|_1 \\ &\leq \|A\|_\infty \cdot \|A\|_1 \cdot \|x^*\|_1 \end{aligned}$$

Matrix Condition Number

- Given a matrix A , the condition number is

$$\kappa(A) = \|A\| \cdot \|A^{-1}\|$$

- Definition: a $n \times n$ matrix A is said to be **nonsingular/invertible** if A^{-1} exists such that $AA^{-1} = A^{-1}A = I$; otherwise, it is **singular**
- Large $\kappa(A)$ implies the matrix is nearly singular

(Matrix) Condition Number

- Working example



If A is a square matrix, show that $\kappa(A) = \underbrace{\|A\|} \cdot \underbrace{\|A^{-1}\|} = \left(\max_{x \neq 0} \frac{\|Ax\|}{\|x\|} \right) \left(\min_{x \neq 0} \frac{\|Ax\|}{\|x\|} \right)^{-1}$

$$\|A^{-1}\| = \max_{x \neq 0} \frac{\|A^{-1}x\|}{\|x\|} \quad \rightarrow \quad z = A^{-1}x \iff x = Az$$

$$= \max_{z \neq 0} \frac{\|z\|}{\|Az\|}$$

$$\kappa(A) = \frac{\max_{\|x\|=1} \|Ax\|}{\min_{\|x\|=1} \|Ax\|}$$

$$= \frac{1}{\min_{\substack{x \neq 0 \\ x \in D}} \frac{\|Ax\|}{\|x\|}}$$

Solving Linear Systems

- For a square matrix A , how to (systematically) solve for $Ax = b$?
 - Transform it into one whose solution is the same but easier to compute
 - Specifically, eliminate x_1 from $n - 1$ equations to get a smaller system $A_2x = b_2$ of size $n - 1$
 - Eventually, reaching the 1 by 1 system $A_nx_n = b_n$ which we know $x_n = b_n/A_n$
 - Working backwards produces x_{n-1}, x_{n-2}, \dots , and eventually x_2 and x_1

$$\begin{array}{c|c|c} A & x & = & b \\ \hline \end{array} \quad \rightarrow \quad \begin{array}{c|c} U & x \\ \hline \end{array} = \tilde{b}$$

Triangular Linear Systems

- What type of linear system is easy to solve?
 - Systems that form **triangular matrices**

$$a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n = a_{1,n+1}$$

$$a_{22}x_2 + \cdots + a_{2n}x_n = a_{2,n+1}$$

$$\ddots \quad \ddots \quad \ddots \quad \vdots \quad \vdots$$

$$a_{nn}x_n = a_{n,n+1}$$

- Back-substitution

$$x_n = \frac{a_{n,n+1}}{a_{nn}},$$

$$x_i = \frac{a_{i,n+1} - \sum_{j=i+1}^n a_{ij}x_j}{a_{ii}}, \quad i = n-1, \dots, 1$$

Elementary Elimination Matrices

- More generally, we can annihilate all entries below k -th in a n -dimensional vector \mathbf{a} by transformation

\downarrow *k -th column*

$$\underline{M_k} \mathbf{a} = \begin{bmatrix} 1 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & 1 & 0 & \cdots & 0 \\ 0 & \cdots & -m_{k+1} & 1 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & -m_n & 0 & \cdots & 1 \end{bmatrix} \begin{bmatrix} a_1 \\ \vdots \\ a_k \\ a_{k+1} \\ \vdots \\ a_n \end{bmatrix} = \begin{bmatrix} a_1 \\ \vdots \\ a_k \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

where $m_i = \frac{a_i}{a_k}$, $i = k + 1, \dots, n$

- Divisor a_k , called **pivot**, must be nonzero

Elementary Elimination Matrices (cont'd)



Example: For a vector $\mathbf{a} = [2, 4, -2]^T$, the elementary elimination matrices \mathbf{M}_1 and \mathbf{M}_2 are (recall: $m_i = \frac{a_i}{a_k}$, $i = k + 1, \dots, n$)

$$\mathbf{M}_1 \mathbf{a} = \begin{bmatrix} 1 & 0 & 0 \\ -2 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} 2 \\ 4 \\ -2 \end{bmatrix} = \begin{bmatrix} 2 \\ 0 \\ 0 \end{bmatrix}$$

$$\mathbf{M}_2 \mathbf{a} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1/2 & 1 \end{bmatrix} \begin{bmatrix} 2 \\ 4 \\ -2 \end{bmatrix} = \begin{bmatrix} 2 \\ 4 \\ 0 \end{bmatrix}$$

Elementary Elimination Matrices (cont'd)

- Matrix \mathbf{M}_k , called elementary elimination matrix, adds multiple of row k to each subsequent row, with multipliers m_i chosen so that result is zero
- \mathbf{M}_k is unit lower triangular and nonsingular
- $\mathbf{M}_k = \mathbf{I} - \underline{\mathbf{m}_k \mathbf{e}_k^T}$, where $\mathbf{m}_k = [0, \dots, 0, m_{k+1}, \dots, m_n]^T$ and \mathbf{e}_k is the k -th column of identity matrix
- $\underline{\mathbf{M}_k^{-1} = \mathbf{I} + \mathbf{m}_k \mathbf{e}_k^T}$, which means $\mathbf{M}_k^{-1} = \mathbf{L}_k$ is the same as \mathbf{M}_k except signs of multipliers are reversed

Elementary Elimination Matrices (cont'd)

- If \mathbf{M}_j , $j > k$, is another elementary elimination matrix, with vectors of multipliers \mathbf{m}_j , then

$$\begin{aligned}\mathbf{M}_k \mathbf{M}_j &= (\mathbf{I} - \mathbf{m}_k \mathbf{e}_k^T)(\mathbf{I} - \mathbf{m}_j \mathbf{e}_j^T) \\ &= \mathbf{I} - \mathbf{m}_k \mathbf{e}_k^T - \mathbf{m}_j \mathbf{e}_j^T + \mathbf{m}_k \mathbf{e}_k^T \mathbf{m}_j \mathbf{e}_j^T \\ &= \mathbf{I} - \mathbf{m}_k \mathbf{e}_k^T - \mathbf{m}_j \mathbf{e}_j^T\end{aligned}$$

where means product is essentially “union”, and similarly for product of inverses, $\mathbf{L}_k \mathbf{L}_j$

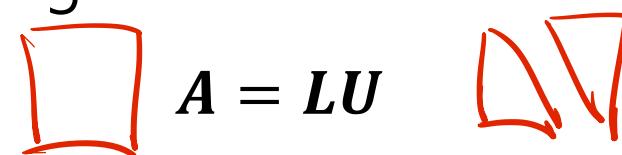
Gaussian Elimination

$$\begin{array}{c|c} \text{Initial Matrix } A & | \\ \hline a_{11} & a_{12} \end{array} \xrightarrow{\text{Row Swap}} \begin{array}{c|c} \text{Pivoted Matrix } A' & | \\ \hline a_{21} & a_{22} \end{array}$$

- To reduce a general linear system of equations $\mathbf{A}\mathbf{x} = \mathbf{b}$ into upper triangular form, do the following
 - Choose M_1 with a_{11} as pivot, annihilate first column of A below first row: the system becomes $M_1\mathbf{A}\mathbf{x} = M_1\mathbf{b}$, but solution remains unchanged
 - Next choose M_2 with a_{22} as pivot, annihilate second column of $M_1\mathbf{A}$ below second row: the system becomes $M_2M_1\mathbf{A}\mathbf{x} = M_2M_1\mathbf{b}$, but the solution still unchanged
 - Process continues for each successive column until all subdiagonal entries have been zeroed, results in upper triangular linear system $M_{n-1} \cdots M_1\mathbf{A}\mathbf{x} = M_{n-1} \cdots M_1\mathbf{b}$, which can be solved via back-substitution
- This process is called **Gaussian elimination**

LU Factorization

- Denote by $M_k^{-1} = L_k$, then $L = M^{-1} = M_1^{-1} \cdots M_{n-1}^{-1} = L_1 \cdots L_{n-1}$ is lower triangular
- By design, $U = \underbrace{M_{n-1} \cdots M_1}_{} A$ is upper triangular
- Therefore, we can factorize matrix A into the product of a lower triangular matrix and an upper triangular matrix


$$A = LU$$

- Thus, Gaussian elimination produces LU factorization of matrix into triangular factors
- Can every square matrix be factorized into an LU decomposition?



LU and Linear System of Equations

- Working example

$$\cancel{\begin{bmatrix} A & B \\ C & D \end{bmatrix}} x = \begin{bmatrix} b \\ c \end{bmatrix}$$

- How would you solve a partitioned linear system of the form

$$\begin{bmatrix} L_1 & O \\ B & L_2 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} b \\ c \end{bmatrix}$$

where L_1 and L_2 are nonsingular lower triangular matrices?

$$\textcircled{1} \quad L_1 x = b$$

$$\textcircled{2} \quad L_2 y = c - Bx$$

Linear Least Square

$$\begin{bmatrix} \cdot & \cdot \\ \cdot & \cdot \end{bmatrix} \times \begin{bmatrix} \cdot \\ \cdot \end{bmatrix} = \begin{bmatrix} \cdot \\ \cdot \end{bmatrix}$$

- Suppose we want to solve a linear system $\mathbf{A}\mathbf{x} = \mathbf{b}$, with \mathbf{A} and \mathbf{b} being known, what does this mean?
 - What does $\mathbf{A}\mathbf{x}$ mean?
 - What does $\mathbf{A}\mathbf{x} = \mathbf{b}$ mean?
- What happens if \mathbf{A} is not a square matrix?
- Instead of solving $\mathbf{A}\mathbf{x} = \mathbf{b}$, we aim to find \mathbf{x} to minimize $\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2$
 - Generally, a perfect fitting may not be possible, and we look for an approximation
 - In data science, this is called Linear Regression

Solving the Linear Least Square

- How to find $\mathbf{x}^* = \operatorname{argmin}_{\mathbf{x}} \|\mathbf{Ax} - \mathbf{b}\|_2$ where $\mathbf{A} \in R^{m \times n}$?
 - The residual error can be written as

$$Err(\mathbf{x}) = (\mathbf{Ax} - \mathbf{b})^T (\mathbf{Ax} - \mathbf{b})$$

- Let $\frac{\partial Err(\mathbf{x})}{\partial \mathbf{x}} = 2\mathbf{A}^T(\mathbf{Ax} - \mathbf{b}) = \mathbf{0}$, which is equivalent to $\underline{(\mathbf{A}^T\mathbf{A})\mathbf{x} = \mathbf{A}^T\mathbf{b}}$
- If $\mathbf{A}^T\mathbf{A}$ is nonsingular, we obtain the solution as

$$\mathbf{x} = (\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T\mathbf{b}$$

- This is known as the **normal equation**

Singularity Issues

$$Ax=b$$

- Working example
 - Show that the following matrix A is singular; if $b = [2 \ 4 \ 6]^T$, how many solutions are there in $Ax = b$?

① singularity
 $\vec{a}_2 = \vec{a}_1 + \vec{a}_3$

$$A = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 2 & 1 \\ 1 & 3 & 2 \end{bmatrix} x = \begin{bmatrix} 2 \\ 4 \\ 6 \end{bmatrix}$$

②

$$\alpha \begin{bmatrix} 0 \\ 2 \\ 0 \end{bmatrix} + (1-\alpha) \begin{bmatrix} 2 \\ 0 \\ 2 \end{bmatrix}$$
$$\alpha \in [0, 1]$$

Solving the Normal Equation

- How to solve $\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{b}$?
- Remember the LU factorization? Since $\mathbf{A}^T \mathbf{A}$ is a square matrix, we can do LU factorization to it and then use back substitution?
- However, consider the following

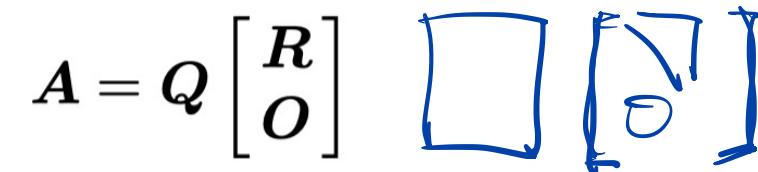
$$\mathbf{A} = \begin{bmatrix} 1 & 1 \\ \epsilon & 0 \\ 0 & \epsilon \end{bmatrix}, \text{ where } 0 < \epsilon < \sqrt{\epsilon_{\text{mach}}}$$

$$\xrightarrow{\hspace{1cm}} \mathbf{A}^T \mathbf{A} = \begin{bmatrix} 1 + \epsilon^2 & 1 \\ 1 & 1 + \epsilon^2 \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

- We want something more robust...

Wisdom from the Past

- When solving a linear system, we decompose the square matrix A into an LU form—we want something similar when A becomes a rectangular matrix
- For a rectangular matrix $A \in R^{m \times n}$ with $m > n$, it can be decomposed as

$$A = Q \begin{bmatrix} R \\ O \end{bmatrix}$$


where Q is an $m \times m$ orthogonal matrix and R is an upper triangular matrix

- This is known as the **QR factorization**

Goodness of Triangles

- Working example
 - What is the Euclidean norm of the minimum residual vector for the following linear least square problem? What is the solution vector for this problem?

$$\begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$

$$\underbrace{\begin{bmatrix} 1 & 1 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}}_A \underbrace{\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}}_x \approx \underbrace{\begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix}}_b$$

$$\underbrace{(x_1+x_2-2)^2 + (x_2-1)^2 + 1}_{\geq 0, \text{ when } x_1=x_2=1} = 0$$

$$\|r\|_2^2 = \|b - Ax\|_2^2$$

$$= \left\| \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_1 + \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_2 - \begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix} \right\|_2^2 = \left\| \begin{bmatrix} x_1+x_2-2 \\ x_2-1 \\ 1 \end{bmatrix} \right\|_2^2$$

QR Factorization for LSQ

$\|b - Ax\|_2$

- If we have a QR factorization of matrix $A \in R^{m \times n}$, i.e.,

$$A = Q \begin{bmatrix} R \\ O \end{bmatrix} \quad \boxed{\cancel{Q} \begin{bmatrix} R \\ O \end{bmatrix} \cancel{R}^{-1} = Q} \quad \boxed{R}$$

$$\boxed{A} \xrightarrow{A^{-1}} \boxed{b}$$

- Then, the equation $Ax = b$ reduces to

$$Q^T A x = \boxed{\begin{bmatrix} R \\ O \end{bmatrix} x} \cong \boxed{\begin{bmatrix} c_1 \\ c_2 \end{bmatrix}} = Q^T b$$

- The following holds

$$\begin{aligned}\|Ax - b\|_2^2 &= \|Q^T(Ax - b)\|_2^2 \\ &= \|Rx - c_1\|_2^2 + \|c_2\|_2^2\end{aligned}$$

QR Factorization for LSQ (Cont'd)

- If we have a QR factorization of matrix A , i.e., $A = QR$, the following holds

$$\|Ax - b\|_2^2 = \|Rx - c_1\|_2^2 + \|c_2\|_2^2$$

- The second term is the **residual error**, which we cannot do anything on it
- The first term, we can choose x to minimize it, and the solution is given by solving $Rx - c_1 = 0$, which can be achieved via back substitution
- In this way, we avoid the cross-product matrix, saving us from round-off error issues
- The remaining question now is... how to perform QR factorization?

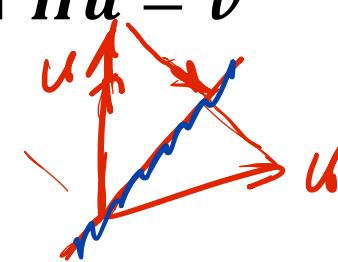
Householder Transformations (cont'd)

- Given a vector $x \in R^n$ with $x^T x = 1$, the **Householder transformation** is

$$H = I - 2xx^T$$

- Reflection property: for any vector $a \in R^n$, Ha reflects a by the hyperplane perpendicular to x
- Let $n \geq 2$ and $u, v \in R^n$ be unit vectors (i.e., $u^T u = v^T v = 1$). Suppose $u \neq v$,

let $x = \frac{u-v}{\|u-v\|_2}$ and construct $H = I - 2xx^T$, then $Hu = v$



Householder Transformations (Cont'd)

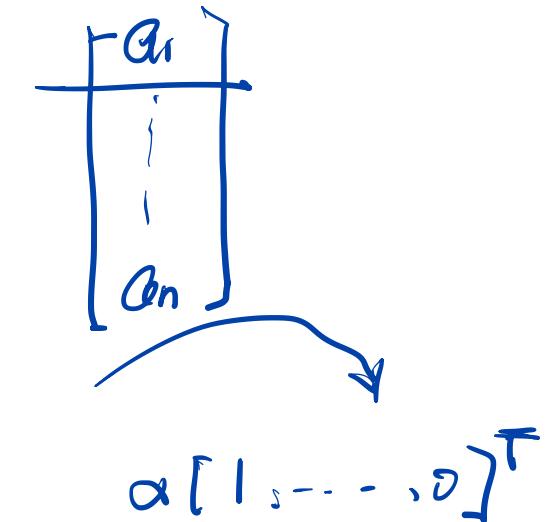
- Suppose we have $\mathbf{a} \in R^n$, and want to annihilate all elements below the first entry while preserving the norm
- Can we leverage some ideas from the Householder transformation?
- Problem: find vector $\mathbf{x} \in R^n$, such that $\mathbf{x}^T \mathbf{x} = 1$ and

$$\mathbf{H}\mathbf{a} = (\mathbf{I} - 2\mathbf{x}\mathbf{x}^T)\mathbf{a} = \alpha \mathbf{e}_1$$

where $\mathbf{e}_1 = (1, 0, \dots, 0)^T$ and $\alpha = \|\mathbf{a}\|_2$

- Solution to this is

$$\mathbf{x} = \frac{\mathbf{a}}{\alpha} \pm \mathbf{e}_1$$



Householder Transformation

- Working example

- Can you identify an eigenvector of H that associates with the eigenvalue -1?; i.e., a vector such that $Hv = -v$?

$$\begin{aligned} H &= I - 2xx^T \\ Hx &= (I - 2xx^T)x \\ &= x - 2x \underbrace{x^Tx}_{=1} = x - 2x = -x \\ Hx &= -x \end{aligned}$$

Householder Transformations (Cont'd)

- For a rectangular matrix $A \in R^{m \times n}$ with $m > n$. Suppose $m = 6$, $n = 5$ and we have computed the following

$$H_1^T H_2^T H_3^T H_4^T H_5^T H_6^T H_2 H_1 A = I$$

Diagram illustrating the Householder transformation process:

The matrix $H_2 H_1 A$ is highlighted in a green box. To its left, a sequence of Householder matrices $H_1^T, H_2^T, \dots, H_6^T$ is shown, with H_2^T and H_1^T circled in blue. The result of applying these transformations to $H_2 H_1 A$ is the identity matrix I .

To the right of the equation, the original matrix A is shown as a 6x5 matrix with entries marked by 'x'. This matrix is transformed into an upper triangular matrix \tilde{A} (circled in green) via the sequence of transformations $H_1^T H_2^T \dots H_6^T$. The matrix \tilde{A} has zeros below the diagonal and non-zero entries above the diagonal, with some entries circled in green.

- Concentrating on the highlighted entries, we determined a matrix $\tilde{H}_3 \in R^{4 \times 4}$ such that

$$\tilde{H}_3 \begin{bmatrix} x \\ x \\ x \\ x \end{bmatrix} = \begin{bmatrix} x \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

Householder Transformations (Cont'd)

- Then, by forming $H_3 = \text{diag}(I_2, \tilde{H}_3)$, we have

$$H_3 H_2 H_1 A = \begin{bmatrix} \times & \times & \times & \times & \times \\ 0 & \times & \times & \times & \times \\ 0 & 0 & \times & \times & \times \\ 0 & 0 & 0 & \times & \times \\ 0 & 0 & 0 & \times & \times \\ 0 & 0 & 0 & \times & \times \end{bmatrix}$$

- More generally, for a given vector $\mathbf{a} = \begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \end{bmatrix}$, where \mathbf{a}_1 is a $(k-1)$ -vector, with $1 \leq k < m$. If we take the Householder vector to be $\mathbf{v} = \begin{bmatrix} \mathbf{0} \\ \mathbf{a}_2 \end{bmatrix} - \alpha \mathbf{e}_k$, where $\alpha = \mp \|\mathbf{a}_2\|_2$, then the resulting Householder transformation annihilates the last $m - k$ components of \mathbf{a} .

Householder Transformations (Cont'd)

- For a rectangular matrix $\mathbf{A} \in R^{m \times n}$ with $m > n$. The QR factorization to this matrix can be written as

$$\mathbf{H}_n \cdots \mathbf{H}_1 \mathbf{A} = \begin{bmatrix} \mathbf{R} \\ \mathbf{O} \end{bmatrix}$$

- The product of successive Householder transformations $\mathbf{H}_n \cdots \mathbf{H}_1$ is itself an orthogonal matrix. Thus, if we take $\mathbf{Q}^T = \mathbf{H}_n \cdots \mathbf{H}_1$, or equivalently, $\mathbf{Q} = \mathbf{H}_1 \cdots \mathbf{H}_n$, then

$$\mathbf{A} = \mathbf{Q} \begin{bmatrix} \mathbf{R} \\ \mathbf{O} \end{bmatrix}$$