

CSC338 Lecture 4

Last time: $A\underline{x} = \underline{b}$ for $A \in \mathbb{R}^{n \times n}$ nonsingular

Use Gauss Elimination and LU Factorization

Question: when does Gauss Elimination fail?

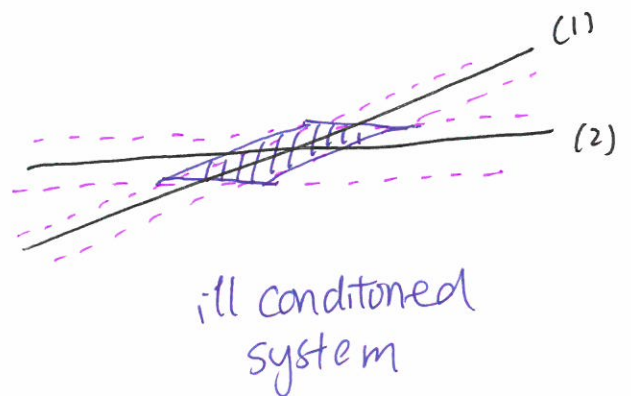
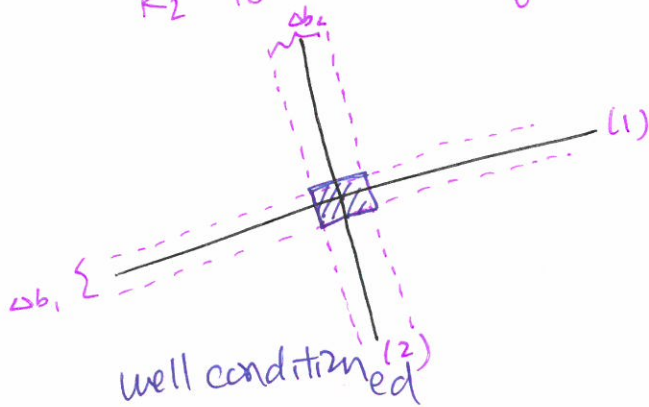
eg// ① $A = \begin{bmatrix} 2 & 1 & 1 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 0 & 2 & 2 \\ 0 & 0 & 2 & 2 \end{bmatrix}$ This element, called the pivot is zero!
 $R_4 \leftarrow R_4 - \frac{2}{0} R_3$ we can't divide by 0!

eg// ② $A = \begin{bmatrix} 2 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 \\ 0 & 0 & 0.0001 & 2 \\ 0 & 0 & 2 & 2 \end{bmatrix}$ This pivot is close to zero
 $R_4 \leftarrow R_4 - \frac{2}{0.0001} R_3$ very large

Examples ① and ② points to issues with the stability of the algorithm

eg// $A = \begin{bmatrix} 4.11 & 5.0 \\ 8.23 & 10.1 \end{bmatrix}$ $\underline{b} = \begin{bmatrix} 2 \\ 4 \end{bmatrix}$ G.E. $\Rightarrow \underline{x} = \begin{bmatrix} 0.487 \\ 0 \end{bmatrix}$ $A\underline{x} = \begin{bmatrix} 2 \\ 4.005 \end{bmatrix}$

This matrix is "almost singular"
 R_2 is almost equal to $2 \times R_1$



(1) - $a_{11} x_1 + a_{12} x_2 = b_1 + \Delta b_1$

(2) - $a_{21} x_1 + a_{22} x_2 = b_2 + \Delta b_2$

eg/ $A = \begin{bmatrix} 0.780 & 0.563 \\ 0.913 & 0.659 \end{bmatrix} \quad \underline{b} = \begin{bmatrix} 0.217 \\ 0.254 \end{bmatrix}$

we obtained 2 solutions computationally:

$\underline{\hat{x}}_1 = \begin{bmatrix} 0.341 \\ -0.087 \end{bmatrix} \quad \underline{\hat{x}}_2 = \begin{bmatrix} 0.999 \\ -1.001 \end{bmatrix}$

which solution is better?

$\underline{r}_1 = \underline{b} - A\underline{\hat{x}}_1 = \begin{bmatrix} +0.000001 \\ 0 \end{bmatrix}$

$\underline{r}_2 = \underline{b} - A\underline{\hat{x}}_2 = \begin{bmatrix} 0.001343 \\ 0.00572 \end{bmatrix}$

but $\underline{x}^* = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$

⇒ Residuals not always a good indicator of accuracy
 (Linearly scaling A, \underline{b} will affect $\|\underline{r}\|$ but not $\|\underline{x}^* - \underline{\hat{x}}\|$)

To be able to discuss conditioning/stability,
 we need a way to talk about the "size" of a vector/matrix.

Def A vector norm $\|\cdot\|$ is a mapping $\mathbb{R}^n \rightarrow \mathbb{R}^{\geq 0}$ such

- that:
1. $\forall \underline{x} \in \mathbb{R}^n, \|\underline{x}\| \geq 0$ with $\|\underline{x}\| = 0$ iff $\underline{x} = \underline{0}$
 2. $\forall \underline{x} \in \mathbb{R}^n, c \in \mathbb{R} \quad \|c\underline{x}\| = |c| \cdot \|\underline{x}\|$
 3. $\forall \underline{x}, \underline{y} \in \mathbb{R}^n \quad \|\underline{x} + \underline{y}\| \leq \|\underline{x}\| + \|\underline{y}\|.$

(triangle inequality)

In this course, we will use the p-norm.

$\|\underline{x}\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{1/p}$

eg/ 1-norm $p=1 \quad \|\underline{x}\|_1 = \sum_{i=1}^n |x_i|$
 2-norm $p=2 \quad \|\underline{x}\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$
 ∞ -norm $p=\infty \quad \|\underline{x}\|_\infty = \max_{1 \leq i \leq n} |x_i|$

~ Euclidean Norm

eg// $\underline{x} = \begin{bmatrix} 2 \\ 3 \\ -4 \\ 1 \end{bmatrix}$

$\|\underline{x}\|_1 = 10$

$\|\underline{x}\|_2 = \sqrt{4+9+16+1} = \dots$

$\|\underline{x}\|_\infty = 4$

Def An induced matrix norm is a matrix norm defined in terms of a vector norm as follows:

$\|A\| \stackrel{\text{def}}{=} \max_{\|\underline{x}\|=1} \|A\underline{x}\| = \max_{\underline{x} \neq 0} \frac{\|A\underline{x}\| \leftarrow \text{vector norm}}{\|\underline{x}\| \leftarrow \text{vector norm}}$

The largest possible value of $\|A\underline{x}\|$ across all possible vectors \underline{x} where $\|\underline{x}\|=1$.

Properties of Matrix Norms

- $\|A\| \geq 0$ with $\|A\|=0$ iff $A=0$ ← the zero matrix.
- $\|cA\| = |c| \|A\|$ for $c \in \mathbb{R}$
- $\|A+B\| \leq \|A\| + \|B\|$
- $\|AB\| \leq \|A\| \|B\|$
- $\|A\underline{x}\| \leq \|A\| \|\underline{x}\|$
- $\|I\| = 1$

~~eg//~~ we can also show: $\|A\|_1 = \max_j \sum_{i=1}^m |a_{ij}| \sim$ max abs column sum

$\|A\|_\infty = \max_i \sum_{j=1}^n |a_{ij}| \sim$ max abs row sum.

eg// $A = \begin{bmatrix} 4 & 0 \\ 0 & 1 \end{bmatrix}$

$\|A\|_1 = 4$
 $\|A\|_\infty = 4$
 $\|A\| = 4$

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

$\|A\|_1 = 5$
 $\|A\|_\infty = 4$

Def The condition number of an $n \times n$ matrix A ④

is defined to be

$$\text{cond}(A) = \begin{cases} \|A\| \cdot \|A^{-1}\| & \text{if } A \text{ is nonsingular} \\ \infty & \text{if } A \text{ is singular} \end{cases}$$

We can show that

$$\|A\| \|A^{-1}\| = \left(\underbrace{\max_{\|\underline{x}\|=1} \|A\underline{x}\|}_{\text{max stretch}} \right) \left(\underbrace{\min_{\|\underline{x}\|=1} \|A\underline{x}\|}_{\text{max shrink}} \right)^{-1}$$

\Rightarrow $\text{cond}(A)$ describes the ratio of the max stretch to max shrinking of unit vectors

Properties of $\text{cond}(A)$

1. $\text{cond}(A) = \|A\| \cdot \|A^{-1}\| \geq \cancel{\|A\| \cdot \|A^{-1}\|} \|A \cdot A^{-1}\| = \|I\| = 1$

2. $\text{cond}(A) \geq 1$

3. $\text{cond}(I) = 1$.

For $\gamma \in \mathbb{R}, \gamma \neq 0$

$$\text{cond}(\gamma A) = \|\gamma A\| \cdot \left\| \frac{1}{\gamma} A^{-1} \right\| = \text{cond}(A)$$

Claim If $\text{cond}(A)$ is small, then $A\underline{x} = \underline{b}$ is well-conditioned.

(small perturbations in A and \underline{b} will not cause large perturbation in the computed solution for \underline{x}).

Specifically, we have

$$\frac{\|\Delta \underline{x}\|}{\|\underline{x}^*\|} \leq \underbrace{\text{cond}(A)}_{\text{we will show this}} \frac{\|\Delta \underline{b}\|}{\|\underline{b}\|}$$

$$\frac{\|\Delta \underline{x}\|}{\|\underline{x}^*\|} \leq \text{cond}(A) \frac{\|\Delta A\|}{\|A\|}$$

in the textbook.

Why did we define $\text{cond}(A) = \|A\| \cdot \|A^{-1}\|$?

Recall the definition of C.N. from lecture 1.

$$\text{C.N.} = \frac{|\Delta y / y|}{|\Delta x / x|}$$

$$\text{or } \underbrace{\left| \frac{\Delta y}{y} \right|}_{\text{relative error of output}} = \text{C.N.} \cdot \underbrace{\left| \frac{\Delta x}{x} \right|}_{\text{relative error of input (problem)}}$$

relative error of output.

relative error of input (problem)

Pf (that $\frac{\|\Delta \underline{x}\|}{\|\underline{x}^*\|} \leq \text{cond}(A) \frac{\|\Delta \underline{b}\|}{\|\underline{b}\|}$)

Consider $A \underline{x} = \underline{b}$.

Let \underline{x}^* be the true solution

$$A \underline{x}^* = \underline{b} \quad (1)$$

$\hat{\underline{x}}$ be the computed solution

$$A \hat{\underline{x}} = \underline{b} + \Delta \underline{b} \quad (2)$$

$$\Delta \underline{x} = \hat{\underline{x}} - \underline{x}^*$$

from (2) — $A \hat{\underline{x}} - \underline{b} = \Delta \underline{b}$

from (1) — $A \underline{x}^* - \underline{b} = \underline{0}$

Subtract $A(\hat{\underline{x}} - \underline{x}^*) = \Delta \underline{b}$ so $A \Delta \underline{x} = \Delta \underline{b}$

So we have $\Delta \underline{x} = A^{-1} \Delta \underline{b} \Rightarrow \|\Delta \underline{x}\| \leq \|A^{-1}\| \cdot \|\Delta \underline{b}\| \quad (3)$

$A \underline{x}^* = \underline{b} \Rightarrow \|A\| \cdot \|\underline{x}^*\| \geq \|\underline{b}\| \quad (4)$

(3)
(4)

$$\frac{\|\Delta \underline{x}\|}{\|A\| \|\underline{x}^*\|} \leq \frac{\|A^{-1}\| \|\Delta \underline{b}\|}{\|\underline{b}\|}$$

nom relative error of the computed solution

$$\left\{ \frac{\|\Delta \underline{x}\|}{\|\underline{x}^*\|} \leq \underbrace{\|A\| \cdot \|A^{-1}\|}_{\text{cond}(A)} \frac{\|\Delta \underline{b}\|}{\|\underline{b}\|} \right\}$$

nom residual of computed solution

Conclusion If $\text{cond}(A)$ is small } property of the problem
 and relative residual is small } property of the computed solution
 then the relative error is small.

But if \underline{b} is small but $\text{cond}(A)$ large, the relative error could still be large.

Improving the Stability of Gauss Elimination

eg/ $\begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix}$ $\begin{bmatrix} 0.0001 & 1 \\ 1 & 1 \end{bmatrix}$

↑ can't divide by 0

↑ multiplication by a large number increase error

Idea: Interchange rows ($R_1 \leftrightarrow R_2$) before eliminating elements below a diagonal, so we maximize the absolute value of the pivot. ~ the value in the diagonal

eg/ $A = \begin{bmatrix} 3 & 2 & 9 \\ 4 & 5 & 1 \\ -5 & 2 & 3 \end{bmatrix} \begin{matrix} \leftarrow R_1 \leftrightarrow R_3 \\ \leftarrow R_2 \leftrightarrow R_3 \end{matrix}$ since $| -5 |$ has the largest abs val in the first column.

We can use a permutation matrix to express the interchange.

$R_1 \leftrightarrow R_3 \Rightarrow P_1 = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix}$ $P_1 A = \begin{bmatrix} -5 & 2 & 3 \\ 4 & 5 & 1 \\ 3 & 2 & 9 \end{bmatrix}$

Def A permutation matrix P has exactly one 1 in each row and each column, and 0 everywhere else.

- Properties
- $P^{-1} = P^T$
 - $\text{cond}(P) = 1$.
 - Multiplying by both sides of $Ax = \underline{b}$ by P does not change the solution $PAx = P\underline{b}$

Gauss Elimination with Partial Pivoting

Before eliminating below a diagonal, permute the matrix A so that the pivot ~~has~~ is the largest magnitude element at or below the diagonal.

Instead of performing

$$M_{n-1} \cdots M_2 M_1 A = U$$

$M_i =$ elementary matrix

We perform

$$M_{n-1} P_{n-1} \cdots M_2 P_2 M_1 P_1 A = U$$

where $P_i =$ permutation matrix

It turns out that the above can be expressed as:

$$\underbrace{(\hat{M}_{n-1} \cdots \hat{M}_2 \hat{M}_1)}_{\text{lower-triangular } L^{-1}} \underbrace{(P_{n-1} \cdots P_2 P_1)}_{\text{permutation matrix } P} A = U \quad \Rightarrow \quad PA = LU$$

Summary To solve $A\underline{x} = \underline{b}$

1. Find P, L, U s.t. $PA = LU$ using G.E + partial pivoting.

$$2. \text{ To solve } A\underline{x} = \underline{b} \Leftrightarrow PA\underline{x} = P\underline{b}$$

$$\Leftrightarrow LU\underline{x} = P\underline{b}$$

$$\text{Solve } \underline{L}\underline{y} = P\underline{b} \quad \text{for } \underline{y}$$

$$\underline{U}\underline{x} = \underline{y} \quad \text{for } \underline{x}.$$

eg// $A = \begin{bmatrix} \epsilon & 1 \\ 1 & 1 \end{bmatrix}$ with $\epsilon \ll \epsilon_{mach}$

To eliminate ~~a_{11}~~
 $R_2 \leftarrow R_2 - \frac{1}{\epsilon} R_1$

without pivoting we get

$$L = \begin{bmatrix} 1 & 0 \\ \frac{1}{\epsilon} & 1 \end{bmatrix} \quad U = \begin{bmatrix} \epsilon & 1 \\ 0 & 1 - \frac{1}{\epsilon} \end{bmatrix} \stackrel{\text{float}}{=} \begin{bmatrix} \epsilon & 1 \\ 0 & -\frac{1}{\epsilon} \end{bmatrix}$$

$$\Rightarrow LU = \begin{bmatrix} 1 & 0 \\ \frac{1}{\epsilon} & 1 \end{bmatrix} \begin{bmatrix} \epsilon & 1 \\ 0 & -\frac{1}{\epsilon} \end{bmatrix} = \begin{bmatrix} \epsilon & 1 \\ 1 & 0 \end{bmatrix} \neq A$$

Complete Pivoting

In complete pivoting (full pivoting) we search for the largest entry in the remaining submatrix.

eg// $\begin{bmatrix} 1 & 2 & 5 & 9 \\ 0 & 3 & -1 & 2 \\ 0 & 4 & 2 & 1 \\ 0 & -1 & -5 & 0 \end{bmatrix}$

Partial pivoting: $R_2 \leftrightarrow R_3$ Pivot = 4.

Full pivoting: $R_2 \leftrightarrow R_4$ Pivot = -5
 $C_2 \leftrightarrow C_3$

Full pivoting changes the solution of the system, but predictably

eg// $\begin{bmatrix} 2 & 0 \\ 0 & 3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$

$\begin{bmatrix} 0 & 2 \\ 3 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$

$x_1 = 1/2$
 $x_2 = 1/3$
 $x_1 = 1/3$
 $x_2 = 1/2$

we need to keep track of the column swaps.

In practice, full pivoting is expensive, and partial pivoting provides good enough stability.

eg // $A = \begin{bmatrix} 1 & 2 \\ k & 1 \end{bmatrix}$, $k \in [0, 2]$ $A^{-1} = \frac{1}{1-2k} \begin{bmatrix} 1 & -k \\ -2 & 1 \end{bmatrix}$ ⑨

What would k have to be for $\text{cond}(A)$ to be large.

Use $\|\cdot\|_1$, we have

$$\text{cond}(A) = \|A\|_1 \cdot \|A^{-1}\|_1,$$

$$= 3 \cdot \frac{1}{1-2k} \cdot 3$$

which is large when $k \rightarrow 0.5$