



Instead of  $A=LU$  (not always possible)

We should use  $PA=LU$

$$L_{n-1} P_{n-2} L_2 P_{2,1} L_1 P_{1,1} A = U \quad \text{an upper triangular}$$

Ex:

$$P = P_{13} \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 1 & 2 & 3 \\ -2 & 4 & 5 \\ 3 & -1 & 0 \end{bmatrix}$$

$$= \begin{bmatrix} 3 & -1 & 0 \\ -2 & 4 & 5 \\ 1 & 2 & 3 \end{bmatrix}$$

$$P_{24} L_1 P_{13} A =$$

$$P_{24} L_1 = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ -\frac{1}{3} & 1 & 0 & 0 \\ \frac{2}{3} & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \end{bmatrix}$$

$$= \tilde{L} P_{24} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 \\ \frac{2}{3} & 0 & 1 & 0 \\ -\frac{1}{3} & 0 & 0 & 1 \end{bmatrix}$$

Thm: After Gaussian elimination with partial column pivoting, we have

$$PA = LU$$

$$P = P_{n-2, i_{n-2}} P_{n-1, i_{n-1}} \dots P_{1, i_1}$$

$L$  is a unit triangular matrix

$$\det(A) = (-1)^m \det(U)$$

$m$  is the number of row exchanges.

Ex:

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

Get ①  $PA = LU$

② Solve  $Ax = b$ .

$$\begin{bmatrix} 1 & -1 & 6 \\ 2 & 0 & 2 \\ 1 & 2 & 4 \end{bmatrix} \xrightarrow{P_{12}} \begin{bmatrix} 2 & 0 & 2 \\ 1 & -1 & 6 \\ 1 & 2 & 4 \end{bmatrix} \begin{matrix} -1 & 0 & -1 \\ 1 & -1 & 6 \\ -1 & 0 & -1 \\ 1 & 2 & 4 \end{matrix}$$

$$\xrightarrow{L_1} \begin{bmatrix} 2 & 0 & 2 \\ \frac{1}{2} & -1 & 5 \\ \frac{1}{2} & 2 & 3 \end{bmatrix} \xrightarrow{P_{23}} \begin{bmatrix} 2 & 0 & 2 \\ \frac{1}{2} & 2 & 3 \\ -\frac{1}{2} & -1 & 5 \end{bmatrix}$$

$$\xrightarrow{L_2} \begin{bmatrix} 2 & 0 & 2 \\ \frac{1}{2} & 2 & 3 \\ -\frac{1}{2} & -1 & 5 \\ & -\frac{1}{2} & \frac{13}{2} \end{bmatrix}$$

$$\det(A) = (-1)^2 \cdot 2 \cdot 2 \cdot \frac{13}{2} = 26$$

$$U = \begin{bmatrix} 2 & 0 & 2 \\ 0 & 2 & 3 \\ 0 & 0 & \frac{13}{2} \end{bmatrix}$$

$$L = \begin{bmatrix} 1 & 0 & 0 \\ \frac{1}{2} & 1 & 0 \\ -\frac{1}{2} & -\frac{1}{2} & 1 \end{bmatrix}$$

$$P = P_{23} P_{12}$$

$$Ax = b$$

$$\underline{P}A x = \underline{P}b$$

$$L\underline{U}x = \underline{P}b$$

$$Ly = \underline{P}b$$

$$Ux = y$$

$$P_{23} P_{12} b \quad b = \begin{bmatrix} 6 \\ 4 \\ 7 \end{bmatrix}$$

$$= P_{23} \begin{bmatrix} 4 \\ 6 \\ 7 \end{bmatrix} = \begin{bmatrix} 4 \\ 7 \\ 6 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 & 0 \\ \frac{1}{2} & 1 & 0 \\ \frac{1}{2} & -\frac{1}{2} & 1 \end{bmatrix} y = \begin{bmatrix} 4 \\ 7 \\ 6 \end{bmatrix}$$

$$y_1 = 4$$

$$y_2 = 7 - \frac{1}{2} \cdot 4 = 5$$

$$y_3 = 6 - \frac{1}{2} y_1 + \frac{1}{2} y_2 = \frac{13}{2}$$

$$Ux = y \quad \begin{bmatrix} 2 & 0 & 2 \\ 0 & 2 & 3 \\ 0 & 0 & \frac{13}{2} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 4 \\ 5 \\ \frac{13}{2} \end{bmatrix}, \quad \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

Other pivoting techniques.

1. Complete pivoting

$$\max_{\substack{1 \leq i \leq n \\ 1 \leq j \leq n}} (a_{ij}^{(1)})$$

Unnecessarily complicated.

2. Scaled pivoting (Useful)

$$\begin{bmatrix} 1 & -1 & 2 & 3 \\ 100 & 200 & -300 & 400 \end{bmatrix} \xrightarrow{\text{row-scale}} \begin{bmatrix} 1 & -1 & 2 & 3 \\ 1 & 2 & -3 & 4 \end{bmatrix} *$$

$$\frac{|a_{k1}|}{s_k} = \max_{1 \leq i \leq n} \frac{|a_{i1}|}{s_i}$$

$$s_i = \sum_{j=1}^n |a_{ij}|$$

Factors affect accuracy

- ① algorithm growth factor  $\delta(n)$
- ② machine precision
- ③ condition of the problem

$x \cdot y$  Well conditioned

$x \cdot y$  may be ill-conditioned if  $x \sim y$ .

Error analysis for solving  $Ax=b$

$$fl(b_i) = b_i(1 + \delta_i), \quad |\delta_i| \leq \epsilon$$

$$fl(\vec{b}) = \vec{b} + \vec{\delta b}, \quad \|\delta b\|_p = \left\| \begin{bmatrix} \delta b_1 \\ \delta b_2 \\ \vdots \\ \delta b_n \end{bmatrix} \right\|_p$$

$$|\delta b_i| \leq |b_i| \epsilon$$

$$\|\delta b\|_p \leq \left\| \begin{bmatrix} |b_1| \epsilon \\ \vdots \\ |b_n| \epsilon \end{bmatrix} \right\|_p = \epsilon \sum \| |b_i| \|_p = \|b\|_p \epsilon$$

$p=1, 2, \infty$

$$\|\delta b\|_p \leq C \|b\|_p \epsilon$$

$$f(A) = A + \delta A$$

$$\|B\|_p = \|B\|_p,$$

$p=1, 2, \infty$

Actually we are

solving  $(A + \delta A)x = b + \delta b$

$$\sup \frac{\|Ax\|}{\|x\|}$$

If  $x^* = A^{-1}b$ ,  $\bar{x} = (A + \delta A)^{-1}(b + \delta b)$

$$\frac{\|x^* - \bar{x}\|}{\|x^*\|} \leq \frac{\left( \frac{\|\delta A\|}{\|A\|} + \frac{\|\delta b\|}{\|b\|} \right)}{\|A\| \|A^{-1}\|}$$

Input error

Sensitivity  $= \text{Cond}(A) \left( \frac{\|\delta A\|}{\|A\|} + \frac{\|\delta b\|}{\|b\|} \right)$

If  $\text{Cond}(A)$  is big, then it is sensitive, ill-conditioned

Otherwise it is called well-conditioned



$$\frac{1}{1+x} = 1 - x + x^2 - x^3 \dots + (-1)^n x^n \dots$$

$$|x| < 1 \quad = \lim_{n \rightarrow \infty} \frac{1-x^{n+1}}{1+x} = \frac{1}{1+x}$$

Banach  
Thm.

$$(1+x)^{-1} = 1 - x + x^2 - x^3 \dots$$

← partial sum

$$\begin{aligned} (A + \delta A)^{-1} &= (A(I + A^{-1}\delta A))^{-1} \\ &= (I + A^{-1}\delta A)^{-1} A^{-1} \\ &= (I + E)^{-1} A^{-1} \end{aligned}$$

Lemma. If  $\|E\| < 1$ , then  $I+E$  is  
invertible, and Banach's  
$$\|(I+E)^{-1}\| \leq \frac{1}{1-\|E\|}$$
Theorem.

$$x_e = A^{-1} b,$$

$$(A + \delta A)x = b + \delta b$$

$$x_e \neq \vec{0}$$

$$\bar{x} = (A + \delta A)^{-1} (b + \delta b)$$

Main theorem.

If  $\|A^{-1} \delta A\| < 1$ , then  $A + \delta A$  is invertible and  $\text{Cond}(A) = \|A\| \|A^{-1}\|$

$$\frac{\|x_e - \bar{x}\|}{\|x_e\|} \leq C \left( \frac{\|\delta A\|}{\|A\|} + \frac{\|\delta b\|}{\|b\|} \right) \text{Cond}(A)$$

$\text{Cond}(A) = \|A\|_2 \|A^{-1}\|_2$

In particular,

GEPP

$A + \delta A_1 + \delta A_2$  — input error from GEPP

$$\frac{\|x_e - \bar{x}\|}{\|x_e\|} \leq C g(n) \text{Cond}(A) \epsilon$$

Banach's Theorem

growth factor  
Worst

$$g(n) = 2^{n-1} \sim n$$

$$\|E\| < 1, \quad \underline{\underline{(I+E)^{-1}}}$$
 exists,  $\frac{1}{1+x} = 1-x$

$$\|(I+E)^{-1}\| \leq \frac{1}{1-\|E\|}$$

Proof

$$B = I - E + E^2 - E^3 + E^4$$

$$+ (-1)^n E^n$$

$$B_n = I - E + E^2 - \dots + (-1)^n E^n$$

$$(I+E)B_n = (I+E)(I - E + E^2 - \dots + (-1)^n E^n)$$

$$= I - E + E - \dots$$

$$= I - \underline{(-1)^n E^{n+1}} \rightarrow I$$

$$\|(-1)^n E^{n+1}\| = \|E^n E\| \leq \|E^n\| \|E\| \leq \|E^{n-1}\| \|E\|$$

$$\lim_{n \rightarrow \infty} B_n = (I+E)^{-1} \leq \|E\|^{n+1} \rightarrow 0$$

$$\|(I+E)^{-1}\| = \|I - E + E^2 - \dots\|$$

$$\leq \|I\| + \|E\| + \|E^2\| + \dots + \|E^n\|$$

$$\leq 1 + \|E\| + \|E\|^2 + \dots + \|E\|^n + \dots$$

$$= \frac{1}{1-\|E\|}$$

$$Ax = b, \quad x_e = A^{-1}b$$

$$(A + \delta A)x = b + \delta b, \quad \bar{x} = (A + \delta A)^{-1}(b + \delta b)$$

$$\begin{aligned} \|\delta x\| &= \|\bar{x} - x_e\| = \|(A + \delta A)^{-1}(b + \delta b) - A^{-1}b\| \\ &= \|(A + \delta A)^{-1}(b + \delta b - (A + \delta A)A^{-1}b)\| \\ &\leq \|(A + \delta A)^{-1}\| \|b + \delta b - b - \delta A x_e\| \\ &\leq \|(A(I + A^{-1}\delta A))^{-1}\| (\|\delta b\| + \|\delta A\| \|x_e\|) \\ &\leq \|(I + A^{-1}\delta A)^{-1}\| \|A^{-1}\| (\|\delta b\| + \|\delta A\| \|x_e\|) \end{aligned}$$

If  $\|A^{-1}\delta A\| < 1$ , then  $A + \delta A$  is invertible  
 $\|A^{-1}\delta A\| \leq \|A^{-1}\| \|\delta A\|$ , this is a  
 stronger condition used by many books.

Ex: Hilbert matrix.

$$\left[ \begin{array}{ccc} \frac{1}{2} & \frac{1}{3} & \\ & \frac{1}{4} & \\ & & \frac{1}{i+j} \end{array} \right] \quad \left[ \begin{array}{ccc} \frac{1}{2} & \frac{1}{3} & \frac{1}{4} \\ \frac{1}{3} & \frac{1}{4} & \frac{1}{5} \\ & & \end{array} \right]$$

$H \quad H^T = H$ , Symmetric positive definite  
 $H^{-1}$  exists  $\lambda_i(H_n) > 0$  SPD  $A = LL^T$   
 Choleski

$$\|\delta x\| \leq \frac{\|A^{-1}\|}{1 - \|A^{-1}\|\|\delta A\|} (\|\delta b\| + \|\delta A\| \|x_e\|)$$

$$\rightarrow \leq \frac{\|A^{-1}\| \|A\|}{1 - \|A^{-1}\|\|\delta A\|} \left( \frac{\|\delta b\|}{\|A\|} + \frac{\|\delta A\|}{\|A\|} \|x_e\| \right)$$

$$\frac{\|\delta x\|}{\|x_e\|} \leq \|A^{-1}\| \|A\|$$

$$Ax_e = b$$

$$\|b\| \leq \|A\| \|x_e\|$$

$$\left( \frac{\|\delta b\|}{\|A\| \|x_e\|} + \frac{\|\delta A\|}{\|A\|} \right) \cdot \frac{1}{1 - \xi} = 1 + \xi + \alpha \xi^2$$

$\xi = \|A^{-1}\| \|\delta A\|$

$$\leq \text{cond}(A) \left( \frac{\|\delta b\|}{\|b\|} + \frac{\|\delta A\|}{\|A\|} \right) + O(\|\delta A\| + \|\delta A\|^2)$$

!