

CS 4220 Feb 26<sup>th</sup>

- Recap on GEPP
- Error Analysis
- Condition Number of Solving Linear Systems

### Alg01 GEPP

$$L = I, P = I$$

for  $j = 1 : n-1$

$$\text{pick } k = \arg \max_{l > j} |A_{lj}|$$

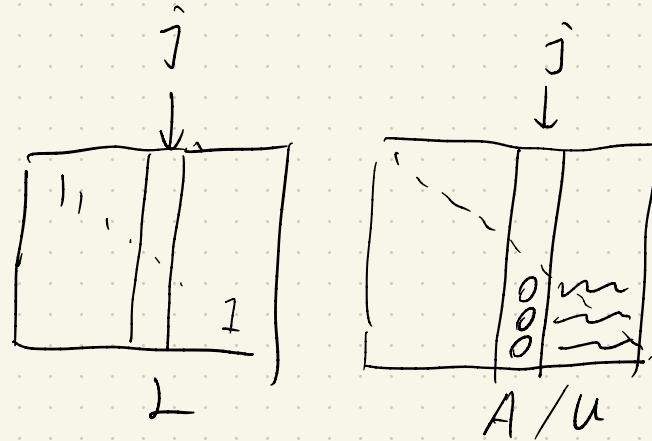
swap row  $k$  with row  $j$  in L.A., update  $P$

$$L(j+1:n, j) = A(j+1:n, j) / A(jj) \leftarrow \text{too small}$$

$$A(j+1:n, j+1:n) = L(j+1:n, j) * A(j, j+1:n)$$

$$U = \text{triu}(A)$$

return  $L, U, P$



$$- Ax = b \quad PAx = Pb$$

$$\textcircled{1} \quad PA = L U \quad \leftarrow \text{GEPP}$$

$$\leftarrow |E| \leq n \cdot \sum |L|/|U|$$

$$\textcircled{2} \quad C = Pb \quad \leftarrow \text{matrix-vector}$$

$$\textcircled{1} + \textcircled{3} + \textcircled{4}$$

$$\textcircled{3} \quad Ly = C \quad \leftarrow \text{forward sub}$$

$$\textcircled{4} \quad Ux = y \quad \leftarrow \text{backward sub.} \quad [\text{Thm}] \text{ After } \textcircled{1} \sim \textcircled{4}, \text{ we get}$$

- Error Analysis

$$A = L U$$

$$A + E = L' U'$$

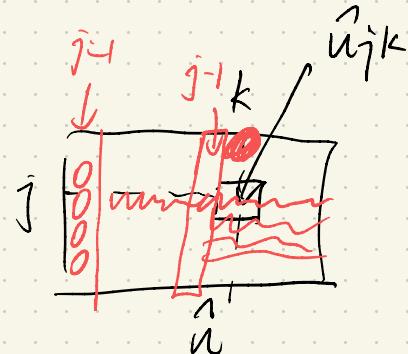
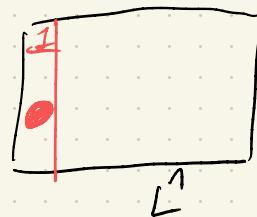
prop.  $|E| \leq n \cdot \epsilon_{\text{mach}} \cdot |L| \cdot |U|$

proof (sketch)

vector  $x$  that solve

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

$$\text{with } |\Delta A| \leq 3 \cdot n \cdot \sum |L|/|U|$$



if  $j > k$ ,  $\hat{u}_{jk} = 0$  at the end.

focus on the case  $j \leq k$

step 0:  $\hat{u}_{jk} = a_{jk}$

step 1:  $\hat{u}_{jk} = \hat{u}_{jk} - \hat{l}_{j2} \hat{u}_{2k}$

step 2:  $\hat{u}_{jk} = \hat{u}_{jk} - \hat{l}_{j2} \hat{u}_{2k}$

:

:

step  $j-1$ :  $\hat{u}_{jk} = \hat{u}_{jk} - \hat{l}_{jj-1} \hat{u}_{j-1k}$

$$\hat{u}_{jk} = \text{fl} \left( a_{jk} - \sum_{i=1}^{j-1} \hat{l}_{ji} \hat{u}_{ik} \right) \quad \begin{matrix} \leftarrow \text{HW1} \\ \text{QG 1a)} \end{matrix}$$

$$= a_{jk}(1+\delta_0) - \sum_{i=1}^{j-1} \hat{l}_{ji} \hat{u}_{ik} (1+\delta_i) \quad |\delta_i| \leq (j-1) \varepsilon_{\text{mach}}$$

$$a_{jk} = \frac{1}{1+\delta_0} \left[ \hat{u}_{jk} \hat{l}_{jj} + \sum_{i=1}^{j-1} \hat{l}_{ji} \hat{u}_{ik} (\delta_i + \delta_0) \right]$$

Lem 1

for  $|\delta| = O(\varepsilon)$

$$\begin{aligned} \frac{1}{1+\delta} &= 1 - \delta + \delta^2 - \dots \\ &= 1 - \delta + O(\varepsilon^2) \end{aligned}$$

Lem 2

for  $|\delta_{1,2}| = O(\varepsilon)$

$$\frac{1+\delta_1}{1+\delta_2} = 1 + \delta_1 - \delta_2 + O(\varepsilon^2)$$

$$\begin{aligned} a_{jk} &= \underbrace{\hat{u}_{jk} \hat{l}_{jj} (1 - \delta_0)}_{\text{red}} + \underbrace{\sum_{i=1}^{j-1} \hat{l}_{ji} \hat{u}_{ik} (\delta_i + \delta_0 - \delta_0)}_{\text{red}} \\ &= (\|\vec{L}\| \|\vec{u}\|)_{jk} + E_{jk} \end{aligned}$$

$$\text{when } E_{jk} = -\hat{u}_{jk} \hat{l}_{jj} \delta_0 + \sum_{i=1}^{j-1} \hat{l}_{ji} \hat{u}_{ik} (\delta_i - \delta_0) + O(\varepsilon^2)$$

and can be bounded by

$$\underline{j} \varepsilon \max(\|\vec{L}\| \|\vec{u}\|)_{jk} + O(\varepsilon^2)$$

if we combine it together with analysis on L and on U,  
we prove the proposition.

- the algo is backstable if

$$\frac{\|\delta A\|}{\|A\|} = O(\varepsilon)$$

$$\|M\| = \max_{i,j} |M_{ij}|$$

$$\|L\| = I \text{ in GEPP}$$

$$\|\delta A\| = O(n \varepsilon \|L\| \|\hat{u}\|)$$

$$\leq O(n \varepsilon (\|L\| \|\hat{u}\|))$$

$$= O(n \varepsilon \|\hat{u}\|)$$

$$\frac{\|\delta A\|}{\|A\|} = O(n \varepsilon \frac{\|A\|}{\|A\|})$$

define growth factor

$$g := \frac{\|\hat{u}\|}{\|A\|} = \frac{\max_{i,j} |A_{ij}|}{\max_{i,j} |A_{ij}|}$$

- in GG

$$A = \begin{bmatrix} 10^{-20} & 1 \\ 1 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 10^{-20} & 1 \end{bmatrix} \begin{bmatrix} 10^{-20} & 0 \\ 0 & 10^{-20} \end{bmatrix}$$

$L$                    $U$

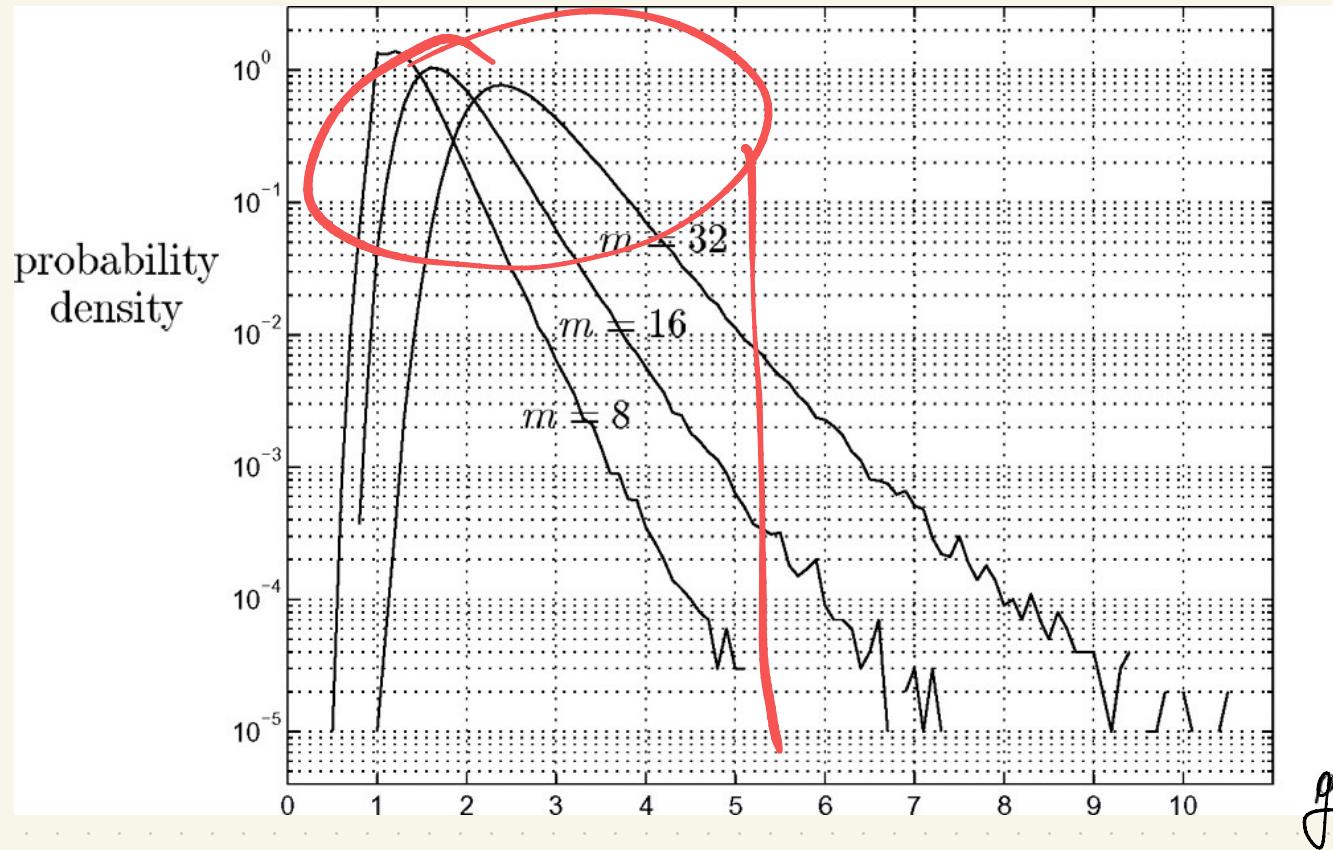
- in GEPP

$$PA = \begin{bmatrix} 1 & 0 & 0 & 1 \\ -1 & 1 & 0 & 1 \\ -1 & -1 & 1 & 1 \\ -1 & -1 & -1 & 1 \end{bmatrix}^m = \begin{bmatrix} 1 & 0 & 0 & 1 \\ -1 & 1 & 0 & 1 \\ -1 & -1 & 1 & 1 \\ -1 & -1 & -1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 2 \\ 0 & 1 & 8 \end{bmatrix} \leftarrow$$

$A$                    $L$                    $U$

Thm  $\|g\| \leq 2^{m-1}$

by def. GEPP is backward stable



g

- Solving Linear Systems

$A \in \mathbb{R}^{n \times n}$ , non-singular

$$(A, b) \xrightarrow{f} x : Ax = b$$

$$(A + \delta A, b) \xrightarrow{f} x + \delta x : (A + \delta A)(x + \delta x) = b \quad \begin{matrix} \text{RE} \\ \text{in output} \end{matrix}$$

$$\begin{aligned} & \Rightarrow Ax + A\delta x + \delta A x + \delta A \delta x = b \\ & \qquad \qquad \qquad \nearrow O(\varepsilon^2) \\ & \qquad \qquad \qquad = b \\ \hline & \qquad \qquad \qquad Ax \end{aligned}$$

$$A\delta x + \delta A x = 0$$

$$\delta x = -A^{-1} \delta A \cdot x$$

Relative error in the output of  $f$ , w.r.t. the relative error in input

$$\frac{\|\delta x\|}{\|x\|} = \frac{\|A^{-1} \delta A \cdot x\|}{\|x\|} \stackrel{\text{triang ineq., tight}}{\leq} \|A^{-1}\| \|\delta A\|$$

$$\begin{aligned} \frac{\|\delta x\|}{\|x\|} &\leq \|A^{-1}\| \|\delta A\| \\ &= \|A^{-1}\| \|A\| \underbrace{\frac{\|\delta A\|}{\|A\|}}_{\text{condition number of } A} \end{aligned}$$

$\kappa(A)$   
condition number of  
 $A$   
and of this problem

- revisit Mat-Vec Mult

$$(A, v) \xrightarrow{f} y : y = Av$$

$$\text{RE } \frac{\|y - \hat{y}\|}{\|y\|} \leq \underbrace{\frac{\|E\|}{\|A\|} \|A^{-1}\| \|A\|}_{\text{cond}(A)}$$

where  $(A + E) \cdot v = \hat{y}$

- GEPP

$$\frac{\|\delta A\|}{\|A\|} = O(n \varepsilon g) \Rightarrow \frac{\|\delta x\|}{\|x\|} \leq \sigma(kA) n \varepsilon g$$

$$k_2(A) = \|A^{-1}\|_2 \|A\|_2$$

$$= \frac{\sigma_1(A)}{\sigma_m(A)}$$