

Welcome to Ma221!. Lecture 12, Sep 20

Gauss. Elim (GE)

for  $i = 1:n-1$

if  $A(i,i) = 0$  but  $A(j,i) \neq 0, j > i$

swap rows  $i$  and  $j$ , record swap  $P$

if all  $A(j,i) = 0$ , signal error

$A(i+1:n, i) = A(i+1:n, i) / A(i, i)$  ... BLAS1

$A(i+1:n, i+1:n) = A(i+1:n, i+1:n) -$   
 $A(i+1:n, i) \cdot A(i, i+1:n)$

when done  $A = \begin{bmatrix} U \\ L \end{bmatrix}$   $P \cdot L \cdot U = \text{original } A$  ... BLAS2

How to "pivot", choose  $A(j, i)$ , or  $A(i, i)$  to put on diagonal

Goal: Backward Stability:

$$P \cdot L \cdot U = A + E \quad \|E\| = O(\epsilon) \cdot \|A\|$$

Not guaranteed just by  $A(i, i) \neq 0$

Ex: single precision  $\epsilon \approx 10^{-7}$

$$A = \begin{bmatrix} 10^{-8} & 1 \\ 1 & 1 \end{bmatrix} \quad A^{-1} \approx \begin{bmatrix} -1 & 1 \\ 1 & -10^{-8} \end{bmatrix}$$

$\kappa(A) \approx 2.6$  very well conditioned

$$L = \begin{bmatrix} 1 & 0 \\ 10^8 & 1 \end{bmatrix} \quad U = \begin{bmatrix} 10^{-8} & 1 \\ 0 & \text{fl}(1 - 10^8 \cdot 1) \end{bmatrix} \approx \begin{bmatrix} 10^{-8} & 1 \\ 0 & -10^8 \end{bmatrix}$$

"forgotten"

$$L \cdot U = \begin{bmatrix} 1 & 0 \\ 10^8 & 1 \end{bmatrix} \cdot \begin{bmatrix} 10^8 & 1 \\ 0 & -10^8 \end{bmatrix} = \begin{bmatrix} 10^8 & 1 \\ 1 & 0 \end{bmatrix}$$

Get same  $L, U$  if  $A(2,2)$  were  $.5, -1, \dots$   
any  $O(1)$  number

$\Rightarrow$  Solving  $Ax=b$  gives same answer  
for all these different  $A$ 's, all wrong

Instead: swap rows 1 & 2, get  $\begin{bmatrix} 1 & 1 \\ 10^8 & 1 \end{bmatrix}$   
proceed with GE, all good to 8 digits

Intuition: want large entry of  $A$  on diagonal  
pick largest entry: "partial pivoting"

Recall: HW 1.10  $C = fl(A \cdot B) = A \cdot B + E$

$$|E| \leq n \cdot \epsilon \cdot |A| \cdot |B|$$

since  $A = P \cdot L \cdot U$ , get similar bound:

Thm: Backward Error for GE:

If  $P, L, U$  from GE:

$$A - E = P \cdot L \cdot U, |E| \leq n \cdot \epsilon \cdot P \cdot |L| \cdot |U|$$

Cor: Solve  $Ax=b$  by 1) GE, 2) forward sub with  $L$   
3) back sub. with  $U$

Computed  $\hat{x}$  satisfies  $(A - F)\hat{x} = b$

$$|F| \leq 3 \cdot n \cdot \epsilon \cdot P \cdot |L| \cdot |U|$$

Proof of Cor: Use HW 1.11 (assume  $P=I$ )

Solve  $Ly=b$ , get  $(L + \delta L)\hat{y} = b$ ,  $|\delta L| \leq n \cdot \epsilon \cdot |L|$

Solve  $Ux = \hat{y}$ , get  $(U + \delta U)\hat{x} = \hat{y}$ ,  $|\delta U| \leq n \cdot \epsilon \cdot |U|$

$$\begin{aligned} b &= (L + \delta L)\hat{y} = (L + \delta L)(U + \delta U)\hat{x} \\ &= (LU + \delta L \cdot U + L \cdot \delta U + \delta L \cdot \delta U)\hat{x} \\ &= (A - E + \delta L \cdot U + L \cdot \delta U + \delta L \cdot \delta U)\hat{x} \quad \text{by Thm} \\ &= (A - F)\hat{x} \end{aligned}$$

$$\begin{aligned} \|F\| &\leq \|E\| + \|\delta L \cdot U\| + \|L \cdot \delta U\| + \|\delta L \cdot \delta U\| \\ &\leq \|E\| + \|\delta L\| \cdot \|U\| + \|L\| \cdot \|\delta U\| + \|\delta L\| \cdot \|\delta U\| \\ &\leq n\epsilon \|L\| \cdot \|U\| + n \cdot \epsilon \|L\| \cdot \|U\| + n \cdot \epsilon \|L\| \cdot \|U\| + n^2 \epsilon^2 \|L\| \cdot \|U\| \\ &\leq 3n\epsilon \|L\| \cdot \|U\| \quad \text{QED of Cor} \end{aligned}$$

for backward stability need  $\|F\| = O(\epsilon) \|A\|$

$$\|F\| \leq n\epsilon \|L\| \cdot \|U\|, \text{ want}$$

$$\|L\| \cdot \|U\| = O(\|A\|) \text{ depends on } P$$

Proof of Thm: Trace through Alg

$U(i,j)$  computed as

$$\begin{aligned} U(i,j) &= A(i,j) - L(i,1) \cdot U(1,j) \\ &\quad - L(i,2) \cdot U(2,j) \\ &\quad \dots \end{aligned}$$

$$= A(i,j) - \sum_{k=1}^{i-1} L(i,k) \cdot U(k,j)$$

also follows from  $A(i,j) = \text{dot product of row } i \text{ of } L \text{ \& col } j \text{ of } U$

use previous analysis of dot products  
when  $i > j$  same idea

$$L(i,j) = \left( A(i,j) - \sum_{k=1}^{j-1} L(i,k) \cdot U(k,j) \right) / U(i,i)$$

again dot product analysis applies QED of Thm

⇒ intuition: choose  $P$  so  $|L(i,j)|$  and  $|U(j,k)|$   
not too large

(1) Standard: Partial pivoting (GEPP)

At each step choose largest entry (in abs. value)  
in each column to put on diagonal

$$|L(k,i)| = |A(k,i) / A(i,i)| \leq 1$$

Thm (easy) with GE,  $|L_{ij}| \leq 1$

$$\text{and } \max_{i,j} (|U(i,j)|) \leq 2^{n-1} \max_{i,j} (|A(i,j)|)$$

Bad news: attainable, lose all precision  
in single, at  $n=24$

Good news: hardly ever happens, so GEPP  
standard alg

Empirical observation:  $\frac{\| |L| \cdot |U| \|}{\|A\|} = g \approx n^{2/3}$

If entries of  $A$  were i.i.d  $N(0,1)$   
true with high probability

(2) Complete Pivoting: Permute rows and columns  
so each  $A(i,i)$  is largest in all remaining  
rows and cols

$$A = P_r \cdot L \cdot U \cdot P_c$$

more stable: Thm:  $g < n^{(\log n)^4}$  vs  $2^{n-1}$

Empirically:  $g < n^{1/2}$  vs  $n^{2/3}$

rarely used, cost =  $O(n^3)$  additional

(3) Can't hit lower bound on latency in parallel case, instead use

Tournament Pivoting

(4) Threshold Pivoting (sparse case)  
tradeoff stability vs. sparsity of  $L, U$

Final Error bound for  $Ax=b$ , get  $(A-f)x=b$

$$\frac{\|x - \hat{x}\|}{\|x\|} \leq \text{condition number} \cdot \text{backward error}$$

$$\leq k(A) \cdot 3 \cdot n \cdot \underbrace{g}_{\text{growth factor}} \cdot \|A\|$$

$\|L\| \cdot \|U\|$

we can estimate  $k(A)$ ,  $g$  in  $O(n^2)$  extra work

What if error bound too large,  
or you want a guarantee?

Try iterative refinement, aka Newton's Method

$$f(x) = Ax - b$$

$$f'(x) = A$$

$$x_{i+1} = x_i - (f'(x_i))^{-1} f(x_i)$$

$$= x_i - A^{-1} (Ax_i - b)$$

$$= A^{-1} b \quad \text{residual}$$

why not immediate convergence?

round off  
residual mostly noise

Use: mixed precision:

most work ( $O(n^3)$  part) in low prec.

rest ( $O(n^2)$  part) in high prec:

computing  $r = Ax_i - b$

Low/high could mean

single/double

half/single

or

bf/lat/b

double/quad (software)

even 3 precisions in use

Do GEPP to solve  $Ax=b$  in low prec

get initial solution  $x(1)$

$i=1$

repeat  $r = A \cdot x(i) - b$  .. in high prec

$O(n^2)$  cost

solve  $Ad=r$  in low prec

$O(n^2)$  cost

update  $x(i+1) = x(i) - d$  in low prec

$O(n)$  cost

until "convergence", depends on goals

(1) Getting a small backward error in

high precision

straight word stopping criterion:

$$\|r\| \leq \epsilon_{\text{high}} \cdot (\|A\| \|x_i\| + \|b\|)$$

(2) Getting small forward error

$$\|x_{\text{comp}} - x_{\text{exact}}\| / \|x_{\text{true}}\| = O(\epsilon)$$

complicated stopping criterion:  
gesvxx in LAPACK

Is it worth doing itref in one precision?  
get  $|E| \leq n \cdot \epsilon |A|$