

Welcome to Ma221! Lec 11, Fall 24

## Recall Sparse Cholesky

Choose ordering (RCM, MD, ND...)

only depends on nonzero locations,  
not values, to minimize work

Build data structures for  $A$ ,  $L$

Perform factorization.

Contrast with GE for general  $A$ :

ordering needs to depend on entries

usual partial pivoting could make  
lots of fill in

Lots of algorithms, examine a few:

- ① Threshold pivoting, among pivot choices at each step pick one within a factor of 2 or 3 of largest, with least fill in: tradeoff stability and speed

## ② Static Pivoting (SuperLU)

- ① reorder and scale  $A$  to make diagonal as large as possible

Thm: for any nonsingular  $A$

$\exists$  perm  $P$  and 2 diagonal  $D_1, D_2$

$$\text{s.t. } B = D_1 A P D_2$$

$$(*) |B(i,i)| = 1 \text{ and } |B(i,j)| \leq 1$$

- ② reorder rows and columns of  $B$  using same techniques as Cholesky, maintains  $(*)$
- ③ During factorization, if a prospective pivot is too small, make it bigger (rare)

Difference between  $B$  and factorized matrix is low rank  $\Rightarrow$  use Sherman-Morrison-Woodbury for fast solver, or GMRES (from Chap 6)

Lots of algorithms + software (see class webpage)

# Structured Matrices

("data sparse")

could be dense, depend on  $O(n)$  data

Many structures (depends on physics, ...)

Common case with common structure

Vandermonde:  $V(i, j) = x_i^{j-1}$

Cauchy  $C(i, j) = \frac{1}{x_i + y_j}$

Toeplitz  $T(i, j) = x_{i-j}$   
constant along diagonals

Hankel  $H(i, j) = x_{i+j}$

Eg:  $Vz = b$  mean  $\sum_{j=1}^n x_i^{j-1} \cdot z_j = b_i$

$\Rightarrow$  polynomial interpolation

$O(n^2)$  using Newton Interp.

similar trick for  $V^T z = b$

Eg Multiplying  $Tz \equiv$  convolution

$\Rightarrow$  use FFT

Eg: Solving  $Cx = b$  arises in  
rational interpolation

Common Structure of all these  $X$ :

$$AX + XB = \text{low rank}$$

for some simple  $A$  and  $B$

Ex: Vandermonde  $V$

$$D = \text{diag}(x_1, \dots, x_n)$$

$$D \cdot V = V \text{ "shifted left"}$$

$$V \cdot P = V \cdot \begin{bmatrix} 0 & 0 & \dots & 1 \\ 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix}, = V \text{ shifted left}$$

$$DV - VP = \text{zero except last column} \\ = \text{rank 1}$$

Ex: Toeplitz  $T$

$$P \cdot T - T \cdot P = T \text{ shifted down} \\ - T \text{ shifted left}$$

= mostly zero, nonzero  
in first row, last column

$$= \text{rank 2}$$

Ex: Cauchy  $C$

$$\text{diag}(x_1, \dots, x_n) \cdot C + C \cdot \text{diag}(y_1, \dots, y_n)$$

$$= \text{all ones} = \text{rank 1}$$

Def: this low rank is called  
"displacement rank"

Thm (Kailath et al) There is an  
 $O(n^2)$  solver if displacement  
rank is  $O(1)$   
(stability not guaranteed)

---

## Chap 3: Least Squares

Ex: polynomial fitting

given sample point  $(y_i, b_i)$   $i=1, \dots, n$   
find "best" polynomial of fixed degree  
to minimize  $\sum_{i=1}^m (p(y_i) - b_i)^2$

minimize  $\|Ax - b\|_2$

$A(i, j) = y_i^{j-1}$ ,  $x$  are coeffs of  $p$   
 $p(y) = x_1 + x_2 y + x_3 y^2 + \dots + x_j y^{j-1}$

Matlab demo: `sin(pi*y/5) + y/5`  
`polyfit(3).m`

Standard Notation:

$\operatorname{argmin}_x \|Ax - b\|_2$   $A^{m \times n}$   $m \geq n$

$m > n$  means over determined  
don't expect  $Ax = b$  exactly

Other variants (all in LAPACK)

Constrained LS  $\operatorname{argmin}_x \|Ax - b\|_2$   
 $x: Bx = y$

where  $\# \text{rows}(B) \leq \# \text{cols}(A) \dots$  so  $Bx = y$  not overdetermined

$\leq \# \text{rows}(A) + \# \text{rows}(B) \dots$  so  $x$  unique

Weighted LS:  $\operatorname{argmin}_x \|y\|_2$  s.t.

$$b = Ax + By$$

if  $B = I$ ,  $y = b - Ax \Rightarrow$  standard LS

if  $B$  square, nonsingular

$$\operatorname{argmin}_x \|B^T(Ax - b)\|_2$$

Underdetermined  $\# \text{rows}(A) < \# \text{cols}(A)$

so  $\operatorname{argmin}_x \|Ax - b\|$  not unique

can also arise if  $A$  not full rank  $\Rightarrow$   
space of solutions (add any  $z: Az = 0$  to  
a solution to get another)

to make solution unique, use

$$\operatorname{argmin}_x \|x\|_2 \quad \text{s.t.} \quad Ax = b$$

# Ridge Regression

$$\operatorname{argmin}_x \|Ax - b\|_2^2 + \lambda \|x\|_2^2$$

$\lambda > 0$  tuning parameter

solution unique if  $\lambda > 0$

# Total Least Squares

$$\operatorname{argmin}_x \|[E, r]\|_2$$

$$(A+E)x = b+r$$

Algorithms for overdetermined LS  
(building blocks for all other cases)

Solve: Normal Equations (NE)

$$A^T A x = A^T b \quad (\text{real case})$$

$A^T A$  s.p.d.  $\Rightarrow$  Cholesky

fastest in dense case

(fewest flops, least comm)

not stable if  $A$  ill-conditioned

Use QR decomposition  $A^{m \times n} = Q R$

$Q^{m \times n}$  orthogonal

$R^{n \times n}$  upper triangular

$$x = R^{-1} Q^T b$$

Gram-Schmidt - unstable if  $A$  ill-conditioned  
(lots of variants, trading off  
speed and stability)

Householder - stable ( $x = A \backslash b$  in Matlab)  
(blocked Householder, to use BLAS3)

(possible to get  $Q$  and  $R$  via NE,  
called CholeskyQR, fast but can be as  
unstable as NE)

SVD: most "complete" solution  
gives cond. number, error bounds,  
works in rank deficient case,  
expensive.

Convert to a square linear system  
with  $A$ , and  $A^T$  in a bigger matrix  
(also for sparsity) (see Q3.3)

---

Normal Equations:

Thm:  $A$  full column rank, then  
solution of  $A^T A x = A^T b$  (NE)  
minimizes  $\|Ax - b\|_2$

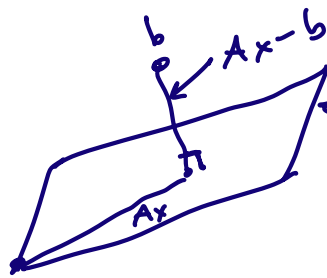
proof: Assume  $x$  satisfies NE, show  
it minimizes  $\|Ax - b\|_2$ : show

$$\|A(x+e) - b\|_2^2 \text{ minimized at } e=0$$



$$\begin{aligned}
&= (A(x+e)-b)^T (A(x+e)-b) \\
&= (Ax + Ae - b)^T (Ax + Ae - b) \\
&= (Ax - b + Ae)^T (Ax - b + Ae) \\
&= (Ax - b)^T (Ax - b) + (Ae)^T (Ae) \\
&\quad + \underbrace{2e^T (A^T (Ax - b))}_{=0 \text{ by NE}}
\end{aligned}$$

$$\begin{aligned}
&= \|Ax - b\|_2^2 + \|Ae\|_2^2 \\
&\geq \|Ax - b\|_2^2, \quad = \text{if } Ae = 0 \\
&\quad \text{since } A \text{ full rank, } \Rightarrow e = 0
\end{aligned}$$



← all vectors  $A \cdot y$

$Ax - b \perp Ay$  for all  $y$

$$(Ay)^T (Ax - b) = 0 \quad \forall y$$

$$y^T (A^T Ax - A^T b) = 0 \quad \forall y$$

$$\Rightarrow A^T Ax = A^T b$$

Cost: form  $A^T A$  + solve  $A^T Ax = b$

in flops cost =  $mn^2$  +  $\frac{n^3}{3}$

in # words know how to minimize both steps  
moved

$$QR: A = QR \quad A^{m \times n}, \quad Q^{m \times n} \quad R^{n \times n}$$

orthonormal columns ▷

solution  $\arg \min_x \|Ax - b\|_2 = R^{-1} Q^T b$

proof 1:  $A = QR = \begin{bmatrix} Q & \hat{Q} \end{bmatrix} \begin{bmatrix} R \\ 0 \end{bmatrix}$  orthogonal

$$\|Ax - b\|_2^2 = \left\| \begin{bmatrix} Q^T \\ \hat{Q}^T \end{bmatrix} (Ax - b) \right\|_2^2$$

$$= \left\| \begin{bmatrix} Q^T \\ \hat{Q}^T \end{bmatrix} (QRx - b) \right\|_2^2$$

$$= \left\| \begin{bmatrix} Q^T QRx - Q^T b \\ \hat{Q}^T QRx - \hat{Q}^T b \end{bmatrix} \right\|_2^2$$

$$= \left\| \begin{bmatrix} Rx - Q^T b \\ 0 - \hat{Q}^T b \end{bmatrix} \right\|_2^2$$

$$= \|Rx - Q^T b\|_2^2 + \|\hat{Q}^T b\|_2^2$$

$$\geq \|\hat{Q}^T b\|_2^2, \quad \text{if } Rx = Q^T b$$

$$\text{or } x = R^{-1} Q^T b$$

proof 2: plug into  $NE$

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

$$= ((QR)^T (QR))^{-1} (QR)^T b$$

$$\begin{aligned}
&= (R^T \underbrace{Q^T Q}_I R)^{-1} R^T Q^T b \\
&= (R^T R)^{-1} R^T Q^T b \\
&= R^{-1} \underbrace{R^T R^T}_I Q^T b \\
&= R^{-1} Q^T b
\end{aligned}$$

Algorithms for  $A=QR$

Classical + Modified Gram-Schmidt  
CGS and MGS

Equate columns of  $A=QR$

$$A(:,i) = \sum_{j=1}^i Q(:,j) R(j,i)$$

since columns of  $Q$  orthogonal

$$Q(:,j)^T \cdot A(:,i) = R(j,i)$$

for  $i=1$  to  $n$

$$tmp = A(:,i)$$

for  $j = 1 : i-1$

cost = 2m ...  $R(j,i) = Q(:,j)^T \cdot A(:,i)$  ... CGS

"  $R(j,i) = Q(:,j)^T \cdot tmp$  ... MGS

"  $tmp = tmp - R(j,i) \cdot Q(:,j)$

end for

$$R(i,i) = \|tmp\|_2$$

$Q(:, i) = \text{tmp} / R(i, i)$   
end for

$$\# \text{ flops} = 2mn^2 + O(mn)$$
$$\sim 2 \cdot \text{cost}(NE) \quad \text{if } m \gg n$$

Householder - stable

MG-S - less stable

CG-S - even less stable

## 2 Metrics for backward stability

want accurate factorization

$$A + E = QR \quad \|E\| = O(\epsilon) \cdot \|A\|$$

also want  $Q$  close to orthogonal

$$\|Q^T Q - I\| = O(\epsilon)$$

Fixes for GS: better but still not <sup>guaranteed</sup> stable

$$\text{MG-S2: MG-S twice} \quad A = QR = (Q_1 R_1) \cdot R$$
$$= Q_1 (R_1 R)$$

$$\text{CG-S2} = \text{CG-S twice}$$