

Last time:  $A \in \mathbb{C}^{m \times n}$  ( $m \geq n$ ),  $b \in \mathbb{C}^n$ ,  $\text{rank}(A) = n$

Least Squares: find  $\hat{x} = \underset{x \in \mathbb{C}^n}{\text{argmin}} \|b - Ax\|_2^2$

Solution:  $\hat{x} = (A^*A)^{-1} A^* b$  (normal equation)  
 $= V \Sigma^\dagger U^* b$  (SVD,  $A = U \Sigma V^*$ )

method	normal eqn	SVD
conditioning	$\approx \kappa_2^2(A)$	$\approx \kappa_2(A)$ ( $\ r\ $ small)
operation count	$\approx 2mn^2$ FLOPs	$\approx 4mn^2$ FLOPs

- Solving least-squares via QR factorization

Computing SVD is expensive. can we work with orthogonal transform but lower cost?

- (reduced) QR factorization

$$A = Q R$$

$\uparrow$   
 $Q \in \mathbb{C}^{m \times n}$   
 $Q^* Q = I_n$   
 orthonormal columns ( $m \geq n$ )

$\leftarrow$   
 $R \in \mathbb{C}^{n \times n}$   
 upper triangular matrix

Not  $Q Q^* = I_m$

when  $\text{rank}(A) = n$ ,  $\Rightarrow \text{rank}(R) = n$

hence  $R(A) = R(Q)$

Let  $Q = [\dot{q}_1 \dots \dot{q}_n]$ ,  $q_i \in \mathbb{C}^m$

projection  
onto  $R(A)$  → Then  $P_A x = P_Q x = \sum_{i=1}^n q_i (q_i^* x) = Q Q^* x$

Hence  $\hat{x} \in \arg \min_x \|Ax - b\|_2^2$

$$\Leftrightarrow r = b - A\hat{x} \perp R(A)$$

$$\Leftrightarrow P_Q r = 0$$

$$\Leftrightarrow QR\hat{x} = QQ^*b$$

$$\Leftrightarrow R\hat{x} = Q^*b$$

$$\Leftrightarrow \hat{x} = R^{-1} Q^*b$$

Actually, from normal equations,

$$\hat{x} = (A^*A)^{-1} A^*b = (R^*R)^{-1} R^* Q^*b = R^{-1} Q^*b$$

- How to compute QR factorization?

Idea 1: convert basis to orthonormal one

Gram-Schmidt: given  $a_1, \dots, a_n \in \mathbb{C}^m$ , produce orthonormal  $q_1, \dots, q_n$  with  $\text{span}\{q_1, \dots, q_n\} = \text{span}\{a_1, \dots, a_n\}$

$$q_1 = \frac{a_1}{\|a_1\|_2}$$

$$q'_2 = a_2 - (q_1^* a_2) q_1, \quad q_2 = \frac{q'_2}{\|q'_2\|_2}$$

$$\dots$$

$$q'_j = a_j - \sum_{i=1}^{j-1} (q_i^* a_j) q_i, \quad q_j = \frac{q'_j}{\|q'_j\|_2} \quad j=3, \dots, n$$

$$\Leftrightarrow a_1 = \|a_1\|_2 q_1$$

$$a_2 = (q_1^* a_2) q_1 + q_2 \|a_2 - (q_1^* a_2) q_1\|_2$$

$$\dots$$

$$a_j = \sum_{i=1}^{j-1} (q_i^* a_j) q_i + q_j \|a_j - \sum_{i=1}^{j-1} (q_i^* a_j) q_i\|_2$$

$$\Leftrightarrow \left| \begin{array}{l} A = QR \\ A = [\hat{a}_1 \dots \hat{a}_n], \quad R_{ij} = \begin{cases} q_i^* a_j, & i < j \\ \|a_j - \sum_{k=1}^{j-1} (q_k^* a_k) q_k\|_2, & i=j \\ 0, & \text{otherwise} \end{cases} \end{array} \right.$$

Implementation:

Classical GS Input:  $A \in \mathbb{C}^{m \times n}$  of rank  $n$

Output:  $Q \in \mathbb{C}^{m \times n}, R \in \mathbb{C}^{n \times n}$

```

for j = 1, ..., n
  for i = 1, ..., j-1      * operation count:
     $R_{ij} = q_i^* a_j \rightarrow (j-1)(2m-1)$ 
  end
   $q'_j = a_j - \sum_{i=1}^{j-1} R_{ij} q_i \rightarrow 2m(j-1)$ 
   $R_{jj} = \|q'_j\|_2, \quad q_j = q'_j / R_{jj} \approx \sum_{j=1}^n 4jm$ 
end

```

$\approx 2mn^2$  FLOPs

• weakness: there is nothing to force orthogonality of  $Q$   
in classical Gram-Schmidt

ex.  $A = \begin{bmatrix} 1 & 1 & 1 \\ \varepsilon & & \\ & \varepsilon & \\ & & \varepsilon \end{bmatrix}$

assume  $\varepsilon$  so small such that  $1/(1+\varepsilon^2) \approx 1$

then  $Q_{GS} = \begin{bmatrix} 1 & 0 & 0 \\ \varepsilon & -1/\sqrt{2} & -1/\sqrt{2} \\ 0 & 1/\sqrt{2} & 0 \\ 0 & 0 & 1/\sqrt{2} \end{bmatrix}$

$q_1 \quad q_2 \quad q_3$

$$q_2^* q_3 = \frac{1}{2}$$

- Remedy: Instead of orthogonalizing  $a_j$  at step  $j$  only,  
can orthogonalize  $c_j$  to  $q_i$  as soon as  $q_i$  is computed

modified GS

Let  $a_k^{(1)} = a_k$ ,  $k = 1, \dots, n$

for  $k = 1, \dots, n$

operation counts  
 $\approx 2mn^2$

$$R_{kk} = \|a_k^{(k)}\|_2, \quad q_k = a_k^{(k)} / R_{kk}$$

for  $j = k+1, \dots, n$

$$R_{kj} = q_k^* a_j^{(k)}, \quad a_j^{(k+1)} = a_j^{(k)} - R_{kj} q_k$$

end

end

classical GS is equivalent to modified GS  $\leftarrow$  exercise

Thm Suppose Modified GS is applied to  $A \in \mathbb{R}^{m \times n}$  of rank  $n$   
yielding  $\hat{Q} \in \mathbb{R}^{m \times n}$ ,  $\hat{R} \in \mathbb{R}^{n \times n}$ ,

$\exists C_i = C_i(m, n)$ , s.t.

$$1) \quad \underline{A + \Delta A_1 = \hat{Q} \hat{R}}, \quad \|\Delta A_1\|_2 \leq C_1 \varepsilon_{\text{mach}} \|A\|_2$$

$$\|\hat{Q}^T \hat{Q} - I\|_2 \leq C_2 \varepsilon_{\text{mach}} K_2(A) / (1 - C_2' \varepsilon_{\text{mach}} K_2(A))$$

2)  $\exists Q \in \mathbb{R}^{m \times n}$  with orthonormal columns s.t.

$$A + \Delta A_2 = Q \hat{R}, \quad \|\Delta A_2\|_2 \leq C_3 \varepsilon_{\text{mach}} \|A\|_2$$

$$\|Q - \hat{Q}\|_2 \leq C_4 \varepsilon_{\text{mach}} K_2(A) / (1 - C_4' \varepsilon_{\text{mach}} K_2(A))$$

Pf: Higham Thm 9.13 

Remark: The theorem states that the departure from orthogonality of  $\hat{Q}$  is bounded by  $O(K_2(A) \epsilon_{mach})$

Remark: Part 2) of the theorem states that  $\hat{R}$  is the exact triangular QR factor of a matrix near to  $A$ , i.e. it is a good  $R$ -factor.

---

• Back to least-squares:

Solve least-squares via QR:

$$\hat{x} = \hat{R}^{-1} \hat{Q}^* b$$

implicitly use  $\hat{Q}^* \hat{Q} = I_n$  but suffer from large  $K_2(A)$

Remark: If we translation the result to perturbation form, we get that the computed  $\hat{x}$  is the exact solution

of the LS problem

$$\|b + \Delta b - (A + \Delta A)y\|_2^2$$

$$\text{where } \|\Delta b\|_2 \lesssim K_2(A) \epsilon_{mach}, \frac{\|\Delta A\|_2}{\|A\|_2} \lesssim \epsilon_{mach}$$

NOT Backward stable

• To resolve this problem, we can apply Modified GS to  $[A \ b]$ , so that  $Q^* b$  is implicitly computed in the last step of GS

$$[A \ b] = [Q, q_{n+1}] \begin{bmatrix} R & z \\ 0 & \rho \end{bmatrix}$$

$$\text{We have } Ax - b = [A \ b] \begin{bmatrix} x \\ -1 \end{bmatrix}$$

$$= [Q, q_{n+1}] \begin{bmatrix} R x - z \\ -\rho \end{bmatrix}$$

$$= Q_1(Rx - z) - \rho q_{n+1}$$

Hence  $\|Ax - b\|_2^2 = \|Rx - z\|_2^2 + \rho^2$

So  $x = R^{-1}z$  is the Least-squares solution

Thm Solving LS via Modified GS for  $[A \ b]$

has forward error as good as a backward stable algo.

- It is possible to perform QR factorization faster than  $2mn^2$  FLOPs if we don't form  $Q$  explicitly

\* Householder :

$$\begin{array}{c} A \\ \begin{pmatrix} x & x & x \\ x & x & x \\ x & x & x \\ x & x & x \end{pmatrix} \end{array} \xrightarrow{Q_1} \begin{array}{c} \begin{pmatrix} x & x & x \\ 0 & x & x \\ 0 & x & x \\ 0 & x & x \end{pmatrix} \end{array} \xrightarrow{Q_2} \begin{array}{c} \begin{pmatrix} x & x & x \\ 0 & x & x \\ 0 & 0 & x \\ 0 & 0 & x \end{pmatrix} \end{array} \xrightarrow{Q_3} \begin{array}{c} \begin{pmatrix} x & x & x \\ 0 & x & x \\ 0 & 0 & x \\ 0 & 0 & 0 \end{pmatrix} \end{array}$$

$\begin{bmatrix} 1 & 0 \\ 0 & H_2 \end{bmatrix}$        $\begin{bmatrix} 1 & 0 \\ 0 & H_3 \end{bmatrix}$

$H_2$

$$Q_3 Q_2 Q_1 A = \begin{bmatrix} R \\ 0 \end{bmatrix} \Rightarrow A = (Q_3 Q_2 Q_1)^* \begin{bmatrix} R \\ 0 \end{bmatrix}$$

$Q_i$  unitary

$$= Q \begin{bmatrix} R \\ 0 \end{bmatrix}$$

← full QR

$Q \in \mathbb{C}^{m \times n}$  unitary

$\begin{bmatrix} R \\ 0 \end{bmatrix} \in \mathbb{C}^{m \times n}$  upper triangular

\* How to choose  $Q_k$ ?

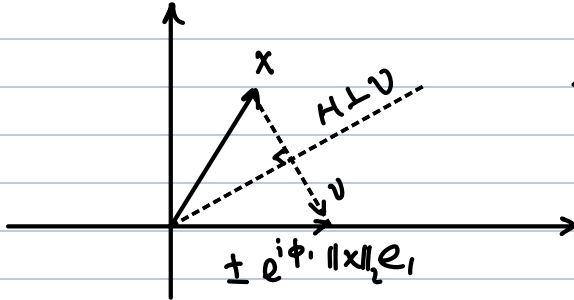
$$\text{Let } Q_k = \begin{bmatrix} \overset{k-1}{\overset{m-(k-1)}}{\mathbf{I}} & 0 \\ 0 & H_k \end{bmatrix} \begin{matrix} k-1 \\ m-(k-1) \end{matrix}$$

keep the first  $k-1$  coordinates unchanged

,  $H_k$  unitary (thus norm preserving)

$x \in \mathbb{C}^{m-(k-1)}$  be  $k$ th, ...,  $m$ th entries of  $k$ th column

Goal:  $H_k x = \pm e^{i\phi_1} \|x\|_2 e_1$   $\begin{pmatrix} x \\ x \\ x \end{pmatrix} \xrightarrow{H_k} \begin{pmatrix} x \\ 0 \\ 0 \end{pmatrix}$



Householder transform:  
reflection across the plane  $H$

$$H_k x = x - 2 P_v x = x - 2 \left( \frac{v^*}{\|v\|_2} x \right) \frac{v}{\|v\|_2}$$

$$= \left( I - \frac{2 v v^*}{\|v\|_2^2} \right) x$$

Householder transform

$\phi_1 = \arg(x_1)$  is added to make  $x, \pm e^{i\phi_1} \|x\|_2 e_1$  symmetric with respect to  $H$ .  
i.e.  $v^* x = -v^* (\pm e^{i\phi_1} \|x\|_2 e_1)$ .  
such that  $H_k x = \pm e^{i\phi_1} \|x\|_2 e_1$ ,  
(otherwise  $H_k$  doesn't reflect  $x$  to  $\pm e^{i\phi_1} \|x\|_2 e_1$ !)

(a) For stability, one may choose  $v = e^{i\phi_1} \|x\|_2 e_1 + x$   
to avoid catastrophic cancellation in computing  $\|v\|_2^2$  when  $\|x\|_2 \approx |x_1|$   
why?

(b) Householder matrices  $H_k$  is never formed in upper triangularizing  $A$ , storage and computations use solely the Householder vector  $v$ .

### Implementation:

Householder QR:

For  $k = 1, \dots, n$

$$x \leftarrow A(k:m, k)$$

$$v_k = \text{sign}(x_1) \|x\|_2 e_1 + x$$

$$v_k \leftarrow v_k / \|v_k\|_2$$

$$A(k:m, k:n) \leftarrow A(k:m, k:n) - 2 v_k (v_k^* A(k:m, k:n))$$

end

### operation count:

$$\approx \sum_{k=1}^n 4(m-k)(n-k) \text{ FLOPs}$$

$$\approx 4 \left( mn^2 - (m+n) \frac{n^2}{2} + \frac{n^3}{3} \right)$$

$$\approx 2mn^2 - \frac{2}{3}n^3$$

subtraction  
(m-k)(n-k)  
FLOPs

outer product

(m-k)(n-k)  
FLOPs

matrix-vector product

2(m-k)(n-k) FLOPs

- Stability of Householder transform

Since the orthogonal transform is performed by Householder vector, the orthogonality of  $Q$  is enforced

Thm Let  $\begin{bmatrix} \hat{R} \\ 0 \end{bmatrix} \in \mathbb{R}^{m \times n}$  be the computed upper triangular QR

factor of  $A \in \mathbb{R}^{m \times n}$ , Let  $Q = (\hat{Q}_n \cdots \hat{Q}_1)^T$  be the exact orthogonal matrix obtained from Householder vectors computed from the algorithm

Then  $A + \Delta A = Q \hat{R}$  with  $\|\Delta A\|_2 \leq C_{m,n} \|A\|_2$

Let  $\hat{Q} = \text{fl}((\hat{Q}_n \cdots \hat{Q}_1)^T)$  be the computed orthogonal matrix

Then  $\hat{Q} = Q(I_m + \Delta I)$  with  $\|\Delta I\|_2 \leq C_{m,n}$

$\hat{Q}$  is very close to orthogonal matrix regardless of  $K(A)$ !

(This a consequence of the backward stability of matrix multiplication)

- Back to solve least-squares

Solve least-squares via QR:

$$\hat{x} = \hat{R}^{-1} \hat{Q}^* b$$

use  $\hat{Q}^* \hat{Q} = I_n$ , but NOT suffer from large  $K_2(A)$

Thm Let  $A \in \mathbb{R}^{m \times n}$  have full rank and that the

least square problem  $\min_x \|Ax - b\|_2$  is solved using Householder QR factorization. The computed solution

$\hat{x}$  is the exact solution to



$$\min_x \|b + \Delta b - (A + \Delta A) \hat{x}\|_2$$

$$\text{where } \|\Delta A\|_2 \leq C_{\min} \|A\|_2, \quad \|\Delta b\|_2 \leq C_{\min} \|b\|_2$$

↑  
Backward stable!

Summary:

Operation count:  
 $O(2mn^2)$

$O(4mn^2)$

normal equation < Householder QR < Modified GS < SVD  
(without computing Q explicitly) (Augmented)

Rounding error:

$O(K_2^2(A) \epsilon_{\text{mach}})$

$O(K_2(A) \epsilon_{\text{mach}})$

normal equation > Householder QR  $\approx$  Modified GS  $\approx$  SVD  
(without computing Q explicitly) (Augmented)