

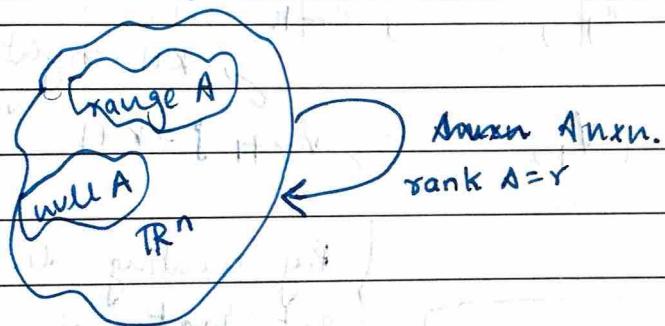
★ → PCA ↔ Eigenvalue decomposition of  $A^T A$ .

$$A^{n \times n} \xrightarrow{A} \mathbb{R}^m$$

Four fundamental subspaces

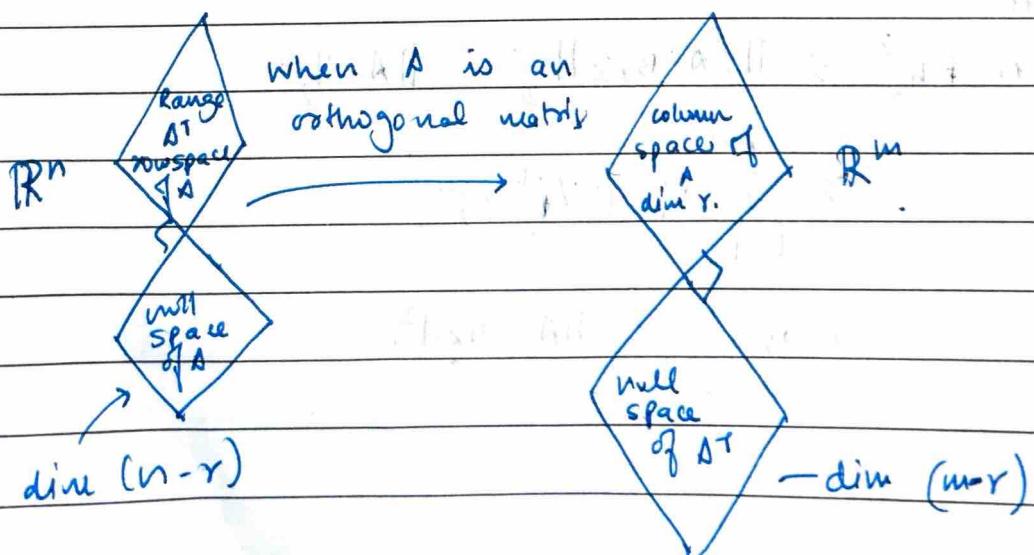
- range  $A$  is subspace of  $\mathbb{R}^m \rightarrow$  (column space of  $A$ )
- null space of  $A$  is subspace of  $\mathbb{R}^n \rightarrow \{x \mid Ax=0\}$
- range  $A^T$  is subspace of  $\mathbb{R}^n \rightarrow$  (column space of  $A^T$ )
- null space of  $A^T$  is subspace of  $\mathbb{R}^m \rightarrow \{y \mid A^T y=0\}$ .

col space of  $A^T = \text{null space of } A$ .



$$\dim(\text{null}(A)) + \dim(\text{range } A) = n$$

$$m-r + r = n.$$



## Projection matrices

$P^2 = P$  • range  $P \cap \text{null } P = \{0\}$

• range  $P^\perp \cap \text{null } P = V$ .

A projection  $P$  is an orthogonal projector  $\iff P = P^*$

range  $P \perp \text{null } P$

Suppose  $P = P^*$ , then let  $x \in \text{null } P$ ,  $y \in \text{range } P$

= range  $(I - P)$

$y = Pw$

$$Py = P^2 w = y$$

want  $\langle x, y \rangle = 0$ .

$$\langle (I - P)x, Py \rangle = y^* P^* (I - P)x$$

$$\begin{aligned} \langle (I - P)x, Py \rangle &= y^* P^* (I - P)x \\ &= y^* (Px - P^2 x) = 0. \end{aligned}$$

[WLOG]:  $q_i$ 's are orthonormal

$$Pq_i = q_i \quad \forall 1 \leq i \leq n$$

(conversely suppose range  $P = \{q_1, \dots, q_n\}$  (say)

& null  $P = \{q_{n+1}, \dots, q_m\}$  say

$$\rightarrow Pq_j = 0 \quad \forall n+1 \leq j \leq m$$

Let

$$Q = \left[ \begin{array}{c|c|c|c} q_1 & \cdots & q_n & | & q_{n+1} & \cdots & q_m \end{array} \right]$$

$$Q^* = \left[ \begin{array}{c} q_1^* \\ \hline q_2^* \\ \vdots \\ q_m^* \end{array} \right]$$

$$PQ = \left[ \begin{array}{c|c|c|c} q_1 & q_2 & \cdots & q_n \\ \hline | & | & \cdots & | \\ 0 & 0 & \cdots & 0 \end{array} \right]$$

$$Q^* PQ =$$

— / —

$$Q^* P Q$$

$$= \begin{bmatrix} q_1^* \\ \vdots \\ q_m^* \end{bmatrix} \begin{bmatrix} q_1 & \cdots & q_n & 0 & \cdots & 0 \end{bmatrix} = \begin{bmatrix} 1 & & & & & \\ & 1 & & & & \\ & & \ddots & & & \\ & & & 0 & & 0 \end{bmatrix}$$

$$Q^* P Q = \text{diag } (1, \underbrace{\dots, 1}_{n \text{ s}}, 0, \dots, 0)$$

$$P = Q \Sigma Q^* \rightarrow \text{SVD of } P.$$

$$Q = \begin{bmatrix} q_1 & \cdots & q_n & 0 & \cdots & 0 \end{bmatrix}$$

$$P = Q \Sigma Q^* = \hat{Q} \hat{\Sigma} \hat{Q}^*$$

$\hat{Q} \rightarrow m \times n$ .

$$P^* = (\hat{Q} \hat{\Sigma} \hat{Q}^*)^* = \hat{Q} \hat{\Sigma}^* \hat{Q}^* = P$$

for any orthogonal projection  $P$

$$P = \hat{Q} \hat{\Sigma} \hat{Q}^* \rightarrow \text{range of } P \text{ is colspace of } \hat{Q}$$

and action of  $P$  on  $V$

projection of  $V$   
on colspace of  $\hat{Q}$

$$P: V \xrightarrow{m} V$$

i.e.  
 $m$        $\dim \text{range } P = n \leq m$

$$P^2 = P$$

(if  $n=m$ , then  $P=I$ )

\*  $\therefore P$  is always rank deficient.

\* Orthogonal projection  $P$  is not an orthogonal matrix

Warning: Orthogonal projections are not orthogonal matrices.

Special case:  $\hat{q} = [q_1] \quad \|q\| = 1$

Then  $P = q_1 q_1^* \rightarrow$  rank 1 matrix.  
(denote by  $P_q$ )

$P_q(v) =$  projection of  $v$  in the direction of  $q$   
(or on the subspace spanned by  $q$ )

$$I - P = I - q q^*$$

is orthogonal projection on the subspace  $q^\perp$   
(rank  $(m-1)$  projection). It is denoted by  $P_{\perp q}$

Summary: The orthogonal projection on a subspace  $\{q_1, \dots, q_n\}$   
where  $q_i$ 's are orthonormal is given by

$$P = \hat{q} \hat{q}^*$$
$$\hat{q} = [q_1 | \dots | q_n]$$

\* Orthogonal projection on any subspace  $A$  with basis

$$S = \{a_1, \dots, a_m\}$$

$$\text{let } A = [a_1 | \dots | a_m]$$

let  $y$  be the orthogonal projection of a vector  $x$  onto vector range ( $A$ ).

$$y - x \in \text{range}(A)$$

let  $P_A$  denote the required orthogonal projection.

$$\text{Then } y = P_A(x)$$

$$y - x = P_A(x) - x \in \text{null } P_A \perp \text{range } P_A$$

$$y - x \perp \text{range } A$$

$$y - x \perp \text{range } A$$

$$y - x \perp a_j \quad \forall 1 \leq j \leq n$$

$$\therefore a_j^* (y - x) = 0 \quad \forall a_j$$

$y \in \text{range } A \Rightarrow y = Av$  for some  $v$ .

$$\therefore a_j^* (Av - x) = 0 \quad \forall a_j$$

$$A^* (Av - x) = 0$$

$$A^* Av - A^* x = 0$$

$$A^* Av = A^* x$$

( $A$  is full rank  $\Rightarrow A^* A$  non singular)

$$v = (A^* A)^{-1} A^* x$$

$$v = A^{-1} (A^* - 1) A^* x$$

$$\boxed{\downarrow = A^{-1} x}$$

$$\boxed{v = (A^* A)^{-1} A^* x}$$

$$y = Pv = \underbrace{A(A^*A)^{-1}A^*}_{P_A} x$$

Summary:

The orthogonal projection on a subspace with basis  $\{a_1, \dots, a_n\}$  which are not orthogonal is given by

$$P_A = A(A^*A)^{-1}A^*, \text{ where } A = [a_1 \dots | a_n]$$

$(I - P)$  is an orthogonal projection on a subspace  $A^\perp$

(orthogonal to  $A$ ).

\* Last time

- Projection  $P: V \rightarrow V$  is a linear transformation satisfying  $P^2 = P$ .
- Orthogonal projection  $P$  :
  - ① range  $P \perp$  null  $P$
  - ②  $P^* = P$  (symmetric)
- Formula for orthogonal projection.
- (i) on a subspace  $W$  spanned by an orthonormal set of vectors  $\{q_1, \dots, q_n\}$  ( $n < m$ ).

$$P_w = \hat{Q}\hat{Q}^*, \text{ where } Q = [q_1 | \dots | q_n].$$

- (ii) Orthogonal projection on 1-dimensional subspace spanned by  $q \in V$  :  $P_q = qq^*$  (rank 1 matrix).  
 $\|q\|=1$

(iii) On a subspace  $A$  spanned by with basis  $\{a_1, \dots, a_n\}$

$$P_A = A (A^* A)^{-1} A^*, \quad A = [a_1 | \dots | a_n]$$

(iv) On a 1-dimensional subspace spanned by  $a \in V$ ,

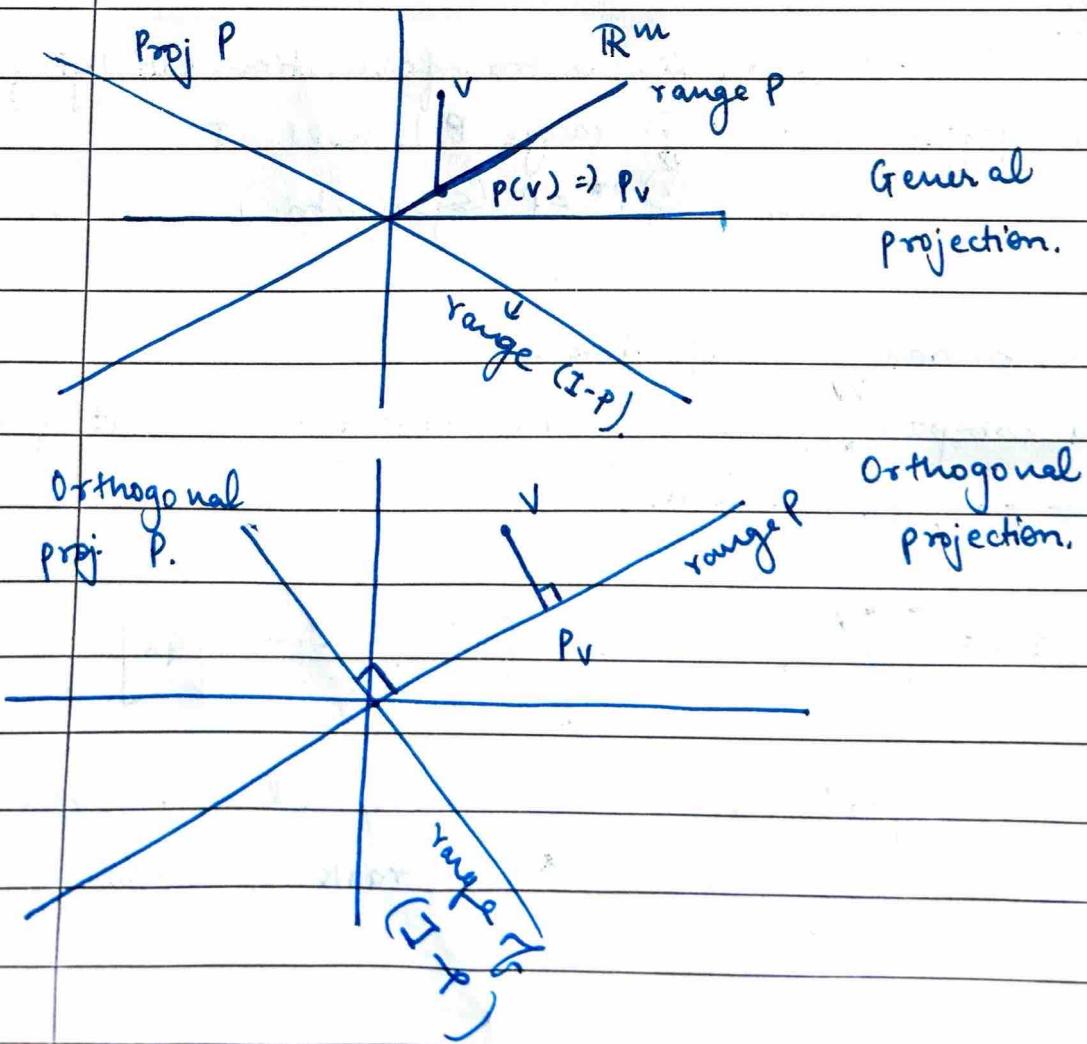
$$\boxed{P_a = \frac{aa^*}{\|a\|}} = \frac{aa^*}{a^*a} \quad \because \|a\| = a^*a.$$

(v)  $I - P_W, I - P_{V^\perp}, I - P_A, I - P_{a^\perp}$

$$P_{W^\perp} \quad P_{V^\perp} \quad P_{A^\perp} \quad P_{a^\perp}.$$



Schematic picture



## \* QR factorization

Defn: Given  $A \in \mathbb{R}^{m \times n}$  ( $m > n$ ), the QR factorization of  $A$  is a factorization of the form

$$\begin{bmatrix} A \\ \end{bmatrix}_{m \times n} = \begin{bmatrix} Q \\ \end{bmatrix}_{m \times m} \begin{bmatrix} R \\ \end{bmatrix}_{m \times n}$$

orthogonal      upper triangular

$$\begin{bmatrix} * & * & * \\ 0 & * & * \\ 0 & 0 & * \\ \end{bmatrix}$$

full QR  
factorization.

When

$$[A]_{m \times n} = [Q]_{m \times n} [R]_{n \times n} \quad (\text{Reduced QR}).$$

Suppose

$$\begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nn} \end{bmatrix} = \begin{bmatrix} q_{11} & \dots & q_{1n} \\ \vdots & \ddots & \vdots \\ q_{n1} & \dots & q_{nn} \end{bmatrix} \begin{bmatrix} r_{11} & \dots & r_{1n} \\ \vdots & \ddots & \vdots \\ 0 & \dots & r_{nn} \end{bmatrix}$$

$$= \begin{bmatrix} r_{11} q_{11} \\ r_{11} q_{21} \\ \vdots \\ r_{11} q_{n1} \end{bmatrix}$$

$$a_{11} = r_{11} q_{11}$$

Gauss Elimination

$$a_1 = r_{11} q_1$$

$$\Rightarrow \boxed{q_1 = a_1 / r_{11}}$$

$$a_2 = r_{12} q_1 + r_{22} q_2$$

$$\therefore q_2 = \frac{a_2 - r_{12} q_1}{r_{22}}$$

Gauss Elimination

$$a_n = r_{1n} q_1 + \dots + r_{nn} q_n$$

$$\text{so, } q_n = \frac{a_n - r_{1n} q_1 - r_{2n} q_2 - \dots - r_{(n-1)n} q_{n-1}}{r_{nn}}$$

Since, columns of  $Q$  are orthonormal,

$$\|q_1\| = 1 \Rightarrow r_{11} = \|a_1\|$$

$$\|q_j\| = 1 \Rightarrow r_{jj} = \left\| a_j - \sum_{i=1}^{j-1} r_{ij} q_i \right\| \quad \left. \begin{array}{l} \text{Graham} \\ \text{Schmidt} \\ \text{orthonormalization.} \end{array} \right\}$$

$$q_i^* a_j$$

Thus, Graham Schmidt orthonormalization process gives rise to QR factorization.

Theorem: Every  $A \in \mathbb{C}^{m \times n}$  (wlog) has a full QR factorization & hence also a reduced QR factorization.

If  $A$  is full rank, then G-S orthonormalization gives QR factorisation.

At each step we can construct  $q_j$ , since  $a_j \neq 0 \forall j$

Recall: In G.S. orthonormalization given a set  $A = \{a_1, \dots, a_n\}$  produce an orthonormal set  $\{q_1, \dots, q_n\}$  satisfying

- ①  $\|q_i\| = 1 \quad \forall i$
- ②  $q_i^* q_j = 0 \quad \forall i \neq j$
- ③  $\text{span}\{a_1, \dots, a_k\} = \text{span}\{q_1, \dots, q_k\} \quad \forall 1 \leq k \leq n$ .

If  $A$  does not have full rank, then for some  $j$ ,

$$\begin{aligned} v_j &= a_j - \sum_{i=1}^{j-1} r_{ij} q_i \\ &= a_j - \sum_{i=1}^{j-1} (q_i^* q_j) a_j = 0. \end{aligned}$$

Then we can choose  $q_j$  arbitrarily satisfying

$$a_j \perp \{q_1, \dots, q_{j-1}\}$$

Then continue with next columns as usual.

$$a_j - \sum (q_i^* a_j) q_i$$

$$a_j = \sum_{i=1}^{j-1} r_{ij} q_i \quad \boxed{a_j - \sum_{i=1}^{j-1} (q_i^* q_j) a_j = r_j}$$

## \* Uniqueness of QR factorization

=> Unique upto the sign of elements of diagonal of R.  
for full rank A.

=> If  $A = QR$  is a reduced QR factorization of A, then for any  $z \in \mathbb{C}$  s.t.  $|z|=1$ , multiplying the  $i^{\text{th}}$  column of Q by z &  $i^{\text{th}}$  row of R by  $z^{-1}$ .

In particular the sign of  $r_{jj}$  can be +ve or -ve.

If A is full rank &  $r_{jj} > 0 \forall j$ , then the QR factorization of A is unique.

## \* Algorithm for Classical Gram Schmidt (CGS)

Only the inner loop.

for  $j=1$  to  $n$

$$v_j = a_j$$

for  $i=1$  to  $j-1$ .

$$r_{ij} = q_i^* a_j$$

$$v_j = v_j - r_{ij} q_i$$

$$r_{jj} = \|v_j\|_2 \quad (\text{2-norm})$$

$$q_i = \frac{v_j}{r_{jj}}$$

\* Modified Gram Schmidt  
for ensuring Numerical stability.

MGS

for  $i=1$  to  $n$ :

$$v_i = q_i$$

$$r_{ii} = \|v_i\|$$

$$q_i = v_i / r_{ii}$$

for  $j = i+1$  to  $n$

$$r_{ij} = q_i * v_j$$

$$v_j = v_j - r_{ij} q_i$$

Operation count  $\approx 2mn^2$ .

\* A high level view of Gram Schmidt orthonormalization

Given  $A = [a_1 | \dots | a_n]$

G.S.  $\leftrightarrow$  right multiplying  $A$  with suitable upper triangular matrices to obtain a matrix with orthogonal columns.

$$(A R_1 R_2) = Q$$

$$A R^{-1} = Q$$

$$A = Q R$$

$$\text{where } R_1 = \begin{bmatrix} 1 & -\frac{\gamma_{12}}{\gamma_{11}} & -\frac{\gamma_{13}}{\gamma_{11}} & \cdots \\ \frac{1}{\gamma_{11}} & 1 & & \\ & & 1 & \\ & & & \ddots \end{bmatrix}$$

$$R_2 = \begin{bmatrix} 1 & & & \\ & 1 & -\frac{\gamma_{23}}{\gamma_{22}} & \\ & \frac{1}{\gamma_{22}} & 1 & \\ & & & \ddots \end{bmatrix}$$

$$R_3 = \begin{bmatrix} 1 & & & \\ & 1 & -\frac{\gamma_{34}}{\gamma_{33}} & \\ & \frac{1}{\gamma_{33}} & 1 & \\ & & & \ddots \end{bmatrix} \text{ so on.}$$

### \* Householder's method for QR factorization

In contrast, Householder's method is to premultiply A with certain unitary matrices  $Q_k$  so that

$(Q_n \dots Q_1) A$  is an upper triangular matrix R

$$(Q_n \dots Q_1) A = R$$

$\underbrace{Q^{-1}}$

$$Q^{-1} A = R$$

$$\boxed{A = QR}$$

\* Householder constructed the matrix  $Q_k$  as follows:

$$Q_k = \begin{bmatrix} I_{k-1} & 0 \\ 0 & F \end{bmatrix} \quad \text{where } F \text{ is a } (n-k+1) \times (n-k+1) \text{ matrix}$$

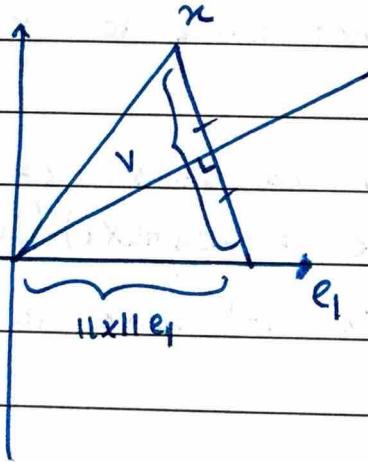
such that  $F$  introduces zeroes below the diagonal in the  $k^{\text{th}}$  column of  $(Q_{k-1} \cdots Q_1) A$

Suppose the  $k^{\text{th}}$  column of  $(Q_{k-1} \cdots Q_1) A$  is

$$\begin{pmatrix} x_1 \\ \vdots \\ x_{n-k+1} \end{pmatrix} \xrightarrow{n} \|x\| e_1 \begin{pmatrix} \|x\| \\ 0 \\ \vdots \\ 0 \end{pmatrix}_{(n-k+1)}$$

$x_1$  is first component.

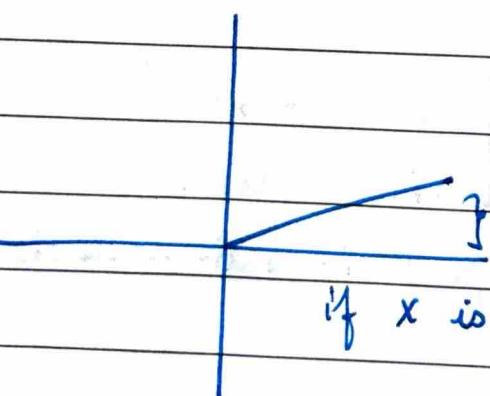
$$k \begin{pmatrix} * \\ * \\ \vdots \\ * \end{pmatrix} \quad \text{vector } x \in \mathbb{C}^{n-k+1}$$



$$P_{\perp v} = I - \frac{vv^*}{v^*v} \quad | \quad v = \|x\| e_1 - x$$

$$P_{\perp v} = I - 2vv^* \quad | \quad \text{Householder reflections.}$$

$$v = -\operatorname{sg}(x_1) \|x\| e_1 - x.$$



If  $x$  is close to  $e_1$

Householder's idea was to choose a vector  $v$  such that the reflection across  $v$  maps  $x$  to  $\|x\|e_1$ .  
 If  $v$  is taken to be ~~length(x)~~  
 $v = \pm \|x\|e_1 - x$

$$v = \pm \|x\|e_1 - x \quad \text{then } f = I - \frac{vv^*}{v^*v}$$

satisfies the purpose.

\* Note:

- ①  $f = f^*$
- ②  $f^2 = I$
- ③  $f$  is unitary (full rank)
- ④ For mathematical purposes we can choose any  $z$   
 s.t.  $|z| \cdot \|x\|e_1 - x$  is not close to  $x$ , where  $z \in \mathbb{C}$   
 s.t.  $|z|=1$ .

But for numerical purposes we choose  $z = -\operatorname{sgn}(x_1)$ , so that  
 $v = -\operatorname{sgn}(x_1) \|x\|e_1 - x = -\operatorname{sgn}(x_1) (\|x\|e_1 + x)$ .

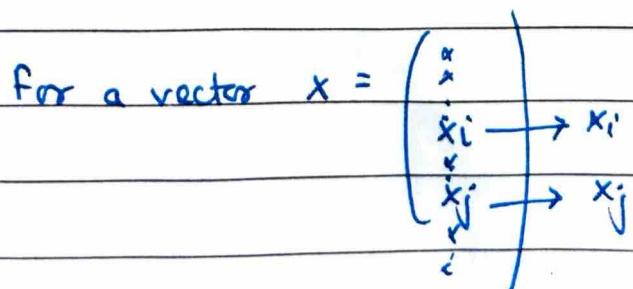
This makes sure that  $v$  is not close to 0. Hence, we avoid cancellation errors.

\*

Givens' Rotations

A Givens rotation  $R(\theta) = \begin{bmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{bmatrix}$

rotates any vector  $x$  through  $\theta$  in anti-clockwise direction.



For a vector  $x$ , the aim is to introduce ' $0$ ' at the  $j^{th}$  spot affecting a change in the  $i^{th}$  spot, without affecting the rest of the entries.

$\begin{matrix} i & & j \\ \vdots & & \vdots \end{matrix}$

let  $R(i, j, \theta) = \begin{bmatrix} 1 & & & & \\ & \ddots & & & \\ & -\cos\theta & & -\sin\theta & \\ & & \sin\theta & \cos\theta & \\ \downarrow & & & & \\ j & & & & 1 \end{bmatrix}$  s.t.

$$\begin{bmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{bmatrix} \begin{bmatrix} x_i \\ x_j \end{bmatrix} = \begin{bmatrix} \sqrt{x_i^2 + x_j^2} \\ 0 \end{bmatrix}$$

$$\cos\theta = \frac{x_i}{\sqrt{x_i^2 + x_j^2}} \quad \text{A} \quad \sin\theta = \frac{-x_j}{\sqrt{x_i^2 + x_j^2}}$$

Rotation matrix preserves length.



### Givens' rotations

For a vector  $x = \begin{pmatrix} x \\ x_i \\ \vdots \\ x_j \\ x \end{pmatrix}$

we can zero out  $x_j$

by affecting a change in  $x_i$

(while not changing the  $\ell_2$  norm of  $x$ )

using

$$R(i, j, \theta) = \begin{bmatrix} 1 & & & & \\ & \ddots & & & \\ & -\cos\theta & & -\sin\theta & \\ & & \sin\theta & \cos\theta & \\ \downarrow & & & & \\ j & & & & 1 \end{bmatrix} \quad \text{s.t.}$$

$$\begin{bmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{bmatrix} \begin{bmatrix} x_i \\ x_j \end{bmatrix} = \begin{bmatrix} \sqrt{x_i^2 + x_j^2} \\ 0 \end{bmatrix}$$

$$\cos\theta = \frac{x_i}{\sqrt{x_i^2 + x_j^2}}$$

$$\sin\theta = \frac{-x_j}{\sqrt{x_i^2 + x_j^2}}$$

e.g. to proceed from.

$$\begin{pmatrix} x & x & x \\ 0 & x & x \\ 0 & 0 & x \\ 0 & 0 & x \\ 0 & 0 & x \end{pmatrix} \text{ to } \begin{pmatrix} x & x & x \\ 0 & x & x \\ 0 & 0 & x \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}$$

we proceed to  $\alpha$ -steps

$$\begin{pmatrix} 1 & x & x \\ 1 & 1 & 1 \\ c & -s & \\ s & c & \end{pmatrix} \begin{pmatrix} x & x & x \\ 0 & x & x \\ 0 & 0 & x \\ 0 & 0 & * \\ 0 & 0 & \alpha \end{pmatrix} = \begin{pmatrix} x & x & x \\ 0 & x & x \\ 0 & 0 & x \\ 0 & 0 & x \\ 0 & 0 & 0 \end{pmatrix}$$

Then apply,

$$\begin{pmatrix} 1 & & \\ 1 & & \\ c & -s & \\ s & c & \end{pmatrix} \begin{pmatrix} x & x & x \\ 0 & x & x \\ 0 & 0 & x \\ 0 & 0 & x \\ 0 & 0 & 0 \end{pmatrix} = \begin{pmatrix} x & x & x \\ 0 & x & x \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}$$

## \* Iterative methods for solving systems of linear equations

- ① Jacobi
- ② Gauss Siedel
- ③ Successive over relaxation (SOR).

To solve,

$$Ax = b$$

$$\frac{+\sqrt{3}}{3} + \frac{7}{\sqrt{3}} = \frac{-3+7}{2\sqrt{3}} = \frac{4}{2\sqrt{3}}$$

— / — / —

Generalities: Assume  $Ax=b$  has a unique solution.

Defn:

\* splitting of  $A$  is a decomposition  $A = S - T$   
 $A = S - T$ , where  $S$  is non-singular.

\* splitting yields an iterative method as follows:

$$Ax = b.$$

$$(S - T)x = b.$$

$$Sx - Tx = b.$$

$$Sx = Tx + b.$$

$$x = S^{-1}Tx + S^{-1}b.$$

This suggest the following iterative method

$$x_{k+1} = Rx_k + c \quad \text{where } R = S^{-1}T \rightarrow \begin{matrix} \text{matrix of the} \\ \text{iterative method.} \end{matrix}$$

$$\& c = S^{-1}b \rightarrow \text{vector}$$

Lemma 1:

Let  $\|\cdot\|$  any induced matrix norm.

If  $\|R\| < 1$ , then  $x_{k+1} = Rx_k + c$  converges for any initial solution  $x_0$ .

Proof:

$$\begin{aligned} x_{k+1} &= Rx_k + c \\ - x &= Rx + c \\ (x_{k+1} - x) &= R(x_k - x) \end{aligned}$$

Now,

$$\begin{aligned} \|x_{k+1} - x\| &= \|R(x_k - x)\| \leq \|R\| \|x_k - x\| \\ &\leq \|R\| \|R\| \|x_{k-1} - x\| \\ &\vdots \\ &\leq \|R\|^{k+1} \|x_0 - x\| \end{aligned}$$

$$\|R\| < 1 \Rightarrow \|R\|^{k+1} \rightarrow 0 \text{ as } k \rightarrow \infty$$

$$\text{i.e. } \|x_{k+1} - x\| \rightarrow 0 \text{ as } k \rightarrow \infty \text{ i.e. } \{x_{k+1}\} \rightarrow x \text{ as } k \rightarrow \infty$$

Lemma 2: Let  $\rho(R)$  denote the spectral radius of  $R$   
 $(\rho(R) = \max \{|\lambda| / \lambda \text{ is an eigenvalue of } R\})$

①  $\rho(R) \leq \|R\|$ , for every induced matrix norm.  
(depends on  $R$  &  $\epsilon$ ).

②  $\forall \epsilon > 0$ ,  $\exists \| \cdot \|_{\star}$  such that  
 $\|R\|_{\star} = \rho(R) + \epsilon.$

Proof : ① let  $\lambda$  be the  $\epsilon$ -value of  $R$  for which  $|\lambda| = \rho(R)$   
let  $x$  be the corresponding eigen vector i.e.  $Rx = \lambda x$   
 $Rx = \lambda x$

$$\|R\| = \max_{g \neq 0} \frac{\|Rg\|}{\|g\|} \geq \frac{\|Rx\|}{\|x\|} = \frac{\|\lambda x\|}{\|x\|} = |\lambda| = \rho(R)$$

② (Lemma 6.5 in Demmel's book).

Theorem: The iterative method  $x_{k+1} = Rx_k + b$   
converges to the solution of  $Ax = b$   $\forall x_0 \leftarrow b$ .  
 $\Leftrightarrow \rho(R) < 1$ .

Proof: If  $\rho(R) \geq 1$ , let  $|\lambda| = \rho(R)$  & let  $Ry = \lambda y$ .

Suppose  $y + x = x_0$ , so  $y = x_0 - x$ .

then

$$Ry = R(x_0 - x) = \lambda(x_0 - x).$$

$$x_{k+1} - x = R(x_k - x) \\ = R^2(x_{k-1} - x)$$

$$\vdots \\ = R^{k+1}(x_0 - x) = \lambda^{k+1}(x_0 - x) \rightarrow 0 \text{ as } k \rightarrow \infty$$

Then in a contradiction to the assumption  
 $\sigma(R) < 1$ .

Conversely, if  $\sigma(R) < 1$ .

choose a norm  $\|\cdot\|_*$  such that

$\|R\|_* < 1$ , then by lemma 1 the iterative method converges.

Questions:

Questions to be considered

① Is the given method convergent?

② Does  $\exists$  a matrix norm in which  $\|R\| < 1$ ? or  
 $\sigma(R) < 1$ ?

② Given 2 iterative methods which one converges faster?

Smaller spectral radius  $\Rightarrow$  faster convergence.

i.e. Which matrix gives smaller spectral radius?

(I)

Jacobi's Method

Assume,  $a_{ii} \neq 0 \forall i$ .

split A as  $D - E - F$  where

$D = \text{diag}(A)$

$-E = \text{lower } \Delta \text{ of } A$

$-F = \text{upper } \Delta \text{ of } A$ .

$$A = \begin{bmatrix} D & -F \\ -E & \end{bmatrix}$$

Now, let  $D = S$ .

and let  $T = E + F$

$$Ax = b$$

$$(D - E - F)x = b.$$

$$Dx - (E + F)x = b.$$

$$x = \underbrace{D^{-1}(E+F)}_{J} x + \underbrace{D^{-1}b}_{C}$$

Jacobi's method

$$x_{k+1} = Jx_k + C.$$

$$\text{where } J = D^{-1}(E+F)$$

$$A \quad C = D^{-1}$$

$$\text{Notice that } J = D^{-1}(E+F)$$

$$= D^{-1}E + D^{-1}F$$

$$= I - D^{-1}D + D^{-1}E + D^{-1}F$$

$$= I + \cancel{D^{-1}(E+F+D)}$$

$$= I - D^{-1}$$

$$= I - D^{-1}(-D+E+F)$$

$$= I - D^{-1}A.$$

$$= I - D^{-1}(D - E - F)$$

$$= I - D^{-1}A.$$

$$\begin{matrix} n \times n & n \times 1 \\ Ax = b. \end{matrix}$$

$$\text{Initial soln} \quad x_0 = \begin{pmatrix} x_1^0 \\ x_2^0 \\ \vdots \\ x_n^0 \end{pmatrix}$$

$$n \times n \quad n \times 1 \quad \left\{ \begin{array}{l} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n = b_1 \\ \vdots \\ a_{ij}x_1 + a_{j2}x_2 + \dots + a_{jn}x_n = b_j \\ \vdots \\ a_{nn}x_1 + a_{n2}x_2 + \dots + a_{nn}x_n = b_n \end{array} \right.$$

Initial soln  $D^{-1} = \begin{pmatrix} Y_{11} & & \\ & \ddots & \\ & & Y_{nn} \end{pmatrix}$

$$C = D^{-1} b = \begin{pmatrix} b_1/a_{11} \\ b_2/a_{22} \\ \vdots \\ b_n/a_{nn} \end{pmatrix}$$

Initial soln  $x_0 = \begin{pmatrix} x_1^0 \\ x_2^0 \\ \vdots \\ x_n^0 \end{pmatrix}$

$$x_1 = (I - D^{-1} A)x_0 + C.$$

$$x_0 - D^{-1} Ax_0$$

$$x_1^1 = x_0^1 + (x_0^1 + \frac{a_{12}x_0^2}{a_{11}} + \dots + \frac{a_{1n}x_0^n}{a_{11}}) + \frac{b_1}{a_{11}}$$

$$\begin{aligned} x_1^1 &= x_1^0 - (x_1^0 + \frac{a_{12}x_1^0}{a_{11}} + \dots + \frac{a_{1n}x_1^0}{a_{11}}) + \frac{b_1}{a_{11}} \\ &= \frac{1}{a_{11}} (b_1 - a_{12}x_1^0 - \dots - a_{1n}x_1^0) \end{aligned}$$

$$x_j^{(k+1)} = \frac{1}{a_{jj}} (b_j - a_{j1}x_1^{(k)} - a_{j2}x_2^{(k)} - \dots - a_{jn}x_n^{(k)})$$

(not present represented by a bracket)

$$x_j^{(1)} = \frac{1}{a_{jj}} (b_j - a_{j1}x_1^{(0)} - a_{j2}\overbrace{x_2^{(0)}}^{\text{(not present represented by a bracket)}} - \dots - a_{jn}x_n^{(0)})$$

$$x_{k+1} = (I - D^{-1}A)x_k + c.$$

$$x_1^{(k+1)} = \frac{1}{a_{11}} (b_1 - a_{12}x_2^{(k)} - \dots - a_{1n}x_n^{(k)})$$

$$x_2^{(k+1)} = \frac{1}{a_{22}} (b_2 - a_{21}x_1^{(k)} - a_{23}\overbrace{x_3^{(k)}}^{\text{(not present represented by a bracket)}} - \dots - a_{2n}x_n^{(k)})$$

$$\vdots$$

$$x_i^{(k+1)} = \frac{1}{a_{ii}} (b_i - a_{i1}x_1^{(k)} - a_{i2}x_2^{(k)} - \dots - \overbrace{a_{ii}x_i^{(k)}}^{\text{(not present represented by a bracket)}} - a_{in}x_n^{(k)})$$

$$x_n^{(k+1)} = \frac{1}{a_{nn}} (b_n - a_{n1}x_1^{(k)} - a_{n2}x_2^{(k)} - \dots - a_{n,n-1}\overbrace{x_{n-1}^{(k)}}^{\text{(not present represented by a bracket)}} - a_{nn}x_n^{(k)})$$

\* Algorithm for one step of Jacobi's method

for  $j = 1$  to  $n$

$$x_j^{(k+1)} = \frac{1}{a_{jj}} \left( b_j - \sum_{\substack{i=1 \\ i \neq j}}^n a_{ji} x_i^{(k)} \right)$$

\* Gauss-Siedel

$$x_1^{(k+1)} = \frac{1}{a_{11}} \left( b_1 - a_{12} x_2^{(k)} - a_{13} x_3^{(k)} - \dots - a_{1n} x_n^{(k)} \right)$$

$$x_2^{(k+1)} = \frac{1}{a_{22}} \left( b_2 - a_{21} x_1^{(k+1)} - \overbrace{a_{22} x_2^{(k)}}^{\text{from previous iteration}} - a_{23} x_3^{(k)} - a_{2n} x_n^{(k)} \right)$$

$$x_j^{(k+1)} = \frac{1}{a_{jj}} \left( b_j - a_{j1} x_1^{(k+1)} - \overbrace{a_{jj} x_j^{(k)}}^{\text{from previous iteration}} - \dots - a_{jn} x_n^{(k)} \right)$$

$$x_n^{(k+1)} = \frac{1}{a_{nn}} \left( b_n - a_{n1} x_1^{(k+1)} - \dots - \overbrace{a_{n,n-1} x_{n-1}^{(k+1)}}^{\text{from previous iteration}} - a_{nn} x_n^{(k+1)} \right)$$

The iterative method corresponds to

$$X^{k+1} = D^{-1} (E X^{k+1} + F X^k + b)$$

$$D X^{k+1} = E X^{k+1} + F X^k + b.$$

$$(D - E) X^{k+1} = F X^k + b.$$

$$X^{k+1} = \underbrace{(D - E)^{-1} F X^k}_{L_1} + \underbrace{(D - E)^{-1} b}_{C}$$

\* Algorithm for one step Gauss-Siedel method

for  $j = 1$  to  $n$        $\leftarrow$  updated terms.

$$x_j^{(k+1)} = \frac{1}{a_{jj}} \left( b_j - \sum_{i=1}^{j-1} a_{ji} x_i^{(k+1)} - \sum_{i=j+1}^n a_{ji} x_i^{(k)} \right)$$

III) Successive over-relaxation (SOR)

$\omega$ : relaxation parameter ( $\neq 0$ )

Splitting of  $A$ .

$$A = \begin{pmatrix} S & T \\ \frac{D-E}{\omega} & \frac{1-\omega}{\omega} D + F \end{pmatrix}$$

If  $\omega > 1$  overrelaxing  
 $\omega < 1$  underrelaxing  
 $\omega = 1$  Gauss Siedel.

The associated iterative method, is

$$X^{k+1} = L_\omega X^k + c,$$

$$\text{where, } L_\omega = \left( \frac{D-E}{\omega} \right)^{-1} \left( \frac{1-\omega}{\omega} D + F \right)$$

$$\text{& } c = \left( \frac{D}{\omega} - E \right)^{-1} b.$$

$$x_1^{(k+1)} = \frac{1}{a_{11}} (a_{11}x_1^{(k)} - w(a_{11}x_1^{(k)} + \dots + a_{1n}x_n^{(k)}) - b_1)$$

$$= (1-w)x_1^{(k)} + \frac{w}{a_{11}} (b_1 - a_{12}x_2^{(k)} - \dots - a_{1n}x_n^{(k)})$$

$$x_2^{(k+1)} = \frac{1}{a_{22}} (a_{22}x_2^{(k)} - w(a_{21}x_1^{(k+1)} + \dots + a_{2n}x_n^{(k+1)}) - b_2)$$

$$= (1-w)x_2^{(k)} + \frac{w}{a_{22}} (b_2 - a_{21}x_1^{(k+1)} - \dots - a_{2n}x_n^{(k)})$$

$$x_n^{(k+1)} = (1-w)x_n^{(k)} + \frac{w}{a_{nn}} (b_n - a_{n1}x_1^{(k+1)} - \dots - a_{n,n-1}x_{n-1}^{(k+1)})$$

Def<sup>n</sup>: The rate of convergence of  $x^{(k+1)} = Rx^{(k)} + c$   
 is.  $\gamma(R) = -\log \rho(R)$  (smaller the value of  $\rho(R)$   
 faster the convergence)

sections 6.5.3 & 6.5.4 of Demmel.

### \* Jacobi's method

- ① # allocations at each step =  $2n$ .
- ② Can compute all components of  $x^{(k+1)}$  in parallel.
- ③ Assumes that the diagonal entries are non-zero.
- ④ Does not always converge, it will converge for an order strictly diagonally dominant matrices.

### \* Strictly diagonally dominant matrices

- ①  $|a_{ii}| > \text{sum of rest of the row entries } + i$ .

### \* Properties of Gauss Siedel method

- 1) # allocations at each step =  $n$ .
- 2) Computation of components of  $x^{(k+1)}$  is not parallelizable.
- 3) Converges if  $A$  is HPD.

\* Common requirement : non-zero diagonal entries  
(which can be attained by rearranging eqns if required).

### \* Least Squares Method (for solving overdetermined $Ax=b$ )

$\begin{cases} A \text{ is full rank} \\ A \text{ is rank deficient.} \end{cases}$

2 cases       $\begin{cases} m \geq n, n \text{ unknowns} \\ m > n. \end{cases}$

If  $m > n$ , then in general the system of eqns  $Ax=b$  does not have a soln.

The aim is then to find  $x$  that minimizes the residual

$$r = Ax - b.$$

The problem is stated as follows:-

given  $A \in \mathbb{C}^{m \times n}$ ,  $m > n$ ,  $b \in \mathbb{C}^m$ , find  $x \in \mathbb{C}^n$   
such that

$$\|b - Ax\|_2 \text{ is minimized.}$$

\* (motivation for a general idea)

\* Polynomial least squares

Given  $n$  distinct pts.  $x_1, \dots, x_m$  & data  $y_1, \dots, y_n$   
at these pts. essentially

$$(x_1, y_1) \dots (x_m, y_m)$$

find a polynomial  $p(x)$  which fits this information in  
a way so that

$$p(x_1) = y_1$$

$$p(x_2) = y_2$$

$$p(x_m) = y_m$$

i.e.  $\sum |p_i(x) - y_i|^2$  is minimized

Suppose,  $p(x) = c_0 + c_1 x + \dots + c_{n-1} x^{n-1}$

$$\begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^{n-1} \\ 1 & x_2 & x_2^2 & \dots & x_2^{n-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^{n-1} \end{bmatrix} \begin{bmatrix} c_0 \\ c_1 \\ \vdots \\ c_{n-1} \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$

$$\begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^{n-1} \\ 1 & x_2 & x_2^2 & \dots & x_2^{n-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^{n-1} \end{bmatrix} \begin{bmatrix} c_0 \\ c_1 \\ \vdots \\ c_{n-1} \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$

## \* least squares problems

In general, we find a pt.  $\hat{x} \in \text{range } A$  s.t.

$$r = b - Ax \text{ is minimized.}$$

Geometrically, the idea is to obtain  $Ax$  ( $\approx x$ )

$Ax = Pb$ , where  $P$  is the orthogonal projection on range  $A$ .

where  $b$  is orthogonal projection

where  $r$  is the orthogonal projection of  $b$  on range of  $A$ .

by orthogonal projection of  $b$  on range  $A$ .

$P$  is orthogonal projection on range  $A$ .

Theorem:

$A \in \mathbb{C}^{m \times n}$ ,  $m > n$

A vector  $x \in \mathbb{C}^n$  minimizes  $\|r\|_2 = \|Ax - b\|_2 \Leftrightarrow r \perp \text{range } A$   
i.e.  $A^*r = 0$ .

Proof:

Note that for  $b \in \mathbb{C}^m$ , we can write  $b$  as  $b_1 + b_2$   
where  $b_1 \in \text{range } A$   $b_2 \in \text{null } A^*$

$$Ax - b_1 \in \text{range } A.$$

$$r = Ax - b = Ax - b_1 + b_2$$

$$\underbrace{r}_{\in \text{range } A} = \underbrace{b_2}_{\in \text{null } A^*}$$

$$\& b_2 \perp Ax - b_1$$

$$\|Ax - b\|_2^2 = \|Ax - b_1\|_2^2 + \|b_2\|_2^2$$

Since  $b_2$  is fixed, to minimize  $\|Ax - b\|_2^2$ , we must minimize  $\|Ax - b_1\|_2^2$ .  $\|Ax - b_1\|_2^2$  is minimized  $\Leftrightarrow Ax = b_1$   
 $Ax = b_1 \Leftrightarrow r = Ax - b = b_2 \in \text{null}(A^*) \Leftrightarrow r \perp \text{range } A$ .

this is the content of the earlier theorem.

$$r \perp \text{range } A \Leftrightarrow x \text{ minimizing } r = b - Ax$$
$$\Leftrightarrow A^* A x = A^* b.$$

Proof: Suppose  $r \perp \text{range } A$ .

let  $x$  be s.t.  $r = b - Ax$  is minimal.

$$A^* r = A^* (b - Ax) = A^* b - A^* Ax = 0.$$

Conversely if  $A^* Ax = A^* b$ ,  $\rightarrow$

$$\text{then } A^* (Ax - b) = 0.$$

$\underbrace{\quad}_{r}$

System of normal eq's for LSP  
 $Ax = b$  ( $A$  is full rank).

$$\Rightarrow r \perp \text{range } A.$$

Remarks: ① To see that  $Ax$  is indeed the orthogonal projection of  $b$  on range  $A$ ,

Consider  $P = A (A^* A)^{-1} A^*$

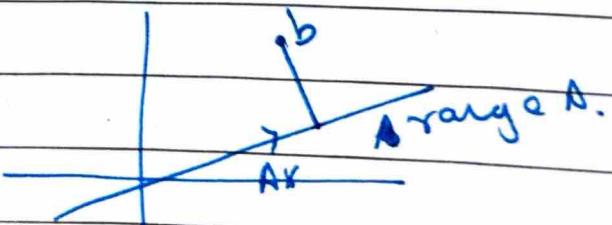
$$\text{Then } Pb = A (A^* A)^{-1} A^* b = A (A^* A)^{-1} A^* Ax$$
$$= Ax.$$

Theorem: The system of normal eq's  $A^* Ax = A^* b$  has a unique solution  $\Leftrightarrow A$  is full rank.

To summarize: If  $A$  has full rank, the soln to the LSP is

unique & is given by

$$x = \underbrace{(A^* A)^{-1}}_{A^+} A^* b = A^+ b, A^+ = \text{pseudo inverse of } A.$$



## \* Computing $x_{LS}$

There are 3 popular ways of computing  $x_{LS}$ :

(i)

### I Solve the normal eq's

when  $A$  is full rank

$A^*A$  is Hermitian Positive Definite

So,  $A^*Ax = A^*b$

can be solved using Cholesky factorization.

① form  $A^*A + A^*b$ .

② find  $R$  (upper triangular)  $A^*A = R^*R$

③ Solve  $R^*w = A^*b$

④ then solve for  $Rx = w$ .

Operation count is dominated by the first 2 steps

$$\frac{1}{3}n^3m^2 + \frac{1}{2}n^3 \approx O\left(\frac{mn^2 + n^3}{3}\right)$$

### II QR factorization

If  $A = QR$ .

then  $P = QG^*$

so,  $y = Pb = QG^*b$ .

Since,  $y \in \text{range } A$ ,  $Ax = y$  for some  $x$ .

$QRx = QG^*b$ .

$Rx = G^*b$

$y \in \text{range } A$  so we have unique solution.

III Using SVD.

$$\text{If } A = U\Sigma V^*$$

$$\text{Then } P = UU^*$$

$$y = Pb = UU^*b.$$

~~Ans~~

$\exists x \text{ s.t. } Ax = y \quad (\because y \text{ belongs to range } A).$

$$U\Sigma V^*x = UU^*b.$$

$$\boxed{\Sigma V^*x = U^*b}$$

### \* Rank deficient LSP's

If  $A$  is rank deficient, then there are infinitely many solutions for the corresponding LSP.

Consider the system,  $Ax = b$ , where  $A \in \mathbb{C}^{m \times n}$ ,  $m \geq n$ .

&  $\text{rank } A = r < n$ .

If  $x$  minimizes  $\|Ax - b\|_2$ , where  $z \in \text{null}(A)$   
also minimizes  $\|Ax - b\|_2$ .

To solve a rank deficient LSP, we consider a "complete orthogonal factorization" i.e. orthogonal matrices  $Q \in \mathbb{C}^{n \times n}$  s.t.

$$Q^T A \tilde{z} = \left[ \begin{array}{c|c} T_{11} & 0 \\ \hline 0 & 0 \end{array} \right] \gamma$$

$\underbrace{\hspace{1cm}}_{r}, \underbrace{\hspace{1cm}}_{n-r}$

If such a factorization is found then the

$$\begin{aligned} \|Ax - b\|_2^2 &= \|A \tilde{z} Z^T x - b\|_2^2 \\ &= \|Q^T A \tilde{z} Z^T x - Q^T b\|_2^2 \\ &= \| \underbrace{c}_{\tilde{z}} \underbrace{d}_{Z^T x} \|_2^2 \end{aligned}$$

$$= \|g^T Z Z^T \Delta x - g^T b\|_2^2$$

*w is first component of*

$$= \|T_w w - b\|_2^2 + \|d\|_2^2$$

$$= \| q^T A Z Z^T x - q^T b \|_2^2$$

$\underbrace{q^T A Z Z^T x}_\text{Term 1} \quad \underbrace{q^T b}_\text{Term 2} \quad \xrightarrow{\text{This}} \quad \text{Term 1}$

$$\| Ax - b \|_2^2 = \| T_{11} w - c \|_2^2 + \| d \|_2^2 \quad q^T b = \begin{pmatrix} c \\ d \end{pmatrix}$$

$$\|Ax - b\|_2^2 = \|T_{11}w - c\|_2^2 + \|d\|_2^2$$

$$\boxed{\text{Intermediate Step}} \Rightarrow \|q^T A Z Z^T x - q^T b\|_2^2 = \|T_{II} w - c + d\|_2^2$$

$$c = \begin{pmatrix} 1 \\ 2 \\ 6 \\ 0 \\ 0 \end{pmatrix} \left\{ \text{first 5 components of } g^T b. \right.$$

$$d = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ x \end{pmatrix} \Big\} \gamma \quad \text{last } n-r \text{ components of } g^T b,$$

If  $x$  is to minimize  $\|Ax - b\|^2$ , then

$$w \text{ must minimize } \|T_1 w - c\|_2^2$$

$$11T_{11}w_2 - c11_2^2 = 0.$$

$$w = T_{11}^{-1} c \rightarrow \text{works.}$$

We can choose  $y$  to be zero so that  $Z^T x = \begin{pmatrix} T_{11}^{-1} c \\ 0 \end{pmatrix}_{n-r}$

$$so \quad x_{ls} = Z \begin{pmatrix} T_{11}^{-1} c \\ 0 \end{pmatrix}$$

★

Rank-deficient LSP's

Last time: If  $A$  has a "complete orthogonal factorisation"

$$A = Q \begin{bmatrix} T_{11} & 0 \\ 0 & 0 \end{bmatrix} Z^T \quad \text{i.e. } Q^T A Z = \begin{bmatrix} T_{11} & 0 \\ 0 & 0 \end{bmatrix}$$

then the LS solution

$$x_{LS} = Z \begin{pmatrix} T_{11}^{-1} c \\ 0 \end{pmatrix} \quad \text{rank } A = r, \quad \text{and } n-r$$

$$g^T b = \begin{pmatrix} c \\ d \end{pmatrix} \quad \text{rank } A = r, \quad \text{and } n-r$$

Theorem: Suppose  $A = U \Sigma V^*$  is the SVD &  $A \in \mathbb{C}^{m \times n}$ ,  $r = \text{rank } A$ .

If  $U = [u_1 | \dots | u_m]$ ,  $V = [v_1 | \dots | v_n]$  &  $b \in \mathbb{C}^m$ ,

then the  $j^{\text{th}}$  component of  $x_{LS}$  is given by,

$$(x_{LS})_j = \sum_{i=1}^r \frac{u_i^T b}{\sigma_i} v_{ji}.$$

$x_{LS}$  minimizes  $\|Ax - b\|_2$ .

This  $x_{LS}$  has the smallest  $\ell^2$ -norm among all minimizers.

Proof :

This is the re-wording of the earlier proof in the case of SVD.

$$\|Ax - b\|_2^2 = \|U \Sigma V^T x - b\|_2^2$$

$$\begin{aligned} \|Ax - b\|_2^2 &= \|U^T Ax - U^T b\|_2^2 = \|U^T A V V^T x - U^T b\|_2^2 \\ &= \|\sum w_i - U^T b\|_2^2 \\ &= \left( \sum_{i=1}^r \sigma_i w_i - u_i^T b \right)^2 \end{aligned}$$

Since we have no control over

$$\sum_{i=1}^m (\sigma_i w_i - v_i^T b)^2 = \sum_{i=1}^r (\sigma_i v_i - v_i^T b)^2 + \sum_{i=r+1}^m (v_i^T b)^2$$

Since we have no control over  $\sum_{i=r+1}^m (v_i^T b)^2$

To minimize  $\|Ax-b\|_2^2$ , we minimize  $\sum_{i=1}^r (\sigma_i w_i - v_i^T b)^2$

Equating  $\sum_{i=1}^r (\sigma_i w_i - v_i^T b)^2 = 0$  gives each  $\sigma_i w_i = v_i^T b$  i.e.

$$w_i = \frac{v_i^T b}{\sigma_i}$$

$$\text{let } w_i = \begin{cases} \frac{v_i^T b}{\sigma_i} & \text{for } 1 \leq i \leq r \\ 0 & \text{for } r+1 \leq i \leq m. \end{cases}$$

$$\text{So } w = \begin{pmatrix} \frac{v_1^T b}{\sigma_1} \\ \vdots \\ \frac{v_r^T b}{\sigma_r} \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

$$\therefore x_{LS} = Vw = V \begin{pmatrix} \frac{v_1^T b}{\sigma_1} \\ \vdots \\ \frac{v_r^T b}{\sigma_r} \\ 0 \\ \vdots \\ 0 \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^r \frac{v_i^T b}{\sigma_i} v_{ii} \\ \sum_{i=1}^r \frac{v_i^T b}{\sigma_i} v_{ri} \\ \vdots \\ \sum_{i=1}^r \frac{v_i^T b}{\sigma_i} v_{mi} \end{pmatrix}$$

(check) As in the case of full rank LSPs, this sol<sup>n</sup> can also be expressed as

$$x_{LS} = A^+ b.$$

where  $A^+ = V \Sigma^+ U^T$

$$\Sigma^+ = \begin{pmatrix} \sqrt{\lambda_1} & & & \\ & \sqrt{\lambda_2} & & \\ & & \ddots & \\ & & & 0 \end{pmatrix}$$

II

### QR with column pivoting

Usual  $A = QR$

Why is column pivoting required?

$$A = [a_1 | a_2 | a_3 | a_4]_{m \times 4} = [q_1 | q_2 | q_3 | q_4]_{m \times 4} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \end{bmatrix}_{4 \times 4}$$

rank  $A = 3$ .

$$[q_1 | q_1 | q_1 + q_2 | q_1 + q_2 + q_3 + q_4]$$

$\uparrow \quad \uparrow \quad \uparrow \quad \uparrow$   
 $a_1 \quad a_2 \quad a_3 \quad a_4$

range (A)

observe that range(CS) is not in the span of any 3 columns of Q.

range { $q_1, q_2, q_3$ }

The QR factorisation gives no info about rank A  
 $(\therefore$  about null(A)).

$$A\pi = QR$$

↳ permutation matrix so that we have linearly independent columns to the left.

The algorithm produces the factorisation

$$Q^T A \pi = \begin{bmatrix} R_{11} & | & R_{12} \\ 0 & | & R_{22} \end{bmatrix} = \begin{bmatrix} R_{11} & | & R_{12} \\ 0 & | & 0 \end{bmatrix} \}^{n-r}$$

permutation.

$R$  is rank  $A$

$Q$  is orthogonal

$R_{11}$  is upper triangular and non-singular.

Denote  $A\pi = [a_{e_1} | \dots | a_{e_n}]$

$$\leftarrow p = [q_1 | \dots | q_n]$$

then  $a_{ek} = \sum_{i=1}^{\min\{r, k\}} r_{ik} q_i \in \text{span}\{q_1, \dots, q_r\}$

Range of  $A = \text{span}\{q_1, \dots, q_r\}$ .

Steps of the algorithm

1) Compute the  $\|a_1\|_2, \dots, \|a_n\|_2$

let  $p$  be the smallest index for which

$$\|a_p\|_2 = \max_{1 \leq i \leq n} \{\|a_i\|_2\}$$

2) let  $a_{e_1} = a_p$ , i.e. apply permutation matrix  $P$ , to exchange columns 1 &  $p$  of  $A$ .

$$Q_1 \begin{pmatrix} A P_1 \\ [a_1 \dots a_n] \end{pmatrix} = \begin{bmatrix} x \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad | \quad au$$

$\nwarrow ap$

3) Apply a suitable Householder matrix to annihilate subdiagonal entries of column 1.

$$A \rightarrow AP_1 \rightarrow Q_1 AP_1 = x_1$$

Now repeat these steps for the  $n \times (n-1)$  lower submatrix of  $A$ .

$$Q_2 \left( (Q_1 A P_1) P_2 \right)$$

$$Q_1 \dots Q_n A P_1 P_2 \dots P_{n-1} = R.$$

$Q^\top A P = R$

★

$\boxed{QRCP}$

$$\begin{array}{c}
 \left[ \begin{array}{|c|c|c|c|} \hline a_1 & \dots & a_p & \dots & a_n \\ \hline \end{array} \right] \xrightarrow{AP_1} \left[ \begin{array}{|c|c|c|c|c|c|} \hline a_{p1} & a_{21} & \dots & a_{11} & \dots & a_{n1} \\ \hline \end{array} \right] \\
 \downarrow \quad \text{first column with max 2 norm} \qquad \qquad \qquad \text{Householder matrix.} \\
 Q_1 AP_1
 \end{array}$$

exchange

$$\left[ \begin{array}{|c|c|c|c|} \hline \|a_{p1}\| & & & \\ \hline 0 & \|a_{p2}\|_2 & \dots & \\ \hline : & 0 & \dots & \\ \hline 0 & 0 & \dots & \\ \hline \end{array} \right] \xleftarrow{Q_2(Q_1 AP_1)P_2} \left[ \begin{array}{|c|c|c|c|c|c|} \hline \|a_{p1}\| & & & & & \\ \hline 0 & & & & & \\ \hline : & & & & & \\ \hline 0 & & & & & \\ \hline \end{array} \right] \qquad \text{Repeat} \\
 \text{for } (m-1) \times (n-1) \text{ submatrix}$$

$Q = \dots Q_1 AP_1 \dots P_r$

$$\left[ \begin{array}{|c|c|} \hline \|a_{p1}\| & \alpha \\ \hline 0 & \|a_{p2}\|_2 \\ \hline \vdots & \vdots \\ \hline 0 & \|a_{pr}\|_1 \\ \hline \hline 0 & R_{22} \\ \hline \end{array} \right] = R$$

$R_{12}^*$

where  $\|a_{p1}\| \geq \|a_{p2}\| \geq \dots \geq \|a_{pr}\|$

Remark Theoretically  $R_{22} = 0$  (since  $\gamma_k \alpha = \gamma_k R = \gamma$ )

But in practice it rarely happens (unless some exceptional cancellation takes place) In any case  $\|R_{22}\|$  is small.

— / —

$x_{LS} = A^T b$   
 LSPs      full rank      rank deficient  $\Rightarrow$ 

- using SVD
- using formulation
- using QR

 $\min \|x\|_2^2 = \|Ax - b\|_2^2$

if  $A$  has a complete ortho fact  
 $Q^T A Z = \begin{bmatrix} T_{11} & 0 \\ 0 & 0 \end{bmatrix}$

(minimal)  $x_{LS} = Q^T \begin{bmatrix} T_{11}^{-1} c \\ 0 \end{bmatrix}$

methods

- using SVD
- using QRCP (today)

 $Q^T b = \begin{bmatrix} c \\ d \end{bmatrix} \} r$

$A: \mathbb{R}^n \rightarrow \mathbb{R}^m$

\* Using QRCP for rank-deficient LSP (<sup>suppose</sup>  $A = QR$ )

aver~~er~~ aver~~er~~ aver~~er~~  $A \in \mathbb{R}^{m \times n}$ .

$$\|Ax - b\|_2^2 = \|Q^T A \begin{bmatrix} y \\ z \end{bmatrix} - \begin{bmatrix} c \\ d \end{bmatrix}\|_2^2$$

$$= \|R_{11}y - (c - R_{12}z)\|_2^2 + \|d\|_2^2$$

where

$$T^T z = \begin{bmatrix} y \\ z \end{bmatrix} \} r$$

$$Q^T b = \underbrace{\begin{bmatrix} c \\ d \end{bmatrix}}_{\{r\}} \quad Q^T b = \begin{bmatrix} c \\ d \end{bmatrix} \} r$$

(basic sol<sup>n</sup>  $x_B$ ).

$$x_{LS} = T \begin{bmatrix} R_{11}^{-1}(c - R_{12}z) \\ z \end{bmatrix}$$

$\rightarrow$  if  $x_{LS}$  with minimal norm is defined then  $z$  should be set to 0.

$$\therefore \mathbf{x}_B = \mathbf{T} \left[ \begin{matrix} \mathbf{R}_{11}^{-1} \mathbf{c} \\ \mathbf{0} \end{matrix} \right] \}_{r \times n}$$

### \* Numerical rank

In theory, if  $\text{rank } A = r$  &  $A = U\Sigma V^*$

then  $\Sigma = \begin{bmatrix} \sigma_1 & & & \\ & \ddots & & \\ & & \sigma_r & \\ & & & 0 \end{bmatrix}$

### \* Eigenvalue / Eigenvector methods

For degree  $> 5$  there is no closed form solution

use iterative methods instead.

The theoretical method of finding roots of the characteristic polynomial doesn't have a closed form solution (or discriminant like solution) for polynomials of deg  $> 5$  (Abel's theorem).

Hence, the need for iterative methods.

### E-value / E-vector methods

Direct  
method

Indirect  
method.

## Direct methods

- Power iteration
- Block power iteration (simultaneous iteration)
- QR iteration
- Rayleigh quotient iteration.

## Iterative methods

- Bisection method
- Divide & Conquer
- Jacobi's method
- Krylov subspace method
- Lanczos iterations.

### Power iteration

Idea  $\{b, Ab, A^2b, \dots\} \rightarrow$  largest e-vector of  $A$   
initial vector.

choose initial vector  $x_0$

i=0 to convergence

$$\begin{cases} y_{i+1} = Ax_i \\ x_{i+1} = y_{i+1} / \|y_{i+1}\|_2 \quad (\text{approx e-vector}) \\ \lambda_{x_{i+1}} = x_{i+1}^T A x_{i+1} \quad (\text{approx -e-value}). \end{cases}$$

Consider the simplest case  $A = \text{diag } (\lambda_1 \dots \lambda_n)$

where  $|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_n|$

16 10:30 - 1:00  
17 19:30 - 1:00.

16  
17  
18  
21

— / —

Note that in this case the e-vectors are  $e_i$ ;  
i.e. columns of the identity matrix

$$\text{let } x_i = \frac{A^i x_0}{\|A^i x_0\|}$$

$$\text{let } x_0 = \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{pmatrix}, \quad e_i \neq 0.$$

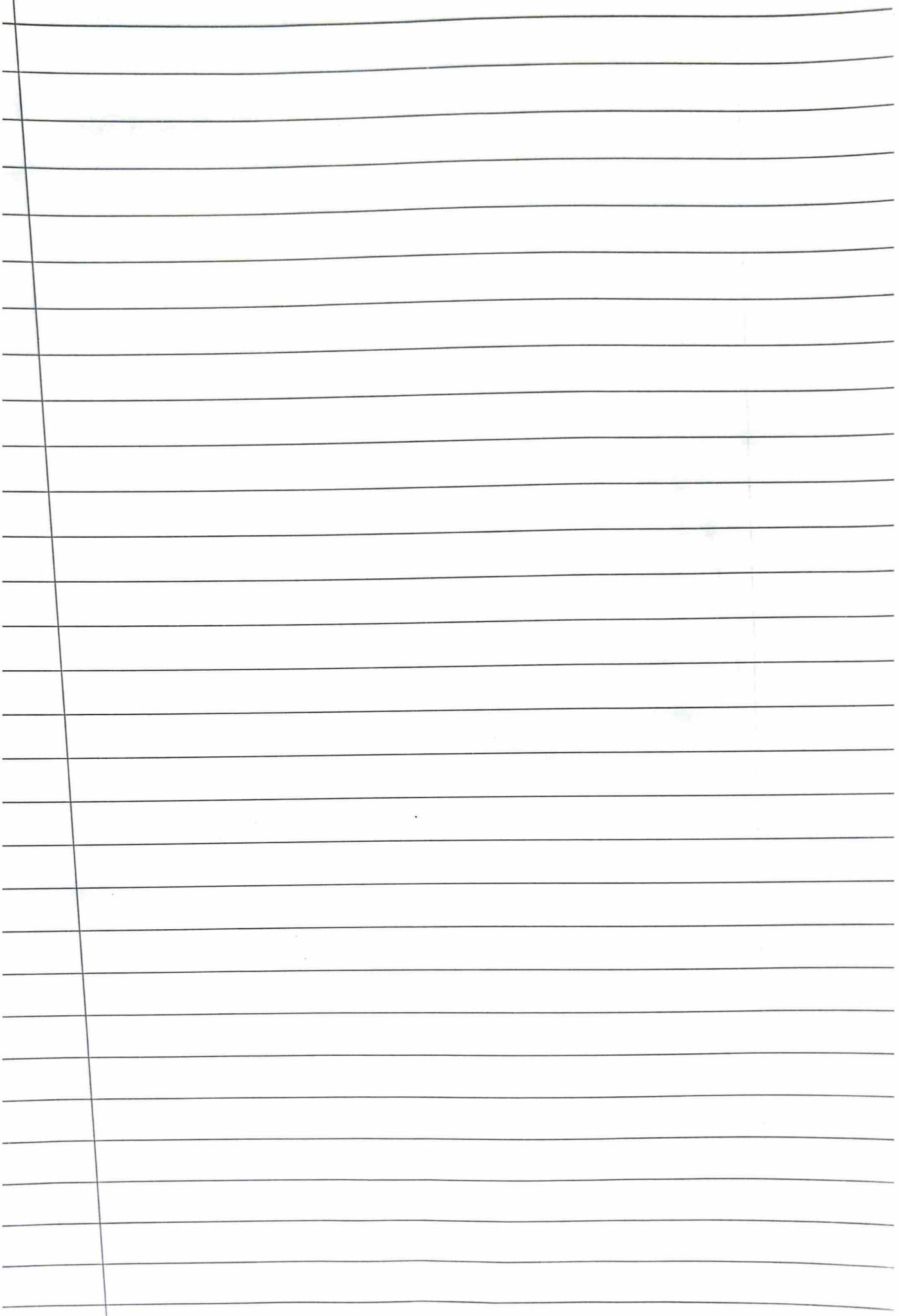
$$\begin{aligned} \text{Then } A^i x_0 &= A^i \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{pmatrix} = \begin{pmatrix} e_1 x_1 \\ e_2 x_2 \\ \vdots \\ e_n x_n \end{pmatrix} \\ &= e_1 \lambda_1^i \begin{pmatrix} \frac{e_1}{\lambda_1} \left( \frac{\lambda_1^i}{\lambda_1} \right) \\ \vdots \\ \frac{e_n}{\lambda_n} \left( \frac{\lambda_n^i}{\lambda_n} \right) \end{pmatrix} \end{aligned}$$

Since each  $|\lambda_i| \leq |\lambda_1|$ ,

$$A^i x_0 \mapsto e_1 \lambda_1^i e_1$$

$\therefore x_i = \frac{A^i x_0}{\|A^i x_0\|} \mapsto \pm e_1 = \text{e-vector corresponding}$   
 $\text{to the largest e-value}$   
 $\lambda_1$ .

— / / —



## \* Analysis of Orthogonal Iteration

1) Note that

$\Delta$  is diagonalizable

$$\text{span}(\mathbf{Q}_k) = \text{span}(\mathbf{z}_k)$$

( $\Delta$  is diagonalizable)

$$= \text{span}(\Delta \mathbf{Q}_{k-1})$$

$$\Rightarrow \Delta = S \Lambda S^{-1}$$

$$= \text{span}(\Delta \mathbf{Q}_{k-2})$$

matrix of  
e-value of  $\Delta$

$$= \text{span}(\Delta^k \mathbf{z}_0)$$

$$= \text{span}(S \Delta^k S^{-1} \mathbf{z}_0).$$

2) By assumption  $|\lambda_1| > \dots > |\lambda_n| > |\lambda_{n+1}| > \dots > |\lambda_m|$

$$\left| \frac{\lambda_i}{\lambda_n} \right| > 1 \text{ for } i < n.$$

$$\left| \frac{\lambda_i}{\lambda_n} \right| < 1 \text{ for } i > n.$$

$$\text{so } \underbrace{\Delta^k S^{-1} \mathbf{z}_0}_{\sim} = \lambda_n^k \begin{bmatrix} [V_k] \\ \vdots \\ [W_k] \end{bmatrix}_{n \times n} \quad \text{where } [W_k] \rightarrow 0 \text{ as } k \rightarrow \infty$$

$$\& [V_k] \rightarrow 0.$$

$\mathbf{z}_0$  an orthogonal matrix

Set  $\mathbf{Q}_0 = \mathbf{z}_0$   $n \times n$

for  $i = 0$  to convergence

$$\mathbf{z}_k = \Delta \mathbf{Q}_{k-1}$$

$$\mathbf{Q}_k \mathbf{R}_k = \mathbf{z}_k$$

$$\lambda_k (\mathbf{Q}_k^* \Delta \mathbf{Q}_k) = \{ \lambda_1^{(k)}, \dots, \lambda_n^{(k)} \}.$$

$$S = \left[ \begin{array}{c|c|c|c} s_1 & \dots & s_n & s_{n+1} & \dots & s_m \end{array} \right]$$

$s_n$        $s_{m-n}$

$$\therefore S \Lambda^k S^{-1} z_0 = \lambda_n^k S \begin{bmatrix} v_k \\ w_k \end{bmatrix}$$

$$S \Lambda^k S^{-1} z_0 = \lambda_n^k [s_n v_k + s_{m-n} w_k]$$

3) Finally,  $\text{span } Q_k = \text{span}(s_n v_k + s_{m-n} w_k)$   
 $\rightarrow \text{span}(s_n v_k)$  as  $k \rightarrow \infty$

= span of first  $n$  e-vectors of  $A$ .

### \* QR iteration

As before, we assume  $|\lambda_1| > \dots > |\lambda_n|$

let

$$z_0 = \left[ \begin{array}{c|c|c} x_1^{(0)} & \dots & x_m^{(0)} \end{array} \right] \quad (\text{initial matrix})$$

orthogonal columns.

Assume that all principal submatrices of  $z_0$  are non-singular.

### Algorithm

$$A \in \mathbb{C}^{m \times m}$$

$$T_0 = Z_0 * A Z_0$$

$K = 1$  to convergence

$$T_{K-1} = Q_K R_K \quad (\text{QR factorization})$$

$$T_K = R_K Q_K$$

end,

$$T_0 = Z_0^* A Z_0 = Q_1 R_1$$

$$T_1 = R_1 Q_1 = Q_2 R_2$$

$$T_2 = R_2 Q_2 = Q_3 R_3.$$

On the one hand,

$$T_{k-1} = Q_{k-1}^* A Q_{k-1}$$

$$Q_0^* A Q_0 = Q_1 R_1$$

### e-value methods

Unsymmetric  
matrices

Symmetric  
matrices

↳ faster iterative methods.

— / —

\*  $\begin{pmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{pmatrix}$   $\xrightarrow{\text{reduction}}_{\text{to a simpler form}} \begin{pmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{pmatrix}$  Hessenberg matrix  $\xrightarrow{\text{e-value method}}$  to obtain e-value & e-vector.

has to be done through similarity transformation.

$$A \neq A^* \quad \begin{pmatrix} x & \dots & \dots & * \\ x & & & \\ x & & & \\ x & & & \\ x & \dots & \dots & -x \end{pmatrix} \xrightarrow{\text{Phase I}} \begin{pmatrix} x & \dots & \dots & x \\ x & \dots & \dots & \\ x & \dots & \dots & \\ 0 & \dots & x & x \end{pmatrix}$$

Phase II

$$\begin{pmatrix} x & \dots & \dots & x \\ & \ddots & & \\ & & \ddots & \\ & & & x \end{pmatrix}$$

$$A = A^* \quad \begin{pmatrix} x & \dots & \dots & x \\ \vdots & & & \\ \vdots & & & \\ x & \dots & \dots & -x \end{pmatrix} \xrightarrow{\text{Phase I}} \begin{pmatrix} x & x & & \\ x & \dots & & \\ 0 & \dots & x & x \\ 0 & \dots & x & x \end{pmatrix} \xrightarrow{\text{Hessenberg matrix (triangular)}} \begin{pmatrix} x & & & \\ & \ddots & & \\ & & \ddots & \\ & & & x \end{pmatrix}$$

Phase II

$$\begin{pmatrix} x & & & \\ & \ddots & & \\ & & \ddots & \\ & & & x \end{pmatrix}$$

\* Conversion to Hessenberg form

Done using Householder matrices.

If we apply it as before

Suppose  $Q_1^* A$  is a Householder matrix such that

$$Q_1^* A = \begin{pmatrix} x & \cdots & x \\ 0 & & \vdots \\ \vdots & & x \\ 0 & \cdots & x \end{pmatrix}$$

To make this a similarity transformation we need right multiplication with  $Q_1$

$$(Q_1^* A) Q_1$$

But this step destroys the zeroes obtained earlier.

Better way

- ① At the first step, choose a householder matrix that leaves the first row of  $A$  unchanged.
  - when  $Q_1^* A$  is computed, it forms linear combinations of rows 2, 3, ..., n (excludes row 1) & introduces zeroes in (3,1), (4,1), ..., (m,1)
  - when  $Q_1$  is multiplied on the right of  $Q_1^* A$ , then it replaces columns of 2, 3, ..., n by their own linear combinations leaving the 1st column unchanged

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix} = \begin{bmatrix} 5+14 \\ 19+16 \end{bmatrix}$$

— / —

$Q_1^* A \leftarrow$

$$\begin{bmatrix} x & \dots & x \\ \vdots & \ddots & | \\ 1 & \dots & 1 \\ x & \dots & x \end{bmatrix} \xrightarrow{Q_1^* A} \begin{bmatrix} x & x & \dots & x \\ x & \vdots & \ddots & | \\ 0 & \vdots & \ddots & | \\ \vdots & \ddots & \ddots & | \\ 0 & \dots & \dots & x \end{bmatrix} \xrightarrow{Q_1^* A Q_1} \begin{bmatrix} x & x & \dots & x \\ x & x & \dots & | \\ 0 & x & \dots & | \\ \vdots & \vdots & \ddots & | \\ 0 & \dots & \dots & x \end{bmatrix}$$

$$\underbrace{Q_{m-2}^* \cdots Q_2^*}_{Q^*} \underbrace{Q_1^* A}_{Q} \underbrace{Q_1 Q_2 \cdots Q_{m-2}}_{Q^*} = Q^* A Q \rightarrow \text{Hessenberg matrix.}$$

$$J - \frac{2vv^*}{v^* v}$$

Algorithm  $A \in \mathbb{R}^{m \times m}$

for  $k=1$  to  $m-2$

$$x = A_{k+1:m, k:m}$$

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

$$v_k = v_k / \|v_k\|$$

$$A_{k+1:m, k:m} = A_{k+1:m, k:m} - 2v_k \left( v_k^* A_{k+1:m, k:m} \right)$$

$$A_{1:m, k+1:m} = A_{1:m, k+1:m} - 2v_k$$

$$A_{1:m, k+1:m} = A_{1:m, k+1:m} - 2(A_{1:m, k+1:m} v_k) v_k^*$$

end.

\* From orthogonal to QR iteration

The QR iteration is a rearrangement of the orthogonal iteration designed to form a sequence of matrices that converges to the Schur form of A.

(Digression)

Recall Schur form of A (decomposition)

Theorem for every square  $n \times n$  matrix A with entries in  $\mathbb{C}$ ,  $\exists$  unitary matrix Q & upper triangular matrix T such that  $A = QTQ^*$

- Observations : ① A & T are similar, diagonal entries of T are e-values of A.  
 ② Columns of Q are the e-vectors of A.

Suppose  $T_k = Q_k^* A Q_k$

i.e.  $\{T_0, T_1, \dots, T_k, \dots\} \rightarrow T = Q^* A Q$

core steps of

Recall the orthogonal iteration

$$Q_0 = Z_0$$

$$A Q_0 = Z_1$$

"  
SIR<sub>1</sub>

$$A Q_1 = Z_2$$

"  
SIR<sub>2</sub>

$$R_k = Q_k^* A Q_{k-1}$$

Let's denote by  $T_k$ :

$$T_{k-1} = Q_{k-1}^* A Q_{k-1} \quad (\text{by def'})$$

$$= Q_{k-1}^* Z_k$$

$$= (Q_{k-1}^* Q_k) R_k. = \hat{Q}_k^* R_k$$

$$T_k = Q_k^* A Q_k = (Q_k^* A Q_{k-1})(Q_{k-1}^* Q_k) = R_k \hat{Q}_k$$

$\therefore T_k$  can be obtained from  $T_{k-1}$  by computing QR factorization of  $T_{k-1}$  & reversing the order of  $Q$  &  $R$ .

QR iteration:  $A \in \mathbb{C}^{n \times n}$

$Z_0$  unitary in  $\mathbb{C}^{n \times n}$

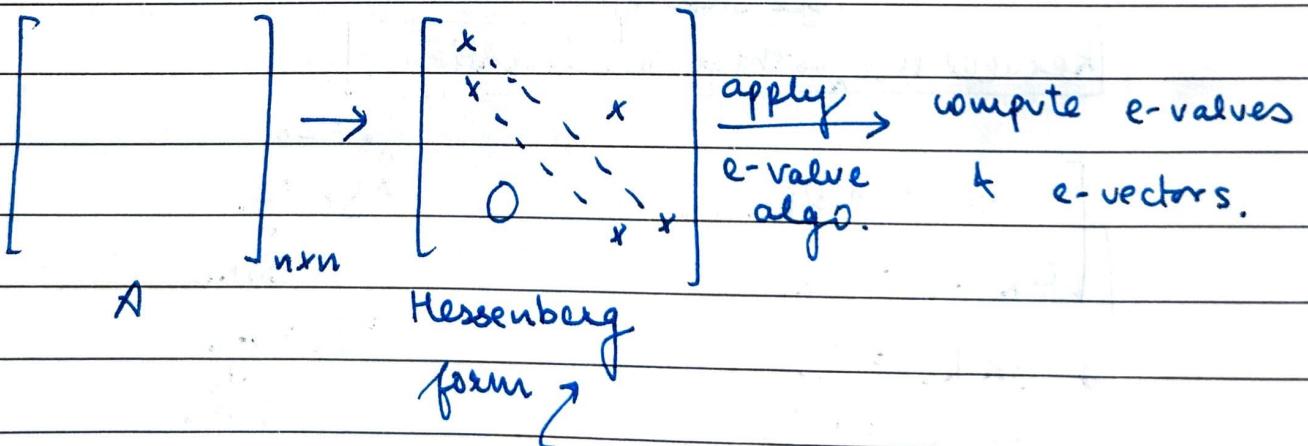
$$T_0 = Z_0^* A Z_0 \quad (\text{principal submatrices of } Z_0 \text{ are invertible}).$$

$$T_{k-1} = Q_{k-1} R_{k-1}$$

$$T_k = R_k Q_k.$$

end.

Last time Any e-value computation is carried out in 2 stages.



Assymmetric (previous methods)

In the symmetric case the Hessenberg form is tridiagonal,

Symmetric case

$$\begin{bmatrix} x & & & \\ x & \ddots & & 0 \\ & & \ddots & \\ 0 & & & x \\ & & & x \end{bmatrix}$$

E-value / E-vector methods for symmetric matrices:

- 1) Bisection
- 2) Rayleigh quotient
- 3) Jacobi's
- 4) Divide & conquer.

\* Stability of Hessenberg matrices

Let  $\tilde{H}$

Let  $\tilde{H}$  be the computed Hessenberg matrix &  $Q$  be the exact unitary matrix corresponding to the  $\tilde{V}_k$ .

Then  $Q\tilde{H}Q^* = A + SA$ , where  $\|SA\| = O(\text{Emach})$ .

$\|A\|$

$\therefore$  The computation of  $\tilde{H}$  is backward stable.

## \* Computing the SVD (Given A)

Mathematically, ① form  $A^*A$

The procedure ② Compute e-value decomposition of  
is unstable.  $A^*A = V\Lambda V^*$

③ let  $\Sigma = \text{diagonal (non-neg. sq. roots of } \Lambda)$   
④ solve  $U\Sigma = AV$  for  $U$ .

## \* A stable method is as follows

(consider the  $2m \times 2n$  Hermitian matrix

If  $A = U\Sigma V^*$ , then  $AV = U\Sigma$   
 $A^*U = V\Sigma$ .

$$\begin{bmatrix} 0 & A^* \\ A & 0 \end{bmatrix}$$

$$\begin{bmatrix} 0 & A^* \\ A & 0 \end{bmatrix} \begin{bmatrix} V & V \\ V & -V \end{bmatrix} = \begin{bmatrix} A^*U & -A^*U \\ AV & AV \end{bmatrix}$$

$$= \begin{bmatrix} V & V \\ V & -V \end{bmatrix} \begin{bmatrix} \Sigma & 0 \\ 0 & -\Sigma \end{bmatrix} \rightarrow \textcircled{A}$$

① is an e-value decomposition of  $H$ .

singular values of  $A$  are the abs. values of the e-values of  $H$ .

For non-sq. matrices a 2-phase method is used.

1st phase  
method - bidiagonalization.

2nd phase : QR method / divide & conquer.

## Golub-Kahan bidiagonalization

$$\begin{array}{c}
 \left[ \begin{array}{cccc} x & x & \cdots & x \\ \vdots & x & \cdots & x \\ 1 & & & \\ 1 & & & \\ x & \cdots & & x \end{array} \right] \xrightarrow[V_1]{\text{(left)}} \left[ \begin{array}{cccc} x & & & x \\ 0 & x & \cdots & x \\ \vdots & & & \\ 1 & & & \\ 0 & x & \cdots & x \end{array} \right] \xrightarrow[V_1]{\text{(right)}} \left[ \begin{array}{cccc} x & x & \cdots & 0 & \cdots & 0 \\ 0 & x & \cdots & x & & \\ \vdots & & & & & \\ 0 & x & \cdots & + & & \end{array} \right]
 \end{array}$$

$$V_2^* (V_1^* \Lambda V_1 V_2) = \text{bidiagonal}$$

math 2

each  $U_i V_i$  is a Householder matrix.

## Algorithm

- ① form H
  - ② compute its. G-K bidiagonalization.
  - ③ Apply e-value methods.

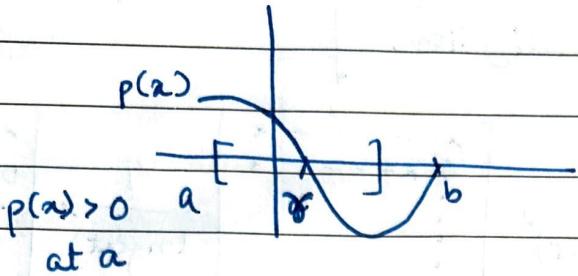
## E-value and E-vector algorithm for symmetric matrices

I] Bisection method - to locate roots of characteristic polynomial of a symmetric matrix A.

If  $p(x)$  has a root at  $r \in [a, b]$  then check if  $p(a) \neq p(b)$  is zero. If yes, then root is found.

If not, let  $c = \frac{at+b}{2}$  & perform similar analysis for

$[a, c] \wedge [c, b]$ . ... continue.



$\Delta p(x) < b$  at  $b$ .

This method is more effective in the case of symmetric matrices because e-values satisfy some nice properties in this case.

let  $A \in \mathbb{R}^{m \times m}$  be symmetric

Then  $A$  can be reduced to tridiagonal form -

$$\begin{pmatrix} a_1 & b_1 & \cdots & 0 \\ b_1 & a_2 & b_2 & \vdots \\ 0 & b_2 & a_3 & \ddots & b_{m-1} \\ & \ddots & \ddots & \ddots & a_m \end{pmatrix} \quad \begin{array}{l} \text{(suppose} \\ b_i \neq 0 \forall i \end{array}$$

① let  $A^{(1)}, \dots, A^{(m)}$  be the principal submatrices of the above tridiag-matrix.

Define  $p_k(x) = \det(A^{(k)} - xI_{k \times k}) = \text{char poly of } A^{(k)}$

These characteristic polynomial satisfy,

$$p_k(x) = (a_k - x) p_{k-1}(x) - b_{k-1}^2 p_{k-2}(x)$$

② The e-values of  $A$  are distinct ( $\lambda_1^{(k)} < \lambda_2^{(k)} < \dots < \lambda_k^{(k)}$ )

(Sturm sequence property) (with the earlier assumptions)

The e-values of  $A^{(k)}$  strictly interlace with e-values

If  $x^{(k+1)}$  is :

$$\lambda_j^{(k+1)} < \lambda_j^{(k)} < \lambda_{j+1}^{(k+1)}$$

$$\begin{array}{ccccccc}
 A^{(k-1)} & \lambda_1^{(k-1)} & < & \lambda_2^{(k-1)} & < & \lambda_3^{(k-1)} & \dots \\
 & \swarrow & & \searrow & & & \\
 A^{(k)} & \lambda_1^{(k)} & < & \lambda_2^{(k)} & < & \lambda_3^{(k)} & \\
 & | & & & & & \\
 A^{(k+1)} & \lambda_1^{(k+1)} & < & \lambda_2^{(k+1)} & < & & \dots
 \end{array}$$

If  $a(\lambda) = \# \text{ of sign changes in the sequence}$   
 $p_0(\lambda), p_1(\lambda), \dots, p_m(\lambda)$  (Sturm sequence)

then  $a(\lambda) = \# \text{ of } \epsilon\text{-values of } A \text{ that are less than } \lambda$ .

Using  $a(\lambda)$  one can find the number of  $\epsilon$ -values of  $A$  in any interval.

$a(0) = \# \text{ of sign changes in the seq.}$

$p_0(0), p_1(0), p_2(0), \dots, p_m(0)$ .

$= \# \text{ of -ve } \epsilon\text{-values of } A$ .

To find ~~# of  $\epsilon$ -values in  $(p, q]$~~   $= a(q) - a(p)$

To find # of  $\epsilon$ -values in  $(p, q] = a(q) - a(p)$ .

Ex find # of  $\epsilon$ -values of  $A$

$$A = \begin{pmatrix} 1 & 1 & 0 & 0 \\ 1 & -2 & 1 & 0 \\ 0 & 1 & 3 & 1 \\ 0 & 0 & 1 & -4 \end{pmatrix} \quad \text{in the interval } [-1, 4].$$

$$\left( \# \text{ of e-values} \right) \text{ in } (-\infty, q) - \left( \# \text{ of e-values} \right) \text{ in } (-\infty, p)$$

}

$$= \# \text{ of -ve e-values} \\ \text{of } A + bI.$$

$$\begin{pmatrix} a_1 & b_1 & & & \\ b_1 & a_2 & & & \\ & b_2 & \ddots & & \\ & & \ddots & b_i & \\ & & & b_i & \ddots \\ & & & & \ddots & a_m \end{pmatrix}$$

If same  $b_i = 0$ ,  
then "divide & conquer"  
(as a strategy).

\*

### Divide and conquer method

Let  $T \in \mathbb{R}^{n \times n}$  be a symmetric tridiagonal matrix with non-zero off diagonal entries.

Let  $n \geq m_2$ , split  $T$  into submatrices

$$[T_1]_{n \times n} \& [T_2]_{(m-n) \times (m-n)}$$

→ subtract  $\beta$ .

$$T = \begin{bmatrix} T_1 & & \\ & \beta & \\ & & T_2 \end{bmatrix} = \begin{bmatrix} \hat{T}_1 & & \\ & \hat{T}_2 & \\ & & \beta \end{bmatrix} + \begin{bmatrix} & & \\ & \beta & \\ & & \beta \end{bmatrix}$$

$\hat{T}_1$  &  $\hat{T}_2$  are  
also tridiagonal  
symmetric.

Rank-1  
collection.

If we have diagonalizations

$$\hat{T}_1 = Q_1 D_1 Q_1^T$$

A

$$\hat{T}_2 = Q_2 D_2 Q_2^T$$

then,

$$T = \begin{bmatrix} \hat{T}_1 & \\ & \hat{T}_2 \end{bmatrix} + \begin{bmatrix} \beta & \beta \\ \beta & \beta \end{bmatrix}$$

$$= \begin{bmatrix} Q_1 & \\ & Q_2 \end{bmatrix} \left( \begin{bmatrix} D_1 & \\ & D_2 \end{bmatrix} + \beta Z Z^T \right) \begin{bmatrix} Q_1^T & \\ & Q_2^T \end{bmatrix}$$

$$\text{where } Z^T = (q_1^T, q_2^T)$$

↓      ↗  
last row      first row  
of  $Q_1$       of  $Q_2$ .

\* What are the e-values of  $D + w w^T$ ? ( $D$  diagonal,  $w \neq 0$ ).  
If  $\lambda$  is an e-value of  $D + w w^T$ , with e-vector  $q$  ( $\neq 0$ )  
then  $(D + w w^T) q = \lambda q$

$$\text{i.e. } (D - \lambda I) q + w w^T q = 0.$$

$$q + (D - \lambda I)^{-1} w w^T q = 0.$$

$$w^T q + w^T (D - \lambda I)^{-1} w w^T q = 0. \quad (w^T q \Rightarrow \text{scalar})$$

i.e.

$$w^T q (1 + w^T (D - \lambda I)^{-1} w) = 0$$

$$\underbrace{\quad}_{= f(\lambda)}$$

(scalar function related to this system)

This is true if

$$f(\lambda) = 0.$$

$$1 + \sum_{i=1}^n \frac{w_i^2}{d_i - \lambda} = 0.$$

Roots of  $f(\lambda)$  are eigen values of  $D + w w^T$  (Answer).

### III] Rayleigh quotient iteration

$A$  is symmetric  $R_A : V \setminus \{0\} \rightarrow \mathbb{R}$

$$x \mapsto \frac{x^T A x}{x^T x}$$

$$x^T x.$$

If  $\lambda$  is an e-value of  $A$  with e-vector  $v$ , then

$$R_A(v) = \lambda.$$

The problem of estimating  $\lambda$  can be framed as a least squares problem as follows-

Given  
 $x \neq 0,$

find a scalar  $\lambda$  that minimizes

$$\|Ax - \lambda x\|_2.$$

The normal eqns for the system,

$$\therefore x^T x \lambda = x^T A x.$$

$$\begin{aligned} Ax &= b \\ \Rightarrow A^T A x &= A^T b \end{aligned}$$

It can be shown that the e-vectors of  $A$  are the stationary pts of  $R_A(x)$ .

As a consequence it can be shown that

$$R_A(x) - R_A(q) = O(\|x - q\|^2)$$

as  $x \rightarrow q$

$\therefore R_A(x)$  is a "quadratically accurate" estimate of  $R_A(q)$

★

## Algorithm

- 1)  $\Rightarrow$  Choose  $v^{(0)}$  with  $\|v^{(0)}\|_2 = 1$
- 2)  $\Rightarrow$   $\lambda^{(0)} = v^{(0)\top} A v^{(0)}$
- 3) for  $k=1$  to convergence
- 4) Solve  $(A - \lambda^{(k-1)} I) w = v^{(k-1)}$  for  $w$ .
- 5)  $v^{(k)} = \frac{w}{\|w\|_2}$
- 6)  $\lambda^{(k)} = R_A(v^{(k)})$
- 7) end.

## Idea

- ① Start with a vector  $v^{(0)}$
- ② Calculate  $R_A(v^{(0)})$  - first approximation.
- ③ Apply inverse iteration to  $\lambda^{(0)}$  to get  $v^{(1)}$ , then apply  $R_A$  to  $v^{(1)}$  & so on.

★

## Jacobi's method

Idea: Try to eliminate/reduce the magnitude of off-diagonal entries so that eventually they become small enough to be declared as zero.

This is done with a sequence of orthogonal similarity transformation  $A \mapsto Q^T A Q$ .

Define  $\text{off}(A) = \sqrt{\sum_{i=1}^n \sum_{j=1, j \neq i}^n |a_{ij}|^2}$  - Frobenius norm of off diagonal entries of  $A$ .

The transformations

$A \mapsto J_1^T A J_1 \rightarrow Q_2^T A_1 Q_2 \mapsto$  should result in a reduction

The matrix  $Q$  used for this is the Givens' rotation.

$$J(p, q, \theta) = \begin{pmatrix} p & q \\ -q & p \end{pmatrix} = \begin{pmatrix} c & s \\ -s & c \end{pmatrix}$$

The basic step of the procedure is as follows -

$$\begin{pmatrix} b_{pp} & b_{pq} \\ b_{qp} & b_{qq} \end{pmatrix} = \begin{pmatrix} c & s \\ -s & c \end{pmatrix} \begin{pmatrix} a_{pp} & a_{pq} \\ a_{qp} & a_{qq} \end{pmatrix} \begin{pmatrix} c & -s \\ s & c \end{pmatrix}$$

Since Frobenius norm is preserved by orthogonal transformation

$$a_{pp}^2 + a_{qq}^2 + 2a_{pq}^2 = b_{pp}^2 + b_{qq}^2 + 2b_{pq}^2 = b_{pp}^2 + b_{qq}^2$$

$$\begin{aligned} \text{Consider } \text{eff}(B)^2 &= \|B\|_F^2 - \sum_{i=1}^n b_{ii}^2 \\ &= \|A\|_F^2 - \sum_{i=1}^n b_{ii}^2 - b_{pp}^2 - b_{qq}^2 + (a_{pp}^2 + a_{qq}^2) \\ &= \|A\|_F^2 - \sum_{i=1}^n a_{ii}^2 + (a_{pp}^2 + a_{qq}^2 - b_{pp}^2 - b_{qq}^2) \\ &= \text{eff}(A)^2 - 2a_{pq}^2 \\ \text{eff}(B) &\leq \text{eff}(A). \end{aligned}$$

$$sc(a-b) + d(c^2 - s^2) = 0$$

$$\therefore \frac{d}{b-a} = \frac{sc}{c^2 - s^2}$$

$$\tan 2\theta.$$

$$\begin{pmatrix} c & -s \\ s & c \end{pmatrix} \begin{pmatrix} a & d \\ d & b \end{pmatrix} \begin{pmatrix} c & s \\ -s & c \end{pmatrix} = \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix}$$

$$\begin{pmatrix} c^2a + s^2b - 2sd & sc(a-b) + d(c^2 - s^2) \\ sc(a-b) + d(c^2 - s^2) & s^2a + c^2b + 2scd \end{pmatrix} = \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix}$$

### Algorithm

Given  $n \times n$  symmetric matrix  $A$  & integers  $p, q$ ,  $1 \leq p < q \leq n$ .  
 Compute  $(c, s)$  s.t. if  $B = J(p, q, \theta)^T A J(p, q, \theta)$   
 then  $b_{pq} = b_{qq} = 0$ .

if  $a_{pq} \neq 0$  then

$$t = (a_{qq} - a_{pp}) / 2a_{pq}$$

if  $t \geq 0$

$$t = \frac{1}{(t + \sqrt{1+t^2})}$$

else

$$t = \frac{1}{(\sqrt{t} - \sqrt{1+t^2})}$$

end

$$c = \frac{1}{\sqrt{1+t^2}}$$

$$s = tc.$$

\* Classical Jacobi algorithm  $\rightarrow$  parameter that user sets  
 conv :  $\text{rff}(V^T A V) \leq \text{tol.} (\|A\|_F)$  (convergence criterial).  
 $V = I_n$ ,  $s = \text{tol.} (\|A\|_F)$ .

while  $\text{rff}(A) > s$

choose  $(p, q)$  so that  $|a_{pq}| = \max_{i \neq j} |a_{ij}|$

$\rightarrow$  obtain  $(c, s)$  for this choice of  $p, q$

$$A = J(p, q, \theta)^T A J(p, q, \theta)$$

$$V = V J(p, q, \theta).$$

\* Iterative methods for e-value/e-vector computation

Krylov subspace methods

$$Ax = b.$$

$\{b, Ab, A^2b, \dots\}$  Krylov sequence.

$$K_1 = \langle b \rangle$$

$$K_2 = \langle b, Ab \rangle$$

$$K_p = \langle b, Ab, \dots, A^{p-1}b \rangle$$

Given a matrix  $A$  and a vector  $b$ , let

$$K = \left[ \begin{array}{c|c|c|c|c} b & Ab & A^2b & \dots & A^{n-1}b \end{array} \right]$$

then  $AK = \left[ \begin{array}{c|c|c|c} Ab & A^2b & \dots & A^n b \end{array} \right]$

If  $K$  is non-singular, let  $c = -K^{-1}A^n b$ .

then  $AK = K \left[ \begin{array}{c|c|c|c|c} e_2 & e_3 & \dots & e_n & -c \end{array} \right] = KC$

where

$$C = \left[ \begin{array}{ccccc} 0 & 0 & \cdots & -c_1 \\ 1 & 0 & \cdots & -c_2 \\ 0 & 1 & \cdots & -c_3 \\ 0 & \cdots & \cdots & -c_n \end{array} \right]$$

Note that  $C$  is a Hessenberg matrix

$C$  is the companion companion matrix of

$$p(x) = x^n + \sum_{i=1}^n c_i x^{i-1}$$

$$C = K^{-1} A K$$

Here,  $K$  may be ill-conditioned : A possible solution is to replace  $K$  with an orthogonal matrix  $Q$  s.t. for all leading  $k$  columns of  $K$  &  $Q$  span the same subspace

$$K = QR$$

Then,

$$C = K^{-1} A K = R^{-1} Q^T A Q R = R^{-1} H R.$$

$$H = R C R^{-1}$$

\  
upper triangular.  $\Rightarrow H$  is upper Hessenberg.

To compute the columns of  $Q$ .

$$\text{Let } Q = [q_1 \mid \dots \mid q_n]$$

$$Q^T A Q = H$$

$$\Rightarrow A Q = Q H$$

Equating the  $j^{th}$  column on both sides -

$$A q_j = \sum_{i=1}^{j+1} h_{ij} q_i$$

Since  $q_i$  are orthogonal  $\forall 1 \leq i \leq j$ ,

$$\therefore q_m^T A q_j = \sum_{i=1}^{j+1} h_{ij} q_m^T q_i$$

$$= h_{mj} \quad \rightarrow \quad \therefore A q_j = \sum_{i=1}^j h_{ij} q_i + h_{j+1,j} q_{j+1}$$

$$h_{j+1,j} q_{j+1} = A q_j - \sum_{i=1}^j h_{ij} q_i$$

\* Arnoldi's algorithm

$$q_1 = \frac{b}{\|b\|_2} \quad k = \# \text{ of columns of } Q(AH) \text{ to be computed.}$$

for  $j = 1 \text{ to } k$

$$z = A q_j$$

for  $i = 1 \text{ to } j$

$$h_{ij} = q_i^T z$$

$$z = z - h_{ij} q_i$$

} like  
MGS.

end

$$h_{j+1,j} = \|z\|_2$$

if  $h_{j+1,j} = 0$ ; quit.

$$q_{j+1} = \frac{z}{h_{j+1,j}}$$

end.

$$Q = [Q_k \mid Q_u]$$

↑ known      ↓ unknown.

$$\begin{aligned} H &= Q^T A Q \\ &= [Q_k \mid Q_u]^T A [Q_k \mid Q_u] \end{aligned}$$

$$= \begin{bmatrix} Q_k^T A Q_k & Q_k^T A Q_u \\ Q_u^T A Q_k & Q_u^T A Q_u \end{bmatrix}$$

$$= \begin{bmatrix} (H_K) & H_{UK} \\ H_{KU} & H_{UU} \end{bmatrix}$$

↑ known.

Last time: Arnoldi's algorithm for generating an orthogonal basis for the Krylov subspace

$$K_k(A, b) = \{b, Ab, A^2b, \dots, A^{k-1}b\}$$

Why Krylov subspaces are useful?

- General dense matrix  $A$ : factorisation methods to solve  $AX=b$  & general eigenvalue methods for  $A$  are effective
- If  $A$  is huge & sparse: need other methods. - iterative methods like Krylov subspace methods.

Solving  $AX=b \leftrightarrow$  optimisation problem.

$$\text{Quadratic } f(x) = \frac{1}{2} x^T Ax - b^T x + c.$$

Using  
matrix derivative

$$f'(x) = \frac{1}{2} A^T x + \frac{1}{2} Ax - b$$

\*\* If  $A$  is symmetric p.d.

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

$$b - Ax_0 = r_0$$

$$K_k(A, r_0) = \{r_0, Ar_0, \dots, A^{k-1}r_0\}$$

gradients are linear combinations of

When  $A$  is symmetric

$$RA(x) = \frac{x^T Ax}{x^T x}$$

tri-diagonal  
matrix

$\leftarrow T \leftrightarrow$  optimizing  $R_T(x)$

$$\nabla R_T \in \langle A, Ar, A^2r, \dots, A^{k-1}r \rangle.$$

## Krylov subspace methods

① Arnoldi's

$$\begin{array}{c} Ax = b \\ \Delta x = \lambda k \end{array}$$

Conjugate

Gradient

Lanczos

GMRES.

\* Lanczos algorithm for extremal eigenvalues of a symmetric p.d. A (large, sparse)

want:  $Q^T A Q = T$        $T \Rightarrow$  tridiagonal

$$AQ = QT$$

$$T = \begin{bmatrix} \alpha_1 & \beta_1 & & & \\ \beta_1 & \ddots & & & \\ & & \ddots & & \beta_{n-1} \\ & & & \ddots & \alpha_n \end{bmatrix}$$

$$Aq_j = \beta_{j-1} q_{j-1} + \alpha_j q_j + \beta_j q_{j+1}$$

$$q_j^T A q_j = q_j^T q_{j-1}$$

$\leftarrow$   $Q$  is orthogonal.

$$q_j^T A q_j = \alpha_j$$

$$q_1^T A q_1 = \alpha_1, \quad q_1 = \frac{a_1}{\|a_1\|}, \quad \beta = 0$$

$$q_1 = \frac{a_1}{\|a_1\|}, \beta = 0$$

for  $j = 1$  to  $K$

$$z = A q_j$$

$$\kappa_j = q_j^T z$$

$$z = z - \kappa_j q_j - \beta_{j-1} q_{j-1}$$

$$\beta_j = \|z\|_2$$

if  $\beta_j = 0$ , quit

$$q_{j+1} = \frac{z}{\beta_j}$$

$$T = Q^T A Q$$

$$= \left[ \begin{array}{c|c} T_K & T_{KU} \\ \hline T_{KU} & T_U \end{array} \right] \left[ \begin{array}{c} f_K \\ f_{n-K} \end{array} \right]$$

$T_{KU}^* = T_{U K}$ .