

## FUNDAMENTALS

MATRIX - VECTOR MULTIPLICATION

$x$ :  $n$  COLUMN VECTOR  $\in \mathbb{C}^n$

$A$ :  $\underset{\substack{m \times n \\ T}}{\text{MATRIX}} \in \mathbb{C}^{m \times n}$

# OF ROWS      # OF COLUMNS.

$$b = Ax$$

$b$ :  $m$  COLUMN VECTOR  $\in \mathbb{C}^m$

$$b_i = \sum_{j=1}^n a_{ij} x_j \quad i = 1, \dots, m$$

THIS IS A LINEAR MAP

$$x \mapsto Ax$$

①

SINCE

$$A(x+y) = Ax + Ay$$

$$A(\alpha x) = \alpha Ax$$

$$x, y \in \mathbb{C}^n, \alpha \in \mathbb{C}$$

WE CAN WRITE M-V MULTIPLICATION SLIGHTLY DIFFERENTLY.

$$b = \sum_{j=1}^n x_j \underbrace{a_j}_{\text{jth column of } A}.$$

$$\begin{bmatrix} b \end{bmatrix} = x_1 \begin{bmatrix} a_1 \end{bmatrix} + x_2 \begin{bmatrix} a_2 \end{bmatrix} + \dots + x_n \begin{bmatrix} a_n \end{bmatrix}$$

②

b is a linear combination of the columns of A.

MATRIX TIMES A MATRIX.

A:  $l \times m$ , C:  $m \times n$ , B:  $l \times n$

$$(1) b_{ij} = \sum_{k=1}^m a_{ik} c_{kj} \quad (AB = AC)$$

$$(2) b_j = Ac_j = \sum_{k=1}^m c_{kj} a_k$$

j<sup>th</sup> column of B      | b<sub>j</sub> is a lin. combination of the columns of A with coefficients c<sub>kj</sub>

(3)

EXAMPLE: OUTER PRODUCT.

$$u \in \mathbb{C}^m$$

$$v \in \mathbb{C}^n$$

$$uv^T = \begin{bmatrix} v_1 & v_2 & \dots & v_n \end{bmatrix}^T u$$

$$= \begin{bmatrix} v_1 u & v_2 u & \dots & v_n u \end{bmatrix}$$

DEF: RANGE OF A,  $\text{RANGE}(A)$  IS THE SET OF VECTORS THAT CAN BE EXPRESSED AS  $Ax$  FOR SOME  $x$ .

(4)

Thm:  $\text{RANGE}(A)$  IS THE  
SPACE SPANNED BY THE  
COLUMNS OF  $A$ .

(5)

(THE LESSER OF  $m$  AND  $n$ )

IF  $m \geq n$  ~~IT~~ THEN A  
MUST HAVE  $n$  LINEARLY  
INDEPENDENT COLUMNS TO  
BE FULL RANK.

DEF: NULL SPACE IF  $A \in \mathbb{C}^{m \times n}$ ,  
NULL( $A$ ), IS THE SET  
OF VECTORS  $x$  THAT  
SATISFY  $Ax = 0$ .

DEF: RANK OF  $A$  IS  
THE DIMENSION OF THE  
COLUMN SPACE OF  $A$ .

A MATRIX  $A$  ( $m \times n$ )  
IS FULL RANK IF IT  
HAS THE MAXIMUM POSSIBLE  
RANK.

Thm: A MATRIX  $A \in \mathbb{C}^{m \times n}$   
WITH  $m \geq n$  HAS FULL RANK  
IF IT MAPS NO TWO  
DISTINCT VECTORS TO THE  
SAME VECTOR.

INVERSE

A NONSINGULAR OR INVERTIBLE  
MATRIX IS A SQUARE MATRIX  
( $m \times m$ ) OF FULL RANK.

⇒ COLUMNS PROVIDE A  
BASIS FOR  $\mathbb{C}^m$ .

CAN WRITE

$$\frac{e_j}{\bar{T}} = \sum_{i=1}^m z_{ij} \frac{a_i}{\bar{T}}$$

UNIT VECTOR  
WITH 1 IN THE  $j^{\text{TH}}$  ENTRY

$i^{\text{TH}}$  COLUMN  
OF  $A$ .

$$\left[ e_1 \mid e_2 \mid \dots \mid e_m \right] = \frac{I}{\bar{T}} = A \sum_{\bar{T}}^{m \times m}$$

IDENTITY  $A^{-1}$ .

## INVERSE

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

$\swarrow$

INVERSE OF  $A$

THEM: For  $A \in \mathbb{C}^{n \times n}$  THE  
FOLLOWING ARE EQUIVALENT.

(a)  $A$  HAS AN INVERSE  $A^{-1}$

(b)  $\text{RANK}(A) = n$

(c)  $\text{RANGE}(A) = \mathbb{C}^n$

(d)  $\text{NULL}(A) = \{0\}$

(e)  $0$  IS NOT AN EIGENVALUE  
OF  $A$

{  
    (f)  $0$  IS NOT A SINGULAR VALUE  
    OF  $A$

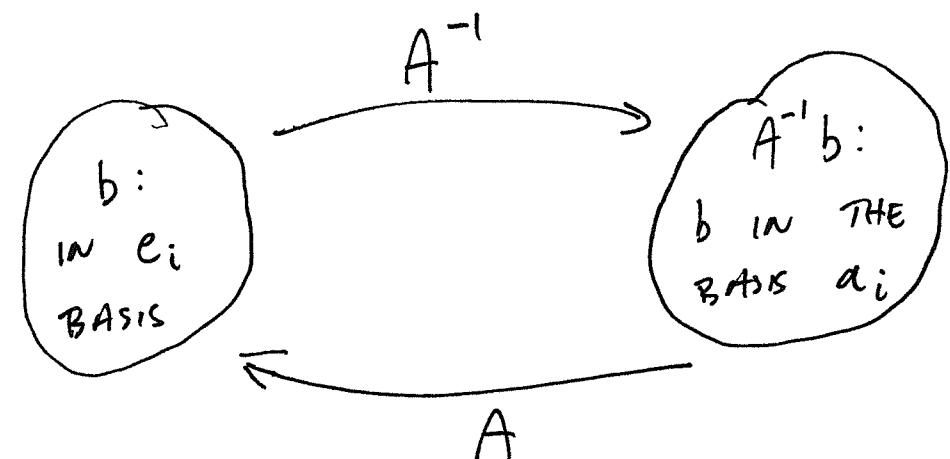
(g)  $\det(A) \neq 0$

①

$$(Ax = b)$$

SINCE  $x = A^{-1}b$  IS  
THE VECTOR OF COEFFICIENTS  
OF THE UNIQUE LINEAR  
EXPANSION OF  $b$  IN THE  
BASIS OF THE COLUMNS OF  
 $A$  MULTIPLICATION BY  $A^{-1}$   
CAN BE VIEWED AS A  
CHANGE OF BASIS.

②



### ORTHOGONAL VECTORS AND MATRICES

(3)

DEF: ADJOINT OR HERMITIAN

CONJUGATE OR  $A \in \mathbb{C}^{m \times n}$

WE WRITE AS  $A^* \in \mathbb{C}^{n \times m}$

WHERE

$$a_{ij}^* = \bar{a}_{ji} \quad \left( \begin{array}{l} \text{COMPLEX CONJUGATE} \\ \text{OF } a_{ji} \end{array} \right)$$

### EXAMPLE

$$A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \\ a_{31} & a_{32} \end{bmatrix}, \quad A^* = \begin{bmatrix} \bar{a}_{11} & \bar{a}_{21} & \bar{a}_{31} \\ \bar{a}_{12} & \bar{a}_{22} & \bar{a}_{32} \end{bmatrix}$$

IF  $A = A^*$ , we SAY  $A$  IS HERMITIAN.

FOR REAL MATRICES.

$$A^* = A^T$$

IF  $A = A^T$  THE MATRIX IS SYMMETRIC.

FOR MATRICES ~~WITH~~  $A, B$  WITH COMPATIBLE DIMENSIONS.

$$(AB)^* = B^* A^*$$

$$\left( \begin{array}{l} \text{FOR INVERTIBLE MATRICES} \\ (AB)^{-1} = B^{-1} A^{-1} \end{array} \right)$$

INNER PRODUCT

$$x, y \in \mathbb{C}^m$$

$$\rightarrow x^* y = \sum_{i=1}^m \bar{x}_i y_i$$

## EUCLIDEAN LENGTH OF A VECTOR

$$\|x\| = \sqrt{x^* x}$$

INNER PRODUCT IS ~~BILINEAR~~  
BILINEAR.

- $(x_1 + x_2)^* y = x_1^* y + x_2^* y$
- $x^*(y_1 + y_2) = x^* y_1 + x^* y_2$
- $(\alpha x)^*(\beta y) = \bar{\alpha} \beta x^* y$

## ORTHOGONAL VECTORS

- $x, y \in \mathbb{C}^n$  ARE ORTHOGONAL  
IF  $x^* y = 0$
- SETS OF VECTORS X AND  
Y ARE ORTHOGONAL IF

(5)

$$x^* y = 0 \quad \forall x \in X, y \in Y.$$

- A SET OF VECTORS S IS ITSELF ORTHOGONAL IF FOR ALL  $x, y \in S$   $x \neq y$   $x^* y = 0$ .
- THE SET S IS ORTHONORMAL IF WE ALSO HAVE  $\|x\| = 1$  FOR EACH  $x \in S$

THM: THE VECTORS IN A  
ORTHOGONAL SET S ARE  
LINEARLY INDEPENDENT.

# COMPONENTS OF A VECTOR

$\{q_1, q_2, \dots, q_n\}$  is an ORTHONORMAL SET

AND  $v$  is AN ARBITRARY VECTOR ( ALL  $\in \mathbb{C}^m$  )

DEFINE:

$$r = v - (q_1^* v)q_1 - \dots - (q_n^* v)q_n$$

$r$  is ORTHOGONAL TO THE SET  $\{q_1, \dots, q_n\}$  since

$$q_i^* r = q_i^* v - (q_1^* v)(q_i^* q_1) - \dots - (q_n^* v)(q_i^* q_n)$$

⑦ SINCE  $q_i^* q_j = 0$  FOR  $i \neq j$  ⑧

$q_i^* q_i = 1$

$q_i^* v = q_i^* v - \underbrace{(q_i^* v)(q_i^* q_i)}_1 = 0$

Thus,

$v = r + \sum_{i=1}^n \underbrace{(q_i^* v) q_i}_{\substack{\text{COEFFICIENT} \\ \text{VECTOR}}$

$= r + \sum_{i=1}^n \underbrace{(q_i q_i^*)}_{{\substack{\text{RANK - 1}}} v$

IF  $\{q_i\}$  IS A BASIS FOR  $\mathbb{C}^n$ , THEN  $n=m$  AND  $r=0$ .  
PROJECTOR.

Thus,

$v = \sum_{i=1}^m \underbrace{(q_i q_i^*)}_{} v$

## UNITARY MATRICES

$Q \in \mathbb{C}^{m \times m}$  is UNITARY

$$\text{IF } Q^* = Q^{-1}.$$

For REAL MATRICES, A MATRIX  
IS ORTHOGONAL IF  $Q^T = Q^{-1}$   
 $(Q \in \mathbb{R}^{m \times m})$

For A UNITARY MATRIX,

$$Q^* Q = I$$

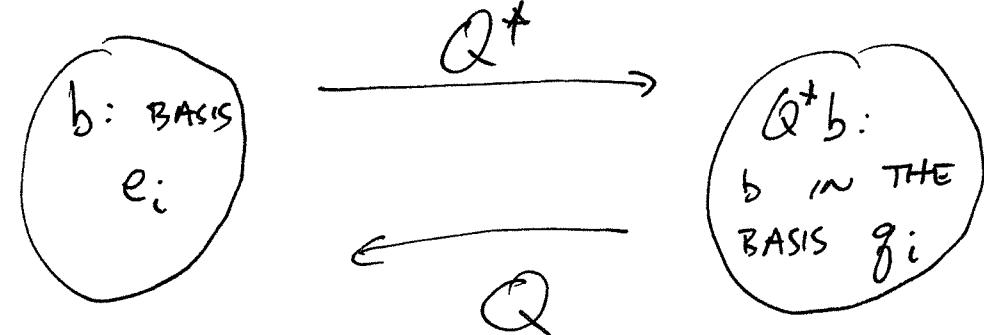
$$\left[ \begin{array}{c} q_1^* \\ q_2^* \\ \vdots \\ q_m^* \end{array} \right] \left[ \begin{array}{|c|c|c|} \hline q_1 & q_2 & \dots & q_m \\ \hline \end{array} \right] = I$$

⑨

THIS HOLDS IF

$$q_i^* q_j = \delta_{ij} \quad (\text{Kronecker Delta})$$

MULTIPLICATION BY  $Q^*$ ,  $Q$



## OTHER PROPERTIES

- $\|Qx\| = \|x\|$

- $(Qx)^* (Qy) = x^* y$

⑩

• HERMITIAN NOT HERMETIAN

• Ann Brew (LIBRARIAN)

ann.brew@imperial.ac.uk

• LEC. 3 4 LEC. 6 (TREFETHEN AND BORG)

### VECTOR NORMS

• MEASURE THE "LENGTH" OF A VECTOR.

A norm is a function

$$\|\cdot\| : \mathbb{C}^m \rightarrow \mathbb{R}$$
 THAT

SATISFIES

(1)  $\|x\| \geq 0$  AND  $\|x\| = 0$  ONLY IF  $x = 0$

(2)  $\|x+y\| \leq \|x\| + \|y\|$  ( $\Delta$ -INEQUALITY)

(3)  $\|\alpha x\| = |\alpha| \|x\|$

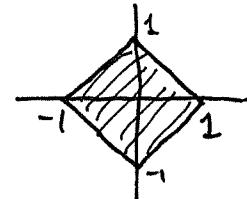
$\forall x, y \in \mathbb{C}^m$  AND  $\alpha \in \mathbb{C}$

①

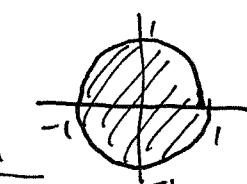
ALREADY SAW THE EUCLIDEAN NORM, BUT THIS IS PART OF A LARGER CLASS OF P-NORMS.

$$\cdot \|x\|_1 = \sum_{i=1}^m |x_i|$$

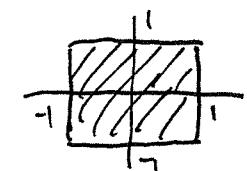
$$\text{for } x \in \mathbb{R}^2 \\ \{x \mid \|x\| \leq 1\}$$



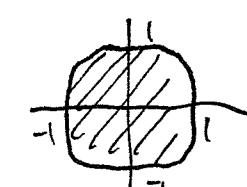
$$\cdot \|x\|_2 = \left( \sum_{i=1}^m |x_i|^2 \right)^{1/2}$$



$$\cdot \|x\|_\infty = \max_{1 \leq i \leq m} |x_i|$$



$$\cdot \|x\|_p = \left( \sum_{i=1}^m |x_i|^p \right)^{1/p}$$



## WEIGHTED NORMS

$$\|x\|_w = \|Wx\|$$

↑  
W IS A DIAGONAL  
MATRIX WITH  $w_{ii} \neq 0$   
A i.

## WEIGHTED 2-NORM

$$\|x\|_w = \left( \sum_{i=1}^m |w_{ii}x_i|^2 \right)^{1/2}$$

## PROJECTORS

A PROJECTOR IS A SQUARE  
MATRIX THAT SATISFIES

$$P^2 = P$$

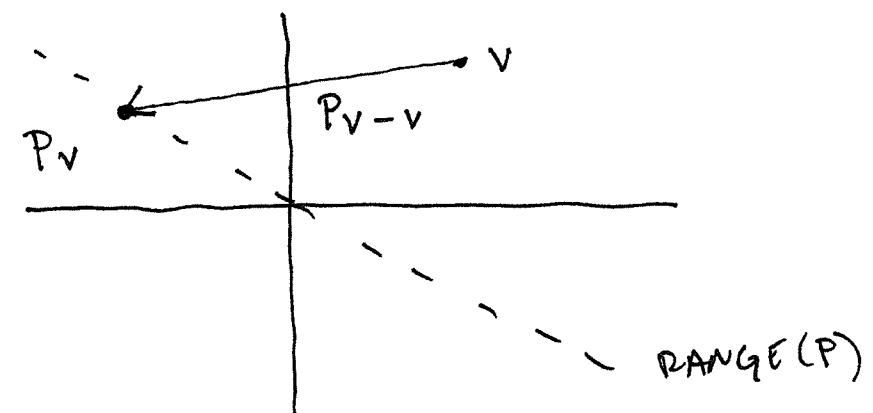
(ALSO SAID TO BE IDIOMPOENT)

⑤

- IF  $v \in \text{RANGE}(P)$   
 $\Rightarrow v = Px$  FOR SOME  $x$

$$Pv = P^2x = Px = v$$

- SUPPOSE THAT  $Pv \neq v$



$$P(Pv - v) = P^2v - Pv = 0$$

$$Pv - v \in \text{NULL}(P).$$

COMPLEMENTARY PROJECTOR.

⑥

IF  $P$  IS A PROJECTOR, (5)  
 $I - P$  IS ALSO A PROJECTOR.

$$(I - P)^2 = I^2 - 2P + P^2 = I - P$$

$I - P$  IS THE COMPLEMENTARY  
PROJECTOR OF  $P$ .

IF  $Pu = 0$  WE HAVE

$$(I - P)u = u$$

$$\Rightarrow \text{RANGE}(I - P) \supseteq \text{NULL}(P)$$

ALSO HAVE  $\text{RANGE}(I - P) \subseteq \text{NULL}(P)$   
 SINCE

$$(I - P)u = u - Pu \in \text{NULL}(P)$$

$$\Rightarrow \text{RANGE}(I - P) = \text{NULL}(P)$$

SINCE WE ~~CAN~~ ~~WE CAN~~ <sup>ALSO</sup> HAVE  $P = I - (I - P)$ , (6)

$$\text{NULL}(I - P) = \text{RANGE}(P)$$

$$\text{Thus, } \text{NULL}(I - P) \cap \text{NULL}(P) = \{0\}$$

$$\text{AND, } \text{RANGE}(P) \cap \text{RANGE}(I - P) = \{0\}$$

A PROJECTOR SEPARATES  $C^n$   
 INTO TWO SPACES.

IN FACT, SUPPOSE SUBSPACES,

$$S_1, S_2 \subseteq C^n \text{ s.t. } S_1 \cap S_2 = \{0\}$$

$$\text{AND } S_1 + S_2 = C^n. \text{ THEN}$$

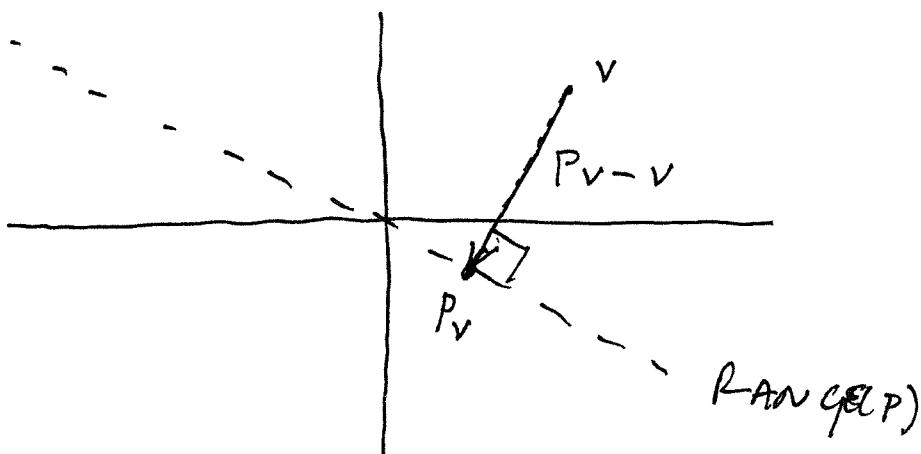
$\exists$  PROJECTOR  $P$  s.t.

$$\text{RANGE}(P) = S_1 \text{ AND } \text{NULL}(P) = S_2$$

P PROJECTS ONTO  $S_1$  ALONG  $S_2$ .<sup>(2)</sup>

ORTHOGONAL PROJECTORS

AN ORTHOGONAL PROJECTOR IS  
ONE THAT PROJECTS ONTO  $S_1$ ,  
ALONG  $S_2$  WHERE  $S_1$  AND  $S_2$   
ARE ORTHOGONAL.



CAN CONSTRUCT ORTHOGONAL  
PROJECTORS FROM SETS OF  
ORTHONORMAL VECTORS.

LET  $\{q_1, \dots, q_n\}$  A SET  
OF n ORTHONORMAL VECTORS  
IN  $\mathbb{C}^n$

LET

$$\hat{Q} = [q_1 \mid q_2 \mid \dots \mid q_n]$$

HAVE SEEN THAT  $v \in \mathbb{C}^n$   
 $v = r + \sum_{i=1}^n (q_i q_i^*) v$   
 COMPONENT OF  
 $v$  IN THE  
 COLUMN SPACE  
 OF  $\hat{Q}$

Thus THE map

$$v \mapsto \underbrace{\sum_{i=1}^n (q_i q_i^*) v}_{y}$$

is AN ORTHOGONAL PROJECTION  
ONTO RANGE ( $\hat{Q}$ ).

$$y = \sum_{i=1}^n (q_i q_i^*) v = \underbrace{\hat{Q} \hat{Q}^*}_{P} v$$

THE COMPLEMENTARY PROJECTOR

$$I - \hat{Q} \hat{Q}^*$$

AN IMPORTANT CASE  $n = 1$

$$P = gg^* \quad (\text{RANK } 1 \text{ ORTHOGONAL PROJECTOR})$$

$$P_\perp = I - gg^* \quad (\text{RANK } n-1)$$

(9)

## QR FACTORIZATION

$$A = Q R \quad (\text{FULL QR})$$

↗   
 $m \times n$    
 ↗   
 $m \times m$  UNITARY MATRIX ( $Q^{-1} = Q^*$ )   
 ↗   
 $m \times n$  UPPER TRIANGULAR MATRIX   
 $r_{ij} = 0 \quad i > j$

$$R = \begin{bmatrix} \text{Hatched} \\ 0 \end{bmatrix}$$

$$A = \overset{\wedge}{Q} \overset{\wedge}{R} \quad (\text{REDUCED QR})$$

↗   
 $m \times n$    
 ↗   
 $n \times n$  UPPER TRIANGULAR MATRIX USE COLUMNS FROM AN ORTHONORMAL SET OF VECTORS.

①

WHAT DO THIS?

②

(1) SOLVE LINEAR SYSTEMS

$$A x = b$$

$m \times n$ , FULL RANK.

(a)  $Q R x = b$  (PERFORM QR)

(b) MULTIPLY BY  $Q^*$

$$\underbrace{Q^* Q}_I R x = \underbrace{Q^* b}_y$$

(c) USE BACKWARDS SUBSTITUTION  
TO SOLVE  $R x = y$

LAST ROW

$$r_{mm} x_m = y_m$$

$$x_m = y_m / r_{mm}$$

PENULTIMATE ROW

$$r_{m-1, m-1} x_{m-1} + r_{m-1, m} x_m = y_{m-1}$$

(2) LEAST SQUARE PROBLEMS.

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

FIND  $x$  s.t.

$\|Ax - b\|_2$  is minimized.

THE SOLUTION IS GIVEN  
BY

$$Ax = \hat{Q}\hat{Q}^* b$$

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

MULTIPLY BY  $\hat{Q}^*$

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

SOLVE FOR  $x$  USING  
BACKWARDS SUBSTITUTION.

(3)

COMPUTING QR FACTORIZATION. (4)

- CLASSICAL GRAM-SCHMIDT.
- MODIFIED GRAM-SCHMIDT.
- HOUSEHOLDER TRIANGULARIZATION.

USING THE MATRIX  $A$  ( $m \times n$ )  
WE CAN CONSTRUCT SUCCESSIVE  
SPACES SPANNED BY ITS  
COLUMNS.

$$\langle a_1 \rangle \subseteq \langle a_1, a_2 \rangle \subseteq \langle a_1, a_2, a_3 \rangle \dots$$

$\langle \dots \rangle$  SPACE SPANNED BY  
THE VECTORS IN THE BRACKETS.

IDEA BEHIND QR:  
CONSTRUCT SEQUENCE OF  
ORTHONORMAL VECTORS  $q_1, q_2, \dots$   
THAT SPAN THE SUCCESSIVE  
SPACES.

$$\left[ \begin{array}{c|c|c|c} a_1 & a_2 & \dots & a_n \end{array} \right] = \underbrace{\left[ \begin{array}{c|c|c|c} q_1 & q_2 & \dots & q_n \end{array} \right]}_{A \text{ } (m \times n)} \times \underbrace{\left[ \begin{array}{cccc} r_{11} & r_{12} & \dots & r_{1n} \\ r_{21} & r_{22} & & \\ 0 & & \ddots & \\ & & & r_{nn} \end{array} \right]}_{\hat{R} \text{ } (n \times n)}$$

THE REDUCED QR FACTORIZATION.

For  $m \geq n$

$$\underbrace{\left[ \begin{array}{c} n \\ m \end{array} \right] A}_{= \left[ \begin{array}{c} n \\ m \end{array} \right] \hat{Q}} \left[ \begin{array}{c} n \\ m \end{array} \right] \hat{R}$$

⑤

FULL QR ( $m \geq n$ )

$$\underbrace{\left[ \begin{array}{c} n \\ m \end{array} \right] A}_{= \left[ \begin{array}{c} m \\ n \end{array} \right] Q} = \left[ \begin{array}{c} m \\ n \end{array} \right] \underbrace{\left[ \begin{array}{c} 1 \\ Q \\ 0 \end{array} \right]}_{Q^{\text{full}}} \left[ \begin{array}{c} n \\ m \end{array} \right] \hat{R}$$

ORTHONORMAL ADDITIONAL VECTORS THAT IF A IS FULL RANK SPAN  $\text{RANGE}(A)^\perp$  {on  $\text{NULL}(A^*)$ } DONE TO MAKE Q UNITARY.

GRAM-SCHMIDT ORTHOGONALIZATION

GIVEN  $a_1, a_2, \dots, a_n$  FIND  $q_1, q_2, \dots, q_n$  AND  $r_{ij}$  BY SUCCESSIVE ORTHOGONALIZATION

FIND  $g_j \in \langle a_1, a_2, \dots, a_j \rangle$  ②

THAT IS ORTHOGONAL TO

$q_1, q_2, \dots, q_{j-1}$

$$1. v_j = a_j - (q_1^* a_j) q_1 - (q_2^* a_j) q_2 - \dots - (q_{j-1}^* a_j) q_{j-1}$$

$$2. g_j = \frac{v_j}{\|v_j\|_2}$$

BY THE FACT THAT  $A = \overset{\text{1}}{Q} \overset{\text{1}}{R}$

WE KNOW

$$q_1 = \frac{a_1}{r_{11}}$$

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

:

$$g_n = \frac{a_n - \sum_{i=1}^{n-1} r_{in} q_i}{r_{nn}}$$
⑧

BY COMPARING THE TWO APPROACHES WE SEE THAT

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

$$\text{AND } \|r_{jj}\| = \|a_j - \sum_{i=1}^{j-1} r_{ij} q_i\|_2$$

CLASSICAL GRAM-SCHMIDT (CGS)  
ALGORITHM

for  $j = 1$  to  $n$

$$v_j = a_j$$

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

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

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

END FOR

$$r_{ij} = \|v_j\|_2$$

$$g_j = v_j / r_{jj}$$

END FOR

CAN USE CGS TO PROVE

EXISTENCE & UNIQUENESS (FULL RANK)  
 $A$   
(TREFETHEN & LAN p. 51)

BUT TURNS OUT TO BE  
UNSTABLE FOR NUMERICAL  
PURPOSES.

## CGS PROJECTORS

GS CAN ALSO BE WRITTEN  
IN TERMS OF PROJECTORS.

$$q_1 = \frac{P_1 a_1}{\|P_1 a_1\|}, \quad q_2 = \frac{P_2 a_2}{\|P_2 a_2\|}$$

$$q_n = \frac{P_n a_n}{\|P_n a_n\|}$$

$P_j$  IS AN ORTHOGONAL PROJECTOR

AN  $m \times m$  MATRIX THAT  
PROJECTS  $C^m$  ORTHOGONALLY  
ONTO  $\langle q_1, q_2, \dots, q_{j-1} \rangle$

①

LET  $\hat{Q}_{j-1} = [q_1 | q_2 | \dots | q_{j-1}]$

THEN

$$P_j = I - \hat{Q}_{j-1} \hat{Q}_{j-1}^*$$

AT EACH STEP, CGS  
COMPUTES

$$v_j = P_j a_j \quad (*)$$

MODIFIED GS (MGS)

DECOMPOSE  $P_j$  INTO  $j-1$   
PROJECTORS OF RANK  $m-1$

$$P_j = P_{\perp q_{j-1}} \cdots P_{\perp q_2} P_{\perp q_1}$$

②

WHERE

$$P_{\perp g_j} = I - g_j g_j^*$$

thus, (\*) CAN BE WRITTEN

AS

$$v_j = P_{\perp g_{j-1}} \cdots P_{\perp g_2} P_{\perp g_1} q_j$$

AT EACH STEP MGS  
APPLIES  $P_{\perp g_i}$  TO ALL  
COLUMNS.

(3)

### MGS ALGORITHM

for  $i = 1$  to  $n$

$$v_i = q_i$$

END FOR

for  $i = 1$  to  $n$

$$r_{ii} = \|v_i\|_2$$

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

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

$$r_{ij} = g_i^* v_j$$

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

END FOR

END FOR

(4)

(6)

MGS CAN BE THOUGHT OF AS TRIANGULAR ORTHOGONALIZATION.

FIRST STEP CAN BE VIEWED AS

$$\left[ \begin{array}{c|c|c|c} v_1 & v_2 & \cdots & v_n \end{array} \right] \xrightarrow{\text{R}_1} \left[ \begin{array}{cccc} 1 & -\frac{r_{12}}{r_{11}} & \cdots & -\frac{r_{1n}}{r_{11}} \\ 0 & 1 & \cdots & 0 \\ \vdots & \ddots & 1 & 0 \\ 0 & 0 & \cdots & 1 \end{array} \right]$$

$$R_1 = \left[ \begin{array}{c|c|c|c} q_{11} & v_2^{(2)} & \cdots & v_n^{(2)} \\ 0 & \ddots & \ddots & \vdots \\ 0 & 0 & \ddots & 1 \end{array} \right]$$

RIGHT MULTIPLY BY  $R_1$

SECOND STEP: RIGHT MULTIPLY BY  $R_2$

$$R_2 = \left[ \begin{array}{ccccc} 1 & 1 & \cdots & \cdots & 1 \\ 0 & \frac{1}{r_{22}} & -\frac{r_{23}}{r_{22}} & \cdots & -\frac{r_{2n}}{r_{22}} \\ 0 & 0 & 1 & \cdots & 0 \\ 0 & 0 & \cdots & \ddots & 1 \end{array} \right]$$

DOING THIS  $n$  TIMES

$\underbrace{A R_1 R_2 \cdots R_n}_R = Q$

TRIANGULAR ORTHOGONALIZATION  
USING TRIANGULAR MATRICES  
TO REDUCE  $A$  TO A  
MATRIX WITH ORTHONORMAL  
COLUMNS.

$R_i$ 'S ARE NEVER FORMED  
EXPLICITLY BUT THEY  
GIVE US A CLEAR  
WAY OF THINKING ABOUT  
WHAT IS HAPPENING AT  
EACH STEP.

(7)

3. HOUSEHOLDER TRIANGULARIZATION

• ORTHOGONAL TRIANGULARIZATION.

TRIANGULARIZE  $A$  USING  
ORTHOGONAL MATRICES.

$$\underbrace{Q_m \dots Q_2 Q_1}_Q A = R$$

Form FULL QR of  $A$ .TRIANGULARIZE BY INTRODUCING  
ZEROS.

$$A \xrightarrow{Q_1} Q_1 A \xrightarrow{Q_2} Q_2 Q_1 A$$

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

(8)

$$Q_3 \rightarrow \begin{bmatrix} x & x & x \\ 0 & x & x \\ 0 & 0 & x \\ 0 & 0 & 0 \end{bmatrix}$$

$$Q_3 Q_2 Q_1 A$$

FIND  $Q_k$  (UNITARY) TO  
INTRODUCE ZEROS BELOW THE  
KTH DIAGONAL WHILE PRESERVING  
THE ZEROS PREVIOUSLY  
INTRODUCED.

APPROACH

$$Q_k = \begin{bmatrix} I & 0 \\ 0 & F \end{bmatrix}$$

I is  $(k-1) \times (k-1)$  IDENTITY

F is  $(m-k+1) \times (m-k+1)$  UNITARY

MATRIX CALL HOUSEHOLDER  
REFLECTOR.

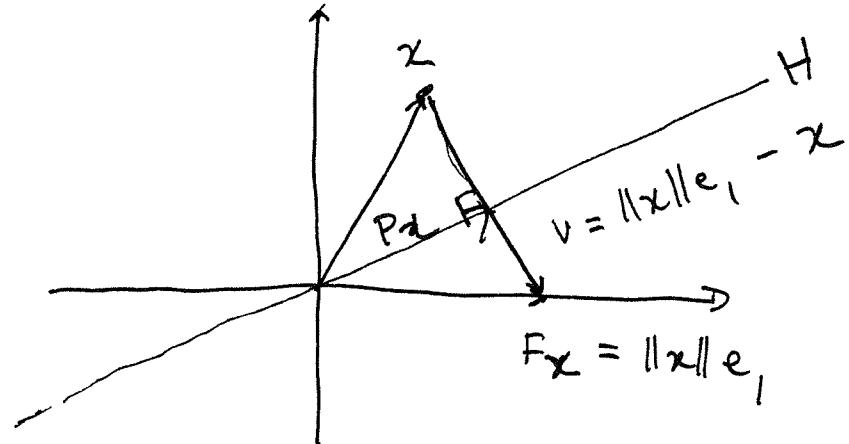
For THE kTH STEP

THE  $k, \dots, m$  ENTRIES

OF THE kTH COLUMN WE  
WRITE AS  $x \in \mathbb{C}^{m-k+1}$

$$x = \begin{bmatrix} x \\ x \\ x \\ x \\ x \end{bmatrix} \xrightarrow{F} Fx = \begin{bmatrix} \|x\| \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} = \|x\|e_1$$

⑦



⑩

F MAPS A VECTOR ON ONE  
SIDE OF H TO ITS  
MIRROR IMAGE.

THE PROJECTOR OF x onto  
H IS

$$P = I - \frac{vv^*}{v^*v}$$

TO REFLECT ABOUT H  
WE NEED TO GO TWICE  
AS FAR IN THE DIRECTION  
OF v.

$$F = I - 2 \frac{VV^*}{V^*V}$$

(11)

FULL RANK AND UNITARY.

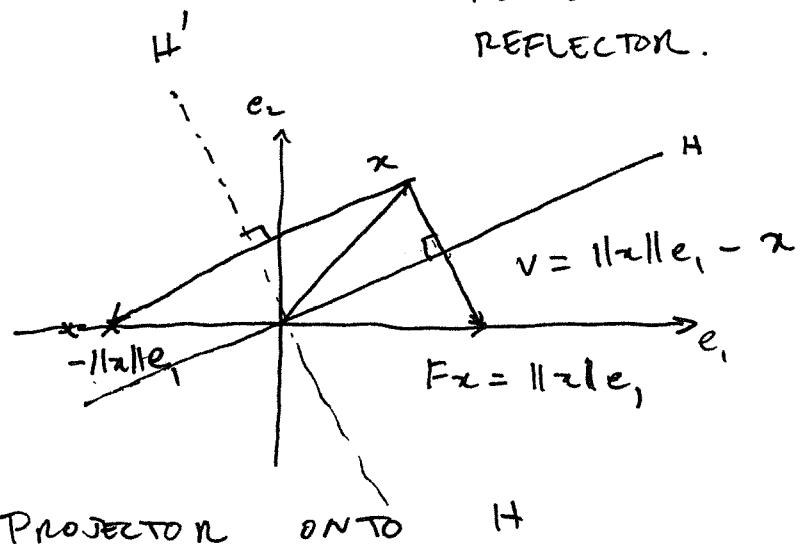
## HOUSE HOLDER TRIANGULARIZATION

①

$$Q_n \cdots Q_2 Q_1 A = R$$

$$Q_k = \begin{bmatrix} I & 0 \\ 0 & F \end{bmatrix}$$

HOUSE HOLDER  
REFLECTOR.



$$P = I - q q^* \quad \text{WHERE} \quad q = \frac{v}{\|v\|}$$

TO GO TWICE THE DISTANCE IN g

$$F = I - 2q q^* = I - 2 \frac{v v^*}{v^* v}$$

② TURNS OUT THERE ARE TWO CHOICES.

WHILE MATHEMATICALLY SPEAKING BOTH CHOICES ARE EQUALLY GOOD, NUMERICALLY WE'D LIKE TO CHOOSE THE REFL. FURTHEST AWAY FROM x.

$$v = -\text{sign}(x_1) \|x\| e_1 - x$$

or

$$v = \text{Sign}(x_1) \|x\| e_1 + x$$

TAKING ALSO  $\text{sign}(0) = 1$

~~Notes~~

## HOUSEHOLDER ALGORITHM

for  $k = 1$  to  $n$

$x = A_{k:m, k} \leftarrow$  VECTOR OF  
THE  $k$  THROUGH  
 $m$  ENTRIES OF  
THE  $k$ TH COLUMN

$$v_k = \text{sign}(x_1) \|x\|_2 e_1 + x \quad \begin{bmatrix} 1 \\ \vdots \\ x_n \end{bmatrix}$$

$$v_k = v_k / \|v_k\|_2$$

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

END. FOR

$A$  IS REDUCED TO UPPER  
TRIANGULAR FORM WITHOUT  
EVER EXPLICITLY COMPUTING  $Q$ .

(3)

HOWEVER, WE KNOW HOW  
TO APPLY  $Q$  (AND  $Q^*$ )  
ONCE THE  $v_k$ 'S ARE KNOWN  
SINCE

$$Q^* = Q_n \dots Q_2 Q_1$$

$$Q = Q_1 Q_2 \dots Q_n$$

RATHER THAN COMPUTING  $Q$   
FROM THE  $Q_k$ 'S THEN APPLYING  
 $Q_1$ , THE SIDE SEQUENCE CAN  
BE APPLIED DIRECTLY

IMPLICIT CALCULATION OF  $Q^* b$

for  $k = 1$  to  $n$

$$b_{k:m} = b_{k:m} - 2v_k(v_k^* b_{k:m})$$

END FOR

(6)

## IMPLICIT CALCULATION OF $Qx^{\top}$

for  $k = n$  down to 1

$$x_{k:m} = \nu_{k:m} - 2\nu_k (\nu_k^T x_{k:m})$$

END FOR

To form  $Q$  we recognize  
THAT

$$Q I = Q$$

AND PERFORM THE IMPLICIT  
CALCULATION

$$Q e_i = g_i \quad \text{For all } i=1,..,m$$

EXAMPLE : LEGENDRE POLYNOMIALS

$$P_n(x) \quad \text{on} \quad [-1, 1]$$

POLYNOMIALS THAT ARE  
ORTHOGONAL WITH RESPECT TO

THE INNER PRODUCT

$$(f, g) = \int_{-1}^1 f(x)g(x) dx$$

For  $P_n(x)$  AND  $P_m(x)$

$$\int_{-1}^1 P_n(x) P_m(x) dx = \frac{2}{2n+1} \delta_{nm}$$

THE FIRST FEW ARE

$$P_0(x) = 1$$

$$P_1(x) = x$$

$$P_2(x) = \frac{3x^2}{2} - \frac{1}{2}$$

$$P_3(x) = \frac{5x^3}{2} - \frac{3}{2}x$$

⋮

THEY CAN BE GENERATED  
BY THE CONTINUOUS VERSION  
OF QR APPLIED TO THE  
"MATRIX" WHOSE "COLUMNS"  
ARE THE MONOMIALS

$$A = \begin{bmatrix} 1 & | & x & | & x^2 & | & x^3 & | \dots & | x^{n-1} \end{bmatrix}$$

USING A CONTINUOUS VERSION  
OF GRAM-SCHMIDT REPLACING

$$q_i^* v_j \text{ BY } \underbrace{\int_{-1}^1 \bar{q}_i(x) v_j(x) dx}_{r_{ij}}$$

② AND  
 $\|v\|_2$  BY  $\left( \int_{-1}^1 \bar{v}(x)v(x) dx \right)^{1/2}$   
 WE OBTAIN

$$A = \begin{bmatrix} q_0(x) & | & q_1(x) & | & \dots & | & q_{n-1}(x) \end{bmatrix} \begin{bmatrix} r_{00} & \dots & r_{0n} \\ \vdots & \ddots & \vdots \\ r_{n0} & \dots & r_{nn} \end{bmatrix}$$

$P_n(x)$  IS RELATED TO  $q_n(x)$   
THROUGH

$$P_n(x) = \frac{q_n(x)}{q_n(1)}$$

# DISCRETE LEGENDRE POLYNOMIALS

①

THE CONTINUOUS FORMULATION  
CAN BE MADE DISCRETE  
BY EVALUATING THE  
MONOMIALS AT  $n$  EQUALLY  
SPACED POINTS.

$$x_i = \frac{2}{m}(i-1) - 1, \quad i=1, \dots, m$$

THIS TURNS A INTO AN  
 $m \times n$  VANDERMONDE MATRIX

$$A = \left[ \begin{array}{|c|c|c|c|} \hline 1 & x_1 & x_1^2 & x_1^{n-1} \\ \hline \vdots & \vdots & \vdots & \vdots \\ \hline 1 & x_m & x_m^2 & x_m^{n-1} \\ \hline \end{array} \right]$$

BY PERFORMING THE  
STANDARD QR, THE  
RESULTING COLUMNS OF Q,  
 $g_k$ , WILL APPROXIMATE  
 $g_k(x_i)$  SINCE THE  
DISCRETE INNER PRODUCT  
IS PROPORTIONAL TO AN  
APPROXIMATION OF THE  
CONTINUOUS ONE

$$\begin{aligned} f^* g &= \sum_{i=1}^m \bar{f}_i g_i \\ &\approx \sum_{i=1}^m \bar{f}(x_i) g(x_i) \\ &\approx \frac{m}{2} \int_{-1}^1 \bar{f}(x) g(x) dx \end{aligned}$$

②

THE APPROXIMATE VALUES OF ③

$P_{\text{approx}}(x_i)$  is

$$P_k(i) = \frac{g_k^{(i)}}{g_k^{(m)}} \quad (\text{MATLAB NOTATION})$$

### LEAST SQUARES PROBLEMS

$$A \in \mathbb{C}^{m \times n} \quad m \geq n$$

$$b \in \mathbb{C}^m$$

FIND  $x \in \mathbb{C}^n$  such THAT

$\|b - Ax\|_2$  IS MINIMIZED  
↓  
 $r$ , RESIDUAL.

### THEOREM

GIVEN THE LEAST SQUARES PROBLEM, A VECTOR

$x$  MINIMIZES  $\|r\|_2$  ④

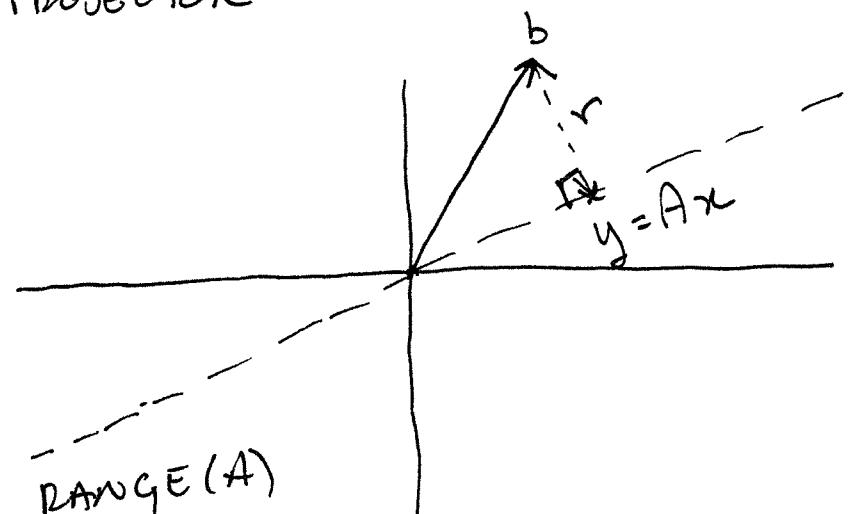
IFF  $r \perp \text{RANGE}(A)$ , THAT IS

$$(i) A^* r = 0$$

$$(ii) A^* A x = A^* b \quad (\text{NORMAL EQUATIONS})$$

$$(iii) P b = Ax$$

WHERE  $P$  IS THE ORTHOGONAL PROJECTION ONTO THE RANGE(A)



IF A IS FULL RANK

THEN  $x$  IS UNIQUE.

USING  $\check{Q}\check{R}$ , WE CAN  
CONSTRUCT THE PROJECTOR  
P AS

$$P = \check{Q} \check{Q}^*$$

SINCE ~~THE~~ THE COLUMNS OF  
 $\check{Q}$  SPAN THE RANGE OF A.

Thus, FROM (iii)

$$Pb = \check{Q} \check{Q}^* b = Ax = \check{Q} \check{R} x$$

LEFT-MULTIPLYING BY  $\check{Q}^*$

WE HAVE

$$\check{R}x = \check{Q}^* b$$

(5)

THIS CAN BE SOLVED USING  
BACKWARDS SUBSTITUTION.  
(6)

EXAMPLE : POLYNOMIAL FITTING.

GIVEN  $x_1, \dots, x_m$  AND  
 $y_1, \dots, y_m$ . FIND

$$p(x) = c_0 + c_1 x + c_2 x^2 + \dots + c_{n-1} x^{n-1}$$

such that

$$\sum_{i=1}^m |p(x_i) - y_i|^2$$

IS MINIMIZED.

WE CAN SET THIS  
UP AS A LEAST SQUARES  
PROBLEM.

$$b = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}, \quad "x" = \begin{bmatrix} c_0 \\ \vdots \\ c_{n-1} \end{bmatrix}$$

(unknowns)

$$A = \begin{bmatrix} 1 & x_1 & & x_1^{n-1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_m & & x_m^{n-1} \end{bmatrix}$$

IF  $x_i$  ARE EQUALLY SPACED, THEN  $\hat{Q}$  ARE RELATED TO THE DISCRETE LEGENDRE POLYNOMIALS ON THE INTERVAL  $[x_1, x_m]$ .

Thus  $\hat{Q}\hat{Q}^* b$  IS THE PROJECTION ONTO THE SPACE

(7)

SPANNED BY THE DISCRETE LEGENDRE POLYNOMIALS.

(8)

## ANALYSING SCHEMES

FOR QR, WE ENCOUNTERED  
THREE DIFFERENT ALGORITHMS

FOR ~~REACH~~ PERFORMING THE  
SAME COMPUTATION.

WHAT MAKES ONE APPROACH  
BETTER THAN THE OTHER?

(1) ONE COULD BE  
FASTER THAN THE  
OTHERS.

(2) ONE METHOD COULD  
BE MORE ROBUST THAN  
THE OTHERS.

①

THE FIRST PROPERTY IS  
CHARACTERIZED BY THE  
OPERATION COUNT.

②

THE SECOND IS RELATED  
TO THE STABILITY OF  
THE ALGORITHM, OR HOW  
SENSITIVE IT IS TO  
PERTURBATIONS.

OPERATION COUNT

• COUNT THE NUMBER  
OF "FLOPS" (FLOATING  
OPERATIONS) THAT  
THE ALGORITHM REQUIRES  
IN THE LIMIT OF  
VERY LARGE MATRICES.

## • THE OPERATIONS

$+, -, \times, \div, \sqrt{\quad}$

ALL COST ONE FLOP

FOR EACH REAL NUMBER

THE OPERATION COUNT  
IS A CLASSICAL MEASURE  
OF THE COST. IN PRACTICE,  
THERE ARE ALSO OTHER  
IMPORTANT FACTORS.

(i) MOVEMENT OF DATA  
IN MEMORY.

(ii) PARALLEL VS. SERIAL  
COMPUTATION.

(3)

WHICH OTHER PROGRAMS RUNNING (4)  
ON THE COMPUTER

.

.

.

OPERATION COUNT FOR MGS.

THE GRAM-SCHMIDT ALGORITHMS  
REQUIRE  $\sim 2mn^2$  FLOPS

TO COMPUTE A QR  
FACTORIZATION OF AN  
 $m \times n$  MATRIX OF REAL  
NUMBERS.

" $n$ " MEANS

$$\lim_{m,n \rightarrow \infty} \frac{\# \text{ OF FLOPS}}{2mn^2} = 1$$

For MGS, when  $m \leq n$   
 ARE LARGE THE WORK  
 IS DOMINATED BY WHAT  
 IS DONE IN THE INNER  
 LOOP

$$(i) r_{ij} = g_i^* v_j$$

INNER PRODUCT OF  
 TWO VECTORS IN  $\mathbb{R}^m$

$m: x$

$m-1: +$

which GIVE  $2m-1$  FLOPS.

$$(ii) v_j = v_j - \underbrace{r_{ij} g_i}_{\text{SCALAR TIMES } A}$$

SUBTRACTION OF  
 TWO VECTORS  
 IN  $\mathbb{R}^m$

$m: x$

$m: -$

(6)

THIS REQUIRES  $2m$  FLOPS.

Thus EACH ITERATION  
 OF THE INNER LOOP  
 REQUIRE  $\sim 4m$  FLOPS.

$$\# \text{ OF FLOPS} \sim \sum_{i=1}^n \sum_{j=i+1}^n 4m$$

$$= 4m \sum_{i=1}^n \sum_{j=i+1}^{n-1} 1$$

$$\approx 4m \sum_{i=1}^n i \left( \int_1^n x dz \right)$$

$$\approx 4m \frac{n^2}{2}$$

$$\approx 2mn^2$$

OPERATION COUNT FOR  
HOUSEHOLDER.

WORK IS DOMINATED  
BY

$$A_{k:m, k:n} = A_{k:m, k:n}$$

$$(3) \quad = (2v_k) \left( v_k^* A_{k:m, k:n} \right)$$

$$(1) \quad v_k^* A_{k:m, k:n}$$

$(n-k)$  INNER PRODUCTS OF  
VECTORS IN  $\mathbb{R}^{m-k}$

$$\# \text{ OF FLOPS} \sim 2(n-k)(m-k)$$

(7)

$$(2) (2v_k) ( )$$

OUTER PRODUCT OF  
TWO VECTORS IN  $\mathbb{R}^{m-k}$   
AND  $\mathbb{R}^{n-k}$

$$\# \text{ OF FLOPS} \sim (m-k)(n-k)$$

(3) SUBTRACTION OF TWO  
 $(m-k) \times (n-k)$  MATRICES

$$\# \text{ OF FLOPS} \sim (m-k)(n-k)$$

TOTAL COST

$$\sim \sum_{k=1}^n 4(n-k)(m-k)$$

$$= 4 \left( \sum_{k=1}^n nm - k(n+m) + k^2 \right)$$

(8)

$$= 4 \sum_{k=1}^n nm - 4 \sum_{k=1}^n k(n+m) + 4 \sum_{k=1}^n k^2$$

⑨

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

$$= 4mn^2 - 2n^3 - 2mn^2 + \frac{4n^3}{3}$$

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

HOUSEHOLDER REQUIRES  
FEWER FLOPS TO COMPUTE  
QR.

## MATRIX NORMS

- CAN THINK OF A MATRIX AS AN  $m \times n$  LENGTH VECTOR.
- APPLY USUAL VECTOR NORMS TO THIS VECTOR.
- MUST SATISFY

$$(1) \|A\| \geq 0 \quad \text{AND} \quad \|A\| = 0 \quad \text{ONLY IF} \quad A=0$$

$$(2) \|A+B\| \leq \|A\| + \|B\|$$

$$(3) \|\alpha A\| = |\alpha| \|A\|$$

FROBENIUS NORM

$$\|A\|_F = \left( \sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2 \right)^{1/2}$$

CAN SHOW:

$$\|AB\|_F \leq \|A\|_F \|B\|_F$$

(1)

- INDUCED MATRIX NORM.

FOR AN  $m \times n$  MATRIX  $A$  AND VECTOR NORMS  $\|\cdot\|_{(n)}$  AND  $\|\cdot\|_{(m)}$  ON THE DOMAIN AND RANGE

$$\|A\|_{(m,n)} = \sup_{\substack{x \in \mathbb{C}^n \\ x \neq 0}} \frac{\|Ax\|_{(m)}}{\|x\|_{(n)}}$$

$$= \sup_{\substack{x \in \mathbb{C}^n \\ \|x\|_{(n)} = 1}} \|Ax\|_{(m)}$$

- MEASURES THE MAXIMUM FACTOR BY WHICH  $A$  WILL "STRETCH" A VECTOR.

(2)

## EXAMPLES

(i)  $\|A\|_2$ ,  $A \in \mathbb{R}^{2 \times 2}$

$$A = \begin{bmatrix} 1 & 2 \\ 0 & 2 \end{bmatrix}$$

ALL UNIT VECTORS IN  $\mathbb{R}^2$   
CAN BE WRITTEN AS

$$x = (\cos \theta, \sin \theta) \quad \theta \in [0, 2\pi)$$

$$Ax = \begin{bmatrix} \cos \theta + 2 \sin \theta \\ 2 \sin \theta \end{bmatrix}$$

$$\|Ax\|_2^2 = \cos^2 \theta + 4 \sin \theta \cos \theta + 8 \sin^2 \theta$$

TO FIND MAX. VALUE.

$$0 = \frac{\partial}{\partial \theta} \|Ax\|_2^2 = 14 \sin \theta \cos \theta + 4(\cos^2 \theta - \sin^2 \theta)$$

③

$$0 = 7 \sin 2\theta + 4 \cos 2\theta$$

$$\Rightarrow \tan 2\theta = -\frac{4}{7}$$

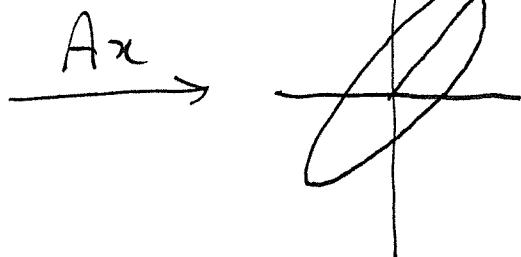
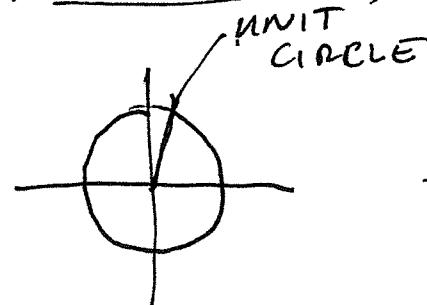
$$2\theta = \tan^{-1}\left(-\frac{4}{7}\right) + n\pi \quad (n \in \mathbb{Z})$$

$$\theta = \frac{1}{2} \tan^{-1}\left(-\frac{4}{7}\right) + n \frac{\pi}{2}$$

MAXIMUM FOR  $n=1$  FOR  
EXAMPLE, WHICH GIVES.

$$\|A\|_2 = 2.9208\dots$$

GRAPHICALLY,



④

(ii) 2 - NORM OF A DIAGONAL MATRIX

$$D = \begin{bmatrix} d_1 & & & \\ & d_2 & & 0 \\ & & \ddots & \\ 0 & & & d_m \end{bmatrix}$$

STRETCHES OR COMPRESSES ALONG THE DIRECTIONS  $e_i$

$$\text{Thus, } \|D\|_2 = \max_i |d_i|.$$

OFTEN DIFFICULT TO FIND EXACT VALUES, SO WE LOOK FOR BOUNDS.

HÖLDER INEQUALITY

$$|x^*y| \leq \|x\|_p \|y\|_q$$

$$\text{For } 1 = \frac{1}{p} + \frac{1}{q}$$

(6)

CAUCHY-SCHWARZ

$$p = q = 2$$

$$|x^*y| \leq \|x\|_2 \|y\|_2$$

EXAMPLE: 2-NORM OF AN OUTER PRODUCT.

$$A = uv^* \text{ SCALAR}$$

$$\begin{aligned} \|A\|_2 &= \|\overbrace{uv^*x}^q\|_2 \\ &= \|v^*x\|_q \|u\|_q \\ &\leq \|u\|_2 \|v\|_2 \|x\|_2 \end{aligned}$$

$$\Rightarrow \|A\|_2 \leq \|u\|_2 \|v\|_2$$

## Bounds for $\|AB\|$

$A: l \times n$

$B: m \times n$

$$\|ABz\|_{(l)} \leq \|A\|_{(l,m)} \|Bz\|_{(m)}$$

$$\leq \|A\|_{(l,m)} \|B\|_{(m,n)} \|z\|_n$$

$$\|AB\|_{(l,n)} \leq \|A\|_{(l,m)} \|B\|_{(m,n)}$$

CONDITIONING AND CONDITION NUMBER

- CONDITIONING PERTAINS TO THE PERTURBATION BEHAVIOUR OF THE MATHEMATICAL PROBLEM.

(2)

- A "PROBLEM" AS A FUNCTION ⑧

$$f: X \rightarrow Y$$

$\uparrow$   
 $\uparrow$

NORMED VECTOR SPACES.

$f$  CAN BE NONLINEAR.

- USUALLY CONCERNED WITH THE BEHAVOUR AT A PARTICULAR VALUE OF  $x \in X$ .

PROBLEM (INSTANCE)

- WELL-CONDITIONED PROBLEM  
ALL SMALL PERTURBATIONS OF  $x$  LEAD TO SMALL CHANGES IN  $f(x)$ .

• ILL-CONDITIONED PROBLEM ①

SMALL CHANGES IN  $x$

GIVE LARGE CHANGES IN  
 $f(x)$ .

MEASURED VIA THE CONDITION  
NUMBER.

## CONDITION NUMBER

ABSOLUTE CONDITION NUMBER

INTRODUCE AN INFINITESIMAL PERTURBATION  $\delta x$  TO  $x$ .

$$\delta f = f(x + \delta x) - f(x)$$

ABSOLUTE CONDITION NUMBER

$$\hat{k} = \sup_{\delta x} \frac{\| \delta f \|}{\| \delta x \|}$$

For DIFFERENTIABLE  $f(x)$   
CAN FIND

JACOBIAN

$$J(x), J_{ij}(x) = \frac{\partial f_i}{\partial x_j}$$

①

For INFINITESIMAL QUANTITIES ②

$$\delta f = J(x) \delta x$$

$$\hat{k} = \sup_{\delta x} \frac{\| J(x) \delta x \|}{\| \delta x \|} \\ = \| J(x) \|$$

RELATIVE CONDITION NUMBER

$$k = \sup_{\delta x} \left( \frac{\| \delta f \|}{\| f(x) \|} \right) \left( \frac{\| \delta x \|}{\| x \|} \right)$$

For DIFFERENTIABLE  $f$

$$k = \sup_{\delta x} \frac{\| J(x) \|}{\| f(x) \|} \left( \frac{\| \delta x \|}{\| x \|} \right)$$

FOR FLOATING POINT NUMBERS ③  
INTRODUCE RELATIVE ERRORS, SO  
K IS USED IN NUMERICAL  
ANALYSIS.

$$K = 1 - 10^{-2} \quad \text{WELL-CONDITIONED}$$

$$K > 10^6 \quad \text{ILL-CONDITIONED.}$$

EXAMPLES (i)  $f: x \mapsto \sqrt{x}$

$$J = \frac{df}{dx} = \frac{1}{2} x^{-\frac{1}{2}}$$

$$K = \frac{\left| \frac{1}{2} x^{-\frac{1}{2}} \right|}{\left| x^{\frac{1}{2}} \right|} = \frac{1}{2}$$

(ii) FINDING THE ROOTS OF  
A POLYNOMIAL GIVEN THE  
COEFFICIENTS.

CONSIDER

$$x^2 - 2x + 1 = (x - 1)^2$$

ROOTS ARE  $x=1$  TWICE.

NOW TAKE

$$x^2 - 2x + 0.9999$$

$$= (x - 0.99)(x - 1.01)$$

ROOTS ARE  $x = 0.99$  AND  $x = 1.01$

RELATIVE CHANGE OF  $10^{-4}$  IN COEFF  
GIVES ~~CH~~ RELATIVE CHANGE OF  $10^{-2}$   
IN THE ROOTS.

$$\text{IN FACT, } \frac{\delta r}{r} = C \sqrt{\frac{\delta c}{c}}$$

THIS GIVES

$$K = \infty$$

THIS PROBLEM IS ILL-CONDITIONED.

④

A MORE COMPLEX EXAMPLE  
WITH WILKINSON POLYNOMIAL

$$p(x) = \prod_{i=1}^{20} (x-i)$$

CONDITION NUMBER OF MATRIX  
VECTOR MULTIPLICATION

Fix  $A \in \mathbb{C}^{m \times n}$ ,  $\exists: x \rightarrow Ax$

$$\kappa = \sup_{\delta x} \frac{\|A(x+\delta x) - Ax\|}{\|Ax\|} \left/ \frac{\|\delta x\|}{\|x\|}\right.$$

$$= \sup_{\delta x} \frac{\|A \delta x\|}{\|\delta x\|} \left/ \frac{\|Ax\|}{\|x\|}\right.$$

$$= \|A\| \left/ \frac{\|Ax\|}{\|x\|}\right.$$

(5)

IF  $A \in \mathbb{C}^{m \times m}$  AND NONSINGULAR (6)

$$\|x\| = \|A^{-1} A x\|$$

$$\leq \|A^{-1}\| \|Ax\|$$

$$\Rightarrow \frac{\|x\|}{\|Ax\|} \leq \|A^{-1}\|$$

$$\Rightarrow \kappa \leq \|A\| \|A^{-1}\|$$

THEOREM: LET  $A \in \mathbb{C}^{m \times n}$   
AND NONSINGULAR AND  
CONSIDER  $Ax=b$ . THE PROBLEM  
OF COMPUTING  $b$  GIVEN  $x$  HAS

$$\kappa = \|A\| \frac{\|x\|}{\|b\|} \leq \|A\| \|A^{-1}\|$$

THE PROBLEM OF COMPUTING  $x$   
GIVEN  $b$  HAS

$$\kappa = \|A^{-1}\| \frac{\|b\|}{\|x\|} \leq \|A^{-1}\| \|A\|.$$

THE CONDITION NUMBER OF  
A

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

CONDITION NUMBER OF A

SYSTEM,  $Ax=b$

•  $b$  IS FIXED

•  $f: A \rightarrow x$

HOW DO PERTURBATIONS IN  $A, \delta A$ ,  
CHANGE  $x$ ?

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

USE THE FACT THAT  $Ax=b$

AND IGNORING  $\delta A \delta x$ , WE

HAVE

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

⑦

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

$$\|\delta x\| \leq \|A^{-1}\| \|\delta A\| \|x\|$$

⑧

CONDITION NUMBER

$$\kappa = \sup_{\delta A} \frac{\|\delta x\|}{\|x\|} \Bigg/ \frac{\|\delta A\|}{\|A\|}$$

$$\leq \|A^{-1}\| \|A\| = \kappa(A)$$

EQUALITY IF

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

CAN ALWAYS BE FOUND.

$$\kappa = \kappa(A).$$

THEOREM: LET  $b$  BE

⑦

FIXED AND CONSIDER  $x = A^{-1}b$ ,

WHERE  $A$  IS SQUARE AND

NON-SINGULAR. THE CONDITION

NUMBER OF THIS PROBLEM

W. R. T.  $\rightarrow$  PERTURBATIONS IN

$A \rightarrow$

$$\kappa = \kappa(A) = \|A\| \|A^{-1}\|.$$

# FLOATING POINT NUMBERS AND ALGEBRAIC

①

(b) GAPS BETWEEN ADJACENT NUMBERS.

②

A COMPUTER MUST USE A DISCRETE REPRESENTATION OF  $\mathbb{R}$ .

(a) THERE MUST BE A LARGEST AND SMALLEST POSITIVE NUMBER.

ON DOUBLE PRECISION MACHINE

$$N_{\max} \approx 1.79 \times 10^{308}$$

$$N_{\min} \approx 2.23 \times 10^{-308}$$

THIS IS TYPICALLY NOT THE ISSUE.

ON DOUBLE PRECISION MACHINE.

[1, 2]

1,  $1 + 2^{-52}$ ,  $1 + 2 \times 2^{-52}$ , ..., 2  
NEXT INTERVAL

[2, 4]

2,  $2 + 2^{-51}$ ,  $2 + 2 \times 2^{-51}$ , ..., 4

IN GENERAL, THE INTERVAL

$[2^j, 2^{j+1}]$  IS REPRESENTED AS  $2^j$  TIMES THE NUMBERS

REPRESENTING THE INTERVAL

[1, 2].

IN THE FLOATING POINT ③  
 REPRESENTATION THE GAPS  
 BETWEEN ~~SUCCESSIVE~~ ADJACENT  
 NUMBERS SCALE WITH THEIR  
 SIZE.

CALL SET OF FLOATING  
 POINT NUMBERS  $F \subset \mathbb{R}$ .

CALL  $\epsilon$  ("MACHINE EPSILON")  
 IS THE RESOLUTION OF THE  
 FLOATING POINT NUMBERS. AND  
 IS HALF THE DISTANCE ~~BETWEEN~~  
 BETWEEN 1 AND THE  
 ADJACENT NUMBER.

FOR DOUBLE PRECISION ④

$$\epsilon = 2^{-53} \approx 1.11 \times 10^{-16}$$

~~WEAKENED~~ AS A RESULT

$$\forall x \in \mathbb{R} \exists x' \in F$$

$$\text{s.t. } |x - x'| \leq \epsilon |x|$$

LET  $f_L : \mathbb{R} \rightarrow F$  BE THE  
 FUNCTION THAT ROUNDS  
 $x \in \mathbb{R}$  TO THE NEAREST  
 FLOATING POINT.

FLOATING POINT AXIOM I (FPA I)

$$\forall x \in \mathbb{R} \exists \epsilon' \text{ with}$$

$$|\epsilon'| \leq \epsilon \text{ s.t. } f_L(x) = x(1 + \epsilon')$$

## FLOATING POINT ARITHMETIC

$+, -, \times, \div$  on  $\mathbb{R}$

HAVE ANALOGUES

$\oplus, \ominus, \otimes, \oslash$  on  $\mathbb{F}$

CONSTRUCTED S.T.

$$x \otimes y = f_L(x \otimes^* y)$$

For  $x, y \in \mathbb{F}$

WHERE  $*$  IS  $+, -, \times$ , OR  $\div$ .

## FUNDAMENTAL AXIOM OF FLOATING POINT ARITHMETIC (FPA II)

$\forall x, y \in \mathbb{F} \exists \epsilon'$  WITH

$(\epsilon') \leq \epsilon$  S.T.

$$x \otimes y = (x * y)(1 + \epsilon')$$

(5)

## STABILITY

(6)

- STABILITY ~~PERTAINS~~ TO THE PERTURBATION BEHAVIOR OF THE ALGORITHM USED TO SOLVE THE PROBLEM ON A COMPUTER.

- ALGORITHM:  $\tilde{f}: X \rightarrow Y$  BETWEEN SAME SPACES AS THE PROBLEM.

- Fix:
- PROBLEM  $f$
  - FLOATING PT. COMPUTER
  - AN ALGORITHM FOR  $f$
  - IMPLEMENTATION OF THE ALGORITHM.

$x \in X$  is rounded  $x' = f_2(x)$  ⑦

THEN SUPPLIED TO THE  
PROGRAM.

THE PROGRAM IS RUN  
AND THE RESULT IS  
 $\tilde{f}(x) \in Y$ .

DESPITE THE COMPLEXITY,  
WE CAN MAKE CLEAN  
STATEMENTS ABOUT  $\tilde{f}(x)$   
USING FTA I & II.

### ACCURACY

. ABSOLUTE ERROR

$$\|\tilde{f}(x) - f(x)\|$$

### • RELATIVE ERROR

$$\frac{\|\tilde{f}(x) - f(x)\|}{\|f(x)\|}$$

AN ALGORITHM IS ACCURATE  
IF FOR EACH  $x \in X$

$$\frac{\|\tilde{f}(x) - f(x)\|}{\|f(x)\|} = O(\epsilon)$$

error is on the order  
of machine epsilon.

MORE PRECISELY,  $\exists$  CONSTANT C  
S.T.  $\forall x \in X$

$$\frac{\|\tilde{f}(x) - f(x)\|}{\|f(x)\|} \leq C\epsilon$$

AS  $\epsilon \rightarrow 0$ .

## STABILITY

An Algorithm  $\tilde{f}$  is  
For problem  $f$  is  
STABLE if For EACH

$x \in X$

$$\frac{\|\tilde{f}(x) - f(\tilde{x})\|}{\|f(\tilde{x})\|} = O(\epsilon)$$

For some  $\tilde{x} \in X$

$$\frac{\|\tilde{x} - x\|}{\|x\|} = O(\epsilon)$$

A STABLE ALGORITHM GIVES  
NEARLY THE RIGHT ANSWER  
TO NEARLY THE RIGHT QUESTION.

⑨

⑩

## BACKWARD STABILITY

THE ALGORITHM IS  
BACKWARD STABLE IF

For EACH  $x \in X$

$$\tilde{f}(x) = f(\tilde{x}) \quad \text{for some } \tilde{x}$$

$\tilde{x} \in X$  with

$$\frac{\|\tilde{x} - x\|}{\|x\|} = O(\epsilon)$$

A BACKWARD STABLE ALGORITHM  
GIVES EXACTLY THE RIGHT  
ANSWER TO NEARLY THE  
RIGHT QUESTION.

## STABILITY

An ALGORITHM  $\tilde{f}$  IS STABLE

IF FOR EACH  $x \in X$

$$\frac{\|\tilde{f}(x) - \tilde{f}(\tilde{x})\|}{\|\tilde{f}(\tilde{x})\|} = O(\epsilon)$$

FOR SOME  $\tilde{x} \in X$

$$\frac{\|\tilde{x} - x\|}{\|x\|} = O(\epsilon)$$

ACCURACY OF A BACKWARD  
STABLE ALGORITHM.

THM: SUPPOSE A BACKWARD  
STABLE ALGORITHM  $\tilde{f}$  IS  
APPLIED TO SOLVE A

①

PROBLEM  $f: X \rightarrow Y$  WITH  
CONDITION NUMBER  $k_f$  ON  
A COMPUTER SATISFYING  
FPA I & II. THEN THE  
RELATIVE ERROR SATISFIES

$$\frac{\|\tilde{f}(x) - f(x)\|}{\|f(x)\|} = O(k_f \epsilon)$$

PROOF: SINCE  $\tilde{f}$  IS BACKWARD  
STABLE,  $\tilde{f}(x) = f(\tilde{x})$  FOR SOME  
 $\tilde{x}$  WITH  $\frac{\|\tilde{x} - x\|}{\|x\|} = O(\epsilon)$ .

DEFINE

$$k_f(x) = \sup_{\frac{\|\delta x\|}{\|x\|} \leq \epsilon} \frac{\|\delta f\|}{\|f\|} / \frac{\|\delta x\|}{\|x\|}$$

Thus,

$$\frac{\|\tilde{f}(x) - f(x)\|}{\|f(x)\|} = \frac{\|f(\tilde{x}) - f(x)\|}{\|f(x)\|}$$

$$\leq k_{\epsilon}(x) \frac{\|\tilde{x} - x\|}{\|x\|}$$

TAKING THE LIMIT AS  $\epsilon \rightarrow 0$

$$\frac{\|\tilde{f}(x) - f(x)\|}{\|f(x)\|} = O(k(x) \delta)$$

### BACKWARDS ERROR ANALYSIS

INVESTIGATE ACCURACY USING  
THE CONDITIONING AND STABILITY.

IF OUR ALGORITHM IS  
BACKWARDS STABLE THEN

(3)

THE ERROR IS REFLECTED  
BY THE CONDITION NUMBER.

### FORWARD ERROR ANALYSIS

KEEP A RUNNING TALLY  
OF THE ~~TO~~ ERROR COMMITTED  
AT EACH STEP OR THE  
ALGORITHM.

STABILITY  $\approx$  FLOATING POINT  
OPERATIONS

$(+, -, \times, \div)$  ARE ALL  
BACKWARDS STABLE.

For SUBTRACTION:

DATA:  $x_1, x_2$

SOLUTION:  $f(x_1, x_2) = x_1 - x_2$

(4)

$$\text{ALGORITHM: } \tilde{f}(x_1, x_2) = f_L(x_1) \odot f_L(x_2)^{(5)}$$

BY FPA I:

$$f_L(x_1) = x_1(1 + \epsilon_1)$$

$$f_L(x_2) = x_2(1 + \epsilon_2)$$

WITH  $| \epsilon_1 |, | \epsilon_2 | \leq \epsilon$

BY FPA II:

$$f_L(x_1) \odot f_L(x_2)$$

$$= (f_L(x_1) - f_L(x_2))(1 + \epsilon_3)$$

WITH  $| \epsilon_3 | \leq \epsilon$

COMBINING THESE RESULTS.

$$f_L(x_1) \odot f_L(x_2)$$

$$= [x_1(1 + \epsilon_1) - x_2(1 + \epsilon_2)]$$

$$+ (1 + \epsilon_3)$$

$$= x_1(1 + \epsilon_1)(1 + \epsilon_3) - x_2(1 + \epsilon_2)(1 + \epsilon_3)$$

$$= \overbrace{x_1}^{\tilde{x}_1}(1 + \epsilon_4) - \overbrace{x_2}^{\tilde{x}_2}(1 + \epsilon_5)$$

$$\text{WITH } | \epsilon_4 |, | \epsilon_5 | \leq 2\epsilon + O(\epsilon^2)$$

$$\text{Thus, } \tilde{f}(x_1, x_2) = \tilde{x}_1 - \tilde{x}_2$$

$$= \tilde{f}(\tilde{x}_1, \tilde{x}_2)$$

$$\frac{|\tilde{x}_{1,2} - x_{1,2}|}{|x_{1,2}|} = O(\epsilon)$$

INNER PRODUCT

$$x, y \in \mathbb{R}^m, \quad \alpha = x^T y$$

ALGORITHM:

$$\tilde{\alpha} = \left( \sum_{i=1}^m f_L(x_i) \otimes f_L(y_i) \right)$$

THIS IS BACKWARDS STABLE.

$$= \left( \sum_{i=1}^m x_i (1 + \epsilon_i) \otimes y_i (1 + \delta_i) \right)$$

WITH  $|\epsilon_i|, |\delta_i| \leq \epsilon$

$$= \left( \sum_{i=1}^m [x_i (1 + \epsilon_i) \otimes y_i (1 + \delta_i)] (1 + \gamma_i) \right)$$

WITH  $|\gamma_i| \leq \epsilon$

(7)

$$= \left( \sum_{i=1}^m \right) [x_i (1 + \epsilon_i) y_i (1 + \tilde{\epsilon}_i)]$$

WITH  $|\tilde{\epsilon}_i| \leq 2\epsilon + O(\epsilon^2)$

(8)

TAKE THE FIRST TWO TERMS

$$x_1 (1 + \epsilon_1) y_1 (1 + \tilde{\epsilon}_1) + x_2 (1 + \epsilon_2) y_2 (1 + \tilde{\epsilon}_2)$$

$$= [x_1 (1 + \epsilon_1) y_1 (1 + \tilde{\epsilon}_1) + x_2 (1 + \epsilon_2) y_2 (1 + \tilde{\epsilon}_2)] (1 + \beta_1)$$

WITH  $|\beta_1| \leq \epsilon$

$$= x_1 (1 + \epsilon_1) y_1 (1 + \hat{\epsilon}_1) + x_2 (1 + \epsilon_2) y_2 (1 + \hat{\epsilon}_2)$$

WITH  $|\hat{\epsilon}_1|, |\hat{\epsilon}_2| \leq 3\epsilon + O(\epsilon^2)$

Following the same process  
with the remaining terms

$$\tilde{\alpha} = \sum_{i=1}^m \underbrace{x_i(1+\epsilon_i)}_{\tilde{x}_i} \underbrace{y_i(1+\tilde{\delta}_i)}_{\tilde{y}_i}$$

$$|\tilde{\delta}_i| \leq (m+1)\epsilon + O(\epsilon^2)$$

WE HAVE THAT

$$\tilde{\alpha}(x, y) = \alpha(\tilde{x}, \tilde{y})$$

For  $\left\| \frac{\tilde{x} - x}{\|x\|} \right\| = O(\epsilon) \text{ Ans}$

$$\frac{\|\tilde{y} - y\|}{\|y\|} = O(\epsilon).$$

(9)

STABLE BUT NOT B. STABLE

$$f(x) = x + 1$$

$$\tilde{f}(x) = f(L(x)) + 1$$

CAN SHOW

$$\tilde{f}(x) = x(1+\epsilon_3) + 1 + \epsilon_2 =$$

UNSTABLE ALGORITHM

USING  $\det(\lambda I - A) = 0$  TO  
FIND THE EIGENVALUES OF  
 $A$ .

WHAT DOES STABILITY MEAN  
IN THE CONTEXT OF LINEAR  
ALGEBRA COMPUTATIONS?

PROBLEM: QR FACTORIZATION OF  
A.

INPUT OR DATA: A

OUTPUT OR SOLUTION: Q, R

WITH  $Q R = A$  AND Q IS  
UNITARY AND R UPPER TRIANGULAR.

ALGORITHM: HOUSEHOLDER ON  
A COMPUTER SATISFYING  
FPA I & II.

INPUT: A, OUTPUT:  $\tilde{Q}, \tilde{R}$

①

IF BACKWARDS STABLE ②  
THEN  $\tilde{Q}, \tilde{R}$  ARE SOLUTIONS  
TO THE PROBLEM WITH INPUT  
 $A + \delta A$  FOR SOME  $\frac{\|\delta A\|}{\|A\|} = O(\epsilon)$

SINCE THIS IS A SOLUTION  
WE KNOW

$$\tilde{Q} \tilde{R} = A + \delta A$$

SO IN OUR EXPERIMENT, WE  
WERE MEASURING  $\frac{\|\delta A\|}{\|A\|}$ .

TYPICALLY THIS ERROR IS  
CALLED BACKWARD ERROR OR  
THE RESIDUAL.

<p><u>THM:</u> LET THE QR FACTORIZATION<sup>(3)</sup></p> <p><math>A = QR</math> OR <math>A \in \mathbb{C}^{n \times n}</math> BE COMPUTED BY HOUSEHOLDER IN A COMPUTER SATISFYING FPA I &amp; II.</p>	<p>TYPICALLY WE WOULD LIKE<sup>(4)</sup> TO FACTOR OUR MATRIX TO DO SOMETHING, E.G. SOLVE A LINEAR SYSTEM.</p>
<p>THEN IF WE HAVE</p> $\tilde{Q}\tilde{R} = A + \delta A \quad \frac{\ \delta A\ }{\ A\ } = O(\epsilon)$ <p>FOR SOME <math>\delta A \in \mathbb{C}^{n \times n}</math></p> <p>WHAT IS MEANT BY <math>\tilde{R}</math> AND <math>\tilde{Q}</math>?</p> <p><math>\tilde{R}</math> IS THE TRIANGULAR MATRIX THAT GETS COMPUTED.</p> <p><math>\tilde{Q}</math> IS THE <u>EXACTLY</u> UNITARY MATRIX FORMED USING THE COMPUTED VECTORS <math>\tilde{v}_k</math>.</p>	<p>IS BACKWARD STABILITY ENOUGH OR DO WE NEED ACCURATE <math>Q</math> AND <math>R</math>?</p> <p>SOLVING <math>Ax = b</math> USING QR ALGORITHM</p> <ol style="list-style-type: none"> <li>1. <math>QR = A</math> DONE USING HOUSEHOLDER.</li> <li>2. <math>y = Q^* b</math> CONSTRUCT <math>Q^* b</math> USING IMPLICIT MULTIPLICATION.</li> </ol>

$$3. \quad x = R^{-1}y$$

SOLVING THE  
TRIANGULAR SYSTEM  
USING BACKWARD  
SUBSTITUTION.

(5)

WE'VE ALREADY DISCUSSED THE  
BACKWARD STABILITY OF STEP 1.

AND THAT IT ~~gives~~ OUTPUTS

$$\tilde{Q}, \tilde{R}.$$

STEP 2 ~~is~~ IS ALSO BACKWARD  
STABLE. USING  $\tilde{Q}$  GIVEN  
BY STEP 1. THIS MEANS.

$$(i) \quad (\tilde{Q} + \delta Q) \tilde{y} = b \quad \| \delta Q \| = O(\epsilon)$$

STEP 3 IS ALSO BACKWARD  
STABLE

$$(ii) \quad (\tilde{R} + \delta R) \tilde{x} = \tilde{y}$$

$$\text{WITH } \frac{\| \delta R \|}{\| R \|} = O(\epsilon).$$

PROVING THESE RESULTS IS VERY  
TEDIOUS!!! SEE LECTURE 17  
IN TREFETHEN & BAU.

THEOREM: THE QR ALGORITHM  
TO SOLVE  $Ax = b$  IS BACKWARD  
STABLE, SATISFYING

$$(A + \Delta A) \tilde{x} = b$$

$$\text{FOR SOME } \frac{\| \Delta A \|}{\| A \|} = O(\epsilon).$$

Proof: Combining (i) + (ii)

(7)

$$b = (\tilde{Q} + \delta Q)(\tilde{R} + \delta R) \tilde{x}$$

$$= [\tilde{Q}\tilde{R} + (\delta Q)\tilde{R} + \tilde{Q}(\delta R) \\ + (\delta Q)(\delta R)] \tilde{x}$$

From step 1. we know

$$\tilde{Q}\tilde{R} = A + \delta A$$

$$b = A\tilde{x} + [\delta A + (\delta Q)\tilde{R} + \tilde{Q}(\delta R) \\ + (\delta Q)(\delta R)]\tilde{x}$$

$\Delta A$

Show EACH TERM IS O( $\epsilon$ ) TO GET RESULT.

(8)

$$\text{SINCE } \tilde{Q}\tilde{R} = A + \delta A$$

$$\Rightarrow \tilde{R} = \tilde{Q}^*(A + \delta A)$$

$$\frac{\|\tilde{R}\|}{\|A\|} \leq \frac{\|\tilde{Q}^*\| \|A + \delta A\|}{\|A\|}$$

$$\leq 1 + \frac{\|\delta A\|}{\|A\|} = O(1)$$

Thus

$$\frac{\|\delta Q \tilde{R}\|}{\|A\|} \leq \|\delta Q\| \frac{\|\tilde{R}\|}{\|A\|} = O(\epsilon)$$

$$\frac{\|\tilde{Q}(\delta R)\|}{\|A\|} \leq \frac{\|\delta R\|}{\|\tilde{R}\|} \|\tilde{Q}\| \frac{\|\tilde{R}\|}{\|A\|} = O(\epsilon)$$

$$\frac{\|\delta Q - \delta R\|}{\|A\|} \leq \|\delta Q\| \left( \frac{\|\delta R\|}{\|\tilde{R}\|} + \frac{\|\tilde{R}\|}{\|A\|} \right) = O(\epsilon^2)$$

(9)

Thus  $\frac{\|\Delta A\|}{\|A\|} = O(\epsilon)$ .

## GAUSSIAN ELIMINATION

### EXAMPLE

$$A = \begin{bmatrix} 2 & 1 & 1 & 0 \\ 4 & 3 & 3 & 1 \\ 8 & 7 & 9 & 5 \\ 6 & 7 & 9 & 8 \end{bmatrix}$$

$$L_1 A = \begin{bmatrix} 1 & & & \\ -2 & 1 & & \\ -4 & & 1 & \\ -3 & & & 1 \end{bmatrix} \quad \left[ \begin{array}{c|c} 1 & \\ \hline -2 & 1 \\ -4 & \\ -3 & \end{array} \right]$$

$$= \begin{bmatrix} 2 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 \\ 3 & 5 & 5 & \\ 4 & 6 & 8 & \end{bmatrix}$$

$$L_2 L_1 A = \begin{bmatrix} 1 & & & \\ -\frac{1}{2} & 1 & & \\ -4 & & 1 & \end{bmatrix} \quad \left[ \begin{array}{c|c} 1 & \\ \hline -\frac{1}{2} & 1 \\ -4 & \end{array} \right]$$

$$\textcircled{1} \quad N = \begin{bmatrix} 2 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 \\ 2 & 2 & 2 & \\ 2 & 4 & & \end{bmatrix}$$

$$L_3 L_2 L_1 A = \begin{bmatrix} 1 & & & \\ & 1 & & \\ & & 1 & \\ & & & -1 \end{bmatrix} \quad \left[ \begin{array}{c|c} 1 & \\ \hline & 1 \\ & \\ & -1 \end{array} \right]$$

$$N = \begin{bmatrix} 2 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 \\ 0 & 0 & 2 & \\ 0 & 0 & 2 & \end{bmatrix} = U$$

WHAT ABOUT  $L$ ? IF  $A = LU$

$$L^{-1} = L_3 L_2 L_1$$

$$\Rightarrow L = L_1^{-1} L_2^{-1} L_3^{-1}$$

$$L_1^{-1} = \begin{bmatrix} 1 & & & \\ -2 & 1 & & \\ -3 & -4 & 1 & \\ -4 & 3 & 1 & 1 \end{bmatrix}^{-1}$$

$$= \begin{bmatrix} 1 & & & \\ 2 & 1 & & \\ 4 & 1 & 1 & \\ 3 & 1 & 1 & \end{bmatrix}$$

For  $L_k^{-1}$  just replace  
FACTOR VALUES BELOW  $\lambda_{kk}$   
 BY  $-1$  TIMES THE VALUE.

$$L = L_1^{-1} L_2^{-1} L_3^{-1}$$

$$= \begin{bmatrix} 1 & & & \\ 2 & 1 & & \\ 4 & 3 & 1 & \\ 3 & 4 & 1 & 1 \end{bmatrix}$$

③

TO GET  $L$ , SIMPLY NEED  
 TO COMPILE ALL THE SUBDIAGONAL  
 ENTRIES OF  $L_k^{-1}$  IN ONE  
 NON-ZERO MATRIX!

RECALL:

MGS WE SAID WAS  
 TRIANGULAR ORTHOGONALIZATION

HOUSEHOLDER WE SAID WAS  
 ORTHOGONAL TRIANGULARIZATION.

GAUSSIAN ELIMINATION IS  
 TRIANGULAR TRIANGULARIZATION.

## GENERAL FORMULA

LET  $x_k$  BE THE  $k^{\text{TH}}$  COLUMN OF THE MATRIX AT THE  $k^{\text{TH}}$  STEP.

$$x_k = \begin{bmatrix} x_{1k} \\ \vdots \\ x_{kk} \\ x_{k+1,k} \\ \vdots \\ x_{mnk} \end{bmatrix}$$

WOULD LIKE  $L_k$  SUCH THAT

$$L_k x_k = \begin{bmatrix} x_{1k} \\ \vdots \\ x_{kk} \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

(5)

TO DO THIS, WE MUST HAVE

$$l_{jk} = \frac{x_{jk}}{x_{kk}} \quad \text{for } j=k+1, \dots, m$$

(6)

SO THAT

$$L_k = \begin{bmatrix} 1 & & & & \\ & 1 & & & \\ & & 1 & & \\ & -l_{k+1,k} & & 1 & \\ & \vdots & & & \\ & -l_{mk} & & & 1 \end{bmatrix}$$

THEN CAN WRITE

$$L_k = I - l_k e_k^*$$

WHERE

$$l_k = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ l_{k+1,k} \\ \vdots \\ l_{m,k} \end{bmatrix}$$

AND  $e_k$  IS THE VECTOR

WITH KTH ENTRY 1 AND

0 ELSEWHERE.

SINCE,  $e_k^* l_k = 0$

$$(I - l_k e_k^*)(I + l_k e_k^*) = I$$

Thus,

$$L_k^{-1} = I + l_k e_k^*$$

(7)

$$\begin{aligned} L_k^{-1} L_{k+1}^{-1} &= (I + l_k e_k^*)(I + l_{k+1} e_{k+1}^*) \\ &= I + l_k e_k^* + l_{k+1} e_{k+1}^* \end{aligned}$$

(8)

$$\text{SINCE } e_k^* l_{k+1} = 0.$$

Thus,

$$L = L_1^{-1} L_2^{-1} \cdots L_{m-1}^{-1}$$

$$= \begin{bmatrix} 1 & & & & \\ l_{21} & 1 & & & \\ \vdots & l_{32} & \ddots & & \\ \vdots & \vdots & \ddots & \ddots & \\ l_{m1} & l_{m2} & \cdots & l_{m,m-1} & 1 \end{bmatrix}$$

# GAUSSIAN ELIMINATION ALGORITHM

$A \in \mathbb{C}^{m \times m}$  (SQUARE)

⑨

$U = A$

$L = I$

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

for  $j = k+1$  to  $m$

$$l_{jk} = \frac{u_{jk}}{u_{kk}}$$

$$(*) u_{j,k:m} = u_{j,k:m} - l_{jk} u_{k,k:m}$$

END FOR

END FOR

## OPERATION COUNT

WORK DOMINATED BY LINE (\*)

CONSISTS OF:

(i)  $m-k+1 : x$

(ii)  $m-k+1 : -$

2( $m-k+1$ ) FLOPS EVERY TIME ⑩  
 (\*) IS REACHED.

$$\# \text{ OF FLOPS} \sim \sum_{k=1}^{m-1} \sum_{j=k+1}^m 2(m-k+1)$$

$$= \sum_{k=1}^{m-1} 2(m-k+1) \underbrace{\sum_{j=k+1}^m 1}_{m-k}$$

~~$\#$~~   $\sum_{k=1}^{m-1} 2(m-k+1)(m-k)$

$$\sim \sum_{k=1}^{m-1} 2m^2 - 4mk + 2k^2$$

$$\sim 2m^3 - \frac{4m^3}{2} + \frac{2}{3}m^3$$

$$\# \text{ OF FLOPS} \sim \frac{2}{3}m^3$$

# GAUSSIAN ELIMINATION WITH PIVOTING

SUPPOSE  $A = \begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix}$

STANDARD GAUSSIAN ELIMINATION WILL FAIL AT THE FIRST STEP.

ALSO, PROBLEMATIC IF

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

SINCE DIVIDING BY  $10^{-20}$  COULD LEAD TO ~~numerical~~ NUMERICAL INSTABILITIES.

## PIVOTS

(2)

AT STEP  $k$  HAVE MATRIX

X

$$\left[ \begin{array}{cccc} x & x & x & x \\ x_{kk} & x & x \\ x & x & x \\ x & x & x \end{array} \right] \xrightarrow{L_k} \left[ \begin{array}{cccc} x & x & x & x \\ x_{kk} & x & x \\ 0 & x & x \\ 0 & x & x \end{array} \right]$$

ENTRY  $x_{kk}$  IS CALLED THE PIVOT.

EVERY ENTRY OF  $X_{k+1:m, k:m}$  IS SUBTRACTED BY THE PRODUCT OF AN ENTRY IN ROW  $\cancel{k}$  AND AN ENTRY IN COLUMN  $k$  DIVIDED BY  $x_{kk}$ .

IF THE GOAL IS TO INTRODUCE ZEROS IN A COLUMN, THE CHOICE OF  $x_{kk}$  IS NOT UNIQUE.

TAKE  $x_{ik}$  AS PIVOT:

$$\begin{bmatrix} x & x & x & x \\ x & x & x & \\ x_{ik} & x & x & \\ x & x & x & \end{bmatrix} \rightarrow \begin{bmatrix} x & x & x & x \\ 0 & x & x & \\ x_{ik} & x & x & \\ 0 & x & x & \end{bmatrix}$$

CAN EVEN TAKE  $x_{ij}$

$$\begin{bmatrix} x & x & x & x \\ x & x & x & \\ x & x & x_{ij} & \\ x & x & x & \end{bmatrix} \rightarrow \begin{bmatrix} x & x & x & x \\ x & x & 0 & \\ x & x & x_{ij} & \\ x & x & 0 & \end{bmatrix}$$

(3)

- FREE TO CHOOSE ANY ENTRY IN  $X_{k:m, k:m}$  AS THE PIVOT (4)
- FOR STABILITY, CHOOSE THE LARGEST ENTRY IN  $X_{k:m, k:m}$ .
- CHOOSING ENTRY OTHER THAN  $x_{kk}$  WILL DESTROY THE INDUCED TRIANGULAR STRUCTURE
- ADD ADDITIONAL STEP OF PERMUTING THE ROWS AND COLUMNS TO MOVE CHOSEN PIVOT INTO  $x_{kk}$ 'S POSITION.

## COMPLETE PIVOTING

- Allow all entries of  $x_{k:m, k:m}$  to be a pivot.

$$\# \text{ OF FLOPS} = \mathcal{O}(m^3)$$

## PARTIAL PIVOTING

- Choose pivot from the  $k$ th column
- ONLY row exchanges are involved.
- # of FLOPS =  $\mathcal{O}(m^2)$

(5)

IDEA:

$$\begin{bmatrix} x & x & x & x \\ x & x & x & x \\ x_{ik} & x & x & x \\ x & x & x & x \end{bmatrix} \xrightarrow{P_k} \begin{bmatrix} x & x & x & x \\ x_{ik} & x & x & x \\ x & x & x & x \\ x & x & x & x \end{bmatrix}$$

SELECT THE PIVOT.

$$\xrightarrow{L_k} \begin{bmatrix} x & x & x & x \\ x_{ik} & x & x & x \\ 0 & x & x & x \\ 0 & x & x & x \end{bmatrix}$$

ELIMINATION STEP.

## EXAMPLE

$$A = \begin{bmatrix} 2 & 1 & 1 & 0 \\ 4 & 3 & 3 & 1 \\ 8 & 7 & 9 & 5 \\ 6 & 7 & 9 & 8 \end{bmatrix}$$

(6)

STEP 1:  $P_1 A$

$$\begin{bmatrix} 1 & & 1 \\ & 1 & \\ & & 1 \end{bmatrix} A = \begin{bmatrix} 8 & 7 & 9 & 5 \\ 4 & 3 & 3 & 1 \\ 2 & 1 & 1 & 0 \\ 6 & 7 & 9 & 8 \end{bmatrix}$$

(7)

EVENTUALLY OBTAIN DECOMPOSITION  
THAT MAY BE WRITTEN AS

$$\begin{bmatrix} 1 & & 1 \\ & 1 & \\ & & 1 \end{bmatrix} A = \begin{bmatrix} 1 & & & \\ \frac{3}{4} & 1 & & \\ \frac{1}{2} & -\frac{2}{7} & 1 & \\ \frac{1}{4} & -\frac{3}{7} & \frac{1}{3} & 1 \end{bmatrix}$$

(8)

STEP 2:  $L_1 P_1 A$

$$\begin{bmatrix} 1 & & & \\ -\frac{1}{2} & 1 & & \\ -\frac{1}{4} & & 1 & \\ -\frac{3}{4} & & & 1 \end{bmatrix} P_1 A = \begin{bmatrix} 8 & 7 & 9 & 5 \\ -\frac{1}{2} & -\frac{3}{2} & -\frac{3}{2} & \\ -\frac{3}{4} & -\frac{5}{4} & -\frac{5}{4} & \\ \frac{7}{4} & \frac{9}{4} & \frac{17}{4} & \end{bmatrix}$$

$P$

$$L \times \begin{bmatrix} 8 & 7 & 9 & 5 \\ \frac{7}{4} & \frac{9}{4} & \frac{17}{4} & \\ -\frac{6}{7} & -\frac{2}{7} & \\ & & \frac{2}{3} \end{bmatrix}$$

STEP 3:  $P_2 L_1 P_1 A$

$$\begin{bmatrix} 1 & & 1 \\ & 1 & \\ & & 1 \end{bmatrix} L_1 P_1 A =$$

$PA = LU$  IS THE  
TRUE LU DECOMPOSITION.

## INTERPRETATION:

1. PERMUTE rows of A according P.

2. APPLY GAUSSIAN ELIMINATION WITHOUT PIVOTING TO PA.

WHY IS THIS THE CASE?

WE DID:

$$L_3 P_3 L_2 P_2 L_1 P_1 A = U$$

THESE OPERATIONS CAN BE REORDERED

$$L_3 P_3 L_2 P_2 L_1 P_1 = (L_3' L_2' L_1') (P_3 P_2 P_1)$$

$$L^{-1} \quad P$$

(9)

WHERE

$$L_3' = L_3 \quad L_2' = P_3 L_2 P_3^{-1}$$

$$L_1' = P_3 P_2 L_1 P_2^{-1} P_3^{-1}$$

PERMUTATION MATRICES PRESERVE LOWER TRIANGULAR FORM.

(10)

## LU DECOMPOSITION

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

HAVE WITH PARTIAL PIVOTING.

$$L_{m-1} P_{m-1} \cdots L_2 P_2 L_1 P_1 A = U$$

THIS CAN ALSO BE WRITTEN

AS

$$(L'_{m-1} \cdots L'_2 L'_1) \underbrace{(P_{m-1} \cdots P_2 P_1)}_P A = U$$

$$\underbrace{L'}_{L^{-1}}$$

$$L'_k = P_{m-1} \cdots P_{k+1} L_k P_{k+1}^{-1} \cdots P_{m-1}^{-1}$$

(1)

Thus

$$PA = LU$$

Algorithm For LU

$$U = A$$

$$L = I$$

$$P = I$$

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

SELECT  $i \geq k$  to } SELECT  
MAXIMIZE  $|u_{ik}|$  }

$u_{k,k:m} \longleftrightarrow u_{i,k:m}$  } EXCHANGE  
 $l_{k,1:k-1} \longleftrightarrow l_{i,1:k-1}$  } OF  
 $P_{k,:} \longleftrightarrow P_{i,:}$  } ROWS.

for  $j = k+1$  to  $m$

$$l_{jk} = u_{jk} / u_{kk}$$

(2)

~~END~~

$$u_{j,k:m} = u_{j,k:m} - l_{jk} u_{k,k:m}$$

END FOR

END FOR.

### OPERATION COUNT

$$\# \text{ OF FLOPS} \approx \frac{2}{3} m^3$$

### STABILITY OF LU

CONSIDER AGAIN

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

STANDARD GAUSSIAN ~~ELIM~~  
 ELIMINATION LEADS TO  
 LARGE ENTRIES OF L AND  
 U.

(3)

THIS AMPLIFICATION IS AT THE HEART OF THE INSTABILITY?

(4)

THEM:  $A = LU$  FOR NONSINGULAR  
 $A \in \mathbb{C}^{m \times m}$  BE COMPUTED BY  
 STANDARD GAUSSIAN ELIMINATION  
 ON A COMPUTER SATISFYING  
 FPA I AND FPA II. IF  
 NO ZERO-PIVOTS ARE ENCOUNTERED  
 THEN COMPUTED  $\tilde{L}$  AND  $\tilde{U}$   
 SATISFY

$$\tilde{L} \tilde{U} = A + \delta A \quad \frac{\|\delta A\|}{\|L\| \|U\|} = O(\epsilon)$$

FOR SOME  $\delta A \in \mathbb{C}^{m \times m}$

BACKWARD STABILITY IF  
 $\|L\| \|U\| = O(\|A\|)$

BUT IF  $\|L\| \|U\| \neq O(\|A\|)$  (5)  
CAN EXPECT BACKWARD INSTABILITY.

PIVOTING

THIS APPLIES REPLACING  
A BY PA.

WITH PARTIAL PIVOTING  
WE KNOW

$$\|L\| = O(1)$$

THEN FOR BACKWARD  
STABILITY, WE MUST HAVE

$$\|U\| = O(\|A\|).$$

THIS CAN BE MEASURED  
BY A GROWTH FACTOR.

(5)

$$\rho = \frac{\max_{i,j} |u_{ij}|}{\max_{i,j} |a_{ij}|} \quad (6)$$

SO THEN ~~WE~~ WE HAVE

$$\|U\| = O(\rho \|A\|)$$

THIS:  $PA = LU$  BE COMPUTED  
USING GAUSSIAN ELIMINATION  
WITH PIVOTING ON A COMPUTER  
SATISFYING FPA I AND FPA II.

THEN  $\tilde{L}$ ,  $\tilde{U}$ , AND  $\tilde{P}$  SATISFY

$$\tilde{L} \tilde{U} = \tilde{P}A + \delta A$$

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

FOR SOME  $\delta A \in \mathbb{C}^{n \times n}$  WHERE  $\rho$   
IS THE GROWTH FACTOR.

EXAMPLE:

$$A = \begin{bmatrix} 1 & & & & 1 \\ -1 & 1 & & & 1 \\ -1 & -1 & 1 & & 1 \\ -1 & -1 & -1 & 1 & 1 \\ -1 & -1 & -1 & -1 & 1 \end{bmatrix}$$

(7)

This is huge for  
large matrix!

(8)

AFTER EACH STEP  $k$ , THE  
 $k+1$  THROUGH  $m$  ENTRIES  
IN THE LAST COLUMN WILL  
BE  $2^k$ .

$$u = \begin{bmatrix} 1 & & & 1 \\ & 1 & & 2 \\ & & 1 & 4 \\ & & & 8 \\ & & & 16 \end{bmatrix}$$

ENTRY  $u_{mn}$  IS  $2^{m-1}$

Thus,  $\rho = 2^{m-1}$

# EIGENVALUES AND EIGENVECTORS ①

$A \in \mathbb{C}^{m \times m}$  (square)

$x \in \mathbb{C}^m$  EIGENVECTOR AND

$\lambda \in \mathbb{C}$  CORRESPONDING EIGENVALUE

IF

$$Ax = \lambda x$$

EIGENSPACE: SET OF THE  
EIGENVECTORS

SPECTRUM: SET OF ALL  
EIGENVALUES.

## WHY COMPUTE THIS

1. E-VALS AND E-VECS

ENCODE INFORMATION ABOUT  
A.

2. ARISE IN SAY, STABILITY  
CALCULATIONS.

3. USE E-VECS AS A  
BASIS TO SOLVE A  
LINEAR SYSTEM.

## EIGENVALUE DECOMPOSITION

$$A = X \Lambda X^{-1}$$

X: NON-SINGULAR  $m \times m$  WITH  
 $x_k$  (EIGENVECTOR k) IN THE  
k<sup>TH</sup> COLUMN.

$\lambda$ : DIAGONAL MATRIX WITH

$A_{kk} = \lambda_k$ , THE  $k^{\text{TH}}$  EIGENVALUE.

THIS DECOMPOSITION IS NOT  
ALWAYS POSSIBLE!

GEOMETRIC MULTIPLICITY OF  $\lambda$   
NUMBER OF E-VECS  
CORRESPONDING TO  $\lambda$   
EIGENVALUE  $\lambda$ .

ALGEBRAIC MULTIPLICITY  
CHARACTERISTIC POLYNOMIAL OF  
 $A$

$$P_A(z) = \det(zI - A)$$

③

THE ROOTS OF  $P_A(z)$   
ARE THE EIGENVALUES.

$$P_A(z) = (z - \lambda_1)(z - \lambda_2) \dots (z - \lambda_m)$$

④

THE ALGEBRAIC MULTIPLICITY  
IS THE NUMBER OF TIMES  
 $\lambda$  IS A ROOT OF  $P_A(z)$ .

A E-VAL IS SIMPLE IF IT  
HAS ALGEBRAIC MULTIPLICITY 1.

THM: THE ALGEBRAIC MULTIPLICITY  
OF  $\lambda$  IS AT LEAST AS  
GREAT AS ITS GEOMETRIC  
MULTIPLICITY.

# DEFECTIVE MATRICES AND EIGENVALUES

(5)

- AN EIGENVALUE IS DEFECTIVE IF ITS ALGEBRAIC MULTIPLICITY EXCEEDS ITS GEOMETRIC MULTIPLICITY.
- A MATRIX IS DEFECTIVE IF IT HAS AT LEAST ONE DEFECTIVE E-VAL.

THM: An  $m \times m$  MATRIX  $A$  IS NON-DEFECTIVE IFF IT HAS AN EIGENVALUE DECOMPOSITION

$$A = X \Lambda X^{-1}$$

# SIMILARITY TRANSFORMATION

(6)

For  $X \in \mathbb{C}^{m \times m}$  NONSINGULAR THEN THE MAP  $A \mapsto X^{-1}AX$  IS A SIMILARITY TRANSFORMATION OF  $A$ .

MATRICES  $A$  &  $B$  ARE SIMILAR IF  $\exists X$  S.T.  
 $B = X^{-1}AX$ .

THM: IF  $X$  IS NON-SINGULAR THEN  $A$  AND  $X^{-1}AX$  HAVE SAME CHARACTERISTIC POLYNOMIAL, EIGENVALUES, AND MULTIPLICITIES.

## SCHUR FACTORIZATION

$$A = Q \underset{\text{UNITARY MATRIX}}{\underset{\swarrow}{T}} Q^*$$

UPPER-TRIANGULAR.

THM: EVERY <sup>SQUARE</sup> MATRIX HAS A SCHUR FACTORIZATION.

SINCE THIS IS A SIMILARITY TRANSFORMATION,  $A \leftrightarrow T$  HAVE THE SAME E-VALS.

SINCE  $T$  IS TRIANGULAR THE E-VALS ARE THE DIAGONAL ENTRIES!

⑦

## UNITARY DIAGONALIZATION

⑧

MATRIX  $A$  IS UNITARY DIAGONALIZABLE IF  $\exists$  UNITARY  $Q$  S.T.

$$A = Q \Lambda Q^*$$

THM: A HAMILTONIAN MATRIX ( $A = A^*$ ) IS THE UNITARY DIAGONALIZABLE AND ITS E-VALS ARE REAL.

## EIGENVALUE ALGORITHMS

~~ALREADY~~

ALREADY DISCUSSED THAT A  
ROOT FINDING  $\Leftrightarrow$  ALGORITHM

FOR  $P_A(z)$  WOULD BE  
UNSTABLE.

THIS CONNECTION WITH ROOT  
FINDING, HOWEVER, TELLS  
US THAT WE SHOULD NOT  
BE ABLE TO FIND THE  
EIGENVALUES IN A FINITE  
NUMBER OF STEPS.

⑨

ANY EIGENVALUE SOLVER  
MUST BE ITERATIVE!

⑩

COMPUTES SCHUR FACTORIZATION  
ITERATIVELY.

GENERATE A SEQUENCE OF  
UNITARY  $Q_j$  SUCH THAT

$$Q_j^k \cdots Q_2^* Q_1^* A Q_1 Q_2 \cdots Q_j$$

CONVERGES TO  $T$  AS  
 $j \rightarrow \infty$ .

IF  $A$  IS HERMITIAN,  $T$   
WILL BE DIAGONAL.

# TWO PHASES OF E-VAL

COMPUTATION

$$\begin{bmatrix} x & x & x & x \\ x & x & x & x \\ x & x & x & x \\ x & x & x & x \end{bmatrix} \xrightarrow[\text{DIRECT}]{1} \begin{bmatrix} x & x & x & x \\ x & x & x & x \\ x & x & x & x \\ x & x & x & x \end{bmatrix}$$

$$2 \swarrow \begin{matrix} H \text{ (HESSENBERG} \\ \text{MATRIX)} \end{matrix}$$

ITERATIVE

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

T

# HESSENBERG MATRIX

- PHASE 1 OF THE EIGENVALUE COMPUTATION.

$$H = \begin{bmatrix} x & x & x & x \\ x & x & x & x \\ x & x & x & \\ x & x & & \end{bmatrix}$$

- For HERMITIAN MATRICES  
THE ASSOCIATED HESSENBERG MATRIX WILL BE TRIDIAGONAL

$$\Delta = \begin{bmatrix} x & x & & & \\ x & x & x & & \\ x & x & x & x & \\ & x & x & x & \end{bmatrix}$$

(1)

TO: FORM  $H$  FROM  $A$  (2)  
WE'LL NEED TO USE  
A SIMILARITY TRANSFORMATION  
TO PRESERVE THE E-VALS.  
WE CAN USE HOUSEHOLDER  
REFLECTORS TO CONSTRUCT  
 $H$  FROM  $A$ .

$$\begin{bmatrix} x & x & x & x \\ x & x & x & x \\ x & x & x & x \\ x & x & x & x \end{bmatrix} \xrightarrow{Q_1^*} \begin{bmatrix} x & x & x & x \\ -\bar{x} & \bar{x} & \bar{x} & \bar{x} \\ 0 & x & x & x \\ 0 & -\bar{x} & \bar{x} & \bar{x} \end{bmatrix}$$

$A$

$Q_1^* A$

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

$Q_1^* A Q_1$

For LARGE  $m$

$$\# \text{ of FLOPS} \sim 4m(m-k)$$

For (ii)

~~This for  $m$  times~~

$$\# \text{ of FLOPS} \sim 4 \sum_{k=1}^{m-2} m(m-k)$$

$$\sim 2m^3$$

$$\text{TOTAL COST} = \frac{10}{3}m^3$$

IF  $A$  IS A HERMITIAN MATRIX, THE WORK IS REDUCED TO  $\sim \frac{4}{3}m^3$  DUE TO

(a) SYMMETRY

(b) SPARSITY.

(3)

This ALGORITHM IS BACKWARD STABLE.

(4)

RAYLEIGH QUOTIENT AND INVERSE ITERATION

• CLASSICAL E-VAL ALGORITHMS.

• FROM THIS POINT ON,  
RESTRICT TO REAL,  
SYMMETRICAL MATRICES.

AND TAKE  $\| \cdot \| = \| \cdot \|_2$

• EIGENVALUES ARE ALL REAL,  $\lambda_1, \dots, \lambda_m$

• ORTHONORMAL E-VECS

$$g_1, \dots, g_m.$$

### RAYLEIGH QUOTIENT

$$r(x) = \frac{x^T A x}{x^T x}$$

NOTICE IF  $x$  IS AN  
E-VEC THEN  $r(x) = \lambda$   
WHERE  $\lambda$  IS THE E-VAL.

VIEW  $x \in \mathbb{R}^m$  AS VARIABLE  
AND  $r: \mathbb{R}^m \rightarrow \mathbb{R}$ .

TAKE PARTIAL DERIVATIVE OF  
r w.r.t.  $x_j$

$$\frac{\partial r}{\partial x_j} = \frac{1}{x^T x} \frac{\partial}{\partial x_j} (x^T A x)$$

(5)

$$\begin{aligned} & - \frac{(x^T A x)}{(x^T x)^2} \frac{\partial}{\partial x_j} (x^T x) \\ &= \frac{2}{x^T x} (Ax - r(x)x)_j \end{aligned} \quad (6)$$

COLLECT THE PARTIAL DERIVATIVES  
FOR  $j = 1, \dots, m$  INTO ONE  
M VECTOR

$$\nabla r(x) = \frac{2}{x^T x} (Ax - r(x)x)$$

Thus,  $\nabla r = 0$  IF  $x$  IS  
AN E-VEC.

E-VECTORS ARE STATIONARY  
POINTS OF  $r(x)$ .

WHEN WE RESTRICT TO (8)  
 $\|x\|=1$ ,  $r(x)$  IS A  
 CONTINUOUS FUNCTION ON THE  
 UNIT SPHERE.

IF  $g_J$  IS AN E-VEC,

SINCE  $\nabla r(g_J) = 0$  AND

$r$  IS SMOOTH, TAYLOR

EXPANSION REVEALS

$$r(x) - r(g_J) = O(\|x - g_J\|^2)$$

AS  $x \rightarrow g_J$

RAYLEIGH QUOTIENT GIVES  
QUADRATICALLY ACCURATE  
 ESTIMATE OF THE E-VAL.

### POWER ITERATION

METHOD TO FIND E-VEC  
~~ASSOCIATE~~ ASSOCIATED WITH  
 THE LARGEST E-VAL.

START WITH  $v^{(0)}$   $\|v^{(0)}\|=1$

AND APPLY A AND NORMALIZE  
 AT EACH ITERATION.

$$v^{(1)} = \frac{Av^{(0)}}{\|Av^{(0)}\|}$$

$$v^{(2)} = \frac{Av^{(1)}}{\|Av^{(1)}\|}$$

REPEATING THE ~~THE~~ PROCESS  $\textcircled{9}$   
 $m-2$  TIMES WE OBTAIN  
 THE HESSENBERG FORM

$$Q_{m-2}^* \cdots Q_2^* Q_1^* A Q_1 Q_2 \cdots Q_{m-2}$$

ALGORITHM: HOUSEHOLDER  
 REDUCTION TO HESSENBERG FORM.

FOR  $k = 1$  TO  $m-2$

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

$$v_k = \text{sign}(z_1) \|z\|_2 e_1 + z$$

$$v_k = v_k / \|v_k\|_2$$

$$\begin{aligned} \text{(i)} \quad A_{k+1:m, k:m} &= A_{k+1:m, k:m} \\ &\quad - 2v_k(v_k^* A_{k+1:m, k:m}) \end{aligned}$$

(12)

$$\begin{aligned} \text{(ii)} \quad A_{1:m, k+1:m} &= A_{1:m, k+1:m} \\ &\quad - 2(A_{1:m, k+1:m} v_k) v_k^* \end{aligned}$$

END FOR

OPERATION COUNT

WORK DOMINATED BY (i)  
 AND (ii)

- OPERATION COUNT FOR (i)  
 IS EXACTLY THAT FOR  
 HOUSEHOLDER TRIANGULARIZATION.

$$\# \text{ OF FLOPS} \sim \frac{4}{3} m^3$$

- WORK FOR (ii)

AND KEEP GOING!

(11)

WRITE  $v^{(0)}$  AS A LIN.  
COMB. OF THE E-VECS OF  
A.

$$v^{(0)} = a_1 g_1 + a_2 g_2 + \dots + a_m g_m$$

SINCE  $v^{(k)}$  IS A MULTIPLE  
OF  $A^k v^{(0)}$ , WE HAVE FOR  
SOME  $c_k$

$$\begin{aligned} v^{(k)} &= c_k A^k v^{(0)} \\ &= c_k (a_1 \lambda_1^k g_1 + a_2 \lambda_2^k g_2 + \dots + a_m \lambda_m^k g_m) \end{aligned}$$

SUPPOSE OUR E-VALS ARE INDEXED SUCH THAT

$$|\lambda_1| \geq |\lambda_2| \geq |\lambda_3| \geq \dots \geq |\lambda_m|$$

$$\begin{aligned} v^{(k)} &= c_k \lambda_1^k (a_1 g_1 + a_2 (\lambda_2/\lambda_1)^k g_2 \\ &\quad + \dots + a_m (\lambda_m/\lambda_1)^k g_m) \end{aligned}$$

THEREFORE,

$$\|v^{(k)} - (a_1 g_1)\| = O\left(\left|\frac{\lambda_2}{\lambda_1}\right|^k\right)$$

AND USING THE RAYLEIGH QUOTIENT  $\lambda^{(k)} = r(v^{(k)})$   
WE HAVE

$$|\lambda^{(k)} - \lambda_1| = O\left(\left|\frac{\lambda_2}{\lambda_1}\right|^{2k}\right)$$

## POWER ITERATION ALGORITHM

①

$v^{(0)}$  = some vector with  $\|v^{(0)}\|=1$ .

for  $k=1, 2, \dots$

$$w = A v^{(k-1)}$$

$$v^{(k)} = w/\|w\|$$

$$\lambda^{(k)} = (v^{(k)})^T A v^{(k)}$$

END FOR ONCE "CONVERGED"

POWER ITERATION IS RESTRICTED

TO FINDING THE ~~LARGEST~~

EIGENVECTOR ASSOCIATED

WITH THE LARGEST EIGENVALUE.

POWER ITERATION HAS A  
LINEAR CONVERGENCE WITH  
RATE BASED ON  $|\lambda_2/\lambda_1|$

## INVERSE ITERATION

②

For  $\mu \in \mathbb{R}$  NOT AN  
E-VAL OF  $A$ ,

- THE E-VECS OF  
THE MATRIX  $(A - \mu I)^{-1}$

ARE THE SAME AS  $A$ .

- THE E-VALS OF  $(A - \mu I)^{-1}$   
ARE  $\{(\lambda_j - \mu)^{-1}\}$  WHERE

$\lambda_j$  ARE THE EVALS OF  
 $A$ .

CAN CHOOSE  $\mu$  CLOSE  
TO  $\lambda_j$ , SUCH THAT

$(\lambda_j - \mu)^{-1}$  IS MUCH LARGER  
THAN  $(\lambda_i - \mu)^{-1} \forall i \neq j$ .

### INVERSE ITERATION ALGORITHM ③

$v^{(0)}$  = some vector  $\|v^{(0)}\|=1$

for  $k=1, 2, \dots$

$$\text{SOLVE } (A - \mu I)w = v^{(k-1)}$$

$$v^{(k)} = w / \|w\|$$

$$\lambda^{(k)} = (v^{(k)})^T A (v^{(k)})$$

THOUGH CONVERGENCE IS  
STILL LINEAR, ITS RATE  
CAN BE CONTROLLED BY  
CHOOSING  $\mu$  TO BE CLOSE  
TO THE E-VAL OR INVERT.

BASED ON THE CONVERGENCE  
OF THE POWER ITERATION,

$$\|v^{(k)} - (\pm g_j)\| = O\left(\left|\frac{\mu - \lambda_j}{\mu - \lambda_k}\right|^k\right) \quad (4)$$

$$|\lambda^{(k)} - \lambda_j| = O\left(\left|\frac{\mu - \lambda_j}{\mu - \lambda_k}\right|^{2k}\right)$$

### RAYLEIGH QUOTIENT ITERATION

COMBINE OUTPUT OF THE  
RAYLEIGH QUOTIENT WITH  
INVERSE ITERATION, I.E. SET

$\mu = \lambda^{(k)}$  AT EACH ITERATION.

### ALGORITHM

$v^{(0)}$  = some vector  $\|v^{(0)}\|=1$

$$\lambda^{(0)} = (v^{(0)})^T A v^{(0)}$$

for  $k=1, 2 \dots$

$$\text{SOLVE } (A - \lambda^{(k-1)} I)w = v^{(k-1)} \quad (5)$$

$$v^{(k)} = w/\|w\|$$

$$\lambda^{(k)} = (v^{(k)})^T A v^{(k)}$$

THIS SIMPLE MODIFICATION  
GREATER ACCELERATES CONVERGENCE

SUPPOSE WE HAVE  $v^{(k)}$

AND  $\lambda^{(k)}$  SUFFICIENTLY CLOSE  
TO  $\lambda_J$  AND  $\lambda_{\bar{J}}$ .

$$\|v^{(k+1)} - g_J\| = O(|\lambda^{(k)} - \lambda_J| \|v^{(k)} - g_J\|)$$

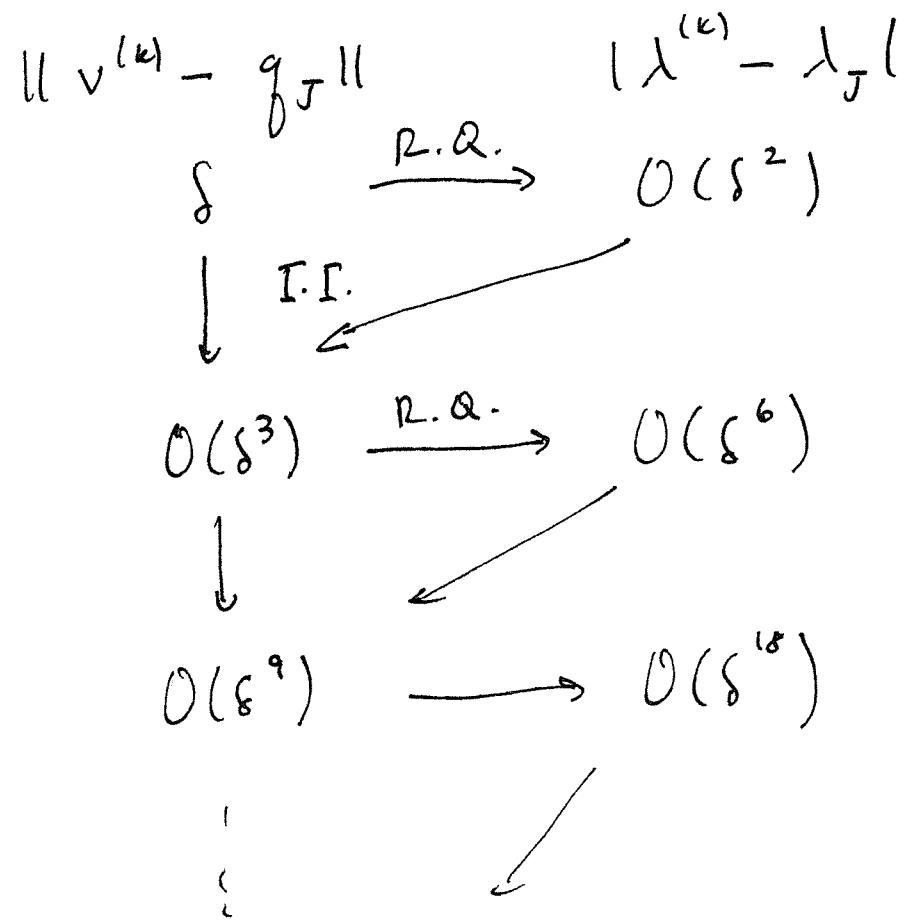
IF  $\|v^{(k)} - g_J\| \leq \delta$  THEN

FROM RAYLEIGH QUOTIENT THAT

$$|\lambda^{(k)} - \lambda_J| \leq O(\delta^2)$$

$$\|v^{(k+1)} - g_J\| \leq O(\delta^3) \quad (6)$$

Thus we see the following  
CONVERGENCE PATTERN EMERGE.



## OPERATION COUNTS

(7)

$\Rightarrow O(m^3)$  AT EACH  
ITERATION.

(8)

### POWER ITERATION

- MATRIX-VEC. MULTIPLICATION  
AT EACH ITERATION

$$\Rightarrow O(m^2)$$

### INVERSE ITERATION

- NEED TO INVERT  
SAME LINEAR SYSTEM  
AT EACH STEP.

- 1  $O(m^3)$  COMPUTATION
- $N_{\text{ITER}} \times O(m^2)$

### KRONECKE RAYLEIGH QUOTIENT

INVOLVES A NEW LINEAR  
SYSTEM AT EVERY ~~STEP~~  
ITERATION.

Now we see the reason  
for "PRE TREATING" A.

IF A IS TRIDIAGONAL  
EACH METHOD IS  $O(m)$ .

THESE METHODS GIVE US  
ONE E-VAL AND ONE  
E-VEC. HOW DO WE  
OBTAIN THE SPECTRUM AND  
EIGEN SPACE OF A?

## QR ALGORITHM

"PURE" QR Algorithm

$$A^{(0)} = A$$

for  $k = 1, 2, \dots$

$$Q^{(k)} R^{(k)} = A^{(k-1)}$$

$$A^{(k)} = R^{(k)} Q^{(k)}.$$

As  $k \rightarrow \infty$   $A^{(k)} \rightarrow T$  (Schur form)

or IF  $A$  is HERMITIAN

$$A^{(k)} \rightarrow \Lambda.$$

(9)

$A^{(k)}$  is similar to  
 $A^{(k-1)}$  SINCE

$$R^{(k)} = (Q^{(k)})^T A^{(k-1)}$$

Thus,

$$A^{(k)} = (Q^{(k)})^T A^{(k-1)} Q^{(k)}.$$

This still seem IMPOSSIBLY SIMPLE!

How/why DOES IT work?

UNNORMALIZED SIMULTANEOUS ITERATION

- APPLY POWER ITERATION TO SEVERAL VECTORS AT ONCE.

(10)

UNNORMALIZED SIMULTANEOUS ITERATION

① CONVERGES TO THE SPACE SPANNED BY THE DOMINANT E-VECS OF A

- APPLY POWER ITERATION TO SEVERAL VECTORS SIMULTANEOUSLY.

- TAKE

$$V^{(0)} = \begin{bmatrix} v_1^{(0)} & | & v_2^{(0)} & | & v_3^{(0)} & | \dots & | & v_n^{(0)} \end{bmatrix}_{m \times n}$$

$$\text{SET } V^{(k)} = A^k V^{(0)}$$

- CAN PROVE THAT THE SPACE

$$\langle A^k v_1^{(0)}, \dots, A^k v_n^{(0)} \rangle$$

$$\langle g_1, \dots, g_n \rangle$$

Thus IF

$$\hat{Q}^{(k)} R^{(k)} = V^{(k)}$$

THEN WE HAVE THAT THE COLUMNS OF  $\hat{Q}^{(k)}$  CONVERGE TO  $g_1, \dots, g_n$ .

SIMULTANEOUS ITERATION

SINCE ALL  $v_j^{(k)}$  CONVERGE TO  $g_1, V^{(k)}$  IS ILL CONDITIONED

REMEDY : COMPUTE

(3)

$$\underline{\hat{Q}}^{(k)} \underline{\hat{R}}^{(k)} = \underline{V}^{(k)}$$

AT EACH  
ITERATION AND USE  
 $\underline{\hat{Q}}^{(k)}$  INSTEAD OF  $\underline{V}^{(k)}$  AT  
THE NEXT ITERATION.

SIMULTANEOUS ITERATION

ALGORITHM

PICK  $\underline{\hat{Q}}^{(0)} \in \mathbb{R}^{m \times n}$

for  $k = 1, 2, \dots$

$$\underline{z} = \underline{A} \underline{\hat{Q}}^{(k-1)}$$

$$\underline{\hat{Q}}^{(k)} \underline{\hat{R}}^{(k)} = \underline{z}$$

IT TURNS OUT SIMULTANEOUS  
ITERATION AND THE QR  
ALGORITHM ARE EQUIVALENT.

TAKE  $n = m$ , REDUCED  $\rightarrow$  FULL  
QR

~~SIMULTANEOUS~~ ITERATION

$$\underline{\hat{Q}}^{(0)} = \underline{I}$$

$$\underline{z} = \underline{A} \underline{\hat{Q}}^{(k-1)}$$

$$\underline{\hat{Q}}^{(k)} \underline{n}^{(k)} = \underline{z}.$$

$$\underline{A}^{(k)} = (\underline{\hat{Q}}^{(k)})^T \underline{A} \underline{\hat{Q}}^{(k)}$$

## "PURE" QR ALGORITHM

$$\underline{A}^{(0)} = A$$

$$\underline{A}^{(k-1)} = \underline{Q}^{(k)} \underline{R}^{(k)}$$

$$\underline{A}^{(k)} = \underline{R}^{(k)} \underline{Q}^{(k)}$$

$$\underline{Q}^{(k)} = \underline{Q}^{(1)} \underline{Q}^{(2)} \dots \underline{Q}^{(k)}$$

For BOTH,

$$\underline{R}^{(k)} = \underline{R}^{(k)} \underline{R}^{(k-1)} \dots \underline{R}^{(1)}$$

THUS: SIMULTANEOUS ITERATION

WITH  $\underline{Q}^{(0)} = I$  AND

QR ALGORITHM GENERATE IDENTICAL  $\Leftrightarrow$  SEQUENCES

$$\underline{R}^{(k)}, \underline{Q}^{(k)} \text{ AND } \underline{A}^{(k)}$$

(5)

WHERE

$$\underline{A}^k = \underline{Q}^{(k)} \underline{R}^{(k)}$$

AND

$$\underline{A}^{(k)} = (\underline{Q}^{(k)})^T \underline{A} \underline{Q}^{(k)}$$

WE SEE NOW WHY QR ALGORITHM WORKS!

- $\underline{Q}^{(k)}$  CONTAINS OUR APPROXIMATION OF THE EIGENVECTORS.

- THE DIAGONAL ENTRIES OF  $\underline{A}^{(k)}$  ARE THE RAYLEIGH QUOTIENTS.

(6)

WE CAN MODIFY THE  
"PURE" QR ALGORITHM TO  
BE "PRACTICAL".

"PRACTICAL" QR ALGORITHM

$$(\underline{Q}^{(0)})^T \underline{A}^{(0)} \underline{Q}^{(0)} = \underline{A}$$

FOR  $k = 1, 2, \dots$

PICK SHIFT  $\mu^{(k)}$

$$\underline{Q}^{(k)} \underline{R}^{(k)} = \underline{A}^{(k-1)} - \mu^{(k)} \underline{I}$$

$$\underline{A}^{(k)} = \underline{R}^{(k)} \underline{Q}^{(k)} + \mu^{(k)} \underline{I}$$

IF OFF-DIAGONAL  $A_{j,j+1}^{(k)}$

IS CLOSE TO ZERO

$$\text{SET } A_{j,j+1} = A_{j+1,j} = 0$$

TO OBTAIN

$$\begin{bmatrix} A_1 & 0 \\ 0 & A_2 \end{bmatrix}$$

(4)

APPLY SAME ALGORITHM (5)  
TO  $A_1, A_2$  SEPARATELY.

CONNECTION TO INVERSE  
ITERATION

$$\underline{A}^k = \underline{Q}^{(k)} \underline{R}^{(k)}$$

TAKE THE INVERSE

$$\underline{A}^{-k} = (\underline{R}^{(k)})^{-1} (\underline{Q}^{(k)})^T$$

USING THE SYMMETRY OF  $\underline{A}^{-1}$ .

$$\underline{A}^{-k} = \underline{Q}^{(k)} (\underline{R}^{(k)})^{-T}$$

TAKE

$$P = \begin{bmatrix} & & 1 \\ & \ddots & \\ 1 & \cdots & \end{bmatrix}$$

$$\text{SINCE } P^2 = \underline{I}$$

$$\left\{ A^{-k} = \left( A^{-1} \right)^k \right\}$$

(9)

$$A^{-k} P = \left[ \underline{Q}^{(k)} P \right] \left[ P \left( \underline{R}^{(k)} \right)^{-T} P \right]$$

This is the simultaneous iteration of  $A^{-1}$  with initial matrix  $P$ .

First column of  $\underline{Q}^{(k)} P$   
is ~~the~~ the last  
column of  $\underline{Q}^{(k)}$

This can be viewed  
as inverse iteration  
applied to  $e_m$ .

THIS CAN THEREFORE BE  
ACCELERATED USING SHIFTS!  
 $\mu^{(k)}$ :

$$A^{(k-1)} - \mu^{(k)} I = \underline{Q}^{(k)} \underline{R}^{(k)}$$

$$A^{(k)} = \underline{R}^{(k)} \underline{Q}^{(k)} + \mu^{(k)} I$$

THIS RETAINS

$$A^{(k)} = (\underline{Q}^{(k)})^T A^{(k-1)} \underline{Q}^{(k)}$$

AND CONSEQUENTLY

$$A^{(k)} = (\underline{Q}^{(k)})^T A \underline{Q}^{(k)}.$$

## CHOICE OF SHIFTS

(1) RAYLEIGH QUOTIENT OF THE LAST COLUMN OF  $\underline{Q}^{(k)}$ .

$$\mu^{(k)} = (\underline{q}_m^{(k)})^T A \underline{q}_m^{(k)}$$

This will give a CUBIC CONVERGENCE TO THE E-VEC AND E-VAL.

This comes FOR FREE SINCE IT IS THE  $m, m$  ENTRY OF  $A^{(k)}$ .

(1)

$$A_{mm}^{(k)} = \underline{e}_m^T A^{(k)} \underline{e}_m$$

(2)

$$= \underline{e}_m^T \underline{Q}^{(k)T} A \underline{Q}^{(k)} \underline{e}_m \\ = \underline{g}_m^{(k)T} A \underline{g}_m^{(k)}$$

SETTING  $\mu^{(k)} = A_{mm}^{(k)}$

IS KNOWN AS THE  
RAYLEIGH      QUOTIENT  
SHIFT.

(2) ~~WILKINSON~~ WILKINSON SHIFT

SET  $B$  ( $2 \times 2$ ) AS THE LOWER RIGHT MOST SUBMATRIX OF  $A^{(k)}$

$$B = \begin{bmatrix} a_{m-1} & b_{m-1} \\ b_{m-1} & a_m \end{bmatrix}$$

(3)

SET  $\mu^{(k)}$  AS THE EIGENVALUE OF  $B$  CLOSEST TO  $a_m$ .

$$\mu^{(k)} = a_m - \frac{\text{sign}(s) b_{m-1}^2}{(|s| + \sqrt{s^2 + b_{m-1}^2})}$$

WHERE

$$s = \frac{a_{m-1} - a_m}{2}$$

### STABILITY

(4)

For  $A \in \mathbb{R}^{m \times m}$  REAL, SYMMETRIC, AND TRIDIAGONAL.

QR ALGORITHM ON A COMPUTER SATISFYING FPA I & II GIVES

$$\tilde{Q} \tilde{\Lambda} \tilde{Q}^T = A + \delta A$$

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

FOR SOME  $\delta A \in \mathbb{R}^{m \times m}$

Question

$$(i) (A^T)^{-1} = (A^{-1})^T$$

THIS IS TRUE

$$(A^{-1} A)^T = I$$

$$A^T (A^{-1})^T = I$$

$$\text{Thus, } (A^{-1})^T = (A^T)^{-1}$$

$$(ii) (A^{-1})^k = (A^k)^{-1}$$

$$(A^{-1} \dots A^{-1}) (A \dots A) = I$$

$\underbrace{\quad}_{k \text{ TIMES}}$

$$(A^{-1})^k A^k = I$$

(5)

Thus,

$$(A^{-1})^k = (A^k)^{-1}.$$

(6)

UNSYMMETRIC E-VAL PROBLEM

MANY OF THE METHODS

USED IN THE SYMMETRIC  
THE E-VAL PROBLEM

CARRY OVER TO THE  
UNSYMMETRIC ONE.

- POWER METHOD. (ITERATION)
- INVERSE ITERATION.
- "PURE" QR
- "PRACTICAL" QR.

BUT, MANY OF THE IMPORTANT  
DETAILS (CHOOSING SHIFTS, ETC)  
BECOME MORE COMPLICATED.

## ITERATIVE METHODS

⑦

- WITH DIRECT METHODS  
OPERATION COUNT TO  
SOLVE A LINEAR  
SYSTEM IS  $\mathcal{O}(m^3)$
- INSTEAD GENERATE  
A SEQUENCE THAT  
CONVERGES TO THE  
SOLUTION (HOPEFULLY)  
VERY QUICKLY.
- EACH ITERATION  
COSTS MATRIX-VECTOR  
MULTIPLICATION  
TOTAL :  $\mathcal{O}(m^2) \times N_{\text{ITER}}$

⑧

- ALLOWS USERS TO SOLVE  
LINEAR SYSTEMS ~~THAT~~  
USING IMPLICIT MATRIX  
VECTOR MULTIPLICATION.
- KRYLOV SUBSPACE  
METHODS.