

# 1 System of Linear Equations, $Ax = b$

## 1.1 $p$ -Norm and Condition Number

Vector  $p$ -Norm:  $\boxed{\|\vec{x}\|_p = \sqrt[p]{\sum_i |x_i|^p}}$

1-Norm :  $\|\vec{x}\|_1 = \sum_i |x_i|$

$\infty$ -Norm :  $\|\vec{x}\|_\infty = \max |x_i|$

- $\|x\|_1 \geq \|x\|_2 \geq \|x\|_\infty$
- $\|x\|_1 \leq \sqrt{n} \|x\|_2 \leq \sqrt{n} \|x\|_\infty$

Matrix  $p$ -Norm:  $\boxed{\|A\|_p = \max_{x \neq 0} \frac{\|Ax\|}{\|x\|}}$

1-Norm :  $\|A\|_1 = \max_j \sum_i |a_{ij}|$

$\infty$ -Norm :  $\|A\|_\infty = \max_i \sum_j |a_{ij}|$

- $\|AB\| \leq \|A\| \cdot \|B\|$
  - $\|Ax\| \leq \|A\| \cdot \|x\|$
- For  $p$ -norms (not necessarily in general)

Function/Vector Condition Number:

$$\begin{aligned} \text{cond}(f(x)) &= \left| \frac{[f(\hat{x}) - f(x)]/f(x)}{[\hat{x} - x]/x} \right| \\ &= \left| \frac{\Delta y/y}{\Delta x/x} \right| = \left| \frac{y' \cdot \Delta x/y}{\Delta x/x} \right| \\ &= \left| \frac{x f'(x)}{f(x)} \right| \end{aligned}$$

Matrix Condition Number:

$\boxed{\text{cond}_p(A) = \|A\|_p \cdot \|A^{-1}\|_p}$  ( $\infty$  if singular)

$$= \frac{\max_{x \neq 0} \|Ax\|_p / \|x\|_p}{\min_{x \neq 0} \|Ax\|_p / \|x\|_p} = \text{cond}_p(\gamma A) \geq 1$$

- Diagonal,  $D$  :  $\text{cond}(D) = \frac{\max |d_i|}{\min |d_i|}$
- $\|z\| = \|A^{-1}y\| \leq \|A^{-1}\| \cdot \|y\|$   
 $\rightarrow \frac{\|z\|}{\|y\|} \leq \max \frac{\|z\|}{\|y\|} \stackrel{?}{=} \|A^{-1}\|$  (optimize)

## 1.2 Error Bounds and Residuals

Error Bound:  $\boxed{\frac{\|\hat{x} - x\|}{\|x\|} \lesssim \text{cond}(A) \epsilon_{\text{mach}}}$   $\rightarrow$  A computed solution is expected to lose about  $\log_{10}(\text{cond}(A))$  digits, so the input data must be more accurate to these digits and the working precision must carry more than these digits.

$$A\hat{x} = b + \Delta b = Ax + A\Delta x$$

- $\|b\| \leq \|A\| \cdot \|x\|$
- $\|\Delta x\| \leq \|A^{-1}\| \cdot \|\Delta b\|$

$$\rightarrow \boxed{\frac{\|\Delta x\|}{\|x\|} \leq \text{cond}(A) \frac{\|\Delta b\|}{\|b\|}}$$

$$A\hat{x} + r = b$$

- $\|\Delta x\| = \|A^{-1}(A\hat{x} - b)\| = \|-A^{-1}r\|$   
 $\leq \|A^{-1}\| \cdot \|r\|$

$$\rightarrow \boxed{\frac{\|\Delta x\|}{\|\hat{x}\|} \leq \text{cond}(A) \frac{\|r\|}{\|A\| \cdot \|\hat{x}\|}}$$

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

- $\|\Delta x\| = \|-A^{-1}(\Delta A)\hat{x}\|$   
 $\leq \|A^{-1}\| \cdot \|\Delta A\| \cdot \|\hat{x}\|$

$$\rightarrow \boxed{\frac{\|\Delta x\|}{\|x\|} \leq \text{cond}(A) \frac{\|\Delta A\|}{\|A\|}}$$

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

- $\|r\| = \|b - A\hat{x}\| = \|\Delta A \cdot \hat{x}\|$   
 $\leq \|\Delta A\| \cdot \|\hat{x}\|$

$$\rightarrow \boxed{\frac{\|r\|}{\|A\| \cdot \|\hat{x}\|} \leq \frac{\|\Delta A\|}{\|A\|}}, \quad \frac{\|\Delta x\|}{\|x\|} \leq \frac{\|A^{-1}\| \cdot \|r\|}{\|\hat{x}\|} \leq \text{cond}(A) \frac{\|\Delta A\|}{\|A\|}$$

$$\left[ A(t)x(t) = b(t) \right] = \left[ (A_0 + \Delta A \cdot t)x(t) = b_0 + \Delta b \cdot t \right]$$

- $x'(t) = \frac{b'(t) - A'(t)x(t)}{A(t)} = A^{-1}(t) \left[ \Delta b - \Delta A \cdot x(t) \right]$
- $x(t) = x_0 + x'(0)t + \mathcal{O}(t^2)$

$$\rightarrow \boxed{\frac{\|x(t) - x_0\|}{\|x_0\|} \leq \text{cond}(A) \left( \frac{\|\Delta b\|}{\|b\|} + \frac{\|\Delta A\|}{\|A\|} \right) |t| + \mathcal{O}(t^2)}$$

## 1.3 Gaussian Elimination with LU/PLU/PLDUQ Decomposition

### Elementary Elimination Matrices, $L_k$

$$\begin{pmatrix} 1 & \dots & 0 & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & 1 & 0 & \dots & 0 \\ 0 & \dots & \frac{-a_{k+1}}{a_k} & 1 & \dots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & \frac{-a_n}{a_k} & 0 & \dots & 1 \end{pmatrix} \begin{pmatrix} a_1 \\ \vdots \\ a_k \\ a_{k+1} \\ \vdots \\ a_n \end{pmatrix} = \begin{pmatrix} a_1 \\ \vdots \\ a_k \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

- $a_k$  is the “pivot”
- is lower triangular
- $\forall i \neq j \quad (L_k^{-1})_{ij} = -(L_k)_{ij}$

Ex :

$$\begin{pmatrix} 1 & 0 & \dots \\ -a_1/a_2 & 1 & \dots \\ \vdots & \vdots & \ddots \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \\ \vdots \end{pmatrix} = \begin{pmatrix} a_1 \\ 0 \\ \vdots \end{pmatrix}$$

### LU/PLU Factorization (w/ partial pivoting)

$$\boxed{A = LU \quad \begin{array}{l} (L \text{ is gen. triang.}) \\ (U \text{ is upp. triang.}) \end{array}} \\ \boxed{L = (\dots L_2 P_2 L_1 P_1)^{-1}}$$

$$\{\dots\}b = (\dots L_2 P_2 L_1 P_1)Ax$$

$$L^{-1}b = (P_1^T L_1^{-1} P_2^T L_2^{-1} \dots)^{-1} Ax \\ = L^{-1}(LU)x = y$$

$$\boxed{b = Ly \quad y = Ux} \\ \text{(forw.-sub.) , (back.-sub.)}$$

- Permutation matrix,  $P_i$ , rowswaps s.t.  $a_k \neq 0$
- $P_i$  rowswaps s.t.  $a_k$  is largest s.t.  $a_{k+i}/a_k \leq 1$  for numerical stability/minimize errors
- Pivoting isn't needed if  $A$  is diag. dom. ( $a_{jj} > \sum_{i \neq j} a_{ij}$ )
- $A$  can be singular

$$\boxed{A = PLU \quad \begin{array}{l} (P \text{ is rowswap permu.}) \\ (L \text{ is unit low. triang.}) \\ (U \text{ is upp. triang.}) \end{array}} \\ \boxed{P = (\dots P_2 P_1)^{-1}}$$

$$\{\dots\}b = (\dots P_2 P_1)Ax$$

$$P^T b = (P_1^T P_2^T \dots)^{-1} Ax \\ = P^T (PLU)x = Ly$$

$$\boxed{P^T b = Ly \quad , \quad y = Ux}$$

$$\boxed{P^T A = LDU \quad (D \text{ is diag.})}$$

- $LDU$  is unique up to  $D$
- $LDU$  is unique if  $L/U$  are unit low./upp. diag., resp.

$$\boxed{P^T A Q^T = LDU \quad \begin{array}{l} (P \text{ is permu. for rows}) \\ (Q \text{ is permu. for cols.}) \end{array}}$$

- “Complete pivoting” search for largest  $a_k$
- Would be most numerically stable
- Expensive, so not really used

$$\text{Error Bound: } \frac{\|r\|}{\|A\|\|x\|} \leq \frac{\|\Delta A\|}{\|A\|} \leq \rho n^2 \epsilon_{\text{mach}} \sim n \epsilon_{\text{mach}} \\ \text{(Wilkinson) (usually)}$$

(growth factor,  $\rho$ , is the largest entry at any point during factorization - usually at  $U$  - divided by the largest entry of  $A$ )

## 1.4 Gaussian-Jordan with MD Decomposition

Elementary Elimination Matrices,  $M_k$

$$\begin{pmatrix} 1 & \dots & \frac{-a_1}{a_k} & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & 1 & 0 & \dots & 0 \\ 0 & \dots & \frac{-a_{k+1}}{a_k} & 1 & \dots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & \frac{-a_n}{a_k} & 0 & \dots & 1 \end{pmatrix} \begin{pmatrix} a_1 \\ \vdots \\ a_{k-1} \\ a_k \\ a_{k+1} \\ \vdots \\ a_n \end{pmatrix} = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ a_k \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

- $a_k$  is the “pivot”
- $\forall i \neq j \quad (M_k^{-1})_{ij} = -(M_k)_{ij}$

MD Factorization (w/ partial pivoting)

$$\boxed{A = MD \quad \begin{array}{l} (M \text{ is elem. elim.}) \\ (D \text{ is diag.}) \end{array}}$$

$$M = (\dots M_2 P_2 M_1 P_1)^{-1}$$

$$\{\dots\}b = (\dots M_2 P_2 M_1 P_1)Ax$$

$$M^{-1}b = (P_1^T M_1^{-1} P_2^T M_2^{-1} \dots)^{-1} Ax$$

$$= M^{-1}(MD)x = y$$

$$\boxed{M^{-1}b = y, \quad y = Dx}$$

(division)

- Permutation matrix,  $P_i$ , rowswaps s.t.  $a_k \neq 0$
- $P_i$  rowswaps cannot ensure numerical stability ( $\leq 1$ )
- Division is  $\mathcal{O}(n)$ , so may be useful for parallel comps.
- Can also find  $A^{-1}$

Finding  $A^{-1}$

$$D^{-1}M^{-1}(A|I) = (I|A^{-1})$$

$$= D^{-1}M^{-1} \left[ \begin{array}{cc|cc} a_{11} & \dots & 1 & 0 \\ \vdots & a_{nn} & 0 & 1 \end{array} \right]$$

$$= \left[ \begin{array}{cc|cc} 1 & 0 & a'_{11} & \dots \\ 0 & 1 & \vdots & a'_{nn} \end{array} \right]$$

## 1.5 Symmetric Matrices

Positive Definite:  $x^T Ax \geq 0$

Cholesky Factorization for Sym., Pos. Def.:  $A = LL^T = LDL^T$

$$\begin{pmatrix} a_{11} & a_{21} & a_{31} & \dots \\ a_{21} & a_{22} & a_{32} & \dots \\ a_{31} & a_{32} & a_{33} & \dots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix} = \begin{pmatrix} l_{11} & 0 & 0 & \dots \\ l_{21} & l_{22} & 0 & \dots \\ l_{31} & l_{32} & l_{33} & \dots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix} \begin{pmatrix} l_{11} & l_{21} & l_{31} & \dots \\ 0 & l_{22} & l_{32} & \dots \\ 0 & 0 & l_{33} & \dots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix} = \begin{pmatrix} l_{11}^2 & \dots & \dots & \dots \\ l_{21}l_{11} & l_{21}^2 + l_{22}^2 & \dots & \dots \\ l_{31}l_{11} & l_{31}l_{21} + l_{32}l_{22} & l_{31}^2 + l_{32}^2 + l_{33}^2 & \dots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix}$$

- Pivoting not needed
- Well defined (always works)
- Only lower triangle needed for storage
- $A = LDL^T$  is sometimes useful, where  $D$  is diag.

Symmetric Indefinite Matrices

- Pivoting Needed:  $PAP^T = LDL^T$
- Ideally,  $D$  is diag., but if not possible, then  $D$  is tridiag. (Aasen) or 1x1/2x2 block diag. (Bunch, Parlett, Kaufmann, etc.)

## 1.6 Banded Matrices

- Similar to normal Gaussian Elim., but less work since more zeroes
- Pivoting means bandwidth will expand no more than double
- Only  $\mathcal{O}(\beta n)$  storage needed

## 1.7 Rank-1 Update with Sherman-Morrison

$$\tilde{A}\tilde{x} = b = (A - uv^T)\tilde{x} \quad \left| \quad \begin{array}{l} \tilde{A}^{-1} = (A - uv^T)^{-1} = A^{-1} + \frac{A^{-1}u}{1 - v^T(A^{-1}u)} v^T A^{-1} \\ \tilde{A}^{-1}b = \quad \quad \quad \tilde{x} = (A^{-1}b) + \frac{A^{-1}u}{1 - v^T(A^{-1}u)} v^T(A^{-1}b) \\ \quad \quad \quad \quad \quad \quad \quad \quad \quad \quad \quad \quad \quad \quad \quad x + \frac{y}{1 - v^T y} v^T x \end{array} \right.$$

General Woodbury Formula:  $\boxed{(A - UV^T)^{-1} = A^{-1} + (A^{-1}U)(I - V^T A^{-1}U)^{-1} v^T A^{-1}}$

- $U$  and  $V$  are general  $n \times k$  matrices
- No guarantee of numerical stability, so caution is needed

## 1.8 Complexity

Explicit Inversion :  $\frac{LUA^{-1} = I}{D^{-1}M^{-1}I = A^{-1}} \rightarrow \mathcal{O}(n^3) \quad , \quad A^{-1}b = x \rightarrow \mathcal{O}(n^2)$

Gaussian Elimination :  $A = LU \rightarrow \mathcal{O}(n^3/3) \quad , \quad LUx = b \rightarrow \mathcal{O}(n^2)$

Gaussian-Jordan :  $A = MD \rightarrow \mathcal{O}(n^3/2) \quad , \quad MDx = b \rightarrow \mathcal{O}(n)$

Symmetric :  $\frac{A = LL^T}{PAP^T = LDL^T} \rightarrow \mathcal{O}(n^3/6) \quad , \quad LL^T x = b \rightarrow \mathcal{O}(n^2)$

Banded :  $A_\beta = LU \rightarrow \mathcal{O}(\beta^2 n) \quad , \quad LUx = b \rightarrow \mathcal{O}(\beta n)$

Sherman-Woodbury :  $\tilde{A} = A - uv^T \rightarrow \mathcal{O}(n^2) \quad , \quad \tilde{x} = \tilde{A}b \rightarrow \mathcal{O}(n^2)$

## 1.9 Diagonal Scaling

Ill-conditioned

$$\begin{pmatrix} 1 & 0 \\ 0 & \epsilon \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 1 \\ \epsilon \end{pmatrix}$$

Well-conditioned

$$\begin{pmatrix} 1 & 0 \\ 0 & 1/\epsilon \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & \epsilon \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1/\epsilon \end{pmatrix} \begin{pmatrix} 1 \\ \epsilon \end{pmatrix}$$

- No general way to correct poor scaling

## 1.10 Iterative Refinement

$$\begin{aligned} r_0 &= b - Ax_0 = A\Delta x_0 \\ r_1 &= b - A(x_0 + \Delta x_0) = b - Ax_1 = A\Delta x_1 \\ r_2 &= b - A(x_1 + \Delta x_1) = b - Ax_2 = A\Delta x_2 \end{aligned}$$

$$\boxed{x = x_0 + \lim_{n=0}^{\infty} \Delta x_n} \quad (\text{terminate when } r_n \text{ is small enough})$$

- Double storage needed to hold original matrix
- $r_n$  usually must be computed with higher precision than  $x_n$
- Useful for badly scaled systems, or making unstable systems stable
- If  $x_n$  is not accurate,  $r_n$  might not need better accuracy

## 2 Least $\|r\|$ Linear Regression/Fit, $Ax + r = b$

- $A = A_{m \times n}$   $\boxed{(m > n, \text{ underdetermined})}$
- $r(y = Ax)$  is cont. & coer.  $\rightarrow \exists \|r(y)\|_{\min}$
- $r(y)$  is strictly convex  $\rightarrow y = Ax$  is unique
- $\boxed{\text{rank}(A) = n}$   $\Rightarrow A(x_1 - x_2) = 0$  (unique  $x$ )  
 $\boxed{(\text{full column rank})} \quad (x_1 - x_2) = 0 \rightarrow x_1 = x_2$

Example - Vandermonde Matrix,  $A$ :

$$Ax = \begin{pmatrix} -\vec{f}(t_1) \\ \vdots \\ -\vec{f}(t_m) \end{pmatrix} \begin{pmatrix} | \\ \vec{x} \\ | \end{pmatrix} = \begin{pmatrix} y(t_1) \\ \vdots \\ y(t_m) \end{pmatrix} = \begin{pmatrix} | \\ \vec{y} \\ | \end{pmatrix} = (x^T A^T)^T, \quad y(t) = \sum_{i=1}^n x_i f_i(t) = \vec{x} \cdot \vec{f}$$

Decompose  $b$ :

$$\begin{aligned} b &= Ax + r \\ &= y + r \\ &= Pb + P_{\perp} b \end{aligned}$$

Projector of  $A$ ,  $P$

$$\text{Projector : } P^2 = P \rightarrow PA = A$$

(Idempotent)                      (Projector of  $A$ )

$$\text{Orthogonal Projector : } P^T = P \rightarrow P_{\perp} A = (I - P)A = 0$$

Minimize residual,  $r$ :

$$\begin{aligned} \nabla \|r\|_2^2 &= 0 \quad \left( \frac{\partial x^2}{\partial x_i} = 0 \right) \\ &= \nabla [(b - Ax)^T (b - Ax)] \\ &= \nabla (b^T b - 2x^T A^T b + x^T A^T A x) \\ 0 &= 2A^T A x - 2A^T b \\ &\downarrow \\ A^T A x &= A^T b \quad (\text{Solvable with Cholesky}) \end{aligned}$$

$$\begin{aligned} \|r\|_2^2 &= \|Pr + P_{\perp} r\|_2^2 = \|b - Ax\|_2^2 \\ &= \|Pr\|_2^2 + \|P_{\perp} r\|_2^2 \\ &= \cancel{\|Pb - Ax\|_2^2} + \|P_{\perp} b\|_2^2 \\ &\downarrow \\ Ax &= Pb \\ A^T A x &= A^T Pb = (P^T A)^T b \\ A^T A x &= A^T b \quad (\text{System of Normal Equations}) \end{aligned}$$

Cross-Product Matrix of  $A$ :  $\boxed{A^T A}$

$$\text{Symmetric : } (A^T A)^T = A^T A$$

$$\begin{aligned} \text{Pos. Def. : } \text{rank}(A) &= n \\ &\rightarrow \langle x | A^T A x \rangle = x^T A^T A x \\ &= (Ax)^T (Ax) \\ &= \|Ax\|^2 \geq 0 \end{aligned}$$

$$\begin{aligned} \text{Nonsingular : } A^T A x &= 0 \\ &\rightarrow \|Ax\|^2 = 0 = Ax \\ &\rightarrow (x = 0) \end{aligned}$$

System of Normal Equations:  $\boxed{A^T A x = A^T b}$

Pseudoinverse,  $A^+$

$$\boxed{x = (A^T A)^{-1} A^T b} \rightarrow \boxed{A^+ \equiv (A^T A)^{-1} A^T}$$

$$\equiv A^+ b \quad A^+ A = I$$

Ortho. Proj.,  $P$

$$\boxed{Ax = A(A^T A)^{-1} A^T b} \rightarrow \boxed{P = A(A^T A)^{-1} A^T}$$

$$= Pb \quad = AA^+$$

## System of Normal Equations Issues:

- Info can be lost forming  $A^T A$ , e.g,  $A = \begin{pmatrix} 1 & 0 \\ \epsilon & 0 \\ 0 & \epsilon \end{pmatrix} \rightarrow A^T A = \begin{pmatrix} 1+\epsilon^2 & 1 \\ 1 & 1+\epsilon^2 \end{pmatrix} \approx \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}$  (singular)

- System of Normal Equations:  $\boxed{\text{cond}(A^T A) = [\text{cond}(A)]^2}$

## 2.1 Error Bounds and Residuals

Error Bound:  $\boxed{\frac{\|\Delta x\|}{\|x\|} \lesssim \text{cond}(A) \epsilon_{\text{mach}}}$   $\rightarrow$  A computed solution is expected to lose about  $\log_{10}(\text{cond}(A))$  digits, so the input data must be more accurate to these digits and the working precision must carry more than these digits.

## Norm and Conditioning:

$$\boxed{\|A\| = \max_{x \neq 0} \left( \frac{\|Ax\|}{\|x\|} = \frac{\|AA^+b\|}{\|A^+b\|} \right)}$$

$$\boxed{\text{cond}(A) = \begin{cases} \|A\|_2 \cdot \|A^+\|_2 & \text{rank}(A) = n \\ \infty & \text{rank}(A) < n \end{cases}}$$

$$A^T A(x + \Delta x) = A^T A(b + \Delta b)$$

$$(A + \Delta A)^T (A + \Delta A)(x + \Delta x) = (A + \Delta A)^T b$$

$$\bullet \quad \|\Delta x\| \leq \|A^+\| \cdot \|\Delta b\|$$

$$\bullet \quad \cancel{A^T Ax} + A^T \Delta Ax + (\Delta A)^T Ax + \cancel{(\Delta A)^T \Delta Ax} = \cancel{A^T b} + (\Delta A)^T b \\ + A^T A \Delta x + \cancel{A^T \Delta A \Delta x} + \cancel{(\Delta A)^T A \Delta x} + \cancel{(\Delta A)^T \Delta A \Delta x}$$

$$\rightarrow \boxed{\frac{\|\Delta x\|}{\|\hat{x}\|} \leq \left( \text{cond}(A) \frac{\|b\|}{\|Ax\|} \right) \frac{\|\Delta b\|}{\|b\|} \\ = \left( \text{cond}(A) \frac{1}{\cos \theta} \right) \frac{\|\Delta b\|}{\|b\|}}$$

$$\bullet \quad \|\Delta x\| = \|(A^T A)^{-1} (\Delta A)^T r - A^+ \Delta Ax\| \\ \leq \|(A^T A)^{-1}\| \cdot \|\Delta A\| \cdot \|r\| + \|A^+\| \cdot \|\Delta A\| \cdot \|x\|$$

- Cond. number is a func. of  $\text{cond}(A)$  and  $b$
- $Pb \approx 0$  or  $\theta \approx 90^\circ$  is highly sensitive

$$\rightarrow \boxed{\frac{\|\Delta x\|}{\|\hat{x}\|} \leq \left( [\text{cond}(A)]^2 \frac{\|r\|}{\|Ax\|} + \text{cond}(A) \right) \frac{\|\Delta A\|}{\|A\|} \\ = \left( [\text{cond}(A)]^2 \tan \theta + \text{cond}(A) \right) \frac{\|\Delta A\|}{\|A\|}}$$

## 2.2 Solving $A^T Ax = A^T b$ with an Augmented Matrix

$$\begin{matrix} r + Ax = b \\ A^T r = 0 \end{matrix} \Rightarrow \begin{pmatrix} I & A \\ A^T & 0 \end{pmatrix} \begin{pmatrix} r \\ x \end{pmatrix} = \begin{pmatrix} b \\ 0 \end{pmatrix} \Rightarrow \begin{pmatrix} \alpha I & A \\ A^T & 0 \end{pmatrix} \begin{pmatrix} r/\alpha \\ x \end{pmatrix} = \begin{pmatrix} b \\ 0 \end{pmatrix}$$

- Solvable with  $LU$  Decomp or Symm. Pos. Def. Methods
- $\alpha$  "controls the relative weights of the two subsystems in choosing pivots from either"
- $\alpha = \max a_{ij}/1000$  (rule of thumb)
- MATLAB uses it for large, sparse systems

## 2.3 QR Decomposition

Motivation:  $Q^T A = \begin{pmatrix} R \\ 0 \end{pmatrix} \rightarrow Q^T Ax + Q^T r = Q^T b \rightarrow \begin{pmatrix} Rx \\ 0 \end{pmatrix} + \begin{pmatrix} r'_1 \\ r'_2 \end{pmatrix} = \begin{pmatrix} b'_1 \\ b'_2 \end{pmatrix} \rightarrow \begin{matrix} \|r'\|^2 = \|b'_1 - Rx\|^2 + \|b'_2\|^2 \\ \downarrow \\ Rx = b'_1, \quad r' = \begin{pmatrix} 0 \\ b'_2 \end{pmatrix} \quad (\text{solve with back-sub}) \end{matrix}$

Orthogonal Matrix,  $Q$

$$Q^T Q = Q Q^T = I$$

QR Factorization

$$A = Q \begin{pmatrix} R \\ 0 \end{pmatrix}$$

Reduced QR Factorization

$$A = Q \begin{pmatrix} R \\ 0 \end{pmatrix} = \begin{pmatrix} Q_{\parallel} & Q_{\perp} \end{pmatrix} \begin{pmatrix} R \\ 0 \end{pmatrix} = Q_{\parallel} R$$

1.  $Q^T$  is a  $\text{span}(A)$  Plane Rotation through  $\mathbb{R}^m$  to  $\text{span}([R \ 0]^T)$

2-norm Preserved ( $Q$  is a rotation/reflection)

- $\|Qv\|^2 = \langle v | Q^T Q v \rangle = \|v\|^2$
- $\|Q^T v\|^2 = \langle v | Q Q^T v \rangle = \|v\|^2$

- $Q^T = H_n \dots H_1$
- $H_i^T H_i = H_i H_i^T = I$
- $A = [a_1 \dots a_n]$
- $I_n = [e_1 \dots e_n]$

$$H_1 a_1 = \alpha_1 e_1 \quad (\|a_1\| = |\alpha_1|)$$

$$H_i \dots H_1 a_i = \sum_j^i c_j e_j = H_n \dots H_1 a_i$$

$$(\|a_i\|^2 = |\alpha_1|^2 = \sum_j^i c_j^2)$$

$$\langle r | a_i \rangle = 0 \quad (1 \leq i \leq n)$$

$$\langle H_i \dots H_1 r | e_j \rangle = 0 \quad (1 \leq j \leq i)$$

→  $Q^T A$  rotates  $A$  until the column vectors are aligned with certain axes described above

2.  $A$  is a Lin. Sum of  $Q_{\parallel}$ 's Orthogonal Column Vectors Given by  $R$

$$\{Q_{\parallel} = Q_{m \times n} \mid \text{span}(Q_{\parallel}) = \text{span}(A)\}$$

$$\rightarrow Q^+ = (Q^T Q)^{-1} Q^T = Q^T$$

$$\rightarrow P = Q_{\parallel} Q_{\parallel}^T$$

$$\rightarrow Q_{\parallel}^T A x = Q_{\parallel}^T P b = \cancel{Q_{\parallel}^T Q_{\parallel}} Q_{\parallel}^T b$$

$$= Q_{\parallel}^T b \quad (\text{System of Orthogonal Equations?})$$

$$A = Q_{\parallel} R = \begin{pmatrix} | & | & | \\ \vec{q}_1 & \dots & \vec{q}_n \\ | & | & | \end{pmatrix} \begin{pmatrix} r_{11} & \dots & r_{1n} \\ 0 & \ddots & \vdots \\ 0 & 0 & r_{nn} \end{pmatrix} = \begin{pmatrix} | & | & | \\ \vec{a}_1 & \dots & \vec{a}_n \\ | & | & | \end{pmatrix}$$

- $\vec{a}_j = \sum_i^j r_{ij} \cdot \vec{q}_i$

→  $R$  transforms the  $Q_{\parallel}$  column vectors about  $\text{span}(A)$ , an  $\mathbb{R}^n$  plane, until they equal the column vectors of  $A$



### 2.3.1 Householder Transformation/Elementary Reflector, $H$

$$\begin{aligned}
 H\vec{a}_1 &= \alpha_1 \vec{e}_1 \quad \begin{array}{l} \|\vec{a}_1\| = |\alpha_1| \\ \text{(rotation)} \end{array} \rightarrow \boxed{H = I - \hat{v}\hat{v}^T = I - \frac{2vv^T}{v^T v}} \quad \bullet \ H = H^T = H^{-1} \\
 &= \boxed{\vec{a}_1 - 2\hat{v}(\hat{v} \cdot \vec{a}_1)} \quad \text{[} v_\perp \text{ bisects } \theta(\vec{a}_1, \vec{e}_1)\text{]} \quad \text{(symmetric and orthogonal)}
 \end{aligned}$$

$$\bullet \ \alpha_1 e_1 = a_1 - (2v_1) \frac{v_1 \cdot a_1}{v_1 \cdot v_1} \Rightarrow v_1 = (a_1 - \alpha e_1) \frac{v_1 \cdot v_1}{2v_1 \cdot a_1} \quad \text{(magnitude doesn't matter)}$$

$$\rightarrow \boxed{v_1 = (a_1 - \alpha e_1)}$$

$$\alpha_1 = \pm \|\vec{a}_1\| \rightarrow \boxed{\alpha_i = -\text{sign}(a_i) \|\vec{a}_i\|} \quad \text{(avoid "cancellation" in finite-calc. of } v \text{ above)}$$

$$H_j \dots H_1 a_i = a_i^j \rightarrow \boxed{v_{j+1} = \begin{pmatrix} 0 \\ \vdots \\ (a_i^j)_i \\ \vdots \\ (a_i^j)_m \end{pmatrix} - \alpha_i e_i}$$

- Store  $v_i$  and  $R$  into  $A$  and an extra  $n$ -vector.
- $Q$  and  $H$  can be computed if needed.
- When column  $i$  is completed, row  $i$  is too.

### 2.3.2 Givens Rotation, $G$

$$\boxed{G = \begin{pmatrix} c & s \\ -s & c \end{pmatrix}} \rightarrow Gx = G \begin{pmatrix} a_1 \\ a_2 \end{pmatrix} = \pm \begin{pmatrix} \|a\| \\ 0 \end{pmatrix}$$

- creates 0's one at a time
- useful for sparse matrices
- When column  $i$  is completed, row  $i$  is too.

$$\rightarrow \boxed{c = \frac{a_1}{\sqrt{a_1^2 + a_2^2}}, \quad s = \frac{a_2}{\sqrt{a_1^2 + a_2^2}}}$$

Avoid squaring any number  $\gg 1$  to prevent overflow/underflow

- $t = \frac{a_2}{a_1} < 1 \rightarrow c = \frac{1}{\sqrt{1+t^2}}, \quad s = c \cdot t$
- $\tau = \frac{a_1}{a_2} < 1 \rightarrow s = \frac{1}{\sqrt{1+\tau^2}}, \quad c = s \cdot \tau$

### 2.3.3 Gram-Schmidt Orthogonalization

Column-Oriented: Build the  $\hat{q}_i$  column vector first, then solve for  $\hat{q}_{i+1}$ , and continue up to  $\hat{q}_n$ .

Classical: Each successive term in the sum for  $q_i$  depends on the original vector,  $a_i$ .

- Extra storage is needed for  $a_i$ .
- Cancellation that causes loss of orthogonality occurs more when ill-conditioned. As a result, performing  $Q_{\parallel}^T b = b'_1$  isn't always best.
- Can't column-pivot, since that depends on rows being completed first.

$$Q_{\parallel}^T = \begin{pmatrix} \hat{q}_1 : q_1 = a_1 \\ \hat{q}_2 : q_2 = a_2 - \hat{q}_1(\hat{q}_1 \cdot a_2) \\ \hat{q}_3 : q_3 = a_3 - \hat{q}_1(\hat{q}_1 \cdot a_3) - \hat{q}_2(\hat{q}_2 \cdot a_3) \\ \vdots \\ \hat{q}_n : q_n = a_n - \sum_j^n \hat{q}_j(\hat{q}_j \cdot a_n) \end{pmatrix}$$

$$R = \begin{pmatrix} \|q_1\| & \hat{q}_1 \cdot a_2 & \hat{q}_1 \cdot a_3 & \dots & \hat{q}_1 \cdot a_n \\ 0 & \|q_2\| & \hat{q}_2 \cdot a_3 & \dots & \hat{q}_2 \cdot a_n \\ 0 & 0 & \|q_3\| & \dots & \hat{q}_3 \cdot a_n \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & \dots & 0 & \|q_n\| \end{pmatrix}$$

Modified: Each successive term in the  $q_i$  sum,  $q_i^{[k+1]}$ , depends on its previous term,  $q_i^{[k]}$ .

$$\left\{ \begin{array}{l} q_i^{[k+1]} \equiv q_i^{[k]} - \hat{q}_k \langle \hat{q}_k | \hat{q}_i^{[k]} \rangle \\ q_i^{[k]} \equiv \hat{q}_k \langle \hat{q}_k | \hat{q}_i^{[k]} \rangle + q_i^{[k+1]} \\ \equiv \sum_{k'=k}^i \hat{q}_{k'} \langle \hat{q}_{k'} | \hat{q}_i^{[k']} \rangle \\ q_1 = x_0 \end{array} \right.$$

- No storage for  $a_i$  is needed.
- Less error in orthog. with finite-precision arithmetic.

$$Q_n^T = \begin{pmatrix} \hat{q}_1 : q_1 = a_1 \\ \hat{q}_2 : q_2 = q_2^{[2]} = \underbrace{a_2 - \hat{q}_1(\hat{q}_1 \cdot a_2)}_{q_2^{[1]}} \\ \hat{q}_3 : q_3 = q_3^{[3]} = \underbrace{a_3 - \hat{q}_1(r_{13}) - \hat{q}_2(\hat{q}_2 \cdot a_3)}_{q_3^{[2]}} \\ \vdots \\ \hat{q}_n : q_n = q_n^{[n]} = a_n - \sum_{m=1}^{n-1} \hat{q}_m(r_{mn}) \end{pmatrix}$$

$$R_n = \begin{pmatrix} \|q_1\| & \langle \hat{q}_1 | \hat{q}_2^{[1]} \rangle & \langle \hat{q}_1 | \hat{q}_3^{[1]} \rangle & \dots & \langle \hat{q}_1 | \hat{q}_n^{[1]} \rangle \\ 0 & \|q_2\| & \langle \hat{q}_2 | \hat{q}_3^{[2]} \rangle & \dots & \langle \hat{q}_2 | \hat{q}_n^{[2]} \rangle \\ 0 & 0 & \|q_3\| & \dots & \langle \hat{q}_3 | \hat{q}_n^{[3]} \rangle \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & \dots & 0 & \|q_n\| \end{pmatrix}$$

Row Oriented: For all  $1 \leq i \leq n$ , solve for  $q_i^{[k]}$  starting first at  $k = 1$ , then continue until  $k = n$ .

- Allows for column pivoting since rows are completed first.
- Cancellation, though still present, is less severe.

Augmented Matrix:

$$(A | b) = (Q_{\parallel} | q_{n+1}) \begin{pmatrix} R & b'_1 \\ 0 & \rho \end{pmatrix}$$

$$\left( \begin{array}{c|c} | & | & | & | \\ a_1 & \dots & a_n & b \end{array} \right) = \left( \begin{array}{c|c} | & | & | & | \\ \hat{q}_1 & \dots & \hat{q}_n & q_{n+1} \end{array} \right) \begin{pmatrix} r_{11} & \dots & r_{1n} & | \\ 0 & \ddots & \vdots & b'_1 \\ \vdots & \ddots & r_{nn} & | \\ 0 & \dots & 0 & \rho \end{pmatrix}$$

- Use Gram-Schmidt QR on this, then solve  $Rx = b'_1$
- This method is preferred numerically to reduce cancelling effects
- Text didn't recommend what  $q_{+1}$  or  $\rho$  should be.
- $\rho$  or  $(q_{n+1})_i$  looks like it should be 0.
- Idk, not much explained.

Reorthogonalizing: Repeating procedure to straighten vectors (usually not needed)

### 2.3.4 Factorization with Column-Pivoting

- Column with largest norm is pivoted to the current column  $i$  to be reduced, and current row  $i$  is completed too.
- Choose the next pivoting column based on norms of the smaller columns from remaining uncompleted submatrix.
- Repeat until the end (rank might be  $n$ ) or if the max norm is smaller than some tolerance (rank might be  $k < n$ )
- Pivoting avoids working with 0's on the diag.

### 2.3.5 Rank Deficiency (or Other) Case

If  $\text{rank}(A) = k < n$  :

$$\begin{array}{cc} (Q^T A P)(P^T x) = Q^T b \\ \begin{pmatrix} R & S \\ 0 & 0' \end{pmatrix} \begin{pmatrix} z \\ 0 \end{pmatrix} = \begin{pmatrix} b'_1 \\ b'_2 \end{pmatrix} \end{array}$$

- $0'$  is approx. 0 since the remaining norms are too small.
- $R = R_{k \times k}$
- $S$  is the remaining columns after  $R$  is completed.
- There are multiple solutions for  $x$ .

- For a quick solution,  $Rz = b'_1, \quad x = P \begin{pmatrix} z \\ 0 \end{pmatrix}$
- For the minimized-norm solution with the smallest  $\|x\|$ ,  $S$  must be annihilated.
- For another method or if underdetermined ( $m < n$ ), something like SVD Decomposition can be used.

## 2.4 Singular Value Decomposition (SVD)

$$\begin{aligned}
 A &= \boxed{U \Sigma V^T} = \left( \begin{array}{c|c|c|c|c} & & & & \\ \hline & & & & \\ & & & & \\ \hline u_1 & \dots & u_k & .. & u_m \\ & & & & \\ & & & & \\ \hline \end{array} \right) \left( \begin{array}{ccccc} \sigma_1 & 0 & \dots & 0 & \\ 0 & \sigma_2 & & \ddots & \vdots \\ 0 & 0 & & \ddots & 0 \\ \vdots & \vdots & & \ddots & \sigma_n \\ 0 & 0 & \dots & 0 & \\ \vdots & \vdots & & \vdots & \\ 0 & 0 & \dots & 0 & \end{array} \right) \left( \begin{array}{c} \text{---} v_1 \text{---} \\ \vdots \\ \text{---} v_k \text{---} \\ \vdots \\ \text{---} v_n \text{---} \end{array} \right) \\
 &= \boxed{U_{\parallel} \Sigma_1 V^T} = (U_{\parallel} \quad U_{\perp}) \begin{pmatrix} \Sigma_1 \\ \mathbf{0} \end{pmatrix} \begin{pmatrix} V_{0\perp}^T \\ V_{0\parallel}^T \end{pmatrix} = \boxed{\sum_i^n \sigma_i \cdot u_i v_i^T}
 \end{aligned}$$

- Underdetermined,  $m < n$  is possible too.
- Analogous to Gaussian-Jordan Diagonalization method.
- $U$  and  $V$  are orthogonal;  $u_i$  and  $v_i$  are the respective “left” and “right” singular vectors.
- Usually, the singular values are ordered such that  $\sigma_1 \geq \sigma_2 \geq \dots$
- $\forall (k < i), \sigma_i = 0 \Rightarrow \text{rank}(A) = k < n$
- $U_{\parallel} = U_{m \times k} : \text{span}(U_{\parallel}) = \text{span}(A), \text{span}(U_{\perp}) = \text{span}(A)^{\perp}$
- $V_{0\perp} = V_{n \times k} : \text{span}(V_{0\parallel}) = \text{null}(A), \text{span}(V_{0\perp}) = \text{null}(A)^{\perp} \quad \text{null}(A) = \{x : Ax = 0\}$

Pseudoinverse:

$$\begin{array}{l} A^+ \equiv V \Sigma^+ U^T \\ \Sigma^+ \equiv \left[ \Sigma^T \text{ and } \sigma_i \rightarrow 1/\sigma_i \quad \forall (\sigma_i \neq 0) \right] \end{array}$$

- $Ax + r = b \rightarrow$

$$\boxed{x = A^+b = (V\Sigma^+U^T)b}$$

- $$x_{\min} = \sum_{\sigma_i \neq 0} \frac{u_i \cdot b}{\sigma_i} v_i$$
 useful for ill-conditioned or rank deficient since small  $\sigma$  can be dropped.

### 2.4.1 Other uses

Euclidean 2-norm :  $\|A\|_2 = \max_{x \neq 0} \frac{\|Ax\|_2}{\|x\|_2} = \sigma_{\max}$

Euclid. Cond. Num. :  $\text{cond}_2(A) = \sigma_{\max}/\sigma_{\min}$

Lower Rank Approx. :  $A \approx A_k = \sum_i^k \sigma_i (u_i v_i^T)$

- Closest rank=  $k$  matrix to  $A$  in the Frobenius norm.
- Frobenius Norm = Euclid. Norm for a “vector” in  $\mathbb{R}^{mn}$ .

Total Least Squares :  $[A \mid y]_{m \times (n+1)} = U \Sigma V_{(n+1) \times (n+1)}^T$

$$\text{rank}([\hat{A} \mid y]) \leq n \rightarrow \sigma_{n+1} = 0 \rightarrow \hat{A} \cdot v_{n+1} = 0$$

$$\left[\widehat{A} \mid y\right] \cdot \begin{bmatrix} x \\ -1 \end{bmatrix} = 0 \rightarrow \begin{bmatrix} x \\ -1 \end{bmatrix} \propto v_{n+1} = \begin{bmatrix} \vec{\nu}_n \\ \nu_{n+1} \end{bmatrix} \rightarrow \boxed{x = \frac{\vec{\nu}_n}{-\nu_{n+1}}}$$

- $\hat{A}$  is an  $A$  with uncertainty, like how  $y$  normally is.

## 2.5 Complexity

Normal, Cholesky

- $A^T A = A'$  costs  $\frac{mn^2}{2}$
- $A' = LL^T$  costs  $\frac{n^3}{6}$
- Rel. Err.  $\propto [\text{cond}(A)]^2$
- Bad if  $\text{cond}(A) \approx 1/\sqrt{\epsilon_{\text{mach}}}$

Householder

- $Q^T A = R$  costs  $mn^2 - \frac{n^3}{3}$
- Rel. Err.  $\propto [\text{cond}(A)]^2 \|r\|_2 + \text{cond}(A)$
- Bad if  $\text{cond}(A) \approx 1/\epsilon_{\text{mach}}$
- More accurate than Cholesky and broadly applicable
- Usable for rank deficient or nearly rank-deficient

Givens

- The normal implementation needs 50% more work than Householder.
- A more complex implementation makes it comparable to Householder.
- Useful if matrix is sparse or zeros need to be maintained.

SVD

- Most expensive cost at  $\propto mn^2 + n^3$ , perhaps 4-10 times or more.
- Robust and reliable.

### 3 Matrix Information

Orthogonal :  $QQ^T = Q^TQ = I$

Unitary :  $UU^\dagger = U^\dagger U = I$

$$\Leftrightarrow U = e^{iH} = (U_h)e^{iD_h}(U_h)^\dagger$$

Normal :  $AA^\dagger = A^\dagger A \Leftrightarrow A = UDU^\dagger$

Symmetric :  $S = S^T = QDQ^T$  ( $D$  is real)

Hermitian :  $H = H^\dagger = UDU^\dagger$  ( $D$  is real)

Similar :  $A(Ty) = \lambda(Ty)$

$$\rightarrow A \sim B = T^{-1}AT$$

Diagonalize :  $T^{-1}AT = D$   $\begin{pmatrix} A \text{ is nondefective} \\ T \text{ is nonsingular} \\ D \text{ is diag.} \end{pmatrix}$

Hessenberg : Triang. but from any diag.

Jordan Form : Nonsing. trans. into a near  
diag. w/ entries in diag. above

•  $Ax = \lambda x$  ,  $\det(A) \neq 0 \Rightarrow A^{-1}x = (1/\lambda)x$

• Shifting :  $(A - \sigma I)x = (\lambda - \sigma)x$

• Simple: Normal [Algebraic] Mult. of 1

• Defective: Geo. Mult. < Alg. Mult.  
(eig. vec. #)

Invariant Subspace:  $\{\mathcal{S} : (A\mathcal{S} \subseteq \mathcal{S}) \equiv (\forall x \in \mathcal{S} \Rightarrow Ax \in \mathcal{S})\}$

#### 3.1 Schur Form

Unitary:  $\boxed{T^\dagger AT = R}$   $\begin{pmatrix} (T \text{ is unitary}) \\ (R \text{ is upp. triang.}) \end{pmatrix}$

Real :  $\boxed{Q^T A Q = R}$   $\begin{pmatrix} (Q \text{ is ortho.}) \\ (R \text{ is block upp. triang.}) \end{pmatrix}$   
(only for real matrix)

•  $A\vec{x}_i = r_{ii}\vec{x}_i$

•  $0 = (R - r_{ii}I) \vec{x}_i$

$$= \begin{pmatrix} R_{11} - r_{ii}I & \vec{u} & R_{13} \\ 0 & 0 & \vec{v}^T \\ \mathcal{O} & 0 & R_{33} - r_{ii}I \end{pmatrix} \vec{x}_i$$

$$\Downarrow$$

$$\vec{x}_i = \begin{pmatrix} \vec{y} \\ -1 \\ 0 \end{pmatrix} : (R_{11} - r_{ii}I) \vec{y} = \vec{u}$$

(nonsingular)

• “Schur form of a real matrix will have complex entries if the matrix has any complex eigenvalues”

•  $R = \begin{pmatrix} R_{11} & \dots & R_{1p} \\ \mathcal{O} & \ddots & \vdots \\ \mathcal{O} & \mathcal{O} & R_{pp} \end{pmatrix}$

•  $\lambda_i(A) = \lambda(R_{ii})$

•  $R_{ii} = \begin{cases} R_{1 \times 1} & \text{(real eigenvalue of } A) \\ R_{2 \times 2} & \text{(complex eigenvalue pairs of } A) \end{cases}$

• All other entries are real

• Reducible: If  $PAP^T = R$  ( $P$  is permu.)

#### Block Upper Triangular Transformation

•  $X_{\parallel} = (\vec{x}_1 \dots \vec{x}_p)$  ,  $Ax_i \in \text{span}(x_i) = \mathcal{S}$

$$\Rightarrow AX_{\parallel} = X_{\parallel}B$$

•  $I_n = X^{-1}X = \begin{pmatrix} Y_{\parallel} \\ Y_{\perp} \end{pmatrix} (X_{\parallel} \ X_{\perp}) = \begin{pmatrix} Y_{\parallel}X_{\parallel} & Y_{\parallel}X_{\perp} \\ Y_{\perp}X_{\parallel} & Y_{\perp}X_{\perp} \end{pmatrix} = \begin{pmatrix} I_p & \mathcal{O} \\ \mathcal{O} & I_{n-p} \end{pmatrix}$

•  $X^{-1}AX = \begin{pmatrix} Y_{\parallel}AX_{\parallel} & Y_{\parallel}AX_{\perp} \\ Y_{\perp}AX_{\parallel} & Y_{\perp}AX_{\perp} \end{pmatrix} = \begin{pmatrix} Y_{\parallel}X_{\parallel}B & Y_{\parallel}AX_{\perp} \\ Y_{\perp}X_{\parallel}B & Y_{\perp}AX_{\perp} \end{pmatrix} = \begin{pmatrix} B & Y_{\parallel}AX_{\perp} \\ \mathcal{O} & Y_{\perp}AX_{\perp} \end{pmatrix} = R$

### 3.2 Upper Hessenberg Transformation (see Householder)

$$\rightarrow H_j^{(k)} = I - \hat{v}\hat{v}^T$$

$$\rightarrow v_{j+1} = \begin{bmatrix} 0 & \dots & 0 & (a_i^j)_k & \dots & (a_i^j)_m \end{bmatrix}^T - \alpha_i e_i$$

- Extra 0s are added to  $v$  to start del. from a diff. row,  $k$ .

## 4 Eigenvalue Equation, $Ax = \lambda x$

### 4.1 Error Bound and Conditioning

$$A + \Delta A = Q(D + \Delta D)Q^{-1}$$

$$(A + \Delta A)(x + \Delta x) = (\lambda + \Delta\lambda)(x + \Delta x)$$

$$\bullet v = (\Delta\lambda I - D)^{-1}(\Delta D)v$$

$$\bullet Ax = \lambda x, \quad y^H A = \lambda y^H$$

$$\bullet \|(\Delta\lambda I - D)^{-1}\|_2^{-1} \leq \|\Delta D\|_2$$

$$\bullet \lambda \text{ is simple} \Rightarrow y^H x \neq 0 \quad (?)$$

$$|\Delta\lambda - \lambda_i| \leq \|Q(\Delta A)Q^{-1}\|_2$$

$$\bullet \cancel{y^H A x} + \cancel{y^H A \Delta x} + y^H (\Delta A)x + y^H (\Delta A)\Delta x$$

$$\approx y^H \lambda x + y^H \lambda \Delta x + y^H (\Delta\lambda)x + y^H (\Delta\lambda)\Delta x$$

$$\rightarrow \boxed{|\Delta\lambda - \lambda_i| \leq \text{cond}(Q) \|\Delta A\|_2}$$

$$\rightarrow \boxed{|\Delta\lambda| \lesssim \frac{\|y\|_2 \cdot \|x\|_2}{|y^H x|} \|\Delta A\|_2 = \frac{1}{\cos \theta} \|\Delta A\|_2}$$

- $AA^\dagger = A^\dagger A \rightarrow \text{cond}(A) = 1$
- Non-simple (multiple) eigenvalue is complicated:
- allows  $y^H x = 0$ , depends on eigenvalue spacings, vector angles, etc.
- Balancing (diagonal rescaling) can improve conditioning

### 4.2 QR Iteration

Power Iteration:

$$A^k x_0 = A^k \sum_{i=1}^n c_i v_i = \sum_{i=1}^n c_i (A^k v_i) = \sum_{i=1}^n c_i (\lambda_i^k v_i)$$

$$= \lambda_1^k \left[ \sum_{i=1}^j c_i v_i + \sum_{i=j+1}^n c_i \left( \frac{\lambda_i}{\lambda_1} \right)^k v_i \right] \quad \begin{matrix} (\lambda_1 = \dots = \lambda_j) \\ (\lambda_1 \geq \dots \geq \lambda_n) \end{matrix}$$

- Eigenvalue converges to  $\lambda_i$  of largest modulus.
- $v_i$  converges to lin. com. if mult. max  $\lambda_i$ .
- Normalize to  $\text{norm}_\infty = 1$  to prevent over/underflow.
- Fails if  $\langle v_i | x_0 \rangle = 0$  (unlikely w/ round. error).
- Real  $A$  and  $x_0$  won't ever converge a complex.
- Convergence Rate:  $C = |\lambda_2/\lambda_1|$

$$\boxed{A^k x_0 = A x_{k-1} = x_k}$$

$$\approx \lambda_1^k \sum_i^j c_i v_i = \lambda_1^k (c_1 v_1) \quad \text{if } j = 1$$

### Power Iteration w/ Shifts:

$$\boxed{A'x_{k-1} = (A - \sigma I)x_{k-1} = x_k} \Rightarrow \boxed{(\lambda'_1 + \sigma) \rightarrow \frac{1}{\lambda'_1 + \sigma}}$$

- Convergence Rate:  $C = \left| \frac{\lambda_2 - \sigma}{\lambda_1 - \sigma} \right| < |\lambda_2/\lambda_1|$
- $x_k \rightarrow v_1$  :  $\sigma = (\lambda_2 + \lambda_n)/2$
- $x_k \rightarrow v_n$  :  $\sigma = (\lambda_1 + \lambda_{n-1})/2$

### Inverse [Power] Iteration:

$$\boxed{(A^{-1})^k x_0 = x_k} \Rightarrow \boxed{(Ax_{k+1} = x_k) \rightarrow (\lambda_1 \rightarrow \frac{1}{\lambda_1})}$$

- Use LU Decomp. or Cholesky to solve for  $x_{k+1}$
- Eigenvalue converges to  $\lambda_i$  of smallest modulus.

### Inverse Iteration w/ Shifts:

$$\boxed{A'x_{k+1} = (A - \sigma I)x_{k+1} = x_k} \Rightarrow \boxed{(\lambda'_1 + \sigma) \rightarrow \frac{1}{\lambda'_1 + \sigma}}$$

- Eigenvalues converges to  $\lambda_i$  closest to  $\sigma$ .
- Using a new shift each iteration requires refactoring each time.

### Rayleigh Quotient (for shifts):

$$\boxed{x\lambda + r = Ax} \Rightarrow \boxed{x^H x \lambda = x^H Ax} \Rightarrow \boxed{\lambda = \frac{\langle x | Ax \rangle}{\langle x | x \rangle}}$$

- Derived with Normal eq. / Least Lin. Regression
- Use  $\lambda$  as  $\sigma$  for shifts

### Deflation (concept):

$$T^{-1}A_k T = \begin{pmatrix} \lambda_1 & b^T \\ 0 & A_{k+1} \end{pmatrix}$$

- Use  $A_{k+1}$  when  $\lambda_1$  is found to good accuracy, then repeat with remaining  $\lambda_i$ .

### Simultaneous/Subspace [Power] Iteration:

$$X_0 = X_{n \times p}, \quad \text{rank}(X_0) = p$$

$$\begin{aligned} \rightarrow & \boxed{X_{k+1} = AX_k = A^{k+1}X_0} \\ \rightarrow & \boxed{\lim_{k \rightarrow \infty} \text{span}(X_k) = \text{span}(A_p)} \end{aligned}$$

- All are found at the same time.
- $\text{span}(A_p)$  is the span of the first  $p$  eigenvectors of largest magnitude.
- Ill-conditioned since all columns of  $X_k$  converge to  $v_1$  (though at different rates so they're still ortho.).
- Normalize to  $\text{norm}_\infty = 1$  to prevent over/underflow.

### Orthogonal [Subspace] Iteration:

$$X = X_{n \times p}, \quad \text{rank}(X) = p$$

$$\begin{aligned} \rightarrow & \boxed{X_k = Q_k R_k} \quad (Q = Q_{n \times p}) \\ \rightarrow & \boxed{X_{k+1} = A Q_k} = Q_k B_{p \times p} \end{aligned}$$

- $\text{span}(Q) = \text{span}(X)$  (see QR Decomp.).
- $B_{p \times p}$  is triag. if  $|\lambda_i| > |\lambda_{i+1}|$ ; else it's block triag. (see Schur Form).
- Orthogonalization is expensive, and convergence may be slow.



### 4.2.1 QR Iteration

$$X_{k+1} \equiv AQ_k \quad (\text{use } X_0 = I \rightarrow X_1 = A)$$

$$Q_{k+1}R_{k+1} = Q_k(Q_k^H A Q_k)$$

$$\rightarrow A_k \equiv Q_k^H X_{k+1} = Q_k^H A Q_k$$

$$= \begin{cases} (Q_k^H Q_{k+1}) R_{k+1} & \equiv Q_{k+1}^{(A)} R_{k+1} \\ Q_k^H (Q_k R_k Q_k^H) Q_k & = R_k Q_k^{(A)} \end{cases}$$

$$\rightarrow A_{k+1} = R_{k+1} Q_{k+1}^{(A)}$$

- $A_k = B_{p \times p}$  (from above), so diag. entries are  $\lambda$ .

- If  $(p = n)$  and  $(X_0 = I)$ , then  $Q_k = Q_1^{(A)} \dots Q_k^{(A)}$  and

$$R^{(k)} = R_k \dots R_1$$

- Since  $(A = Q_1 R_1)$ , then by induction,  $A^k = Q_k R^{(k)}$
- Still, orthogonalization is  $\mathcal{O}(n^3)$  expensive and convergence may be slow.

### Inverse QR Iteration

$$\left( A_k^{-H} = Q_{k+1}^{(A)} R_{k+1}^{-H} \right), \quad \left( A_{k+1}^{-H} = R_{k+1}^{-H} Q_{k+1}^{(A)} \right)$$

$$\rightarrow (A^{-H})^k = Q_k^{(A)} R_k^{-H}$$

- $R^H$  is low. triang.;  $Q^{(A)}$  is built backwards from  $v_n$  to  $v_1$ .
- Columns of  $Q^{(A)}$  for inv. QR iter. of  $A^H$  are the same as QR iter. of  $A$ .
- Means QR iter. of  $A$  is an implicit inv. iter., so shifts are recommended from  $v_n$  to  $v_1$  (see below).

### QR Iteration w/ Shifts & Deflation

$$\text{Rayleigh Quotient: } \sigma^{(n)} = \frac{\langle q_n | A | q_n \rangle}{\langle q_n | q_n \rangle} = \langle q_n | A | q_n \rangle = (A_k)_{nn}$$

$$\rightarrow Q_{k+1}^{(A)} R_{k+1} = A_k - \sigma_k^{(n)} I$$

$$\rightarrow A_{k+1} = R_{k+1} Q_{k+1}^{(A)} + \sigma_k^{(n)} I$$

$$\rightarrow \text{Deflation: } \lim_{k \rightarrow \infty} A_k = \begin{pmatrix} A'_{k+1} & b^T \\ 0 & \sigma_k^{(n)} \end{pmatrix}$$

- Diag. entries of  $A_k$  are automatically Rayleigh Quotients.
- $\sigma$  is the last diag. entry, which corresponds to  $v_n$  (see above).
- If  $\sigma_k = \lambda_n \rightarrow$  the last row of  $Q_{k+1}^{(A)} R_{k+1}$  and  $R_{k+1}$  is 0  $\rightarrow$  the last row of  $R_{k+1} Q_{k+1}^{(A)}$  is 0  $\rightarrow$  Deflation (see right).
- Fails if  $\sigma$  is halfway between two  $\lambda_i$  and favors neither. Or cancellation occurs (rare). Also might require complex arithmetic. Other more robust shifts are available. Convergence is only a few iterations, but cost to factor is still  $\mathcal{O}(n^3)$ .

### Hessenberg QR Iteration [w/ Shifts]

(see Hessenberg Transformation,  $H$ )

$$H = H^{(1)} \Rightarrow H^H A H \text{ isn't Hessenberg at all}$$

$$H = H^{(2)} \Rightarrow H^H A H \text{ is Hessenberg at subdiag.}$$

$$\rightarrow A_1 = H^H A H = H A H$$

$$\rightarrow A_k [-\sigma I] = R_k Q_k^{(A)} \quad (k > 1)$$

$$= R_k \left( A_{k-1} [-\sigma I] \right) R_k^{-1} \text{ is Hessenberg at subdiag.}$$

- Transformation to Hessenberg is done once and costs  $\mathcal{O}(n^3)$ .
- Gives Rotations to factorize a Hessenberg matrix at each iter. costs  $\mathcal{O}(n^2)$ .
- Less iter. needed since already near triangular.
- If  $A = A^H$ , then the transform is tri-diag that costs  $\mathcal{O}(n)$  and  $A_k$  approaches diag.
- Still expensive if  $n$  is large.
- Excessive storage if large and sparse.
- No advantage if only a few  $\lambda_i$  are needed.

### 4.3 Krylov Subspace

Characteristic Polynomial:  $[Ax = \lambda x] \Rightarrow [p_n(\lambda) = 0]$

$$p_n(\lambda) = \left( \sum_{i=0}^{n-1} c_i \lambda^i \right) + \lambda^n = c_0 + c_1 \lambda + \dots + c_{n-1} \lambda^{n-1} + \lambda^n$$

Companion Matrix:  $[p_n(\lambda) = 0] \Rightarrow C_n x = \lambda x$ ,  $C_n = \begin{pmatrix} 0 & 0 & \dots & 0 & -c_0 \\ 1 & 0 & \dots & 0 & -c_1 \\ 0 & 1 & \dots & 0 & -c_2 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & 1 & -c_{n-1} \end{pmatrix} = (e_2 \ e_3 \ \dots \ e_n \ -\vec{c})$

Krylov Matrix:  $K_k = K_{n \times k} = (x_0 \ x_1 \ \dots \ x_{k-1}) = (x_0 \ Ax_0 \ \dots \ A^{k-1}x_0)$  ( $x_0 \neq 0$  vector)

Krylov Subspace:  $[K_k = \text{span}(K_k)] \rightarrow [\mathcal{K}_n = \text{span}(K_n) = \text{span}(A)]$

$$\left. \begin{aligned} AK_n &= (x_1 \ x_2 \ \dots \ x_n) \\ &= K_n (e_2 \ e_3 \ \dots \ \vec{a}) \quad (\vec{a} = K_n^{-1}x_n) \\ &= K_n C_n \end{aligned} \right\} \Rightarrow [K_n^{-1}AK_n = C_n]$$

- Assume  $\exists K^{-1}$
- $C_n$  is Hessenberg
- Successive columns of  $K_n$  converge to  $v_1$ , so problem is ill-conditioned

#### 4.3.1 Arnoldi Iteration

$$K_n = Q_n R_n$$

Arnoldi iteration is simply modified, column-oriented Gram-Schmidt, the benefit being that  $Q_n$  is built by column-vector.  $R_n$  is also built by column when building  $Q_n$ .

$$\left\{ \begin{aligned} q_n^{[i+1]} &\equiv q_n^{[i]} - \hat{q}_i \langle \hat{q}_i | \hat{q}_n^{[i]} \rangle \\ q_n^{[i]} &\equiv \hat{q}_i \langle \hat{q}_i | \hat{q}_n^{[i]} \rangle + q_n^{[i+1]} \\ &\equiv \sum_{i'=i}^n \hat{q}_{i'} \langle \hat{q}_{i'} | \hat{q}_n^{[i]} \rangle \end{aligned} \right\} \left\{ \begin{aligned} q_n &\equiv q_n^{[n]} \\ q_n^{[1]} &\equiv A \hat{q}_{n-1} \\ &= \sum_{i=1}^n \hat{q}_i \langle \hat{q}_i | \hat{q}_n^{[1]} \rangle \\ q_1 &= x_0 \end{aligned} \right.$$

$$Q_n^T = \begin{pmatrix} \hat{q}_1 : q_1 = x_0 \\ \hat{q}_2 : q_2 = q_2^{[2]} = \underbrace{A \hat{q}_1}_{q_2^{[1]}} - \hat{q}_1 \langle \hat{q}_1 | A \hat{q}_1 \rangle \\ \hat{q}_3 : q_3 = q_3^{[3]} = \underbrace{A \hat{q}_2 - \hat{q}_1 \langle \hat{q}_1 | A \hat{q}_2 \rangle}_{q_3^{[2]}} - \hat{q}_2 \langle \hat{q}_2 | q_3^{[2]} \rangle \\ \vdots \\ \hat{q}_k : q_k = q_k^{[k]} = A \hat{q}_{k-1} - \sum_{j=1}^{k-1} \hat{q}_j \langle \hat{q}_j | A \hat{q}_{k-1} \rangle \end{pmatrix}$$

$$R_n = \begin{pmatrix} \|q_1\| & \langle \hat{q}_1 | A \hat{q}_1 \rangle & \langle \hat{q}_1 | q_3^{[1]} \rangle & \dots & \langle \hat{q}_1 | \hat{q}_n^{[1]} \rangle \\ 0 & \|q_2\| & \langle \hat{q}_2 | q_3^{[2]} \rangle & \dots & \langle \hat{q}_2 | \hat{q}_n^{[2]} \rangle \\ 0 & 0 & \|q_3\| & \dots & \langle \hat{q}_3 | \hat{q}_n^{[3]} \rangle \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & \dots & 0 & \|q_n\| \end{pmatrix}$$

$$Q_n^{-1} A Q_n = R_n C_n R_n^{-1} \equiv H \quad \left( \begin{array}{c} \text{Still} \\ \text{Hessenberg} \end{array} \right)$$

$$A Q_n = Q_n H$$

Instead of building  $R_n$  from  $K_n = Q_n R_n$ , build  $H$  from  $A Q_n = Q_n H$  (also by column as  $Q_n$  is built).

$$A \hat{q}_j = \sum_{i=1}^n \hat{q}_i H_{ij} = \sum_{i=1}^n \hat{q}_i \langle \hat{q}_i | A \hat{q}_j \rangle$$

$$\langle \hat{q}_i | A | \hat{q}_j \rangle = H_{ij} = \langle \hat{q}_i | \hat{q}_{j+1}^{[i]} \rangle$$

$$H = \begin{pmatrix} \langle \hat{q}_1 | A | \hat{q}_1 \rangle & \langle \hat{q}_1 | A | \hat{q}_2 \rangle & \langle \hat{q}_1 | A | \hat{q}_3 \rangle & \dots & \langle \hat{q}_1 | \hat{q}_n^{[1]} \rangle \\ \|q_2\| & \langle \hat{q}_2 | q_3^{[2]} \rangle & \langle \hat{q}_2 | \hat{q}_4^{[2]} \rangle & \dots & \langle \hat{q}_2 | \hat{q}_n^{[2]} \rangle \\ 0 & \|q_3\| & \langle \hat{q}_3 | \hat{q}_4^{[3]} \rangle & \dots & \langle \hat{q}_3 | \hat{q}_n^{[3]} \rangle \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & \dots & 0 & \|q_n\| & \langle \hat{q}_n | \hat{q}_{n+1}^{[n]} \rangle \end{pmatrix}$$

### 4.3.2 Rayleigh-Ritz Procedure

$$\begin{array}{l}
 Q_n = [Q_k \ U_k] \\
 Q_k = [\hat{q}_1 \ \hat{q}_2 \ \dots \ \hat{q}_k] \\
 U_k = [\hat{q}_{k+1} \ \dots \ \hat{q}_n]
 \end{array}
 \Rightarrow
 Q_n^H A Q_n = \begin{bmatrix} Q_k^H \\ U_k^H \end{bmatrix} A [Q_k \ U_k]$$

$$= \begin{bmatrix} Q_k^H A Q_k & Q_k^H A U_k \\ U_k^H A Q_k & U_k^H A U_k \end{bmatrix} = H = \begin{bmatrix} H_k & ||| \\ \tilde{H}_k & ||| \end{bmatrix}$$

- $H_k = Q_k^H A Q_k$   
(Is hessenberg)  
(~ Rayleigh Quot.)
- $\tilde{H}_k = \begin{bmatrix} \mathbf{0}^T & \|q_{k+1}\| \\ \mathcal{O} & \mathbf{0} \end{bmatrix}$

- $Q_k$  are the first  $k$  vectors of  $Q_n$  found through Arnoldi Iteration, called Arnoldi vectors.
- The eigenvalues of  $H_k$  are called Ritz values, and its eigenvectors premultiplied by  $Q_k$  are called Ritz vectors.
- At iteration  $H_{k=n} = H$ , then the Ritz values/vectors are equal to the eigenvalues/vectors of  $A$ .
- The Ritz values/vectors of  $H_k$  converge to  $A$ 's eigenvalues/vectors as  $k \rightarrow n$ , and may be found from QR/Orthogonal iteration.
- Extreme eigenvalues are found quickest; shifting + inverting can be used to find interior eigenvalues.
- ( $\|q_{k+1}\| = 0 \rightarrow \tilde{H}_k = 0$ )  $\rightarrow H$  is block triag.,  $\mathcal{K}_k$  is an invar. subspace, and the Ritz values/vectors ARE  $A$ 's eigenvalues/vectors.
- For each  $k$ th iteration, multiplying  $Aq$  costs  $\mathcal{O}(n^2)$ , orthogonalizing costs  $\mathcal{O}(kn)$ , and computing the Ritz values/vectors costs  $\mathcal{O}(k^3)$ . Storage is also needed for  $Q_k$  and  $H_k$ . As this increases with  $k$ , it's better to do a few iterations, then restart the Arnoldi process with a better  $x_0$  based on already computed info.

### 4.3.3 Lanczos Iteration (when $A = A^H$ )

$$\begin{array}{l}
 A = A^H \Rightarrow Q_n^H A Q_n = H_n \text{ hessen. and symm. = tridiag.} \\
 T_k \equiv H_k \text{ (common notation in practice)}
 \end{array}$$

$$Q_n^T = \begin{pmatrix} \hat{q}_1 & : & q_1 & = & x_0 \\ \hat{q}_2 & : & q_2 & = & A\hat{q}_1 - \hat{q}_1 \langle \hat{q}_1 | A \hat{q}_1 \rangle - \overbrace{\hat{q}_0 \langle \hat{q}_1 | \hat{q}_1 \rangle}^{\equiv 0} \\ \hat{q}_3 & : & q_3 & = & A\hat{q}_2 - \hat{q}_2 \langle \hat{q}_2 | A \hat{q}_2 \rangle - \hat{q}_1 \langle q_2 | q_2 \rangle QR \\ \vdots & & & & \\ \hat{q}_{i+1} & : & q_{i+1} & = & A\hat{q}_i - \hat{q}_i \langle \hat{q}_i | A \hat{q}_i \rangle - \hat{q}_{i-1} \langle q_i | q_i \rangle \end{pmatrix}$$

Similar to Arnoldi Iteration in building  $Q_n$  and  $H$  by column and using the Rayleigh-Ritz procedure to estimate the eigenvalues/vectors, but much more simplified with an only three-term reoccurrence for  $Q_n$ , and a symm. tridiag.  $H$  for storage and computation of the Ritz values/vectors.

$$H = T = \begin{pmatrix} \langle \hat{q}_1 | A | \hat{q}_1 \rangle & \langle \hat{q}_2 | \hat{q}_1 \rangle & 0 & 0 \\ \|q_2\| & \langle \hat{q}_2 | A | \hat{q}_2 \rangle & \ddots & 0 \\ 0 & \ddots & \ddots & \|q_n\| \\ 0 & 0 & \|q_n\| & \langle q_n | A | q_n \rangle \end{pmatrix}$$

- Only the diagonal and subdiagonal of  $T$  need to be stored.
- Extreme eigenvalues are found quickest; shifting + inverting can be used to find interior eigenvalues.
- ( $\|q_{k+1}\| = 0 \rightarrow \tilde{T}_k = 0$ )  $\rightarrow T$  is block triag.,  $\mathcal{K}_k$  is an invar. subspace, and the Ritz values/vectors ARE  $A$ 's eigenvalues/vectors.
- Finite-arithmetic means the vectors are more likely to lose orthogonality. Reorthogonalizing them might be expensive, so sometimes it is ignored if the Ritz values are close enough. Not reorthogonalizing might lead to repeated eigenvalues.