



Mathematical  
Institute

# Numerical Linear Algebra

YUJI NAKATSUKASA

*Mathematical Institute, Oxford*

**Computational Techniques**

InFoMM Centre for Doctoral Training

Michaelmas Term 2019, Week 5

Oxford  
Mathematics



# Plan

- ▶ Day 1: SVD and optimality
- ▶ Day 2: LU, QR and linear systems, least-squares
- ▶ Day 3: QR alg, SVD alg (and Krylov)
- ▶ Day 4: Krylov and RandSVD

# References

- ▶ Golub-Van Loan (12): Matrix Computations
  - ▶ classic, encyclopedic
- ▶ Trefethen-Bau (97): Numerical Linear Algebra
  - ▶ covers essentials, beautiful exposition
- ▶ J. Demmel Applied (97): Numerical Linear Algebra
  - ▶ impressive content, some niche
- ▶ N. J. Higham (02), Accuracy and Stability of Algorithms
  - ▶ bible for stability, conditioning
- ▶ Horn and Johnson (12), Matrix Analysis (& topics (86))
  - ▶ amazing theoretical treatise, little numerical treatment

## Linear algebra review

For  $A \in \mathbb{R}^{n \times n}$ , (or  $\mathbb{C}^{n \times n}$ ; hardly makes difference)

The following are equivalent (how many can you name?):

1.  $A$  is nonsingular.

## Linear algebra review

For  $A \in \mathbb{R}^{n \times n}$ , (or  $\mathbb{C}^{n \times n}$ ; hardly makes difference)

The following are equivalent (how many can you name?):

1.  $A$  is nonsingular.
2.  $A$  is invertible:  $A^{-1}$  exists.
3. The map  $A : \mathbb{R}^n \rightarrow \mathbb{R}^n$  is a bijection.
4. all  $n$  eigenvalues of  $A$  are nonzero.
5. all  $n$  singular values of  $A$  are positive.
6.  $\text{rank}(A) = n$ .
7. the rows of  $A$  are linearly independent.
8. the columns of  $A$  are linearly independent.
9.  $Ax = b$  has a solution for every  $b \in \mathbb{C}^n$ .
10.  $A$  has no nonzero null vector. Neither does  $A^T$ .
11.  $A^*A$  is positive definite (not just semidefinite).
12.  $\det(A) \neq 0$ .
13.  $A^{-1}$  exists such that  $A^{-1}A = AA^{-1} = I_n$ .
14. ...

# Important matrices

For square matrices,

- ▶ Symmetric:  $A_{ij} = A_{ji}$  (Hermitian:  $A_{ij} = \bar{A}_{ji}$ )
  - ▶ symmetric positive (semi)definite  $A \succ (\succeq) 0$ : symmetric and positive eigenvalues
- ▶ Orthogonal:  $AA^T = A^T A = I$  (Unitary:  $AA^* = A^* A = I$ )  $\rightarrow$  note  $A^T A = I$  implies  $AA^T = I$
- ▶ Skew-symmetric:  $A_{ij} = -A_{ji}$  (skew-Hermitian:  $A_{ij} = -\bar{A}_{ji}$ )
- ▶ Normal:  $A^T A = AA^T$
- ▶ Tridiagonal:  $A_{ij} = 0$  if  $|i - j| > 1$
- ▶ Triangular:  $A_{ij} = 0$  if  $i > j$

For (possibly nonsquare) matrices  $A \in \mathbb{C}^{m \times n}$ ,  $m \geq n$

- ▶ Hessenberg:  $A_{ij} = 0$  if  $i > j + 1$
- ▶ “orthonormal”:  $A^* A = I_n$ ,
- ▶ sparse: most elements are zero

other structures: Hankel, Toeplitz, circulant, symplectic, ...

# SVD: the most important result in NLA

- ▶ **Symmetric eigenvalue decomposition:**  $A = V\Lambda V^T$   
for symmetric  $A \in \mathbb{R}^{n \times n}$ , where  $V^T V = I_n$ ,  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ .
- ▶ **Singular Value Decomposition (SVD):**  $A = U\Sigma V^T$   
for any  $A \in \mathbb{R}^{m \times n}$ ,  $m \geq n$ . Here  $U^T U = V^T V = I_n$ ,  
 $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_n)$ ,  $\sigma_1 \geq \sigma_2 \geq \dots \sigma_n$ .

SVD proof:

# SVD: the most important result in NLA

- ▶ **Symmetric eigenvalue decomposition:**  $A = V\Lambda V^T$   
for symmetric  $A \in \mathbb{R}^{n \times n}$ , where  $V^T V = I_n$ ,  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ .
- ▶ **Singular Value Decomposition (SVD):**  $A = U\Sigma V^T$   
for any  $A \in \mathbb{R}^{m \times n}$ ,  $m \geq n$ . Here  $U^T U = V^T V = I_n$ ,  
 $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_n)$ ,  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n$ .

SVD proof: Take Gram matrix  $A^T A$  and its eigendecomposition  $A^T A = V\Lambda V^T$ .  $\Lambda$  is nonnegative, and  $(AV)^T(AV)$  is diagonal, so  $AV = U\Sigma$  for some orthonormal  $U$ . Right-multiply  $V^T$ .



# SVD: the most important result in NLA

- ▶ **Symmetric eigenvalue decomposition:**  $A = V\Lambda V^T$   
for symmetric  $A \in \mathbb{R}^{n \times n}$ , where  $V^T V = I_n$ ,  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ .
- ▶ **Singular Value Decomposition (SVD):**  $A = U\Sigma V^T$   
for any  $A \in \mathbb{R}^{m \times n}$ ,  $m \geq n$ . Here  $U^T U = V^T V = I_n$ ,  
 $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_n)$ ,  $\sigma_1 \geq \sigma_2 \geq \dots \sigma_n$ .

SVD proof: Take Gram matrix  $A^T A$  and its eigendecomposition  $A^T A = V\Lambda V^T$ .  $\Lambda$  is nonnegative, and  $(AV)^T(AV)$  is diagonal, so  $AV = U\Sigma$  for some orthonormal  $U$ . Right-multiply  $V^T$ .

SVD and eigendecomposition

- ▶  $V$  eigvecs of  $A^T A$
- ▶  $U$  eigvecs (for nonzero eigvals) of  $AA^T$  (up to sign)
- ▶  $\sigma_i = \sqrt{\lambda_i(A^T A)}$
- ▶ Jordan-Wielandt matrix  $\begin{bmatrix} 0 & A \\ A^T & 0 \end{bmatrix}$ : eigvals  $\pm\sigma_i(A)$ , and  $m - n$  copies of 0. Eigvec matrix is  $\begin{bmatrix} U & U & U_0 \\ V & -V & 0 \end{bmatrix}$ ,  $A^T U_0 = 0$

## Vector norms

For vectors  $x = [x_1, \dots, x_n]^T \in \mathbb{C}^n$

- ▶  $p$ -norm  $\|x\|_p = (|x_1|^p + |x_2|^p + \dots + |x_n|^p)^{1/p}$ 
  - ▶ Euclidean norm=2-norm  $\|x\|_2 = \sqrt{|x_1|^2 + |x_2|^2 + \dots + |x_n|^2}$
  - ▶ 1-norm  $\|x\|_1 = |x_1| + |x_2| + \dots + |x_n|$
  - ▶  $\infty$ -norm  $\|x\|_\infty = \max_i |x_i|$

Inequalities: For  $x \in \mathbb{C}^n$ ,

- ▶  $\frac{1}{\sqrt{n}}\|x\|_2 \leq \|x\|_\infty \leq \|x\|_2$
- ▶  $\frac{1}{\sqrt{n}}\|x\|_1 \leq \|x\|_2 \leq \|x\|_1$
- ▶  $\frac{1}{n}\|x\|_1 \leq \|x\|_\infty \leq \|x\|_1$

$\|\cdot\|_2$  is **unitarily invariant** as  $\|Ux\|_2 = \|x\|_2$  for any unitary  $U$  and any  $x \in \mathbb{C}^n$ .

# Matrix norms

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

- ▶  $p$ -norm  $\|A\|_p = \max_x \frac{\|Ax\|_p}{\|x\|_p}$ 
  - ▶ **2-norm**=spectral norm(=Euclidean norm)  $\|A\|_2 = \sigma_{\max}(A)$   
(largest singular value)
  - ▶ 1-norm  $\|A\|_1 = \max_i \sum_{j=1}^n |A_{ji}|$
  - ▶  $\infty$ -norm  $\|A\|_\infty = \max_i \sum_{j=1}^n |A_{ij}|$
- ▶ **Frobenius norm**  $\|A\|_F = \sqrt{\sum_i \sum_j |A_{ij}|^2}$   
(2-norm of vectorization)
- ▶ **trace norm**=**nuclear norm**  $\|A\|_* = \sum_{i=1}^{\min(m,n)} \sigma_i(A)$

Red: **unitarily invariant** norms  $\|A\|_* = \|UA\|_*$

Inequalities: For  $A \in \mathbb{C}^{m \times n}$ ,

- ▶  $\frac{1}{\sqrt{n}} \|A\|_\infty \leq \|A\|_2 \leq \sqrt{m} \|A\|_\infty$
- ▶  $\frac{1}{\sqrt{m}} \|A\|_1 \leq \|A\|_2 \leq \sqrt{n} \|A\|_1$
- ▶  $\|A\|_2 \leq \|A\|_F \leq \sqrt{\min(m,n)} \|A\|_2$

## Optimal low-rank approximation by SVD

Truncated SVD:  $A_r = U_r \Sigma_r V_r^T$ ,  $\Sigma_r = \text{diag}(\sigma_1, \dots, \sigma_r)$

$$\|A - A_r\|_2 = \sigma_{r+1} = \min_{\text{rank}(B)=r} \|A - B\|_2$$

- Storage savings: if  $\sigma_{r+1} \ll \sigma_1$ ,  $A \approx A_r$  with

$$A \approx A_r = U \Sigma V^T$$

- Optimality holds for any unitarily invariant norm

## SVD optimality proof in 2-norm

Truncated SVD:  $A_r = U_r \Sigma_r V_r^T$ ,  $\Sigma_r = \text{diag}(\sigma_1, \dots, \sigma_r)$

$$\|A - A_r\|_2 = \sigma_{r+1} = \min_{\text{rank}(B)=r} \|A - B\|_2$$

## SVD optimality proof in 2-norm

Truncated SVD:  $A_r = U_r \Sigma_r V_r^T$ ,  $\Sigma_r = \text{diag}(\sigma_1, \dots, \sigma_r)$

$$\|A - A_r\|_2 = \sigma_{r+1} = \min_{\text{rank}(B)=r} \|A - B\|_2$$

- Since  $\text{rank}(B) \leq r$ , we can write  $B = B_1 B_2^T$  where  $B_1, B_2$  have  $r$  columns.

## SVD optimality proof in 2-norm

Truncated SVD:  $A_r = U_r \Sigma_r V_r^T$ ,  $\Sigma_r = \text{diag}(\sigma_1, \dots, \sigma_r)$

$$\|A - A_r\|_2 = \sigma_{r+1} = \min_{\text{rank}(B)=r} \|A - B\|_2$$

- ▶ Since  $\text{rank}(B) \leq r$ , we can write  $B = B_1 B_2^T$  where  $B_1, B_2$  have  $r$  columns.
- ▶ There exists orthogonal  $W \in \mathbb{C}^{n \times (n-r)}$  s.t.  $BW = 0$ . Then  $\|A - B\|_2 \geq \|(A - B)W\|_2 = \|AW\|_2 = \|U \Sigma (V^* W)\|_2$ .

## SVD optimality proof in 2-norm

Truncated SVD:  $A_r = U_r \Sigma_r V_r^T$ ,  $\Sigma_r = \text{diag}(\sigma_1, \dots, \sigma_r)$

$$\|A - A_r\|_2 = \sigma_{r+1} = \min_{\text{rank}(B)=r} \|A - B\|_2$$

- ▶ Since  $\text{rank}(B) \leq r$ , we can write  $B = B_1 B_2^T$  where  $B_1, B_2$  have  $r$  columns.
- ▶ There exists orthogonal  $W \in \mathbb{C}^{n \times (n-r)}$  s.t.  $BW = 0$ . Then  $\|A - B\|_2 \geq \|(A - B)W\|_2 = \|AW\|_2 = \|U\Sigma(V^*W)\|_2$ .
- ▶ Now since  $W$  is  $(n - r)$ -dimensional, there is an intersection between  $W$  and  $[v_1, \dots, v_{r+1}]$ , the  $(r + 1)$ -dimensional subspace spanned by the leading  $r + 1$  left singular vectors ( $[W, v_1, \dots, v_{r+1}]\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 0$  has a solution; then  $Wx_1$  is such a vector).



## SVD optimality proof in 2-norm

Truncated SVD:  $A_r = U_r \Sigma_r V_r^T$ ,  $\Sigma_r = \text{diag}(\sigma_1, \dots, \sigma_r)$

$$\|A - A_r\|_2 = \sigma_{r+1} = \min_{\text{rank}(B)=r} \|A - B\|_2$$

- ▶ Since  $\text{rank}(B) \leq r$ , we can write  $B = B_1 B_2^T$  where  $B_1, B_2$  have  $r$  columns.
- ▶ There exists orthogonal  $W \in \mathbb{C}^{n \times (n-r)}$  s.t.  $BW = 0$ . Then  $\|A - B\|_2 \geq \|(A - B)W\|_2 = \|AW\|_2 = \|U\Sigma(V^*W)\|_2$ .
- ▶ Now since  $W$  is  $(n - r)$ -dimensional, there is an intersection between  $W$  and  $[v_1, \dots, v_{r+1}]$ , the  $(r + 1)$ -dimensional subspace spanned by the leading  $r + 1$  left singular vectors ( $[W, v_1, \dots, v_{r+1}]\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 0$  has a solution; then  $Wx_1$  is such a vector).
- ▶ Then scale  $x_1$  to have unit norm, and  $\|U\Sigma V^*Wx_1\|_2 = \|U\Sigma_{r+1}y_1\|_2$ , where  $\|y_1\|_2 = 1$  and  $\Sigma_{r+1}$  is the leading  $r + 1$  part of  $\Sigma$ . Then  $\|U\Sigma_{r+1}y_1\|_2 \geq \sigma_{r+1}$  can be verified directly.

# Matrix decompositions

- ▶ **SVD**  $A = U\Sigma V^T$
- ▶ Eigenvalue decomposition  $A = X\Lambda X^{-1}$ 
  - ▶ **Normal**:  $X$  unitary  $X^*X = I$
  - ▶ **Symmetric**:  $X$  unitary and  $\Lambda$  real
- ▶ Jordan decomposition:  $A = XJX^{-1}$ ,  
$$J = \text{diag}\left(\begin{bmatrix} \lambda_i & 1 & & \\ & \lambda_i & \ddots & \\ & & \ddots & 1 \\ & & & \lambda_i \end{bmatrix}\right)$$
- ▶ **Schur** decomposition  $A = QTQ^*$ :  $T$  upper triangular
- ▶ **QR**: orthonormal,  $U$ : upper triangular
- ▶ **LU**:  $L$ : lower triangular,  $U$ : upper triangular

**Red**: Orthogonal decompositions, stable computation available

# Numerical stability

For computational task  $Y = f(X)$  and computed approximant  $\hat{Y}$ ,

- ▶ Ideally, error  $\|Y - \hat{Y}\|/\|Y\| = O(u)$ : seldom true  
( $u$ : unit roundoff,  $\approx 10^{-16}$  in standard double precision)
- ▶ Good alg. has **Backward stability**  $\hat{Y} = f(X + \Delta X)$ ,  
 $\frac{\|X - \hat{X}\|}{\|X\|} = O(u)$  “exact solution of slightly wrong input”
- ▶ Forward stability  $\|Y - \hat{Y}\|/\|Y\| = O(\kappa(f)u)$  “error is as small as backward stable alg.” (sometimes used to mean small error; we follow Higham’s book [2002])
- ▶ Most important condition number:

$$\kappa_2(A) = \frac{\sigma_{\max}(A)}{\sigma_{\min}(A)} (\geq 1)$$

e.g. for linear systems. A backward stable soln for  $Ax = b$ ,  
i.e.,  $(A + \Delta A)\hat{x} = (b + \Delta b)$  satisfies

$$\frac{\|\hat{x} - x\|}{\|x\|} \lesssim u\kappa_2(A)$$

# LU decomposition

$$A = LU \in \mathbb{R}^{n \times n}$$

$L$ : lower triangular,  $U$ : upper triangular

- ▶ Cost  $\frac{2}{3}n^3$  flops
- ▶ For  $Ax = b$ ,
  - ▶ first solve  $Ly = b$ , then  $Ux = y$ .
  - ▶ triangular solve is always backward stable: e.g.  $(L + \Delta L)\hat{y} = b$  (see Higham's book)
- ▶ **Pivoting** crucial for numerical stability:  $PA = LU$ , where  $P$ : permutation matrix. Then stability means  $\hat{L}\hat{U} = PA + \Delta A$ 
  - ▶ Even with pivoting, unstable examples exist, but still always stable in practice and used everywhere!
- ▶ Special case where  $A \succ 0$  positive definite:  $A = R^T R$ , **Cholesky** factorization, ALWAYS stable,  $\frac{1}{3}n^3$  flops

# QR decomposition

$$A = QR \in \mathbb{R}^{m \times n}$$

$Q \in \mathbb{R}^{m \times n}$ : orthonormal,  $R \in \mathbb{R}^{n \times n}$ : upper triangular

- ▶ Many algorithms available: Gram-Schmidt, Householder, CholeskyQR, ...
- ▶ For  $Ax = b$ ,
  - ▶ solve  $Rx = Q^T b$ : always stable! But LU used for speed, as QR needs  $\frac{4}{3}n^3$  flops
- ▶ Pivoting  $A = QRP$  not needed for numerical stability
  - ▶ but pivoting gives rank-revealing QR

# Householder QR factorization

Householder reflectors:

$$Q = I - 2vv^T, \quad v = \frac{x - \|x\|_2 e}{\|x - \|x\|_2 e\|_2}$$

satisfies  $Qv = e$

$\Rightarrow$  find  $Q_1$  s.t.  $Q_1 A(:, 1) = e = [1, 0, \dots]^T \|A(:, 1)\|_2$ , repeat to get  $Q_n \cdots Q_2 Q_1 A = R$  upper triangular, then

$$A = (Q_1^T \cdots Q_{n-1}^T Q_n^T) R = QR$$

Properties

- ▶ Cost  $\frac{4}{3}n^3$  flops with Householder-QR (twice that of LU)
- ▶ Unconditionally backward stable:  $\hat{Q}\hat{R} = A + \Delta A$ ,  
 $\|\hat{Q}^T \hat{Q} - I\|_2 = O(u)$
- ▶  $Q_i$  orthogonal+symmetric, eigvals 1 ( $n - 1$  copies) and  $-1$

## Least-squares problem

$$\min_x \|Ax - b\|_2, \quad A \in \mathbb{R}^{m \times n}, m \geq n$$

## Least-squares problem

$$\min_x \|Ax - b\|_2, \quad A \in \mathbb{R}^{m \times n}, m \geq n$$

Let  $A = [Q \ Q_\perp] \begin{bmatrix} R \\ 0 \end{bmatrix}$  be 'full' QR factorization (useful for theory, seldom used in computation). Then

$$\|Ax - b\|_2 = \|Q^T(Ax - b)\|_2 = \left\| \begin{bmatrix} R \\ 0 \end{bmatrix} x - \begin{bmatrix} Q^T b \\ Q_\perp^T b \end{bmatrix} \right\|_2$$

so  $x = R^{-1}Q^T b$  is solution. This also gives algorithm:



## Least-squares problem

$$\min_x \|Ax - b\|_2, \quad A \in \mathbb{R}^{m \times n}, m \geq n$$

Let  $A = [Q \ Q_\perp] \begin{bmatrix} R \\ 0 \end{bmatrix}$  be 'full' QR factorization (useful for theory, seldom used in computation). Then

$$\|Ax - b\|_2 = \|Q^T(Ax - b)\|_2 = \left\| \begin{bmatrix} R \\ 0 \end{bmatrix} x - \begin{bmatrix} Q^T b \\ Q_\perp^T b \end{bmatrix} \right\|_2$$

so  $x = R^{-1}Q^T b$  is solution. This also gives algorithm:

1. Compute (thin) QR factorization  $A = QR$
  2. Solve linear system  $Rx = Q^T b$ .
- ▶ This is backward stable: computed  $\hat{x}$  solution for  $\min_x \|(A + \Delta A)x + (b + \Delta b)\|_2$  (see Higham's book Ch.20)
  - ▶ Mathematically,  $x$  satisfies normal equation  $(A^T A)x = A^T b$ , but this is NOT backward stable

# Power method for eigenproblems

$x := \text{random vector}$ ,  $x = Ax$ ,  $x = \frac{x}{\|x\|}$ ,  $\lambda = x^T Ax$ , repeat

- ▶ Basis for QR algorithm, Krylov methods (Lanczos, Arnoldi,...)
- ▶ Convergence analysis: let  $x_0 = \sum_{i=1}^n c_i v_i$ ,  $Av_i = \lambda_i v_i$  with  $|\lambda_1| > |\lambda_2| > \dots$ . Then after  $k$  iterations,

$$x = C \sum_{i=1}^n \left( \frac{\lambda_i}{\lambda_1} \right)^k c_i v_i \rightarrow C c_1 v_1 \quad \text{as } k \rightarrow \infty$$

- ▶ Converges **geometrically**  $(\lambda, x) \rightarrow (\lambda_1, x_1)$  with **linear rate**  $\frac{|\lambda_2|}{|\lambda_1|}$
- ▶ What does this imply about  $A^n = QR$  as  $n \rightarrow \infty$ ?

# Power method for eigenproblems

$x := \text{random vector}$ ,  $x = Ax$ ,  $x = \frac{x}{\|x\|}$ ,  $\lambda = x^T Ax$ , repeat

- ▶ Basis for QR algorithm, Krylov methods (Lanczos, Arnoldi,...)
- ▶ Convergence analysis: let  $x_0 = \sum_{i=1}^n c_i v_i$ ,  $Av_i = \lambda_i v_i$  with  $|\lambda_1| > |\lambda_2| > \dots$ . Then after  $k$  iterations,

$$x = C \sum_{i=1}^n \left( \frac{\lambda_i}{\lambda_1} \right)^k c_i v_i \rightarrow C c_1 v_1 \quad \text{as } k \rightarrow \infty$$

- ▶ Converges **geometrically**  $(\lambda, x) \rightarrow (\lambda_1, x_1)$  with **linear rate**  
 $\frac{|\lambda_2|}{|\lambda_1|}$
- ▶ What does this imply about  $A^n = QR$  as  $n \rightarrow \infty$ ? First vector of  $Q \rightarrow v_1$

# Power method for eigenproblems

$x := \text{random vector}$ ,  $x = Ax$ ,  $x = \frac{x}{\|x\|}$ ,  $\lambda = x^T Ax$ , repeat

- ▶ Basis for QR algorithm, Krylov methods (Lanczos, Arnoldi,...)
- ▶ Convergence analysis: let  $x_0 = \sum_{i=1}^n c_i v_i$ ,  $Av_i = \lambda_i v_i$  with  $|\lambda_1| > |\lambda_2| > \dots$ . Then after  $k$  iterations,

$$x = C \sum_{i=1}^n \left( \frac{\lambda_i}{\lambda_1} \right)^k c_i v_i \rightarrow C c_1 v_1 \quad \text{as } k \rightarrow \infty$$

- ▶ Converges **geometrically**  $(\lambda, x) \rightarrow (\lambda_1, x_1)$  with **linear rate**  $\frac{|\lambda_2|}{|\lambda_1|}$
- ▶ What does this imply about  $A^n = QR$  as  $n \rightarrow \infty$ ?

Inverse power method:  $x := (A - \mu I)x$ ,  $x = x/\|x\|$

- ▶ Converges with improved **linear rate**  $\frac{|\lambda_{\sigma(2)} - \mu|}{|\lambda_{\sigma(1)} - \mu|}$  to eigval closest to  $\mu$  ( $\sigma$ : permutation)
- ▶  $\mu$  can change adaptively with the iterations. The choice  $\mu := x^T Ax$  gives Rayleigh quotient iteration, with **quadratic** convergence (cubic if  $A$  symmetric)

# QR algorithm for eigenproblems

$$A_1 = Q_1 R_1, A_2 = R_1 Q_1, A_2 = Q_2 R_2, A_3 = R_2 Q_2, \dots$$

- ▶ Basically:  $QR \rightarrow RQ \rightarrow QR \rightarrow RQ \rightarrow \dots$  **triangular**
- ▶ Fundamental work by Francis (61,62) and Kublanovskaya (63)
- ▶ Truly **Magical** algorithm!
  - ▶ backward stable, as based on orthogonal transforms
  - ▶ always converges, but global proof unavailable(!)
  - ▶ uses 'inverse power method' (rational funcs) without inversions

Two techniques to speed up from  $> O(n^4)$  to  $O(n^3)$

- ▶ Initial reduction to **Hessenberg** form via Unitary transform

$$A = \begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \end{bmatrix} \xrightarrow{Q_1} \begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \end{bmatrix} \xrightarrow{Q_2} \begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \end{bmatrix} \xrightarrow{Q_3} \begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \end{bmatrix}$$

- ▶ **shift**  $A_k - s_k I = Q_k R_k, A_{k+1} = R_k Q_k + s_k I$ , repeat  
(effective choice:  $s_k = A_k(n, n)$ )

# QR algorithm for eigenproblems

$$A_1 = Q_1 R_1, A_2 = R_1 Q_1, A_2 = Q_2 R_2, A_3 = R_2 Q_2, \dots$$

- ▶ Basically:  $QR \rightarrow RQ \rightarrow QR \rightarrow RQ \rightarrow \dots$  **triangular**
- ▶ Fundamental work by Francis (61,62) and Kublanovskaya (63)
- ▶ Truly **Magical** algorithm!
  - ▶ backward stable, as based on orthogonal transforms
  - ▶ always converges, but global proof unavailable(!)
  - ▶ uses 'inverse power method' (rational funcs) without inversions

Two techniques to speed up from  $> O(n^4)$  to  $O(n^3)$

- ▶ Initial reduction to **Hessenberg** form via Unitary transform

$$A = \begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \end{bmatrix} \xrightarrow{Q_1} \begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \end{bmatrix} \xrightarrow{Q_2} \begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \end{bmatrix} \xrightarrow{Q_3} \begin{bmatrix} * & * & * & * & * \\ \star & * & * & * & * \\ * & \star & * & * & * \\ * & * & \star & * & * \\ * & * & * & \star & * \end{bmatrix}$$

- ▶ **shift**  $A_k - s_k I = Q_k R_k, A_{k+1} = R_k Q_k + s_k I$ , repeat  
(effective choice:  $s_k = A_k(n, n)$ )

## QR algorithm and power method

QR algorithm:  $A_k = Q_k R_k$ ,  $A_{k+1} = R_k Q_k$

$$A^k = (Q_1 \cdots Q_k)(R_k \cdots R_1) = Q^{(k)} R^{(k)}.$$

Proof by induction: Suppose  $A^{k-1} = Q^{(k-1)} R^{(k-1)}$ . Then  $A_k = R_{k-1} Q_{k-1} = (Q^{(k-1)})^* A Q^{(k-1)}$ , and

$$(Q^{(k-1)})^* A Q^{(k-1)} = Q_k R_k.$$

Then  $A Q^{(k-1)} = Q^{(k-1)} Q_k R_k$ , and so

$$A^k = A Q^{(k-1)} R^{(k-1)} = Q^{(k-1)} Q_k R_k R^{(k-1)} = Q^{(k)} R^{(k)} \square$$

Now take inverse:  $A^{-k} = (R^{(k)})^{-1} (Q^{(k)})^*$ ,

Conjugate transpose:  $(A^{-k})^* = Q^{(k)} (R^{(k)})^{-*}$

$\Rightarrow$  QR factorization of matrix with eigvals  $r(\lambda_i) = \lambda_i^{-k}$

$\Rightarrow$  Connection also with (unshifted) **inverse** power method

NB no matrix inverse performed

## QR algorithm with shifts and shifted inverse power method

1.  $A_k - s_k I = Q_k R_k$  (QR factorization)
2.  $A_{k+1} = R_k Q_k + s_k I$ ,  $k \leftarrow k + 1$ , repeat.



# QR algorithm with shifts and shifted inverse power method

1.  $A_k - s_k I = Q_k R_k$  (QR factorization)
2.  $A_{k+1} = R_k Q_k + s_k I$ ,  $k \leftarrow k + 1$ , repeat.

$$\prod_{i=1}^k (A - s_i I) = Q^{(k)} R^{(k)} (= (Q_1 \cdots Q_k)(R_k \cdots R_1))$$

Proof: Suppose true for  $k - 1$ . Then QR alg. computes

$(Q^{(k-1)})^*(A - s_k I)Q^{(k-1)} = Q_k R_k$ , so

$(A - s_k I)Q^{(k-1)} = Q^{(k-1)}Q_k R_k$ , hence

$$\prod_{i=1}^k (A - s_i I) = (A - s_k I)Q^{(k-1)}R^{(k-1)} = Q^{(k-1)}Q_k R_k R^{(k-1)} = Q^{(k)} R^{(k)}.$$

Inverse conjugate transpose:  $\prod_{i=1}^k (A - s_i I)^{-*} = Q^{(k)}(R^{(k)})^{-*}$

$\Rightarrow$  QR factorization of matrix with eigvals  $r(\lambda_j) = \prod_{i=1}^k \frac{1}{\lambda_j - s_i}$

$\Rightarrow$  Connection with **shifted inverse** power method, hence

**rational approximation**

# QR algorithm for symmetric $A$

- Initial reduction to Hessenberg form  $\rightarrow$  tridiagonal

$$A = \begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \end{bmatrix} \xrightarrow{Q_1} \begin{bmatrix} * & * & & & \\ * & * & * & * & * \\ & * & * & * & * \\ & & * & * & * \\ * & * & * & * & * \end{bmatrix} \xrightarrow{Q_2} \begin{bmatrix} * & * & & & \\ * & * & * & & \\ & * & * & * & * \\ & & * & * & * \\ & & & * & * \end{bmatrix} \xrightarrow{Q_3} \begin{bmatrix} * & * & & & \\ * & * & * & & \\ & * & * & * & \\ & & * & * & * \\ & & & * & * \end{bmatrix}$$

- QR steps for tridiagonal:  $O(n)$  instead of  $O(n^2)$
- powerful alternatives available for tridiagonal eigenproblem (QR, divide-conquer ([Gu-Eisenstat 95]), HODLR ([Kressner-Susnjara 19], exploit low-rank structure),...)
- Cost:  $\frac{4}{3}n^3$  flops for eigvals,  $\approx 10n^3$  for eigvecs
- Since I am speaking, also mention spectral divide-and-conquer (w/ Freund, Higham): all about **rational approximation**

$$A = \begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \end{bmatrix} \xrightarrow{V_1} \begin{bmatrix} * & * & * & & \\ * & * & * & & \\ * & * & * & & \\ & * & * & * & * \\ & & * & * & * \end{bmatrix} \xrightarrow{V_2} \begin{bmatrix} * & & & & \\ & * & * & & \\ & * & * & & \\ & & * & * & \\ & & & * & * \end{bmatrix} \xrightarrow{V_3} \begin{bmatrix} * & & & & \\ & * & & & \\ & & * & & \\ & & & * & \\ & & & & * \end{bmatrix} = \Lambda.$$

# Golub-Kahan for SVD

Apply Householder reflectors from left and right (different ones) to **bidiagonalize**

$$A \xrightarrow{H_{L,1}} \begin{bmatrix} \star & \star & \star & \star \\ & \star & \star & \star \\ \star & \star & \star & \\ \star & \star & \star & \\ \star & \star & \star & \end{bmatrix} \xrightarrow{H_{R,1}} \begin{bmatrix} \star & \star & & \\ & \star & \star & \star \\ \star & \star & \star & \\ \star & \star & \star & \\ \star & \star & \star & \end{bmatrix} \xrightarrow{H_{L,2}} \begin{bmatrix} \star & \star & & \\ & \star & \star & \star \\ & \star & \star & \\ \star & \star & \star & \\ & \star & \star & \end{bmatrix} \xrightarrow{H_{R,2}} \begin{bmatrix} \star & \star & & \\ & \star & \star & \\ & \star & \star & \star \\ & \star & \star & \\ \star & \star & \star & \end{bmatrix} \xrightarrow{H_{L,3}} \begin{bmatrix} \star & \star & & \\ & \star & \star & \\ & \star & \star & \star \\ & \star & \star & \\ & \star & \star & \end{bmatrix} \xrightarrow{H_{R,3}} \begin{bmatrix} \star & \star & & \\ & \star & \star & \\ & \star & \star & \star \\ & \star & \star & \\ & \star & \star & \end{bmatrix} \xrightarrow{H_{L,4}} B,$$

- ▶ Once bidiagonalized,
  - ▶ Mathematically, QR on  $B^T B$
  - ▶ More elegant: dqds algorithm [Fernando-Parlett 1994]
- ▶ Cost:  $\approx 4mn^2$  flops for singvals  $\Sigma$ ,  $\approx 20mn^2$  flops for singvecs  $U, V$

## Polynomial rootfinding $p(x) = 0$ via eigenvalues

►  $p(x) = x^n + a_{n-1}x^{n-1} + \cdots + a_1x + a_0$

$p(\lambda) = 0 \Leftrightarrow \lambda$  eigenvalue of

$$C = \begin{bmatrix} -a_{n-1} & -a_{n-2} & \cdots & -a_1 & -a_0 \\ 1 & & & & \\ & 1 & & & \\ & & \ddots & & \\ & & & 1 & 0 \end{bmatrix} \in \mathbb{C}^{n \times n}$$

►  $\tilde{p}(x) = T_n(x) + a_{n-1}T_{n-1}(x) + \cdots + a_1T_1(x) + a_0T_0(x)$ ,  
 $T_i(x)$ : Chebyshev polynomial.  $\tilde{p}(\lambda) = 0 \Leftrightarrow \lambda$  eigenvalue of

$$\tilde{C} = \frac{1}{2} \begin{bmatrix} -a_{n-1} & 1 - a_{n-2} & -a_{n-3} & \cdots & -a_0 \\ 1 & 0 & 1 & & \\ & \ddots & \ddots & \ddots & \\ & & 1 & 0 & 1 \\ & & & 2 & 0 \end{bmatrix} \in \mathbb{C}^{n \times n}$$

Powerful approach for nonlinear problems: **approximate with polynomial, and solve eigenproblem**

## Exotic but tractable eigenvalue problems

- ▶ Standard eigenvalue problem  $Ax = \lambda x$ 
  - ▶  $A$  symmetric: tridiagonal QR algorithm
  - ▶ nonsymmetric: Hessenberg QR algorithm
  - ▶ QZ algorithm: QR applied to  $B^{-1}A$  implicitly
- ▶ Polynomial eigenvalue problem  $P(\lambda)x = 0$ , e.g.  
 $P(\lambda) = \lambda^2 A + \lambda B + C$ 
  - ▶ usually linearization + QZ
$$\begin{bmatrix} -B & -C \\ I & 0 \end{bmatrix} \begin{bmatrix} \lambda x \\ x \end{bmatrix} = \lambda \begin{bmatrix} A & 0 \\ 0 & I \end{bmatrix} \begin{bmatrix} \lambda x \\ x \end{bmatrix}$$

All these are tractable!  $O(n^3)$  cost (or  $O((nd)^3)$ ,  $d$ : degree) or less

- ▶ Nonlinear eigenvalue problem:  $F(\lambda)x = 0$ , e.g.  
 $F(\lambda) = \exp(\lambda)A + \log(\lambda)B + C$ 
  - ▶ usually: local approximation  $F(\lambda) \approx P(\lambda)$  + linearization + QZ

# Krylov subspace methods

Idea: find solution  $\hat{x}$  in subspace

$$\mathcal{K}_k(A, b) := \text{span}([b, Ab, A^2b, \dots, A^{k-1}b]).$$

equivalent to  $\hat{x} = p_{k-1}(A)b$ ,  $p_{k-1}$  : polynomial of degree  $k - 1$

Why should it work?

- ▶ For eigenproblems  $Ax = \lambda x$ , when looking for dominant eigvals
  - ▶ or better yet, work with  $(A - \mu I)^{-1}$ : shift-invert Arnoldi
- ▶ For linear systems  $Ax = b$ , fast convergence roughly when  $\kappa_2(A) = O(1)$
- ▶ Very useful to understand from **polynomial(/rational) approximation** viewpoint:  $x \approx \hat{x} = p_{k-1}(A)b$

## Arnoldi process

Denote by  $Q \in \mathbb{C}^{n \times k}$  orthonormal that spans ( $k = 1, 2, \dots$ )

$$\mathcal{K}_k(A, b) := \text{span}([b, Ab, A^2b, \dots, A^{k-1}b])$$

(note  $Q(:, 1) = b/\|b\|_2$ ). Then consider matrix  $AQ$ . Careful consideration reveals identity

$$AQ = QH + q_{k+1}[0, \dots, 0, h_{k+1,k}],$$

where  $H$  is **upper Hessenberg**

- ▶  $i$ th column yields  $Aq_i = \sum_{j=1}^i H_{ji}q_j + H_{i+1,i}q_{j+1}$ : obtain  $H_{ji}$  by orthogonalization (Householder or Gram-Schmidt)
- ▶ far superior (in stability) to forming matrix  $[b, Ab, A^2b, \dots, A^{k-1}b]$  then computing QR

# Lanczos

When  $A$  symmetric, Arnoldi simplifies to

$$AQ = QT + q_{k+1}[0, \dots, 0, h_{k+1,k}],$$

where  $T$  is **tridiagonal**

- ▶ three-term recurrence, orthogonalize necessary only against past two vecs  $q_i, q_{i-1}$

Lanczos' algorithm for symmetric eigenproblem:

- ▶ Conceptually, find  $Q$  and do **Rayleigh-Ritz**: compute eigvals of  $Q^T A Q$ : projection method
- ▶ We actually have  $Q^T A Q = T$ , tridiagonal eigenproblem



## GMRES for $Ax = b$

Conceptually, solve

$$\min_{x \in \mathcal{K}_k(A, b)} \|Ax - b\|_2$$

Given  $AQ = QH + q_{k+1}h_{k+1,k}e_{k+1}$  where

$e_{k+1} := [0, \dots, 0, h_{k+1,k}]$ , equivalent to (same trick as in least-squares)

$$\begin{aligned}\min_y \|AQy - b\|_2 &= \min_y \|(QH + q_{k+1}h_{k+1,k}e_{k+1})y - b\|_2 \\ &= \min_y \left\| \begin{bmatrix} H \\ h_{k+1,k}e_{k+1} \end{bmatrix} y - \begin{bmatrix} Q^T \\ q_{k+1}^T \end{bmatrix} b \right\|_2\end{aligned}$$

Solve final problem via QR (Givens rotations)+triangular solve,  
 $O(n^2)$

## GMRES convergence

Recall that  $x \in \mathcal{K}_k(A, b) \Rightarrow x = p_{k-1}(A)b$ . Hence GMRES solution is

$$\begin{aligned}\min_{x \in \mathcal{K}_k(A, b)} \|Ax - b\|_2 &= \min_{p_{k-1} \in \mathcal{P}_{k-1}} \|Ap_{k-1}(A)b - b\|_2 \\ &= \min_{\tilde{p} \in \mathcal{P}_k, \tilde{p}(0)=0} \|(\tilde{p}(A) - I)b\|_2 \\ &= \min_{p \in \mathcal{P}_k, p(0)=1} \|p(A)b\|_2\end{aligned}$$

If  $A$  diagonalizable  $A = X\Lambda X^{-1}$ ,

$$\begin{aligned}\|p(A)\|_2 &= \|Xp(\Lambda)X^{-1}\|_2 \leq \|X\|_2 \|X^{-1}\|_2 \|p(\Lambda)\|_2 \\ &= \kappa_2(X) \max_{z \in \lambda(A)} |p(z)|\end{aligned}$$

Interpretation: find polynomial s.t.  $p(0) = 1$  and  $|p(\lambda_i)|$  small (demo)

## CG, MINRES

When  $A$  symmetric, Lanczos gives  $AQ = QT + q_{k+1}[0, \dots, 0, 1]$ ,  
 $T$ : tridiagonal

- ▶ CG: when  $A \succ 0$ , solve  $Q^T(AQy - b) = 0, x = Qy$   
→ “Galerkin orthogonality”: residual orthogonal to  $Q$ 
  - ▶ three-term recurrence reduces cost to  $O(k)$   $A$ -matmuls
  - ▶ minimizes  $A$ -norm of error  $x_k = \operatorname{argmin}_{x \in Q} \|x - x_*\|_A$

$$\begin{aligned}(x - x_*)^T A(x - x_*) &= (Qy - x_*)^T A(Qy - x_*) \\ &= y^T (Q^T A Q) y - 2b^T Qy + b^T x_*,\end{aligned}$$

minimizer is  $y = (Q^T A Q)^{-1} Q^T b$ , so  $Q^T(AQy - b) = 0$

- ▶ MINRES: symmetric (indefinite) version of GMRES
  - ▶ again, three-term recurrence,  $O(k)$   $A$ -matmuls

# Preconditioning for GMRES, CG etc

$$Ax = b$$

Instead find  $M \approx A^{-1}$  and solve

$$MAx = Mb$$

Desiderata of  $M$ :

- ▶  $M$  simple enough s.t. lin. systems  $Mx = b$  easy
- ▶  $MA$  has clustered eigenvalues
- ▶ Finding effective preconditioners is never-ending research topic  
Andy Wathen is our Oxford expert!

# Randomized SVD

[Halko, Martinsson, Tropp 2011]

- (i) Generate a random matrix  $\Omega \in \mathbb{R}^{n \times (r+\ell)}$ , where  $\ell$  is a small integer (say 5).
- (ii) Compute  $A\Omega$  and its QR factorization  $A\Omega = QR$ .
- (iii) Compute  $Q^T A$  and its SVD  $Q^T A = \tilde{U} \hat{\Sigma} \hat{V}^T$ .
- (iv) Take  $\hat{U} = Q\tilde{U}$ , and let  $\hat{U}_r, \hat{\Sigma}_r, \hat{V}_r$  be the leading  $r$  parts of  $\hat{U}, \hat{\Sigma}, \hat{V}$  respectively. Output  $\hat{A}_r = \hat{U}_r \hat{\Sigma}_r \hat{V}_r^T$  as a rank- $r$  approximant to  $A$ .

Approximation quality:

$$\mathbb{E}[\|A - \hat{A}_r\|_F] \leq \sqrt{1 + \frac{r}{\ell-1}} \|A - A_r\|_F$$

Optimal up to  $\sqrt{1 + \frac{r}{\ell-1}} = O(1)$

# Understanding randomized SVD

Want  $\|A - QQ^T A\|_F / \|A - A_r\|_F = O(1)$ , where  $A\Omega = QR$ ,

$$A\Omega = U \begin{bmatrix} \Sigma_1 V_1^T \Omega \\ \Sigma_2 V_2^T \Omega \end{bmatrix}, \quad \Sigma_1 = \text{diag}(\sigma_1, \dots, \sigma_r)$$

- ▶  $V_1^T \Omega$ :  $r \times (r + \ell)$  **rectangular** Gaussian, so *well-conditioned*, so w.h.p.  $\|(V_1^T \Omega)^\dagger\| = O(1)$  ( $X^\dagger$ : pseudoinverse)

- ▶ Hence

$$A\Omega(\Sigma_1 V_1^T \Omega)^\dagger = U \begin{bmatrix} I \\ F \end{bmatrix}, \quad F = \Sigma_2 V_2^T \Omega (V_1^T \Omega)^\dagger \Sigma_1^{-1}$$

- ▶ Note  $\|(I - QQ^T)A\|_F^2 = \|A\|_F^2 - \|Q^T A\|_F^2$

- ▶ Take  $\tilde{Q} := A\Omega(\Sigma_1 V_1^T \Omega)^\dagger = U \begin{bmatrix} I \\ F \end{bmatrix} (I + F^T F)^{-1/2}$ . Then

$$\begin{aligned} \|Q^T A\|_F^2 &\geq \|\tilde{Q}^T A\|_F^2 \geq \|(I + F^T F)^{-1/2} \Sigma_1\|_F^2 = \text{Tr}(\Sigma_1 (I + F^T F)^{-1} \Sigma_1) \\ &\geq \text{Tr}(\Sigma_1 (I - F^T F) \Sigma_1) = \|\Sigma_1\|_F^2 - \|\Sigma_2 V_2^T \Omega (V_1^T \Omega)^\dagger\|_F^2 \\ &= \|\Sigma_1\|_F^2 - \|\Sigma_2 M\|_F^2, \quad \|M\| = O(1) \end{aligned}$$

- ▶ Thus  $\|(I - QQ^T)A\|_F^2 \leq \|\Sigma_2 M\|_F^2 = O(\|\Sigma_2\|_F^2)$

## Important (N)LA topics not treated

- ▶ **tensors**
- ▶ FFT (values $\leftrightarrow$ coefficients) for polynomials
- ▶ sparse direct solvers
- ▶ multigrid
- ▶ functions of matrices
- ▶ generalized, polynomial eigenvalue problems
- ▶ perturbation theory
- ▶ compressed sensing
- ▶ model order reduction
- ▶ communication-avoiding algorithms
- ▶ differential equations, optimisation, machine learning,... LA is everywhere in applied maths!

# Eigenvalue perturbation theory

- ▶ first-order eigenvalue perturbation:  $Ax = \lambda x$ ,  $y^T A = \lambda y^T$ ,  $\lambda$  simple then

$$\lambda(A + \epsilon E) = \epsilon \frac{y^T E x}{y^T x} + O(\epsilon^2)$$

- ▶ Weyl's theorem:  $A, E$  symmetric then

$$\lambda_i(A) + \lambda_i(A + E) \leq \lambda_i(A) + \|E\|_2$$

- ▶ Davis-Kahan:  $A$  symmetric,  $Ax_i = \lambda_i x_i$ ,  $(A + E)\hat{x}_i = \hat{\lambda}_i \hat{x}_i$ ,  
 $gap_i := \min_j |\lambda_i - \lambda_j|$

$$\sin \angle(x_i, \hat{x}_i) \leq \frac{\|E\|}{gap_i}$$



## Sherman–Morrison(-Woodbury) formula

$A$ :  $n \times n$  invertible,  $U, V$ :  $n \times k$ ,  $k < (\ll) n$ . Then

$$(A + UV^T)^{-1} = A^{-1} - A^{-1}U(I_k + V^T A^{-1}U)^{-1}V^T A^{-1}$$

► Low-rank update of  $A^{-1}$

Similar updates possible for QR, (SVD), ...