

CS 598 EVS: Tensor Computations

Matrix Computations Background

Edgar Solomonik

University of Illinois, Urbana-Champaign

Matrices and Tensors

- ▶ What is a matrix?
 - ▶ A collection of numbers arranged into an array of dimensions $m \times n$, e.g., $M \in \mathbb{R}^{m \times n}$
 - ▶ A linear operator $f(x) = Mx$
 - ▶ A bilinear form $x^T M y$

- ▶ What is a tensor?
 - ▶ A collection of numbers arranged into an array of a particular order, with dimensions $l \times m \times n \times \dots$, e.g., $\mathcal{T} \in \mathbb{R}^{l \times m \times n}$ is order 3
 - ▶ A multilinear operator $z = f(x, y)$

$$z_i = \sum_{j,k} t_{ijk} x_j y_k$$

- ▶ A multilinear form $\sum_{i,j,k} t_{ijk} x_i y_j z_k$

Matrix Norms

- ▶ **Properties of matrix norms:**

$$\|\mathbf{A}\| \geq 0$$

$$\|\mathbf{A}\| = 0 \iff \mathbf{A} = \mathbf{0}$$

$$\|\alpha\mathbf{A}\| = |\alpha| \cdot \|\mathbf{A}\|$$

$$\|\mathbf{A} + \mathbf{B}\| \leq \|\mathbf{A}\| + \|\mathbf{B}\| \quad (\textit{triangle inequality})$$

- ▶ **Frobenius norm:**

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

- ▶ **Operator/induced/subordinate matrix norms:**

For any vector norm $\|\cdot\|_p$, the induced matrix norm is

$$\|\mathbf{A}\|_p = \max_{\mathbf{x} \neq \mathbf{0}} \|\mathbf{A}\mathbf{x}\|_p / \|\mathbf{x}\|_p = \max_{\|\mathbf{x}\|_p=1} \|\mathbf{A}\mathbf{x}\|_p$$

Existence of SVD

- ▶ Consider any maximizer $\mathbf{x}_1 \in \mathbb{R}^n$ with $\|\mathbf{x}_1\|_2 = 1$ to $\|\mathbf{A}\mathbf{x}_1\|_2$

Let $\mathbf{y}_1 = \mathbf{A}\mathbf{x}_1 / \|\mathbf{A}\mathbf{x}_1\|_2$ and $\sigma_1 = \mathbf{y}_1^T \mathbf{A}\mathbf{x}_1 = \|\mathbf{A}\mathbf{x}_1\|_2$, then consider any maximizer \mathbf{x}_2 of

$$\|(\mathbf{A} - \sigma_1 \mathbf{y}_1 \mathbf{x}_1^T) \mathbf{x}_2\|_2.$$

We can see that $\mathbf{x}_1 \perp \mathbf{x}_2$ since, otherwise, we have $\mathbf{x}_2 = \alpha \mathbf{x}_1 + \tilde{\mathbf{x}}_2$ with $\tilde{\mathbf{x}}_2 \perp \mathbf{x}_1$ and $\|\tilde{\mathbf{x}}_2\|_2 < \|\mathbf{x}_2\|_2$ and

$$\|(\mathbf{A} - \sigma_1 \mathbf{y}_1 \mathbf{x}_1^T)(\alpha \mathbf{x}_1 + \tilde{\mathbf{x}}_2)\|_2 = \|(\mathbf{A} - \sigma_1 \mathbf{y}_1 \mathbf{x}_1^T) \tilde{\mathbf{x}}_2\|_2.$$

Hence we have a contradiction, since

$$\|(\mathbf{A} - \sigma_1 \mathbf{y}_1 \mathbf{x}_1^T) \mathbf{x}_2\|_2 < (1 / \|\tilde{\mathbf{x}}_2\|_2) \|(\mathbf{A} - \sigma_1 \mathbf{y}_1 \mathbf{x}_1^T) \tilde{\mathbf{x}}_2\|_2.$$

More generally, we can see that any maximizer \mathbf{x}_{i+1} to

$$\|(\mathbf{A} - [\mathbf{y}_1 \ \cdots \ \mathbf{y}_i] \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_i \end{bmatrix} [\mathbf{x}_1 \ \cdots \ \mathbf{x}_i]^T) \mathbf{x}_{i+1}\|_2$$

is orthogonal to $\mathbf{x}_1, \dots, \mathbf{x}_i$ and similar for \mathbf{y}_{i+1} .

Singular Value Decomposition

- ▶ The singular value decomposition (SVD)

We can express any matrix \mathbf{A} as

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$$

where \mathbf{U} and \mathbf{V} are orthogonal, and $\mathbf{\Sigma}$ is square nonnegative and diagonal,

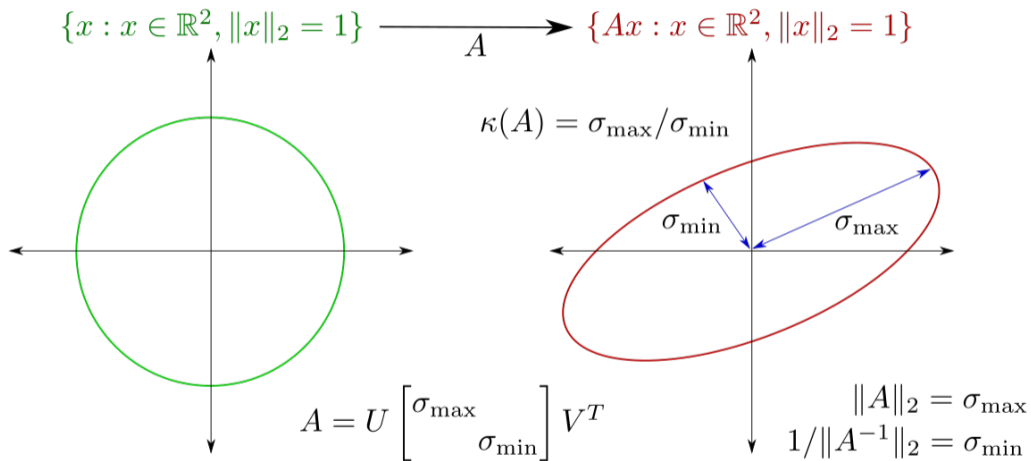
$$\mathbf{\Sigma} = \begin{bmatrix} \sigma_{max} & & & \\ & \ddots & & \\ & & \ddots & \\ & & & \sigma_{min} \end{bmatrix}$$

The columns of \mathbf{U} and \mathbf{V} are left and right singular vectors of \mathbf{A} , i.e.,

$$\mathbf{A}\mathbf{v}_i = \sigma_i\mathbf{u}_i, \quad \mathbf{u}_i^T \mathbf{A} = \sigma_i\mathbf{v}_i^T$$

- ▶ Condition number in terms of singular values
 - ▶ *We have that $\|\mathbf{A}\|_2 = \sigma_{max}$ and if \mathbf{A}^{-1} exists, $\|\mathbf{A}^{-1}\|_2 = 1/\sigma_{min}$*
 - ▶ *Consequently, $\kappa(\mathbf{A}) = \|\mathbf{A}\|_2\|\mathbf{A}^{-1}\|_2 = \sigma_{max}/\sigma_{min}$*

Visualization of Matrix Conditioning



Matrix Condition Number

- ▶ The matrix condition number $\kappa(\mathbf{A})$ is the ratio between the max and min distance from the surface to the center of the unit ball (norm-1 vectors) transformed by \mathbf{A} :
 - ▶ *The max distance to center is given by the vector maximizing $\max_{\|\mathbf{x}\|=1} \|\mathbf{A}\mathbf{x}\|_2$.*
 - ▶ *The min distance to center is given by the vector minimizing $\min_{\|\mathbf{x}\|=1} \|\mathbf{A}\mathbf{x}\|_2 = 1/(\max_{\|\mathbf{x}\|=1} \|\mathbf{A}^{-1}\mathbf{x}\|_2)$.*
 - ▶ *Thus, we have that $\kappa(\mathbf{A}) = \|\mathbf{A}\|_2 \|\mathbf{A}^{-1}\|_2$*
- ▶ The matrix condition number bounds the worst-case amplification of error in a matrix-vector product: *Consider $\mathbf{y} + \delta\mathbf{y} = \mathbf{A}(\mathbf{x} + \delta\mathbf{x})$, assume $\|\mathbf{x}\|_2 = 1$*
 - ▶ *In the worst case, $\|\mathbf{y}\|_2$ is minimized, that is $\|\mathbf{y}\|_2 = 1/\|\mathbf{A}^{-1}\|_2$*
 - ▶ *In the worst case, $\|\delta\mathbf{y}\|_2$ is maximized, that is $\|\delta\mathbf{y}\|_2 = \|\mathbf{A}\|_2 \|\delta\mathbf{y}\|_2$*
 - ▶ *So $\|\delta\mathbf{y}\|_2 / \|\mathbf{y}\|_2$ is at most $\kappa(\mathbf{A}) \|\delta\mathbf{x}\|_2 / \|\mathbf{x}\|_2$*

Linear Systems

- ▶ Given a square matrix $A \in \mathbb{R}^{n \times n}$ with rank n , consider solving $Ax = b$ given b
- ▶ The SVD allows explicit inversion of A

$$A^{-1} = V\Sigma^{-1}U^T$$

- ▶ However, Gaussian elimination is more computationally efficient
 - ▶ *Can factorize arbitrary A as $A = PLU$ for permutation matrix P and triangular L, U*
 - ▶ *For symmetric A LDLT factorization is $A = PLDL^T P^T$, where D has diagonal entries of 2-by-2 anti-diagonal symmetric blocks*
 - ▶ *If positive definite, Cholesky requires no pivoting/permutation*
 - ▶ *Suffices to solve linear systems in $O(n^2)$ cost using triangular solve*
- ▶ Given a factorization of A , solving a system with $A + uv^T$ has cost $O(n^2)$ via the Sherman-Morrison-Woodbury formula

Linear Least Squares

- ▶ Find $\mathbf{x}^* = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{Ax} - \mathbf{b}\|_2$ where $\mathbf{A} \in \mathbb{R}^{m \times n}$:

Since $m \geq n$, the minimizer generally does not attain a zero residual $\mathbf{Ax} - \mathbf{b}$. We can rewrite the optimization problem constraint via

$$\mathbf{x}^* = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{Ax} - \mathbf{b}\|_2^2 = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \left[(\mathbf{Ax} - \mathbf{b})^T (\mathbf{Ax} - \mathbf{b}) \right]$$

- ▶ Given the SVD $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ we have $\mathbf{x}^* = \underbrace{\mathbf{V}\mathbf{\Sigma}^\dagger\mathbf{U}^T}_{\mathbf{A}^\dagger} \mathbf{b}$, where $\mathbf{\Sigma}^\dagger$ contains the reciprocal of all nonzeros in $\mathbf{\Sigma}$, and more generally \dagger denotes pseudoinverse:
 - ▶ *The minimizer satisfies $\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T \mathbf{x}^* \cong \mathbf{b}$ and consequently also satisfies*

$$\mathbf{\Sigma}\mathbf{y}^* \cong \mathbf{d} \quad \text{where } \mathbf{y}^* = \mathbf{V}^T \mathbf{x}^* \text{ and } \mathbf{d} = \mathbf{U}^T \mathbf{b}.$$

- ▶ *The minimizer of the reduced problem is $\mathbf{y}^* = \mathbf{\Sigma}^\dagger \mathbf{d}$, so $y_i = d_i/\sigma_i$ for $i \in \{1, \dots, n\}$ and $y_i = 0$ for $i \in \{n+1, \dots, m\}$.*

Normal Equations

Demo: Normal equations vs Pseudoinverse

Demo: Issues with the normal equations

- ▶ *Normal equations* are given by solving $\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{b}$:

If $\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{b}$ then

$$(\mathbf{U} \mathbf{\Sigma} \mathbf{V}^T)^T \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T \mathbf{x} = (\mathbf{U} \mathbf{\Sigma} \mathbf{V}^T)^T \mathbf{b}$$

$$\mathbf{\Sigma}^T \mathbf{\Sigma} \mathbf{V}^T \mathbf{x} = \mathbf{\Sigma}^T \mathbf{U}^T \mathbf{b}$$

$$\mathbf{V}^T \mathbf{x} = (\mathbf{\Sigma}^T \mathbf{\Sigma})^{-1} \mathbf{\Sigma}^T \mathbf{U}^T \mathbf{b} = \mathbf{\Sigma}^\dagger \mathbf{U}^T \mathbf{b}$$

$$\mathbf{x} = \mathbf{V} \mathbf{\Sigma}^\dagger \mathbf{U}^T \mathbf{b} = \mathbf{x}^*$$

- ▶ However, solving the normal equations is a more ill-conditioned problem than the original least squares algorithm

Generally we have $\kappa(\mathbf{A}^T \mathbf{A}) = \kappa(\mathbf{A})^2$ (the singular values of $\mathbf{A}^T \mathbf{A}$ are the squares of those in \mathbf{A}). Consequently, solving the least squares problem via the normal equations may be unstable because it involves solving a problem that has worse conditioning than the initial least squares problem.

Solving the Normal Equations

- ▶ If \mathbf{A} is full-rank, then $\mathbf{A}^T \mathbf{A}$ is symmetric positive definite (SPD):
 - ▶ *Symmetry is easy to check* $(\mathbf{A}^T \mathbf{A})^T = \mathbf{A}^T \mathbf{A}$.
 - ▶ *\mathbf{A} being full-rank implies $\sigma_{\min} > 0$ and further if $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ we have*

$$\mathbf{A}^T \mathbf{A} = \mathbf{V}^T \mathbf{\Sigma}^2 \mathbf{V}$$

which implies that rows of \mathbf{V} are the eigenvectors of $\mathbf{A}^T \mathbf{A}$ with eigenvalues $\mathbf{\Sigma}^2$ since $\mathbf{A}^T \mathbf{A} \mathbf{V}^T = \mathbf{V}^T \mathbf{\Sigma}^2$.

- ▶ Since $\mathbf{A}^T \mathbf{A}$ is SPD we can use Cholesky factorization, to factorize it and solve linear systems:

$$\mathbf{A}^T \mathbf{A} = \mathbf{L}\mathbf{L}^T$$

QR Factorization

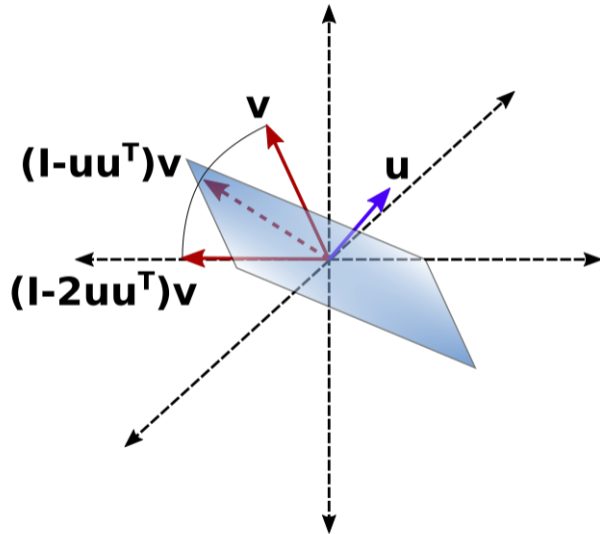
- ▶ If A is full-rank there exists an orthogonal matrix Q and a unique upper-triangular matrix R with a positive diagonal such that $A = QR$
 - ▶ Given $A^T A = LL^T$, we can take $R = L^T$ and obtain $Q = AL^{-T}$, since $\underbrace{L^{-1}A^T}_{Q^T} \underbrace{AL^{-T}}_Q = I$ implies that Q has orthonormal columns.
- ▶ A reduced QR factorization (unique part of general QR) is defined so that $Q \in \mathbb{R}^{m \times n}$ has orthonormal columns and R is square and upper-triangular. A full QR factorization gives $Q \in \mathbb{R}^{m \times m}$ and $R \in \mathbb{R}^{m \times n}$, but since R is upper triangular, the latter $m - n$ columns of Q are only constrained so as to keep Q orthogonal. The **reduced QR** factorization is given by taking the first n columns Q and \hat{Q} the upper-triangular block of R , \hat{R} giving $A = \hat{Q}\hat{R}$.
- ▶ We can solve the normal equations (and consequently the linear least squares problem) via reduced QR as follows

$$A^T A x = A^T b \quad \Rightarrow \quad \hat{R}^T \underbrace{\hat{Q}^T \hat{Q}}_I \hat{R} x = \hat{R}^T \hat{Q}^T b \quad \Rightarrow \quad \hat{R} x = \hat{Q}^T b$$

Computing the QR Factorization

- ▶ The Cholesky-QR algorithm uses the normal equations to obtain the QR factorization
 - ▶ Compute $A^T A = LL^T$, take $R = L^T$, and solve for Q triangular linear systems $LQ^T = A^T$
 - ▶ If A is $m \times n$, forming $A^T A$ has cost mn^2 , computing Cholesky factorization has cost $(2/3)n^3$, and solving the triangular systems (if Q is needed) costs mn^2 , yielding total cost $2mn^2 + (2/3)n^3$
 - ▶ However, this method is unstable since $A^T A$ is ill-conditioned. This is addressible by iterating on the computed (nearly-orthogonal) Q factor (CholeskyQR2).
- ▶ Orthogonalization-based methods are most efficient and stable for QR factorization of dense matrices
 - ▶ Apply a sequence of orthogonal transformations Q_1, \dots, Q_k to reduce A to triangular form $(Q_1 \cdots Q_k)^T A = R$
 - ▶ Householder QR uses rank-1 perturbations of the identity matrix (reflectors) $Q_i = I - 2u_i u_i^T$ to zero-out each sub-column of A
 - ▶ Givens rotations zero-out a single entry at a time
 - ▶ Both approaches have cost $O(mn^2)$ with similar constant to Cholesky-QR

Householder orthogonalization



Eigenvalue Decomposition

- ▶ If a matrix A is diagonalizable, it has an *eigenvalue decomposition*

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

where X are the right eigenvectors, X^{-1} are the left eigenvectors and D are eigenvalues

$$AX = [Ax_1 \quad \cdots \quad Ax_n] = XD = [d_{11}x_1 \quad \cdots \quad d_{nn}x_n].$$

- ▶ If A is symmetric, its right and left singular vectors are the same, and consequently are its eigenvectors.
- ▶ More generally, any *normal* matrix, $A^H A = A A^H$, has unitary eigenvectors.
- ▶ A and B are *similar*, if there exist Z such that $A = ZBZ^{-1}$
 - ▶ Normal matrices are *unitarily similar* ($Z^{-1} = Z^H$) to diagonal matrices
 - ▶ Symmetric real matrices are *orthogonally similar* ($Z^{-1} = Z^T$) to real diagonal matrices
 - ▶ Hermitian matrices are unitarily similar to real diagonal matrices

Similarity of Matrices

Invertible similarity transformations $Y = XAX^{-1}$

<i>matrix (A)</i>	<i>reduced form (Y)</i>
arbitrary	bidiagonal
diagonalizable	diagonal

Unitary similarity transformations $Y = UAU^H$

<i>matrix (A)</i>	<i>reduced form (Y)</i>
arbitrary	triangular (Schur)
normal	diagonal
Hermitian	real diagonal

Orthogonal similarity transformations $Y = QAQ^T$

<i>matrix (A)</i>	<i>reduced form (Y)</i>
real	Hessenberg
real symmetric	real diagonal
real SPD	real positive diagonal

Field of Values

- ▶ For any square matrix \mathbf{A} and vector \mathbf{x} the *Rayleigh quotient* is

$$\rho_{\mathbf{A}}(\mathbf{x}) = \frac{\mathbf{x}^H \mathbf{A} \mathbf{x}}{\mathbf{x}^H \mathbf{x}}$$

- ▶ Its magnitude is bounded by the singular values as

$$1/\|\mathbf{A}\|_2^{-1} \leq |\rho_{\mathbf{A}}(\mathbf{x})| \leq \|\mathbf{A}\|_2$$

- ▶ If \mathbf{x} is an eigenvector of \mathbf{A} , so $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$ or $\mathbf{x}^H \mathbf{A} = \lambda\mathbf{x}^H$, then

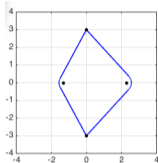
$$\rho_{\mathbf{A}}(\mathbf{x}) = \lambda$$

- ▶ The set $\mathcal{F}_{\mathbf{A}} = \{\rho_{\mathbf{A}}(\mathbf{x}) : \mathbf{x} \in \mathbb{C}^n, \mathbf{x} \neq 0\}$ is the *field of values* of \mathbf{A}

Field of Values and Eigenvalues

- ▶ Clearly any eigenvalue λ of \mathbf{A} is in \mathcal{F}_A

- ▶ For the matrix $\mathbf{A} = \begin{bmatrix} & 3 & \\ -3 & & \\ & & 3 \\ & 1 & 1 \end{bmatrix}$, \mathcal{F}_A is¹



- ▶ The field of values of a normal matrix is easy to characterize
 - ▶ If \mathbf{A} is normal, \mathcal{F}_A is the convex hull of the eigenvalues.
 - ▶ If \mathbf{A} is Hermitian and positive definite, $\mathcal{F}_A = [\sigma_{min}, \sigma_{max}]$
- ▶ In general, eigenvectors are obtained from critical points of the Rayleigh quotient on the unit circle

$$\mathcal{L}_A(\mathbf{x}, \lambda) = \mathbf{x}^H \mathbf{A} \mathbf{x} + \lambda(1 - \mathbf{x}^H \mathbf{x})$$

$$\nabla \mathcal{L}_A(\mathbf{x}, \lambda) = \begin{bmatrix} 2\mathbf{A}\mathbf{x} - 2\lambda\mathbf{x} \\ 1 - \mathbf{x}^H \mathbf{x} \end{bmatrix} = 0,$$

¹Credit to <https://www.chebfun.org/examples/linalg/FieldOfValues.html>

Singular Vectors as Critical Points

- ▶ Like eigenvectors, we can also derive singular vectors from an optimization (critical point) perspective
 - ▶ *Again, consider the critical points of the Lagrangian function of an optimization problem on the unit-sphere,*

$$\mathcal{L}_A(\mathbf{u}, \mathbf{v}, \lambda_1, \lambda_2) = 2\mathbf{u}^H \mathbf{A} \mathbf{v} + \lambda_1(1 - \mathbf{u}^H \mathbf{u}) + \lambda_2(1 - \mathbf{v}^H \mathbf{v})$$

$$\nabla \mathcal{L}_A(\mathbf{u}, \mathbf{v}, \lambda_1, \lambda_2) = \begin{bmatrix} 2\mathbf{A} \mathbf{v} - 2\lambda_1 \mathbf{u} \\ 2\mathbf{A}^H \mathbf{u} - 2\lambda_2 \mathbf{v} \\ 1 - \mathbf{u}^H \mathbf{u} \\ 1 - \mathbf{v}^H \mathbf{v} \end{bmatrix} = 0,$$

- ▶ *At a critical point, we can see that $\lambda_1 = \lambda_2$, since $\mathbf{u}^H \mathbf{A} \mathbf{v} = \lambda_1 = \lambda_2$.*

Matrix Functions

- ▶ Consider a polynomial p , for a diagonalizable matrix $\mathbf{A} = \mathbf{X}\mathbf{D}\mathbf{X}^{-1}$,

$$p(\mathbf{A}) = \mathbf{X}p(\mathbf{D})\mathbf{X}^{-1}$$

$$\begin{aligned} p(\mathbf{A}) &= \sum_{i=0}^{\deg(p)} c_i \mathbf{A}^i = \sum_{i=0}^{\deg(p)} c_i \prod_{j=1}^i \mathbf{X}\mathbf{D}\mathbf{X}^{-1} \\ &= \sum_{i=0}^{\deg(p)} c_i \mathbf{X}\mathbf{D}^i \mathbf{X}^{-1} = \mathbf{X} \left(\sum_{i=0}^{\deg(p)} c_i \mathbf{D}^i \right) \mathbf{X}^{-1} \end{aligned}$$

- ▶ The above definition readily extends to other analytic functions f , but non-diagonalizable matrices require a more sophisticated definition

Crouzeix's conjecture

- ▶ So far, we have seen close connections between the matrix 2-norm and singular values, and between the Rayleigh quotient and the eigenvalues
- ▶ An important open problem in numerical analysis that relates the norm with the Rayleigh quotient is Crouzeix's conjecture
 - ▶ *For any polynomial p and complex matrix \mathbf{A} ,*

$$\|p(\mathbf{A})\|_2 \leq 2 \max_{z \in \mathcal{F}_A} |p(z)|$$

- ▶ *The conjecture is known to hold for some subclasses of matrices and with constant 1.08 instead of 2 (Crouzeix's theorem)*
- ▶ *If valid, the bound of 2 is tight, including for $p(\mathbf{A}) = \mathbf{A}$, by choosing $\mathbf{A} = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}$*

Computing Eigenvalue and Singular Value Decompositions

- ▶ Direct methods for eigenvalue problems start by reducing the matrix to upper-Hessenberg form
 - ▶ *Seek a sequence of unitary similarity transformations*
 $H = Q_k \cdots Q_1 A Q_1^T \cdots Q_k^T$ so that H is zero below the first subdiagonal (upper-Hessenberg)
 - ▶ *Can pick each Q_i as a Householder transformation acting on the last $n - i$ rows*
 - ▶ *$O(n^3)$ cost to reduce to upper-Hessenberg or tridiagonal if symmetric*
 - ▶ *To obtain singular vectors, can work with $A^T A$ or perform 'bidiagonal reduction'*
 - ▶ *If matrix is sparse, fill may be introduced*
- ▶ Iterative methods are generally based on products with the matrix
 - ▶ *Power iteration converges to the largest eigenvalue eigenvectors of A*
 - ▶ *Convergence rate is linear and depends on ratio of two largest eigenvalues*
 - ▶ *Integrating diagonal shifts and inversion yields other methods: inverse iteration, Rayleigh-quotient iteration*
 - ▶ *Most iterative methods involve only products with A or a related matrix*

Introduction to Krylov Subspace Methods

- ▶ *Krylov subspace methods* work with information contained in the $n \times k$ matrix

$$\mathbf{K}_k = [\mathbf{x}_0 \quad \mathbf{A}\mathbf{x}_0 \quad \cdots \quad \mathbf{A}^{k-1}\mathbf{x}_0]$$

We seek to best use the information from the matrix vector product results (columns of \mathbf{K}_k) to solve eigenvalue problems.

- ▶ Assuming \mathbf{K}_n is invertible, the matrix $\mathbf{K}_n^{-1}\mathbf{A}\mathbf{K}_n$ is a *companion matrix* \mathbf{C} :
Letting $\mathbf{k}_n^{(i)} = \mathbf{A}^{i-1}\mathbf{x}$, we observe that

$$\mathbf{A}\mathbf{K}_n = \begin{bmatrix} \mathbf{A}\mathbf{k}_n^{(1)} & \cdots & \mathbf{A}\mathbf{k}_n^{(n-1)} & \mathbf{A}\mathbf{k}_n^{(n)} \end{bmatrix} = \begin{bmatrix} \mathbf{k}_n^{(2)} & \cdots & \mathbf{k}_n^{(n)} & \mathbf{A}\mathbf{k}_n^{(n)} \end{bmatrix},$$

therefore premultiplying by \mathbf{K}_n^{-1} transforms the first $n - 1$ columns of $\mathbf{A}\mathbf{K}_n$ into the last $n - 1$ columns of \mathbf{I} ,

$$\begin{aligned} \mathbf{K}_n^{-1}\mathbf{A}\mathbf{K}_n &= \begin{bmatrix} \mathbf{K}_n^{-1}\mathbf{k}_n^{(2)} & \cdots & \mathbf{K}_n^{-1}\mathbf{k}_n^{(n)} & \mathbf{K}_n^{-1}\mathbf{A}\mathbf{k}_n^{(n)} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{e}_2 & \cdots & \mathbf{e}_n & \mathbf{K}_n^{-1}\mathbf{A}\mathbf{k}_n^{(n)} \end{bmatrix} \end{aligned}$$

Krylov Subspaces

- ▶ Given $Q_k R_k = K_k$, we obtain an orthonormal basis for the Krylov subspace,

$$\mathcal{K}_k(\mathbf{A}, \mathbf{x}_0) = \text{span}(\mathbf{Q}_k) = \{p(\mathbf{A})\mathbf{x}_0 : \deg(p) < k\},$$

where p is any polynomial of degree less than k .

- ▶ The Krylov subspace includes the $k - 1$ approximate dominant eigenvectors generated by $k - 1$ steps of power iteration:
 - ▶ *The approximation obtained from $k - 1$ steps of power iteration starting from \mathbf{x}_0 is given by the Rayleigh-quotient of $\mathbf{y} = \mathbf{A}^k \mathbf{x}_0$.*
 - ▶ *This vector is within the Krylov subspace, $\mathbf{y} \in \mathcal{K}_k(\mathbf{A}, \mathbf{x}_0)$.*
 - ▶ *Consequently, Krylov subspace methods will generally obtain strictly better approximations of the dominant eigenpair than power iteration.*

Rayleigh-Ritz Procedure

- ▶ The eigenvalues/eigenvectors of \mathbf{H}_k are the *Ritz values/vectors*:

$$\mathbf{H}_k = \mathbf{X} \mathbf{D} \mathbf{X}^{-1}$$

eigenvalue approximations based on Ritz vectors \mathbf{X} are given by $\mathbf{Q}_k \mathbf{X}$.

- ▶ The Ritz vectors and values are the *ideal approximations* of the actual eigenvalues and eigenvectors based on only \mathbf{H}_k and \mathbf{Q}_k :

Assuming \mathbf{A} is a symmetric matrix with positive eigenvalues, the largest Ritz value $\lambda_{\max}(\mathbf{H}_k)$ will be the maximum Rayleigh quotient of any vector in $\mathcal{K}_k = \text{span}(\mathbf{Q}_k)$,

$$\max_{\mathbf{x} \in \text{span}(\mathbf{Q}_k)} \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{x}} = \max_{\mathbf{y} \neq 0} \frac{\mathbf{y}^T \mathbf{Q}_k^T \mathbf{A} \mathbf{Q}_k \mathbf{y}}{\mathbf{y}^T \mathbf{y}} = \max_{\mathbf{y} \neq 0} \frac{\mathbf{y}^T \mathbf{H}_k \mathbf{y}}{\mathbf{y}^T \mathbf{y}} = \lambda_{\max}(\mathbf{H}_k),$$

which is the best approximation to $\lambda_{\max}(\mathbf{A}) = \max_{\mathbf{x} \neq 0} \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{x}}$ available in \mathcal{K}_k . The quality of the approximation can also be shown to be optimal for other eigenvalues/eigenvectors.

Arnoldi Iteration

- ▶ Arnoldi iteration computes the i th column of \mathbf{H}_n , \mathbf{h}_i and the i th column of \mathbf{Q}_n directly using the recurrence $\mathbf{A}\mathbf{q}_i = \mathbf{Q}_n\mathbf{h}_i = \sum_{j=1}^{i+1} h_{ji}\mathbf{q}_j$

- ▶ Note that

$$\mathbf{q}_i^T \mathbf{A}\mathbf{q}_j = \mathbf{q}_i^T (\mathbf{Q}_n \mathbf{H}_n \mathbf{Q}_n^T) \mathbf{q}_j = \mathbf{e}_i^T \mathbf{H}_n \mathbf{e}_j = h_{ij}.$$

- ▶ The Arnoldi algorithm computes \mathbf{q}_{i+1} from $\mathbf{q}_1, \dots, \mathbf{q}_i$ by first computing $\mathbf{u}_i = \mathbf{A}\mathbf{q}_i$ then orthogonalizing,

$$\mathbf{q}_{i+1} h_{i+1,i} = \mathbf{u}_i - \sum_{j=1}^i \mathbf{q}_j h_{ji}, \quad h_{ji} = \mathbf{q}_j^T \mathbf{u}_i$$

then computing the norm of the vector to obtain $h_{i+1,i}$, yielding the i th column of \mathbf{H}_n .

Multidimensional Optimization

- ▶ Minimize $f(\mathbf{x})$
 - ▶ *In the context of constrained optimization, also have equality and or inequality constraints, e.g., $\mathbf{Ax} = \mathbf{b}$ or $\mathbf{x} > 0$*
 - ▶ *Unconstrained local optimality holds if $\nabla f(\mathbf{x}^*) = 0$ and $H_f(\mathbf{x}^*)$ is positive semi-definite*
 - ▶ *Reduces to solving nonlinear equations via optimality condition*
 - ▶ *Unconstrained local optimality conditions are looser, need the gradient to be zero or positive in all unconstrained directions at \mathbf{x}^**
 - ▶ *The condition $\nabla f(\mathbf{x}^*) = 0$ implies poor conditioning, perturbations that change the function value in the k th digit can change the solution in the $(k/2)$ th digit*
- ▶ Quadratic optimization $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \mathbf{Ax} - \mathbf{b}^T \mathbf{x}$
 - ▶ *Quadratic optimization problems can provide local approximations to general nonlinear optimization problems via Newton's method (where \mathbf{A} is the Hessian and \mathbf{b}^T is the gradient)*
 - ▶ *Equivalent to solving linear system $\mathbf{Ax} = \mathbf{b}$ by optimality condition*
 - ▶ *Accordingly, conditioning relative to perturbation in \mathbf{b} is $\kappa(\mathbf{A})$*

Basic Multidimensional Optimization Methods

- ▶ Steepest descent: minimize f in the direction of the negative gradient:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)$$

such that $f(\mathbf{x}_{k+1}) = \min_{\alpha_k} f(\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k))$, i.e. perform a line search (solve 1D optimization problem) in the direction of the negative gradient.

- ▶ Given quadratic optimization problem $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{x}$ where \mathbf{A} is symmetric positive definite, the error $\mathbf{e}_k = \mathbf{x}_k - \mathbf{x}^*$ satisfies

$$\|\mathbf{e}_{k+1}\|_{\mathbf{A}} = \mathbf{e}_{k+1}^T \mathbf{A} \mathbf{e}_{k+1} = \frac{\sigma_{\max}(\mathbf{A}) - \sigma_{\min}(\mathbf{A})}{\sigma_{\max}(\mathbf{A}) + \sigma_{\min}(\mathbf{A})} \|\mathbf{e}_k\|_{\mathbf{A}}$$

- ▶ When sufficiently close to a local minima, general nonlinear optimization problems are described by such an SPD quadratic problem.
- ▶ Convergence rate depends on the conditioning of \mathbf{A} , since

$$\frac{\sigma_{\max}(\mathbf{A}) - \sigma_{\min}(\mathbf{A})}{\sigma_{\max}(\mathbf{A}) + \sigma_{\min}(\mathbf{A})} = \frac{\kappa(\mathbf{A}) - 1}{\kappa(\mathbf{A}) + 1}.$$

Gradient Methods with Extrapolation

- ▶ We can improve the constant in the linear rate of convergence of steepest descent by leveraging *extrapolation methods*, which consider two previous iterates (maintain *momentum* in the direction $\mathbf{x}_k - \mathbf{x}_{k-1}$):

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k) + \beta_k (\mathbf{x}_k - \mathbf{x}_{k-1})$$

- ▶ The *heavy ball method*, which uses constant $\alpha_k = \alpha$ and $\beta_k = \beta$, achieves better convergence than steepest descent:

$$\|\mathbf{e}_{k+1}\|_{\mathbf{A}} = \frac{\sqrt{\kappa(\mathbf{A})} - 1}{\sqrt{\kappa(\mathbf{A})} + 1} \|\mathbf{e}_k\|_{\mathbf{A}}$$

Nesterov's gradient optimization method is another instance of an extrapolation method that provides further improved optimality guarantees.

Conjugate Gradient Method

- ▶ The *conjugate gradient method* is capable of making the optimal (for a quadratic objective) choice of α_k and β_k at each iteration of an extrapolation method:

$$(\alpha_k, \beta_k) = \underset{\alpha_k, \beta_k}{\operatorname{argmin}} \left[f\left(\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k) + \beta_k (\mathbf{x}_k - \mathbf{x}_{k-1})\right) \right]$$

- ▶ *For SPD quadratic programming problems, conjugate gradient is an optimal first order method, converging in n iterations.*
- ▶ *It implicitly computes Lanczos iteration, searching along A -orthogonal directions at each step.*
- ▶ *Parallel tangents* implementation of the method proceeds as follows
 1. *Perform a step of steepest descent to generate $\hat{\mathbf{x}}_k$ from \mathbf{x}_k .*
 2. *Generate \mathbf{x}_{k+1} by minimizing over the line passing through \mathbf{x}_{k-1} and $\hat{\mathbf{x}}_k$.*

The method is equivalent to CG for a quadratic objective function.

Krylov Optimization

- ▶ Conjugate gradient (CG) finds the minimizer of $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \mathbf{A}\mathbf{x} - \mathbf{b}^T \mathbf{x}$ (which satisfies optimality condition $\mathbf{A}\mathbf{x} = \mathbf{b}$) within the Krylov subspace of \mathbf{A} :
 - ▶ It constructs Krylov subspace $\mathcal{K}_k(\mathbf{A}, \mathbf{b}) = \text{span}(\mathbf{b}, \mathbf{A}\mathbf{b}, \dots, \mathbf{A}^{k-1}\mathbf{b})$.
 - ▶ At the k th step conjugate gradient yields iterate

$$\mathbf{x}_k = \|\mathbf{b}\|_2 \mathbf{Q}_k \mathbf{T}_k^{-1} \mathbf{e}_1,$$

where \mathbf{Q}_k is an orthogonal basis for Krylov subspace $\mathcal{K}_k(\mathbf{A}, \mathbf{b})$ and $\mathbf{T}_k = \mathbf{Q}_k^T \mathbf{A} \mathbf{Q}_k$.

- ▶ This choice of \mathbf{x}_k minimizes $f(\mathbf{x})$ since

$$\begin{aligned} \min_{\mathbf{x} \in \mathcal{K}_k(\mathbf{A}, \mathbf{b})} f(\mathbf{x}) &= \min_{\mathbf{y} \in \mathbb{R}^k} f(\mathbf{Q}_k \mathbf{y}) \\ &= \min_{\mathbf{y} \in \mathbb{R}^k} \mathbf{y}^T \mathbf{Q}_k^T \mathbf{A} \mathbf{Q}_k \mathbf{y} - \mathbf{b}^T \mathbf{Q}_k \mathbf{y} \\ &= \min_{\mathbf{y} \in \mathbb{R}^k} \mathbf{y}^T \mathbf{T}_k \mathbf{y} - \|\mathbf{b}\|_2 \mathbf{e}_1^T \mathbf{y} \end{aligned}$$

is minimized by $\mathbf{y} = \|\mathbf{b}\|_2 \mathbf{T}_k^{-1} \mathbf{e}_1$.

Conjugate Gradient Method: Optimized Form

After initialization $\mathbf{x}_0 = 0$, $\mathbf{r}_0 = \mathbf{b}$, $\mathbf{p}_0 = \mathbf{r}_0$, the k th iteration of CG computes

$$\mathbf{q}_k = \mathbf{A}\mathbf{p}_k$$

$$\alpha_k = \frac{\mathbf{r}_k^T \mathbf{r}_k}{\mathbf{q}_k^T \mathbf{p}_k}$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k$$

$$\mathbf{r}_{k+1} = \mathbf{r}_k - \alpha_k \mathbf{q}_k$$

At this point if the residual norm ($\|\mathbf{r}_{k+1}\|$) is small, terminate, otherwise prepare for next iteration,

$$\mathbf{p}_{k+1} = \mathbf{r}_{k+1} + \frac{\mathbf{r}_{k+1}^T \mathbf{r}_{k+1}}{\mathbf{r}_k^T \mathbf{r}_k} \mathbf{p}_k$$

See Jonathan Shewchuk 1994 notes on CG or James Demmel's book for the derivation of this form of the algorithm.

Conjugate Gradient Convergence Analysis

- ▶ In previous discussion, we assumed \mathbf{K}_n is invertible, which may not be the case if \mathbf{A} has $k < n$ distinct eigenvalues, however, then CG converges in $k - 1$ iterations (in exact arithmetic)
 - ▶ *To prove this, we can analyze the ‘minimizing’ polynomials in the Krylov subspace in terms of the (real and positive) eigenvalues of \mathbf{A}*
 - ▶ *The approximate solution x_k obtained by CG after $k - 1$ iterations is given by minimizing $z \in \mathcal{K}_k(\mathbf{A}, \mathbf{b})$, which means $z = \rho_{k-1}(\mathbf{A})\mathbf{b}$ for some polynomial ρ_{k-1} of degree $k - 1$*
 - ▶ *Now, consider the residual*

$$\mathbf{A}\mathbf{x} - \mathbf{b} = (\mathbf{A}\rho_{k-1}(\mathbf{A}) - \mathbf{I})\mathbf{b}$$

- ▶ *Choosing ρ_{k-1} as a polynomial interpolant so that $\rho_{k-1}(\lambda) = 1/\lambda$ for $\lambda \in \lambda(\mathbf{A})$, results in a zero residual since then $\rho_{k-1}(\mathbf{A}) = \mathbf{A}^{-1}$.*

Round-off Error in Conjugate Gradient

- ▶ CG provides strong convergence guarantees for SPD matrices in exact arithmetic
 - ▶ Classically, CG was viewed as a direct method, since its guaranteed to convergence in n iterations
 - ▶ In practice, round-off error prevents CG from achieving this for realistic matrices, so CG was actually abandoned for a while due to being viewed as unstable
 - ▶ Later, it was realized that CG is highly competitive as an iterative method
- ▶ Due to round-off CG may stagnate / have plateaus in convergence
 - ▶ A formal analysis of round-off error² reveals that CG with round-off is equivalent to exact CG on a matrix of larger dimension, whose eigenvalues are clustered around those of A
 - ▶ Using this view, CG convergence plateaus may be explained by the polynomial q_k developing more and more zeros near the same eigenvalue of A

²A. Greenbaum and Z. Strakos, SIMAX 1992

Preconditioning

- ▶ Convergence of iterative methods for $\mathbf{Ax} = \mathbf{b}$ depends on $\kappa(\mathbf{A})$, the goal of a preconditioner \mathbf{M} is to obtain \mathbf{x} by solving

$$\mathbf{M}^{-1}\mathbf{Ax} = \mathbf{M}^{-1}\mathbf{b}$$

with $\kappa(\mathbf{M}^{-1}\mathbf{A}) < \kappa(\mathbf{A})$

- ▶ need not form $\mathbf{M}^{-1}\mathbf{A}$ but only compute matrix-vector products $\mathbf{M}^{-1}(\mathbf{Ax})$
- ▶ want $\mathbf{M}^{-1}\mathbf{x}$ to be easy to compute (easier than $\mathbf{A}^{-1}\mathbf{x}$)
- ▶ so generally one extracts some $\mathbf{M} \approx \mathbf{A}$ that is easy to solve linear systems with
- ▶ however, $\mathbf{M} \approx \mathbf{A}$ may be insufficient/unnecessary, primary goal is to improve conditioning to accelerate iterative methods, i.e., want $\kappa(\mathbf{M}^{-1}\mathbf{A}) \ll \kappa(\mathbf{A})$
- ▶ Common preconditioners select parts of \mathbf{A} or perform inexact factorization
 - ▶ (block-)Jacobi preconditioner takes \mathbf{M} to be (block-)diagonal of \mathbf{A}
 - ▶ incomplete LU (ILU) preconditioners compute $\mathbf{M} = \mathbf{LU} \approx \mathbf{A}$ (+pivoting)
 - ▶ ILU variants constraint sparsity of \mathbf{L} and \mathbf{U} factors during factorization to be the same or not much more than that of \mathbf{A}
 - ▶ good problem-specific preconditioners are often available in practice and theory, applying also to problems beyond linear systems (eigenvalue problems, optimization, approximate graph algorithms)

Newton's Method

- ▶ Newton's method in n dimensions is given by finding minima of n -dimensional quadratic approximation using the gradient and Hessian of f :

$$f(\mathbf{x}_k + \mathbf{s}) \approx \hat{f}(\mathbf{s}) = f(\mathbf{x}_k) + \mathbf{s}^T \nabla f(\mathbf{x}_k) + \frac{1}{2} \mathbf{s}^T \mathbf{H}_f(\mathbf{x}_k) \mathbf{s}.$$

The minima of this function can be determined by identifying critical points

$$\mathbf{0} = \nabla \hat{f}(\mathbf{s}) = \nabla f(\mathbf{x}_k) + \mathbf{H}_f(\mathbf{x}_k) \mathbf{s},$$

thus to determine \mathbf{s} we solve the linear system,

$$\mathbf{H}_f(\mathbf{x}_k) \mathbf{s} = -\nabla f(\mathbf{x}_k).$$

Assuming invertibility of the Hessian, we can write the Newton's method iteration as

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \underbrace{\mathbf{H}_f(\mathbf{x}_k)^{-1} \nabla f(\mathbf{x}_k)}_{\mathbf{s}}.$$

Quadratic convergence follows by fixed point function analysis, beyond smoothness, a sufficient assumption is that $\mathbf{H}_f(\mathbf{x}^)$ is SPD.*

Nonlinear Least Squares

- ▶ An important special case of multidimensional optimization is *nonlinear least squares*, the problem of fitting a nonlinear function $f_{\mathbf{x}}(t)$ so that $f_{\mathbf{x}}(t_i) \approx y_i$:
For example, consider fitting $f_{[x_1, x_2]}(t) = x_1 \sin(x_2 t)$ so that

$$\begin{bmatrix} f_{[x_1, x_2]}(1.5) \\ f_{[x_1, x_2]}(1.9) \\ f_{[x_1, x_2]}(3.2) \end{bmatrix} \approx \begin{bmatrix} -1.2 \\ 4.5 \\ 7.3 \end{bmatrix}.$$

- ▶ We can cast nonlinear least squares as an optimization problem to minimize residual error and solve it by Newton's method:

Define residual vector function $\mathbf{r}(\mathbf{x})$ so that $r_i(\mathbf{x}) = y_i - f_{\mathbf{x}}(t_i)$ and minimize

$$\phi(\mathbf{x}) = \frac{1}{2} \|\mathbf{r}(\mathbf{x})\|_2^2 = \frac{1}{2} \mathbf{r}(\mathbf{x})^T \mathbf{r}(\mathbf{x}).$$

Now the gradient is $\nabla \phi(\mathbf{x}) = \mathbf{J}_{\mathbf{r}}^T(\mathbf{x}) \mathbf{r}(\mathbf{x})$ and the Hessian is

$$\mathbf{H}_{\phi}(\mathbf{x}) = \mathbf{J}_{\mathbf{r}}^T(\mathbf{x}) \mathbf{J}_{\mathbf{r}}(\mathbf{x}) + \sum_{i=1}^m r_i(\mathbf{x}) \mathbf{H}_{r_i}(\mathbf{x}).$$

Gauss-Newton Method

- ▶ The Hessian for nonlinear least squares problems has the form:

$$\mathbf{H}_\phi(\mathbf{x}) = \mathbf{J}_r^T(\mathbf{x})\mathbf{J}_r(\mathbf{x}) + \sum_{i=1}^m r_i(\mathbf{x})\mathbf{H}_{r_i}(\mathbf{x}).$$

The second term is small when the residual function $r(\mathbf{x})$ is small, so approximate

$$\mathbf{H}_\phi(\mathbf{x}) \approx \hat{\mathbf{H}}_\phi(\mathbf{x}) = \mathbf{J}_r^T(\mathbf{x})\mathbf{J}_r(\mathbf{x}).$$

- ▶ The *Gauss-Newton* method is Newton iteration with an approximate Hessian:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \hat{\mathbf{H}}_\phi(\mathbf{x}_k)^{-1}\nabla f(\mathbf{x}_k) = \mathbf{x}_k - (\mathbf{J}_r^T(\mathbf{x}_k)\mathbf{J}_r(\mathbf{x}_k))^{-1}\mathbf{J}_r^T(\mathbf{x}_k)\mathbf{r}(\mathbf{x}_k).$$

Recognizing the normal equations, we interpret the Gauss-Newton method as solving linear least squares problems $\mathbf{J}_r(\mathbf{x}_k)\mathbf{s}_k \cong \mathbf{r}(\mathbf{x}_k)$, $\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{s}_k$.

Low Rank Matrix Approximation

- ▶ Given a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ seek rank $r < m, n$ approximation
 - ▶ Given by matrices $\mathbf{U} \in \mathbb{R}^{m \times r}$ and $\mathbf{V} \in \mathbb{R}^{n \times r}$ so

$$\mathbf{A} \approx \mathbf{UV}^T$$

- ▶ Reduces memory footprint and cost of applying \mathbf{A} from mn to $mr + nr$
 - ▶ This factorization is nonunique, $\mathbf{UV}^T = (\mathbf{UM})(\mathbf{VM}^{-T})^T$
- ▶ Eckart-Young (optimal low-rank approximation by SVD) theorem
 - ▶ Truncated SVD approximates \mathbf{A} as

$$\mathbf{A} \approx \tilde{\mathbf{A}} = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T$$

where $\sigma_1, \dots, \sigma_r$ are the largest r singular values, while \mathbf{u}_i and \mathbf{v}_i are the associated left and right singular vectors

- ▶ Eckart-Young theorem demonstrates that the truncated SVD minimizes

$$\underbrace{\|\mathbf{A} - \tilde{\mathbf{A}}\|_2}_{\sigma_{r+1}} \quad \text{and} \quad \underbrace{\|\mathbf{A} - \tilde{\mathbf{A}}\|_F^2}_{\sum_{i=r+1}^{\min(m,n)} \sigma_i^2}$$

Rank Revealing Matrix Factorizations

- ▶ Computing the SVD
 - ▶ *Can compute full SVD with $O(mn \min(m, n))$ cost via bidiagonalization*
 - ▶ *unconditionally stable and accurate*
 - ▶ *inefficient for low r or if A is sparse*
 - ▶ *Given any low-rank approximation composed of U and V , compute QR of each and SVD of product of R factors to obtain SVD with total cost $O((m+n)r^2)$*
- ▶ QR with column pivoting
 - ▶ *By selecting columns of largest norm in the trailing matrix during QR factorization, we obtain a pivoted factorization with permutation matrix P*

$$AP = QR$$

- ▶ *Truncating this factorization can be done after applying r Householder reflectors (or another QR algorithm on r columns), with cost $O((m+n)r)$*
- ▶ *Approximation is somewhat suboptimal in theory, but in practice almost always as accurate as truncated SVD*

Simultaneous and Orthogonal Iteration

- ▶ Orthogonal iteration computing many eigenvectors at once in an iterative way

- ▶ Initialize $\mathbf{X}_0 \in \mathbb{R}^{n \times k}$ to be random, orthogonalize it to obtain \mathbf{Q}_0 , then iterate via

$$\mathbf{Q}_{i+1}\mathbf{R}_{i+1} = \mathbf{A}\mathbf{Q}_i$$

- ▶ For random starting guess, with high probability, $\lim_{i \rightarrow \infty} \text{span}(\mathbf{X}_i) = \mathbb{S}$ where \mathbb{S} is the subspace spanned by the k eigenvectors of \mathbf{A} with the largest eigenvalues in magnitude.
 - ▶ Can use this to compute the right singular vectors of matrix \mathbf{M} by using $\mathbf{A} = \mathbf{M}^T \mathbf{M}$ (no need to form \mathbf{A} , just multiply \mathbf{Q}_i by \mathbf{M}^T then \mathbf{M}).
 - ▶ QR has cost $O(nk^2)$ while product has cost $O(n^2k)$ (or more generally, k products with \mathbf{A}) per iteration.
 - ▶ QR iteration performs orthogonal iteration implicitly when $n = k$

Orthogonal Iteration Convergence

- ▶ If \mathbf{A} has distinct eigenvalues and \mathbf{R}_i has positive decreasing diagonal, the j th column of \mathbf{Q}_i converges to the j th Schur vector of \mathbf{A} linearly with rate $\max(|\lambda_{j+1}/\lambda_j|, |\lambda_j/\lambda_{j-1}|)$.
 - ▶ *Convergence of the first column of \mathbf{Q}_i follows by analogy to power iteration*
 - ▶ *Span of first j columns of \mathbf{Q}_i converges to the span of the first j Schur vectors with rate $|\lambda_{j+1}/\lambda_j|$*
 - ▶ *Hence orthogonal iteration converges similarly to k instances of inverse iteration with shifts chosen near the k largest magnitude eigenvalues*
 - ▶ *Block-Krylov methods, which consider $\text{span}\{\mathbf{X}_0, \mathbf{A}\mathbf{X}_0, \dots, \mathbf{A}^{k-1}\mathbf{X}_0\}$ provide some improvement over orthogonal iteration for low rank SVD (see works by Ming Gu and others)*

Randomized SVD

- ▶ Orthogonal iteration for SVD can also be viewed as a randomized algorithm
 - ▶ Suppose that we have an exact low-rank factorization $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ with $\mathbf{\Sigma} \in \mathbb{R}^{r \times r}$
 - ▶ If $\mathbf{Q}^{(0)}$ is a random orthogonal matrix, so is $\mathbf{V}^T \mathbf{Q}^{(0)}$
 - ▶ Consequently, $\mathbf{A}\mathbf{Q}^{(0)}$ is a set of r random linear combinations of columns of $\mathbf{U}\mathbf{\Sigma}$
 - ▶ Further, after the QR $\mathbf{Q}^{(1)}\mathbf{R}^{(1)} = \mathbf{A}\mathbf{Q}^{(0)}$,

$$\mathbf{U}\mathbf{U}^T = \mathbf{Q}^{(1)}\mathbf{Q}^{(1)T}$$

holds with probability 1 (suffices to have $\mathbf{A}\mathbf{Q}^{(0)}$ full rank)

- ▶ Consequently, we can compute SVD of $\mathbf{Q}^{(1)T} \mathbf{A}$ (with cost $O(nr^2)$) and recover \mathbf{U} by premultiplying the computed left singular vectors by $\mathbf{Q}^{(1)}$
- ▶ When \mathbf{A} is not exactly low-rank, span of leading singular vectors can be captured by oversampling (e.g., selecting each $\mathbf{Q}^{(i)}$ to have $r + 10$ columns)
- ▶ Initial guess $\mathbf{Q}^{(0)}$ need not be orthogonal (Gaussian random performs well, structured pseudo-random enables $O(mn \log n)$ complexity for one-shot randomized SVD), but better accuracy is obtained with orthogonality

Generalized Nyström Algorithm

- ▶ The generalized Nyström algorithm provides an efficient way of computing a low-rank factorization given an approximation of its span³
 - ▶ *Given matrices $S_1 \in \mathbb{R}^{k \times n}$ and $S_2 \in \mathbb{R}^{k \times m}$ the rank k factorization of a matrix $A \in \mathbb{R}^{m \times n}$ is obtained via*

$$\hat{A}_k = AS_1^T (S_2 AS_1^T)^\dagger S_2 A$$

- ▶ *The truncated SVD is recovered if S_1 and S_2 contain the largest eigenvectors*
- ▶ *Generally, we expect $S_2 AS_1^T$ to be full rank, otherwise factorization is rank-deficient*
- ▶ *If $S_2 AS_1^T$ is invertible, $\forall \mathbf{u}$, $AS_1^T \mathbf{u} = \hat{A}_k S_1^T \mathbf{u}$*
- ▶ *The skeleton decomposition is obtained by choosing both S_1 and S_2 to be sampling matrices (each row being a unit vector)*
- ▶ *Instead, S_1 and S_2 may be chosen as random ‘sketch matrices’*
- ▶ *The interpolative decomposition is obtained by choosing either of the two to be a sampling matrix.*

³Nakatsukasa, Yuji, Fast and stable randomized low-rank matrix approximation, 2020.

Analysis of Generalized Nyström Algorithm

- ▶ Consider $\mathbf{F}_1 = \mathbf{A}\mathbf{S}_1^T$ and $\mathbf{F}_2 = \mathbf{A}\mathbf{S}_2^T$, derive the minimizer \mathbf{Z} to

$$\|\mathbf{A} - \mathbf{F}_1\mathbf{Z}\mathbf{F}_2^T\|_F$$

$$\text{vec}(\mathbf{A} - \mathbf{F}_1\mathbf{Z}\mathbf{F}_2^T) = \text{vec}(\mathbf{A}) - \text{vec}(\mathbf{F}_1 \otimes \mathbf{F}_2) \text{vec}(\mathbf{Z})$$

$$\text{vec}(\mathbf{Z}) = \text{vec}(\mathbf{F}_1 \otimes \mathbf{F}_2)^+ \text{vec}(\mathbf{A})$$

$$= \text{vec}(\mathbf{F}_1^+ \otimes \mathbf{F}_2^+) \text{vec}(\mathbf{A})$$

$$= \mathbf{F}_1^+ \mathbf{A} (\mathbf{F}_2^+)^T$$

- ▶ The generalized Nyström algorithm may be interpreted as applying a two-sided oblique projection of \mathbf{A}

- ▶ *Optimal solution above is given by orthogonal projections $\mathbf{F}_1\mathbf{F}_1^+$ and $\mathbf{F}_2\mathbf{F}_2^+$*
- ▶ *Generalized Nyström approximation instead uses the oblique projections*

$$\mathbf{P}_1 = \mathbf{A}\mathbf{S}_1(\mathbf{S}_2\mathbf{A}\mathbf{S}_1^T)^+ \mathbf{S}_2, \mathbf{P}_2 = \mathbf{S}_1^T(\mathbf{S}_2\mathbf{A}\mathbf{S}_1^T)^+ \mathbf{S}_2\mathbf{A}$$

where $\mathbf{P}_1\mathbf{P}_1 = \mathbf{P}_1$ and $\mathbf{P}_2\mathbf{P}_2 = \mathbf{P}_2$, while the approximation obtained via $\mathbf{P}_1\mathbf{A}\mathbf{P}_2$