# CS 598 EVS: Tensor Computations
## Matrix Computations Background

Edgar Solomonik

University of Illinois at Urbana-Champaign

# Conditioning

- **Absolute Condition Number**:

- **(Relative) Condition Number**:
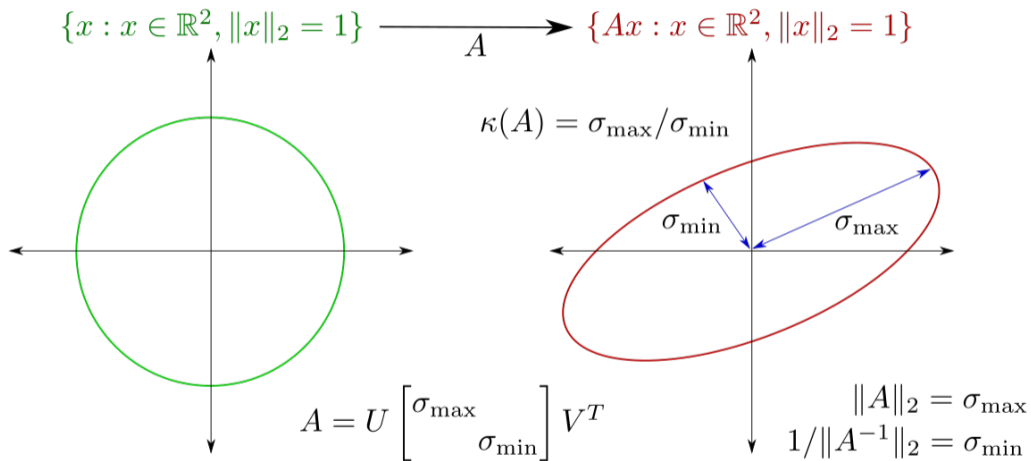
# Posedness and Conditioning

- **What is the condition number of an ill-posed problem?**

# Matrix Condition Number

- The matrix condition number $\kappa(\boldsymbol{A})$ is the ratio between the max and min distance from the surface to the center of the unit ball (norm-1 vectors) transformed by $\boldsymbol{A}$:

- The matrix condition number bounds the worst-case amplification of error in a matrix-vector product:

# Singular Value Decomposition

- The singular value decomposition (SVD)

- Condition number in terms of singular values

# Visualization of Matrix Conditioning

# Linear Least Squares

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

- Given the SVD $A = U\Sigma V^T$ we have $x^\star = \underbrace{V\Sigma^\dagger U^T}_{A^\dagger} b$, where $\Sigma^\dagger$ contains the reciprocal of all nonzeros in $\Sigma$, and more generally † denotes pseudoinverse:

# Normal Equations

- *Normal equations* are given by solving $A^T A x = A^T b$:

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

# Solving the Normal Equations

- If $A$ is full-rank, then $A^T A$ is symmetric positive definite (SPD):

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

# 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 $A = QR$

- 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

- We can solve the normal equations (and consequently the linear least squares problem) via reduced QR as follows

# Computing the QR Factorization

- The Cholesky-QR algorithm uses the normal equations to obtain the QR factorization

- Orthogonalization-based methods are most efficient and stable for QR factorization of dense matrices

# Eigenvalue Decomposition

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

- $A$ and $B$ are *similar*, if there exist $Z$ such that $A = ZBZ^{-1}$

# Similarity of Matrices

| matrix | similarity | reduced form |
|---:|---|---|
| SPD | | |
| real symmetric | | |
| Hermitian | | |
| normal | | |
| real | | |
| diagonalizable | | |
| arbitrary | | |

# Rayleigh Quotient

▸ For any vector $x$ that is close to an eigenvector, the *Rayleigh quotient* provides an estimate of the associated eigenvalue of $A$:

# Introduction to Krylov Subspace Methods

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

$$K_k = \begin{bmatrix} x_0 & Ax_0 & \cdots & A^{k-1}x_0 \end{bmatrix}$$

starting vectors

$\text{span}(k_k)$

$\min_{x \in \text{span}(k_k)} \left\{ \frac{x^T A x}{x^T x}, \|Ax - b\|_2 \right\}$

- $A$ is similar to *companion matrix* $C = K_n^{-1} A K_n$:

# Krylov Subspaces

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

$$\mathcal{K}_k(A, x_0) = span(Q_k) = \{p(A)x_0 : deg(p) < k\},$$

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:

Power iteration

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

$$\text{compute } \rho = \frac{x^{(k)^T} A x^{(k)}}{x^{(k)^T} x^{(k)}}$$

$$x^{(k)} = x^{(k)} / \rho$$

# Krylov Subspace Methods

$$AX = XD$$

- The $k \times k$ matrix $H_k = Q_k^T A Q_k$ minimizes $||AQ_k - Q_k H_k||_2$:

Ritz vectors/values

(eigvec/eigvals of $H_k$)

approximate of A

$AQ_k \approx Q_v H_k$

- $H_k$ is upper-Hessenberg, because the companion matrix $C_n$ is upper-Hessenberg:

if A is symmetric

# Rayleigh–Ritz Procedure

- The eigenvalues/eigenvectors of $H_k$ are the *Ritz values/vectors*:

- The Ritz vectors and values are the *ideal approximations* of the actual eigenvalues and eigenvectors based on only $H_k$ and $Q_k$:

# Low Rank Matrix Approximation

▸ Given a matrix $\boldsymbol{A} \in \mathbb{R}^{m \times n}$ seek rank $r < m, n$ approximation



$$f_A(x) = Ax$$

$$\text{rank}(\tilde{A}) = r$$

$$\tilde{A} - A$$

▸ Eckart-Young (optimal low-rank approximation by SVD) theorem



$$\min_{\substack{\tilde{A} \\ \text{rank}(\tilde{A})=r}} \| \tilde{A} - A \|_F \quad \text{or} \quad \| U\tilde{A} - A \|_2 \qquad m \geq n$$

$$\sigma_i \geq \sigma_{i+1}$$

$$\| X \|_F^2 = \sum_i \sigma_i(A)^2$$

# Rank Revealing Matrix Factorizations

- Computing the SVD

diagonalize $\underline{A^T A}$ $\rightarrow$ ergvecs of $A^T A \hat{=}$ singvecs of A

$\overline{A V_{\cdot k}} = \bigcirc{U_1} \underline{S_1}$

$\underline{U S \cdot V^T}$

$\square = Q \bigcirc{R}$ $\quad U_R \underline{S V^T}$

$\bigcirc{Q U_R}$ $\quad$ $S$ $\quad \boxed{V^T}$

$u$

- QR with column pivoting

Golub-Kahan bidiagonalization $\underline{O(m^3)}$

$Q$ $\square$ $\xrightarrow{Q^T}$ $\square$ $\quad u$ $\square$ $\xrightarrow{V^T}$ $\square$

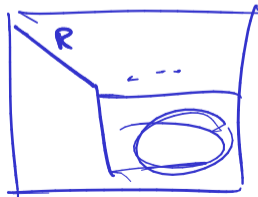# QR with column pivoting



largest norm

after $k$ steps

$\Rightarrow$ $Q_1$

$Q_1 \cdots Q_n$

$R$

# Orthogonal Iteration

▸ For sparse matrices, QR factorization creates fill, so must revert to iterative methods

$R = L$ factor of $A^T A$

$A$ is symmetric

$B_k = A Q_k$

$Q_{k+1} R = B_k$

$\text{span}(Q_\infty) = \text{span}($ of leading $k$ singular vectors of $A)$

$(A^T A)$  $A^T (A Q_k)$

Orthogonal iteration interleaves deflation and power iteration

# Randomized SVD

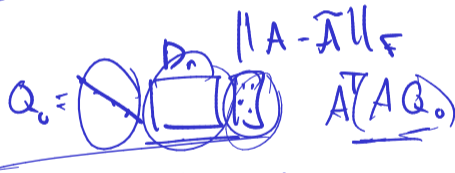$QR - col-piv$     $O(mnr)$

▸ Orthogonal iteration for SVD can also be viewed as a randomized algorithm

pick $Q_0$ to be random & orthogonal

$AQ_0$ if $A \to$ low rank (rank $r$)

$O(mn \cdot \log n)$

$U SV^t Q_0$

$Q_0 =$     $\|A - \hat{A}\|_F$

$D_n$     $A^T(AQ_0)$

$V^T Q_0$ as a set of random normalized linear combinations of cols of $US$

$Q_1 R = QR(AQ_0)$     $span(Q_1) = span(u)$     $u_1 \sigma_1 \cdots u_r \sigma_r$

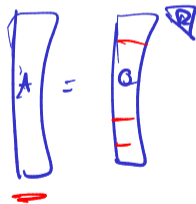$(USU^+ + E) \, Q_0$

$n \overline{\phantom{nnn}}$
$\quad n+10$

# Generalized Nyström Algorithm

▸ The generalized Nyström algorithm provides an efficient way of computing a sketched low-rank factorization

$$\widetilde{A} = A S_1 \left( S_2^T A S_1 \right)^+ S_2^T A$$

where $S_1$ and $S_2$ are sketching matrices

- Gaussian random

- FFT structure SRFT

- leverage-score sampling

# Multidimensional Optimization

- Minimize $f(\boldsymbol{x})$

- Quadratic optimization $f(\boldsymbol{x}) = \frac{1}{2}\boldsymbol{x}^T\boldsymbol{A}\boldsymbol{x} - \boldsymbol{b}^T\boldsymbol{x}$

# Basic Multidimensional Optimization Methods

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

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

    $$||\boldsymbol{e}_{k+1}||_{\boldsymbol{A}} =$$

    - 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 $\boldsymbol{A}$, since

# 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 $x_k - x_{k-1}$):

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

# Conjugate Gradient Method

- The *conjugate gradient method* is capable of making the optimal (for a quadratic objective) choice of $\alpha_k$ and $\beta_k$ at each iteration of an extrapolation method:

- *Parallel tangents* implementation of the method proceeds as follows

# Krylov Optimization

- Conjugate gradient (CG) finds the minimizer of $f(x) = \frac{1}{2}x^T A x - b^T x$ (which satisfies optimality condition $Ax = b$) within the Krylov subspace of $A$:

## CG and Krylov Optimization

The solution at the $k$th step, $\boldsymbol{y}_k = ||\boldsymbol{b}||_2 \boldsymbol{T}_k^{-1} \boldsymbol{e}_1$ is obtained by CG from $\boldsymbol{y}_{k+1}$ with a single matrix-vector product with $\boldsymbol{A}$ and vector operations with $O(n)$ cost

# Preconditioning

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

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

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

- Common preconditioners select parts of $A$ or perform inexact factorization

# Conjugate Gradient Convergence Analysis

- In previous discussion, we assumed $K_n$ is invertible, which may not be the case if $A$ has $m < n$ distinct eigenvalues, however, in exact arithmetic CG converges in $m - 1$ iterations[1]

---

[1] This derivation follows *Applied Numerical Linear Algebra* by James Demmel, Section 6.6.4

# Conjugate Gradient Convergence Analysis (II)

- Using $z = \rho_{k-1}(A)Ax$, we can simplify $\phi(z) = (x - z)^T A(x - z)$ as

- We can bound the objective based on the eigenvalues of $A = Q\Lambda Q^T$ using the identity $p(A) = Qp(\Lambda)Q^T$,

# Conjugate Gradient Convergence Analysis (III)

- Using our bound on the square of the residual norm $\phi(z)$, we can see why CG converges after $m - 1$ iterations if there are only $m < n$ distinct eigenvalues

- To see that the residual goes to $0$, we find a suitable polynomial in $\mathcal{Q}_m$ (the set of polynomials $q_m$ of degree $m$ with $q_m(0) = 1$)

# Round-off Error in Conjugate Gradient

- CG provides strong convergence guarantees for SPD matrices in exact arithmetic

- Due to round-off CG may stagnate / have plateaus in convergence

# Graph and Matrix Duality

- graphs have have a natural correspondence with sparse matrices

- matrix-based representations of graphs can be used to devise algorithms

# Graph Partitioning from Eigenvectors

- The Laplacian matrix provides a model of interactions on a graph that is useful in many contexts

- The second-smallest-eigenvalue eigenvector of the Laplacian (the Fiedler vector), gives a good partitioning of the graph

# 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 $f$:

# Nonlinear Least Squares

- An important special case of multidimensional optimization is *nonlinear least squares*, the problem of fitting a nonlinear function $f_{\boldsymbol{x}}(t)$ so that $f_{\boldsymbol{x}}(t_i) \approx y_i$:

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

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

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

# Constrained Optimization Problems

- We now return to the general case of *constrained* optimization problems:

$$\min_{\boldsymbol{x}} f(\boldsymbol{x}) \quad \text{subject to} \quad \boldsymbol{g}(\boldsymbol{x}) = \boldsymbol{0} \quad \text{and} \quad \boldsymbol{h}(\boldsymbol{x}) \leqslant \boldsymbol{0}$$

- Generally, we will seek to reduce constrained optimization problems to a series of simpler optimization problems: