# 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 $\boldsymbol{A}$:
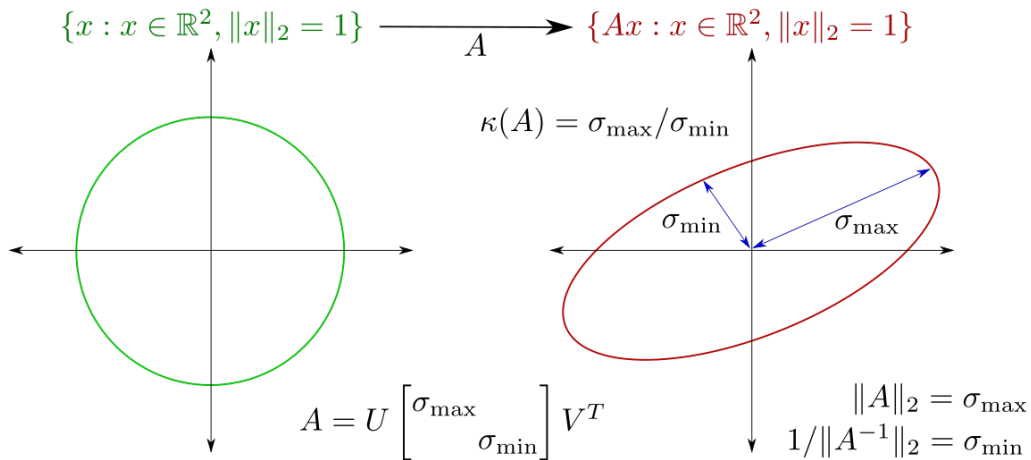
- The matrix condition number bounds the worst-case amplification of error in a matrix-vector product:

# 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 $\dagger$ 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

$$\boldsymbol{K}_k = \begin{bmatrix} \boldsymbol{x_0} & \boldsymbol{A}\boldsymbol{x_0} & \cdots & \boldsymbol{A}^{k-1}\boldsymbol{x_0} \end{bmatrix}$$

- $\boldsymbol{A}$ is similar to *companion matrix* $\boldsymbol{C} = \boldsymbol{K}_n^{-1}\boldsymbol{A}\boldsymbol{K}_n$:

# Krylov Subspaces

▸ Given $\boldsymbol{Q}_k \boldsymbol{R}_k = \boldsymbol{K}_k$, we obtain an orthonormal basis for the Krylov subspace,

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

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:

# Krylov Subspace Methods

- The $k \times k$ matrix $\boldsymbol{H}_k = \boldsymbol{Q}_k^T \boldsymbol{A} \boldsymbol{Q}_k$ minimizes $||\boldsymbol{A}\boldsymbol{Q}_k - \boldsymbol{Q}_k\boldsymbol{H}_k||_2$:

- $\boldsymbol{H}_k$ is upper-Hessenberg, because the companion matrix $\boldsymbol{C}_n$ is upper-Hessenberg:

# 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 $A \in \mathbb{R}^{m \times n}$ seek rank $r < m, n$ approximation

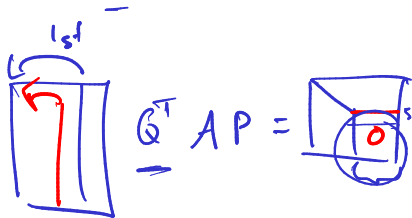- Eckart-Young (optimal low-rank approximation by SVD) theorem

# Rank Revealing Matrix Factorizations

‣ Computing the SVD



‣ QR with column pivoting

$$AP = QR$$

$$Q^T AP = $$

# Orthogonal Iteration

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

$$Q_{k+1} \leftarrow QR(\overset{\text{symmetrize}}{A}Q_k)$$

$$\underbrace{\phantom{AAAA}}_{A^T A}$$

- Orthogonal iteration interleaves deflation and power iteration

# Randomized SVD

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

# Generalized Nyström Algorithm

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

$S$ – sketch matrix, $n \times k$ ← sample size

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

• Gaussian random ← will densify

$O(n \cdot k)$

• SRFT $\quad S = \sqrt{} D_n P$

$O(n \log n)$

↑ diagonal $\quad$ ↑ random FFT permutation

efficient if sketched data dense

• for sparse $A$, want to sample, Countsketch
• leverage score sampling

$S_1 S_2$

$$\underbrace{A S_2^T}_{n \times k} \underbrace{(S_1 A S_2^T)}_{k \times k}^+ \underbrace{S_1 A}_{k \times n} = \tilde{A}$$

# Multidimensional Optimization

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

$$x \in \mathbb{R}^n, \quad f: \mathbb{R}^n \to \mathbb{R}$$

- constrained or unconstrained

  equality

  inequality

Constrained & nonlinear

Lagrangian | interior point

unconstrained
nonlinear optimization

$\nabla f(x^*) = 0$
or KKT

$\downarrow$ Newton's method

nonlinear solve $\downarrow$ Newton's method

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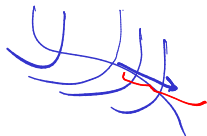
  linear system

$$\nabla f(x^*) = 0$$

$$A x = b$$

$A$ is SPD

# Basic Multidimensional Optimization Methods

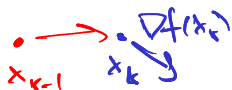- Steepest descent: minimize $f$ in the direction of the negative gradient:



- 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

$$\|\boldsymbol{e}_{k+1}\|_A = \boldsymbol{e}_{k+1}^T \boldsymbol{A} \, \boldsymbol{e}_{k+1} = \|\boldsymbol{e}_k\|_A \, \frac{\kappa(A) - 1}{\kappa(A) + 1}$$

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

$$x_{k+1} - x_k = \alpha_k \nabla f(x_k) + \beta_k (x_k - x_{k-1})$$

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

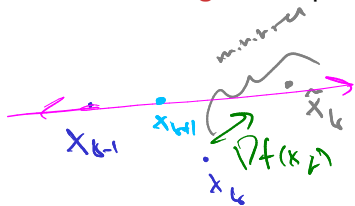$$\| e_{k+1} \|_A = \| e_k \|_A \frac{\sqrt{\kappa(A)} - 1}{\sqrt{\kappa(A)} + 1}$$

Nestrov's method

# 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:

- if $A$ is $n \times n$, CG converges in $n$ iterations

- $x_{k+1} - x_k$ is $A$-orthogonal to prior directions

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



1st minimized $\quad x_k + \alpha \nabla f(x_k) \Rightarrow \bar{x}_k$

2nd minimized $\quad x_{k-1} \perp \beta (\bar{x}_k - x_{k-1})$

# 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 $A$:

. each iteration of CG involves 1 matvec with $A$
  and vector ops

$$x_k \in \min_{\substack{x \in K_k(A,b) \\ \underbrace{\quad}_{span(Q_k)}}} f(x) = \min_{\substack{y \in R^k \\ x = Q_k y}} \frac{1}{2} y^T \underbrace{Q_k^T A Q_k}_{T_k} y \overbrace{\quad}^{\text{is tridiagonal}} - y Q_k^T b$$

$$\underbrace{\qquad\qquad\qquad\qquad\qquad}_{\varphi(y)}$$

$e_1 \|b\|_2$

$$0 = \nabla \varphi(y) = \underbrace{Q_k^+ A Q_k}_{T_k} y - Q_k^T b$$

$$T_k \, y = e_1 \|b\|_2$$

# 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

$$T_{k+1} = \begin{bmatrix} & & \\ & T_k & \\ & & \end{bmatrix}$$

$$y_{k+1} = ||b||_2 \, T_{k+1}^{-1} e_1$$

$$O(k)$$

$$\text{ran}\left( T_{k+1} - \begin{bmatrix} T_k & \\ & T_{k+1}(k+1,k+1) \end{bmatrix} \right) = 2$$

$$\left( M - u v^\top \right)^{-1} = \left( M^{-1} + \frac{M^{-1} u v^\top M^{-1}}{1 - v^\top M^{-1} u} \right) \qquad O(k)$$

# Preconditioning

$$\min \| {}_{=}Ax - {}^{-1}b \|_{M^{-1}}$$

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

often, pick $M \approx A$
so that $M^{-1}$ is easy
to obtain/apply

never form $M^{-1}A$
only apply $A$ and solve with $M$

$\underset{\sim}{M^{-1}A}$

$$A = \begin{bmatrix} u_1 & \cdots & u_n \end{bmatrix} \begin{smallmatrix} \sigma_{max} \\ & \sigma_{min} \end{smallmatrix} \begin{bmatrix} v_1 & \cdots & v_n \end{bmatrix}^{\top} \to A - u_1 \sigma_{max} v_1^{\top} + u_1 v_1^{\top}$$

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

  • incomplete LU (LU without fill)

# 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

- 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:

# Lagrangian Duality

- The Lagrangian function with constraints $g(x) = 0$ and $h(x) \leqslant 0$ is

- The Lagrangian dual problem is an unconstrained optimization problem:

$$\max_{\lambda} q(\lambda), \quad q(\lambda) = \begin{cases} \min_x \mathcal{L}(x, \lambda) & \text{if } \lambda \geqslant 0 \\ -\infty & \text{otherwise} \end{cases}$$

The unconstrained optimality condition $\nabla q(\lambda^*) = 0$, implies

# Sequential Quadratic Programming

- *Sequential quadratic programming (SQP)* reduces a nonlinear equality constrained problem to a sequence of constrained quadratic programs via a Taylor expansion of the Lagrangian function $\mathcal{L}_f(\boldsymbol{x}, \boldsymbol{\lambda}) = f(\boldsymbol{x}) + \boldsymbol{\lambda}^T \boldsymbol{g}(\boldsymbol{x})$:

- SQP ignores the constant term $\mathcal{L}_f(\boldsymbol{x}_k, \boldsymbol{\lambda}_k)$ and minimizes $\boldsymbol{s}$ while treating $\boldsymbol{\delta}$ as a Lagrange multiplier:

# Interior Point Methods

- Barrier functions provide an effective way of working with inequality constraints $h(x) \leqslant 0$:

- Interior point methods additionally incorporate Lagrangian optimization