

CS 598 EVS: Tensor Computations

Matrix Computations Background

Edgar Solomonik

University of Illinois at Urbana-Champaign

Conditioning of Problems - sensitivity to perturbations

- ▶ **Absolute Condition Number:**

Problem: $f(x)$
 ↑
 input

$$\kappa_{\text{abs.}}(f, x) = f'(x) = \lim_{\delta x \rightarrow 0} \frac{|f(x + \delta x) - f(x)|}{|\delta x|}$$

- ▶ **(Relative) Condition Number:**

$$\kappa(f, x) = \lim_{\delta x \rightarrow 0} \frac{|f(x + \delta x) - f(x)| / |f(x)|}{|\delta x| / |x|} = \frac{|f'(x)| |x|}{|f(x)|}$$

$$\kappa(f) = \max_{x \in D} \kappa(f, x)$$

Posedness and Conditioning

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

$$f'(x) = \infty$$

well-posed problem

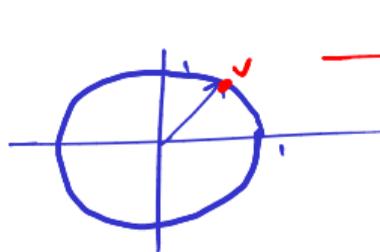
- unique
- continuously varying solution

$\kappa(f)$ to be finite

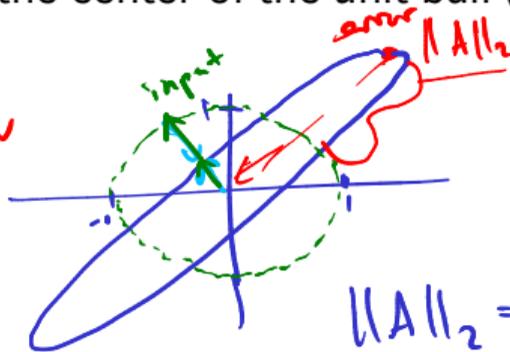
Matrix Condition Number

$$f(x) = Ax \quad \|Ax\| \leq \|A\|_F \|B\|_F$$

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



$$A \rightarrow Av$$



$$\kappa(A) = \|A\|_2 \cdot \|x\|_2$$

$$\|A\|_F = \| \text{vecc}(A) \|_2 = \sqrt{\sum_{i,j} a_{ij}^2}$$

$$\|A\|_2 = \max_{\|x\|_2=1} \|Ax\|_2$$

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

$$\begin{aligned}
 & x \rightarrow x + \delta x \\
 & \leq \frac{\|A\|_2 \|\delta x\|_2}{\|Ax\|} \leq \underbrace{\|A\|_2 \|A^{-1}\|_2}_{\kappa(A)} \frac{\|A \delta x\|}{\|x\|} \leq \dots
 \end{aligned}$$

Singular Value Decomposition

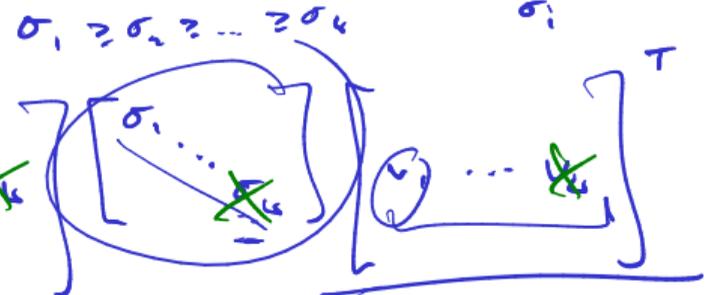
- ▶ The singular value decomposition (SVD)

$$A = U S V^T = \begin{bmatrix} u_1 & \dots & u_k & \dots & u_{m-k} \end{bmatrix} \begin{bmatrix} \sigma_1 & & & & \\ & \dots & & & \\ & & \sigma_k & & \\ & & & \dots & \\ & & & & \sigma_{m-k} & & & \end{bmatrix} \begin{bmatrix} v_1 & \dots & v_k & \dots & v_{n-k} \end{bmatrix}^T$$

\downarrow orthogonal / positive
 \uparrow orthogonal, $U^T U = I$
 \uparrow orthogonal, $V^T V = I$

U and V contain singular vectors u_i and v_i .
 $A v_i = \sigma_i u_i$

$$A v_i = U S V^T v_i = U \underbrace{S e_i}_{\sigma_i} = u_i \sigma_i$$



vectors of A

$$A^T u_i = \sigma_i v_i$$

- ▶ Condition number in terms of singular values

$$\|A\|_2 = \sigma_{\max}(A)$$

$$\kappa(A) = \frac{\sigma_{\max}(A)}{\sigma_{\min}(A)}$$

$$\|A^{-1}\|_2 = \frac{1}{\sigma_{\min}(A)}$$

$$\|A^+\|_2 = \frac{1}{\sigma_{\min}(A)}$$

$$S = \begin{bmatrix} \sigma_1 & & \\ & \dots & \\ & & \sigma_k & \\ & & & \dots & \\ & & & & \sigma_{m-k} & \\ & & & & & \dots & \end{bmatrix}$$

$\alpha_1, \dots, \alpha_k$

$$A(\alpha_1 v_1 + \dots + \alpha_k v_k) = \sigma_1 \alpha_1 u_1 + \dots + \sigma_k \alpha_k u_k$$

$$A^+ A = I$$

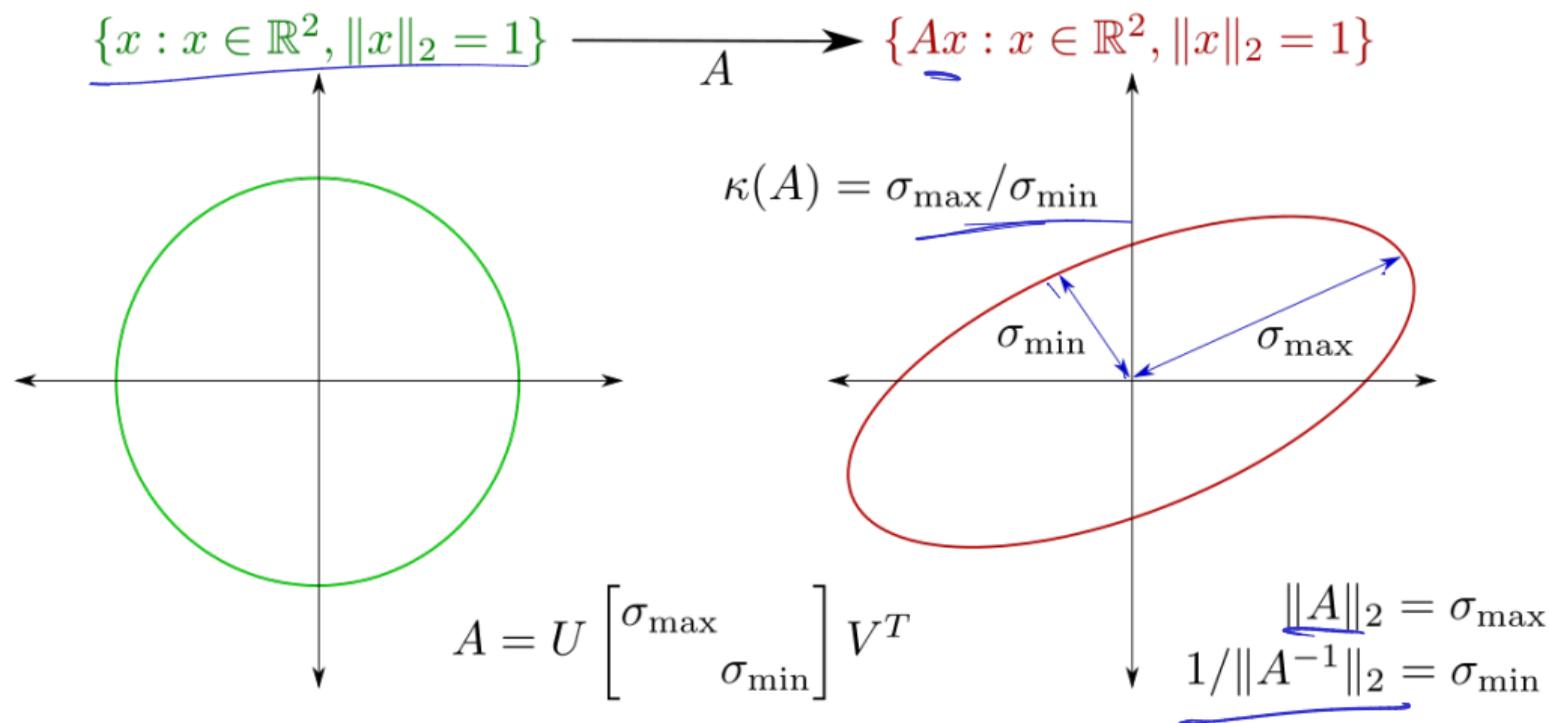
$$\frac{x = A^+ b}{A x \approx b} \quad \| \cdot \|_2$$

$$A = U S V^T$$

$$A^+ = V S^+ U^T \rightarrow S^{-1}$$

$$A^+ = (A^+ A)^+ A^+$$

Visualization of Matrix Conditioning



Normal Equations

Demo: Normal equations vs Pseudoinverse

Demo: Issues with the normal equations

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

$A^T A$ should be SPD

$$\lambda(A^T A) > 0$$

$$A^+ = (A^T A)^{-1} A^T$$

$$A^T A = A_1^T A_1 \oplus A_2^T A_2$$

$$A = \begin{bmatrix} \Delta_1 \\ \Gamma \end{bmatrix} \begin{bmatrix} \Sigma_1 \\ \Pi \end{bmatrix}$$

$$\underbrace{x^T}_{y^T} \underbrace{A^T A}_{\lambda x} \underbrace{x}_{y} = \lambda \underbrace{x^T x}_{> 0}$$

$$A^T A = V S U^T U S V^T = V S^2 V^T$$

$$\kappa(A^T A) = \kappa(A)^2$$

$$\begin{bmatrix} \sigma_1^2 & & \\ & \dots & \\ & & \sigma_n^2 \end{bmatrix}$$

- However, solving the normal equations is a more ill-conditioned problem than 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:

A is SPD

$$A = A^T$$

$$\lambda(A) > 0$$

$Ax = b$ $LL^T x = b$ $Ly = b$
 $\Delta y = b$

\Rightarrow $\begin{bmatrix} \hat{n} \\ A \end{bmatrix} = \left(\begin{bmatrix} L \\ \end{bmatrix} \begin{bmatrix} L^T \end{bmatrix} \right)$ $\frac{O(n^3)}{\frac{2}{3}n^3}$ $L^T x = y$
 $\Delta x = y$

$PA = \begin{bmatrix} L \\ \end{bmatrix} \begin{bmatrix} U \end{bmatrix}$

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

$$\underline{A}x \cong b \qquad \underline{QR}x \cong b \qquad Rx = \underline{Q^T}b$$
$$x = \underline{R^{-1}}\underline{Q^T}b$$

- ▶ 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

$$Ax = \lambda x$$

$$X = \begin{bmatrix} x_1 & \dots & x_n \end{bmatrix}$$

$$AX = XD = \begin{bmatrix} \lambda_1 & & \\ & \dots & \\ & & \lambda_n \end{bmatrix}$$

$$X^{-1}AX = D$$

$$A = XD X^{-1}$$

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

equivalents of A and B are the same

$$Ax = \lambda x$$

$$ZBZ^{-1}x = \lambda x$$

$$BZ^{-1}x = \lambda Z^{-1}x$$

$$\underbrace{\quad}_y = \lambda \underbrace{\quad}_y$$

Eigenvalue Decomposition

$$A = \underbrace{U}_{\text{orthogonal}} \underbrace{D}_{\text{orthogonal}} \underbrace{U^T}_{\text{orthogonal}} \quad U^T U = I$$

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

— orthogonal
— orthogonal
— positive diagonal
— diagonal matrix
— tridiagonal matrix

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



Similarity of Matrices

<i>matrix</i>	<i>similarity</i>	<i>reduced form</i>
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

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

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

Krylov Subspaces

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

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

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 $\mathbf{H}_k = \mathbf{Q}_k^T \mathbf{A} \mathbf{Q}_k$ minimizes $\|\mathbf{A} \mathbf{Q}_k - \mathbf{Q}_k \mathbf{H}_k\|_2$:

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

Rayleigh-Ritz Procedure

- ▶ The eigenvalues/eigenvectors of \mathbf{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 \mathbf{H}_k and \mathbf{Q}_k :

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

Multidimensional Optimization

- ▶ Minimize $f(\mathbf{x})$

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

Basic Multidimensional Optimization Methods

- ▶ Steepest descent: minimize f 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

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

- ▶ 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

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 $\mathbf{x}_k - \mathbf{x}_{k-1}$):

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

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} :

CG and Krylov Optimization

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

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

- ▶ Common preconditioners select parts of \mathbf{A} or perform inexact factorization

Conjugate Gradient Convergence Analysis

- ▶ In previous discussion, we assumed \mathbf{K}_n is invertible, which may not be the case if \mathbf{A} has $m < n$ distinct eigenvalues, however, in exact arithmetic CG converges in $m - 1$ iterations¹

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

Conjugate Gradient Convergence Analysis (III)

- ▶ Using our bound on the square of the residual norm $\phi(\mathbf{z})$, 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 $q_m(0) = 1$)

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 $f_{\mathbf{x}}(t)$ so that $f_{\mathbf{x}}(t_i) \approx y_i$:

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