

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 $\kappa(\mathbf{A})$ 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} :

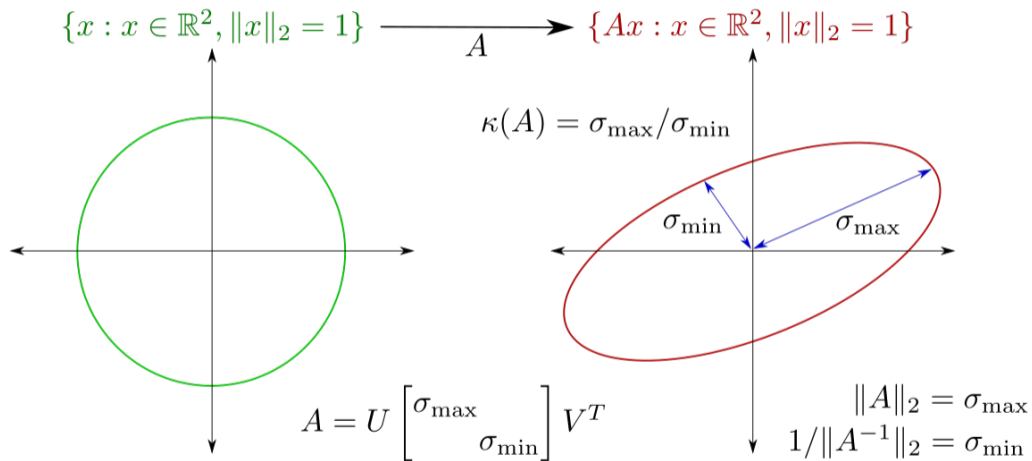
- ▶ The matrix condition number bounds the worst-case amplification of error in a matrix-vector product:

Singular Value Decomposition

- ▶ The singular value decomposition (SVD)

- ▶ Condition number in terms of singular values

Visualization of Matrix Conditioning



Linear Least Squares

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

- ▶ 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:

Normal Equations

Demo: Normal equations vs Pseudoinverse

Demo: Issues with the normal equations

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

- ▶ 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:

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

<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 :

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

Orthogonal Iteration

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

Multidimensional Optimization

- ▶ Minimize $f(\mathbf{x})$ \mathbb{R}^n
nonlinear

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

Basic Multidimensional Optimization Methods

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

line search

- ▶ 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:

$$\lim_{k \rightarrow \infty} \frac{\|e_{k+1}\|}{\|e_k\|} \leq \frac{1}{e}$$

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:

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

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 (II)

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

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

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 $q_m(0) = 1$)

Graph and Matrix Duality

- ▶ graphs have have a natural correspondence with sparse matrices

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

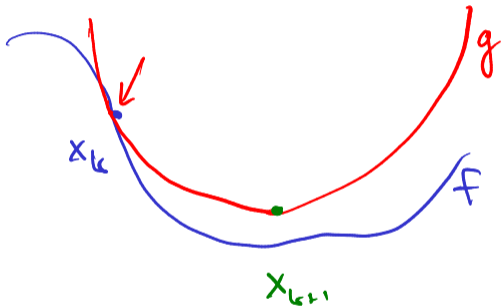
Graph Partitioning from Eigenvectors

- ▶ 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 :



$$g(x) = \frac{1}{2}x^T A x - x^T b$$

$$f(x_k + x) = f(x_k) + \nabla f(x_k)^T x_k$$

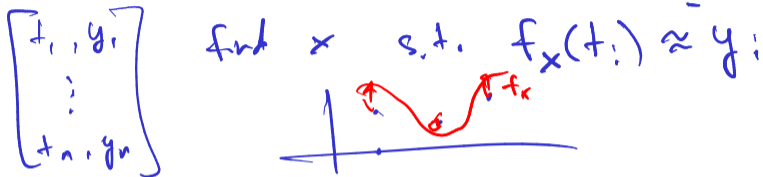
$$+ \frac{1}{2} x^T H_f(x_k) x + \dots$$

$$A x = b$$

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

Nonlinear Least Squares

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



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

$$r_i(x) = f_x(t_i) - y_i$$
$$\min_x e(x), \quad e(x) = \frac{1}{2} r^T(x) r(x)$$

$$\underline{H_e(x_0)} s_k = -\nabla e(x_0)$$

Gauss-Newton Method

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

$$\nabla \text{vec} \mathcal{H} = \int_r^T(x) r(x)$$

$$H_k(x) = \underbrace{\int_r^T(x) J_r(x)}_{\text{small}} + \underbrace{\sum_{i=1}^n w_{r_i(x)}(x) r_i(x)}_{\text{small}}$$

where r_i is small

$$\int_r^T(x) J_r(x) s_k = \int_r^T(x) r(x)$$

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

normal equations

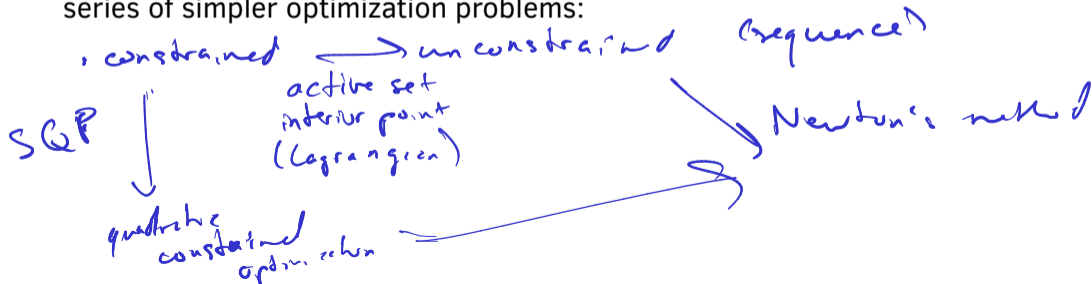
$$\underbrace{J_r(x_k)}_{\text{circle}} s_k \approx r(x_k)$$

Constrained Optimization Problems

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

$$\min_x f(x) \quad \text{subject to} \quad \underbrace{g(x) = 0}_{\begin{bmatrix} g_1(x) \\ \vdots \\ g_m(x) \end{bmatrix} = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}} \quad \text{and} \quad \underbrace{h(x) \leq 0}_{\begin{bmatrix} h_1(x) \\ \vdots \\ h_p(x) \end{bmatrix} \leq \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}}$$

- ▶ 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 $h(x) \leq 0$ is

$$\mathcal{L}(x, \lambda) = \underbrace{f(x)}_{\text{unconstrained min.}} + \underbrace{\begin{bmatrix} \lambda^T h(x) \\ g(x) \end{bmatrix}}_{\text{critical with } \mathcal{L}}$$

unconstrained min. on a critical point of \mathcal{L}

- ▶ The Lagrangian dual problem is an unconstrained optimization problem:

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

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

$$\begin{aligned} \max(\lambda^*, h(x)) &= \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} && \text{ineq is exactly satisfied} \\ \max(\lambda^*, g(x)) &= \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} && \text{or } \lambda = 0 \\ &&& \Rightarrow g(x) = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} \end{aligned}$$

Sequential Quadratic Programming

ignore h.c.t.

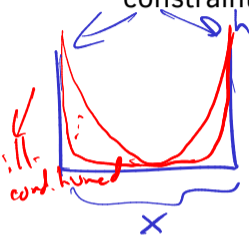
- ▶ **Sequential quadratic programming (SQP)** reduces a nonlinear equality constrained problem to a sequence of constrained quadratic programs via a Taylor expansion of the Lagrangian function $\mathcal{L}_f(x, \lambda) = f(x) + \lambda^T g(x)$:

$$\mathcal{L}_f(x_k + s, \lambda_k + u) \approx f(x_k) + s^T \nabla f(x_k) + \frac{1}{2} s^T H_f(x_k) s + \lambda_k^T g(x_k) + u^T g(x_k) + \frac{1}{2} s^T \sum_i \lambda_i H_{g_i}(x_k) s$$

- ▶ SQP ignores the constant term $\mathcal{L}_f(x_k, \lambda_k)$ and minimizes s while treating δ as a Lagrange multiplier:

Interior Point Methods

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



$$f(x) \rightarrow \frac{f(x) - \mu \log(-h(x))}{1}$$

$$\rightarrow \frac{f(x) + \mu \sum \frac{1}{h_i(x)}}{1}$$

$$\mu_1 > \mu_2 > \dots > \mu_r$$



- ▶ Interior point methods additionally incorporate Lagrangian optimization

Karush-Kuhn-Tucker (KKT) conditions

Consider the linear-constrained Quadratic program (QP): Its Lagrangian

function may be used to derive an interior point method The first-order

optimality (KKT) conditions are

Primal-dual Interior Point Method (IPM)

Solve perturbed KKT conditions after introducing slack variables $s \in \mathbb{R}^{m_2}$

Interior Point Method (IPM): KKT system

Newton's method applied to KKT equations results in linear systems

These linear systems become ill-conditioned as the interior point method approaches convergence