

CS 598 EVS: Tensor Computations

Matrix Computations Background

Edgar Solomonik

University of Illinois at Urbana-Champaign

Conditioning

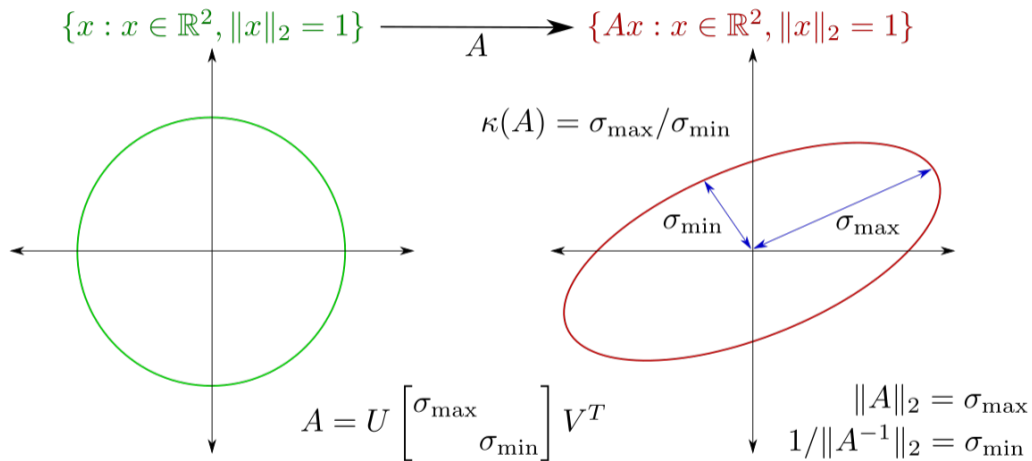
- ▶ **Absolute Condition Number:**

- ▶ **(Relative) Condition Number:**

Posedness and Conditioning

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

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

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

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:

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 (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 $q_m(0) = 1$)

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:

Constrained Optimization Problems

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

$$\min_x f(x) \quad \text{subject to} \quad g(x) = \mathbf{0} \quad \text{and} \quad h(x) \leq \mathbf{0}$$

equality *inequality*

$$Ax = b \qquad x \geq 0$$

- ▶ 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 $h(x) \leq 0$ is

$$\mathcal{L}(x, \lambda) = \underline{f(x)} + \lambda^T \begin{bmatrix} h(x) \\ g(x) \end{bmatrix}$$

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

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

strong duality

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

$$\max \left(\lambda^*, \begin{bmatrix} h(x^*) \\ g(x^*) \end{bmatrix} \right) = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix}$$

$$f^* = q^* \\ \mathcal{L}(x^*, \lambda^*) = f^* = q^*$$

Optimality and Complementarity Slackness Condition

Consider the inequality-constrained optimization problem, $\mathbf{h}(\mathbf{x}) \leq \mathbf{0}$,

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) - \boldsymbol{\lambda}^T \mathbf{h}(\mathbf{x})$$

- ▶ The pair \mathbf{x}^* and $\boldsymbol{\lambda}^*$ are a primal-dual optimal solution \mathbf{x}^* is feasible, $\boldsymbol{\lambda}^* \geq \mathbf{0}$, and strong duality holds, $f(\mathbf{x}^*) = q(\boldsymbol{\lambda}^*) = \mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*)$

$$\mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*) = f(\mathbf{x}^*) - \underbrace{\boldsymbol{\lambda}^{*T} \mathbf{h}(\mathbf{x}^*)}_0$$

$$\lambda_i = 0 \text{ or } h_i(\mathbf{x}^*) = 0$$

Sequential Quadratic Programming

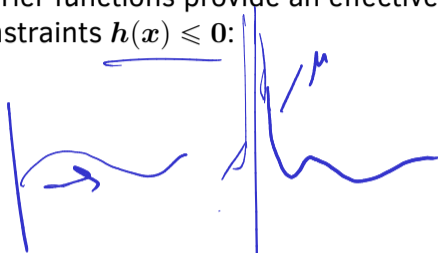
- ▶ *Sequential quadratic programming (SQP)* reduces a nonlinear equality constrained problem to a sequence of constrained quadratic programs via a Taylor expansion of the Lagrangian function $\mathcal{L}_f(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \boldsymbol{\lambda}^T \mathbf{g}(\mathbf{x})$:

$$H_{\mathcal{L}_f} = \begin{matrix} \mathbf{x} & \boldsymbol{\lambda} \\ \begin{bmatrix} H_f(\mathbf{x}) & \mathbf{J}_g^T(\mathbf{x}) \\ \mathbf{J}_g(\mathbf{x}) & \mathbf{0} \end{bmatrix} \end{matrix}$$

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

Interior Point Methods

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



- ▶ Interior point methods additionally incorporate Lagrangian optimization

Karush-Kuhn-Tucker (KKT) conditions

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

$$f(x) = \frac{1}{2} x^T H x + x^T b \quad \Bigg| \quad \mathcal{L}(x, \lambda, u) = \frac{1}{2} x^T H x + x^T c - \lambda (Ax - b) - u (Cx - d)$$

$Ax = b \quad Cx \geq d$ f

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

optimality (KKT) conditions are

$$\nabla \mathcal{L}_f = 0$$

$$\lambda, u \geq 0$$

$$Ax - b = 0$$

$$Cx - d \geq 0 \Rightarrow Cx - d - s = 0$$

$$\underline{s \geq 0}$$

$$u^T (Cx - d) = 0$$

Primal-dual Interior Point Method (IPM)

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

μ -barrier parameter

$$\nabla \mathcal{L} f = 0 \mid s, \lambda, \nu \geq 0 \mid Ax - b = 0 \mid Cx - d - s = 0$$

$$s_i \cdot \nu_i = \mu \frac{\langle s, \nu \rangle}{\dim(s)}$$

as $\mu \rightarrow 0$

δ^* for KKT
gives δ^* for
original problem



Interior Point Method (IPM): KKT system

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

$$\underbrace{J_{\mathcal{D}f(x)}^T}_{H_c} s_k = -\mathcal{D}f(x_k)$$

Newton for minimizing $f(x)$ = Newton for solving NLS with $\mathcal{D}f(x) = 0$

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

$$\mathbb{R}^M \begin{bmatrix} x & \lambda & u \\ -H & A^T & C \\ A & & \\ C & & \end{bmatrix} \begin{bmatrix} dx \\ d\lambda \\ du \end{bmatrix} = \mathcal{D}^{(k)}$$
$$\mathcal{D}^{(k)} = \begin{bmatrix} S U^{-1} \\ \dots \\ s_n / u_n \end{bmatrix}$$

$$f(x) = \frac{1}{2} x^T U x$$

$$Ax = b$$

$$x \geq 0$$

$$\Rightarrow Cx \geq 0$$

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

$$U \in \mathbb{R}^{n \times n}$$

with $C \in \mathbb{I}$

$$C \in \mathbb{R}^{n \times n}$$