

CS 598: Provably Efficient Algorithms for Numerical and Combinatorial Problems

Part 1: Background Numerical Analysis Tools

Edgar Solomonik

University of Illinois at Urbana-Champaign

Matrices and Tensors

- ▶ What is a matrix?

- ▶ What is a tensor?

Matrix and Tensor Decompositions

- ▶ What is a matrix factorization?
- ▶ What is a tensor decomposition?

Graphs and Hypergraphs

- ▶ What is a graph?

- ▶ What is a hypergraph?

Numerical and Combinatorial Problems

- ▶ What numerical problems will we look at?
- ▶ What combinatorial problems will we look at?
- ▶ What applications do these have?

Provably Efficient Algorithms

- ▶ What makes an algorithm provably efficient?

Provably Efficient Parallel Schedules

- ▶ What makes a parallel schedule good
- ▶ What makes a parallel schedule for a given algorithm optimal?

Error Analysis

- ▶ Forward Error:

- ▶ Backward Error:

Conditioning

- ▶ Conditioning measures the worst-case sensitivity of the output with respect to perturbations of the input
- ▶ The absolute condition number is a property of the problem, which measures its sensitivity to perturbations in input
- ▶ The relative condition number considers relative perturbations in input and output, so that

Rounding Error in Floating Point Operations

- ▶ Addition and Subtraction

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 transformed by $\kappa(\mathbf{A})$:
- ▶ 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

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

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,

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

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 Hessenberg, because the companion matrix \mathbf{C}_k is 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 :

General Multidimensional Optimization

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

- ▶ Given quadratic optimization problem $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \mathbf{A}\mathbf{x} + \mathbf{c}^T \mathbf{x}$ 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:

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 finds the minimizer of $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \mathbf{A}\mathbf{x} + \mathbf{c}^T \mathbf{x}$ (which satisfies optimality condition $\mathbf{A}\mathbf{x} = -\mathbf{c}$) within the Krylov subspace of \mathbf{A} :

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:

Gauss-Newton Method

- ▶ The Hessian for nonlinear least squares problems has the form:
- ▶ The *Gauss-Newton* method is Newton iteration with an approximate Hessian:

Tensors

► A *tensor* $\mathcal{T} \in \mathbb{R}^{n_1 \times \cdots \times n_d}$ has

► Order d tensors represent d -dimensional arrays

Reshaping Tensors

When using tensors, it is often necessary to transition between high-order and low-order representations of the same object

- ▶ Recall for a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ its *unfolding* is given by

$$\mathbf{v} = \text{vec}(\mathbf{A}) \Rightarrow$$

- ▶ A tensor $\mathcal{T} \in \mathbb{R}^{n_1 \times \dots \times n_d}$ can be fully unfolded the same way

$$\mathbf{v} = \text{vec}(\mathcal{T}) \Rightarrow$$

- ▶ Often we also want to *fold* tensors into higher-order ones
- ▶ Generally, we can *reshape* (fold or unfold) any tensor

$$\mathcal{U} = o_{n_1 \times \dots \times n_d}(\mathcal{V}) \Rightarrow$$

Canonical Polyadic (CP) Decomposition

- ▶ A rank R *CP decomposition* of an $s \times s \times s \times s$ tensor is
- ▶ We can represent the CP using the following *tensor diagram*:
- ▶ Finding an approximate tensor decomposition corresponds to a nonlinear least squares problem: