

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?
 - ▶ A collection of numbers arranged into an array of dimensions $m \times n$, e.g., $M \in \mathbb{R}^{m \times n}$
 - ▶ A linear operator $f(x) = Mx$
 - ▶ A bilinear form $x^T M y$

- ▶ What is a tensor?
 - ▶ A collection of numbers arranged into an array of a particular order, with dimensions $l \times m \times n \times \dots$, e.g., $\mathcal{T} \in \mathbb{R}^{l \times m \times n}$ is order 3
 - ▶ A multilinear operator $z = f(x, y)$

$$z_i = \sum_{j,k} t_{ijk} x_j y_k$$

- ▶ A multilinear form $\sum_{i,j,k} t_{ijk} x_i y_j z_k$

Matrix and Tensor Decompositions

- ▶ What is a matrix factorization?
 - ▶ *A decomposition of a matrix in terms of other matrices with desirable properties*
 - ▶ *$M = QR = USV^T = XDX^{-1}$ are examples of factorizations, where Q, U, V are orthogonal, R is upper-triangular, while S and D are diagonal*
 - ▶ *Factorizations enable compression of M , solution to linear systems and least squares problems with M , and computation of eigenvalues of M*
- ▶ What is a tensor decomposition?
 - ▶ *A decomposition of a tensors in terms of other tensors with desirable properties*
 - ▶ *For example, the canonical polyadic (CP) decomposition of an order 3 tensor is*

$$t_{ijk} = \sum_{r=1}^R a_{ir} b_{jr} c_{kr}$$

- ▶ *Factorizations enable compression of \mathcal{T} , may uncover semantic structure of \mathcal{T} , or may permit tensor network algorithms that represent \mathcal{T} implicitly*

Graphs and Hypergraphs

- ▶ What is a graph?
 - ▶ A set of vertices and edges $G = (V, E)$, with $V = \{1, \dots, n\}$ and $E \subseteq V \times V$ possibly with a weight function $w : E \rightarrow \mathbb{R}$
 - ▶ An adjacency matrix A , where for each $e = (i, j) \in E$, $a_{ij} = w(e)$
- ▶ What is a hypergraph?
 - ▶ A set of vertices and hyperedges $G = (V, H)$, where each $h \in H$ is a subset of vertices, $h \subseteq V$, possibly with a weight function w for hyperedges
 - ▶ A set of tensors $\mathcal{T}^{(1)}, \dots, \mathcal{T}^{(n)}$ where $\mathcal{T}^{(k)}$ is order k and stores hyperedges of order k , e.g. if $h = (i, j, k, l) \in H$ then $\mathcal{T}_{ijkl}^{(4)} = w(h)$

Numerical and Combinatorial Problems

- ▶ What numerical problems will we look at?
 - ▶ *Matrix factorizations via direct and approximate methods*
 - ▶ *Iterative methods for linear systems and eigenvalue problems*
 - ▶ *Discrete convolution*
 - ▶ *Tensor decomposition, tensor completion, tensor networks, and associated numerical optimization problems*
- ▶ What combinatorial problems will we look at?
 - ▶ *Sorting and partitioning*
 - ▶ *Graph problems: shortest paths, connectivity, minimal spanning tree*
- ▶ What applications do these have?
 - ▶ *Data mining and data compression (matrix/tensor decomposition)*
 - ▶ *Numerical methods for PDEs (linear systems and eigenvalue problems)*
 - ▶ *Machine learning, e.g. collaborative filtering (tensor completion, convolution)*
 - ▶ *Computational quantum chemistry and physics (tensor networks)*
 - ▶ *Graph analytics*

Provably Efficient Algorithms

- ▶ What makes an algorithm provably efficient?
 - ▶ *Polynomial time with respect to size of input*
 - ▶ *Low asymptotic worst-case execution time*
 - ▶ *Low asymptotic average-case execution time*
 - ▶ *Low (leading order) constant factors in execution time*
 - ▶ *Numerical stability (small errors made during execution are not amplified)*
 - ▶ *Approximation quality / error bounds (for inexact algorithms)*
 - ▶ *The algorithm is parallelizable (low depth)*
 - ▶ *The algorithm exhibits data-reuse / requires little communication*
 - ▶ *The algorithm requires little synchronization between threads or processes*
 - ▶ *The algorithm is cache oblivious*
 - ▶ *The runtime of an implementation of the algorithm is lower than comparable implementations of alternative algorithms on representative problem instances*

Provably Efficient Parallel Schedules

- ▶ What makes a parallel schedule good
 - ▶ *Low parallel execution time given an infinite number of processors*
 - ▶ *Low parallel execution time given any fixed number of processors*
 - ▶ *Different processors work on different data*
 - ▶ *Low communication and synchronization costs*
 - ▶ *Few cache misses*
 - ▶ *Low interprocessor communication volume*
 - ▶ *Few messages communicated*
 - ▶ *Low critical path costs*
- ▶ What makes a parallel schedule for a given algorithm optimal?
 - ▶ *Its execution time with an infinite number of processors is equal to the depth of the algorithm*
 - ▶ *Its execution time decreases linearly with the number of processors until it equals the depth*
 - ▶ *Its communication and synchronization costs match communication lower bounds for the algorithm*

Error Analysis

- ▶ Forward Error:

Forward error is the computational error of an algorithm

- ▶ Absolute: $\hat{f}(x) - f(x)$
- ▶ Relative: $(\hat{f}(x) - f(x))/f(x)$

- ▶ Backward Error:

Backward error analysis enables us to measure computational error with respect to data propagation error

- ▶ An algorithm is *backward stable* if its a solution to a nearby problem
- ▶ If the computed solution $\hat{f}(x) = f(\hat{x})$ then

$$\text{backward error} = \hat{x} - x$$

- ▶ More precisely, we want the nearest \hat{x} to x with $\hat{f}(x) = f(\hat{x})$

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

For scalar problem f with input x it is simply the derivative of f at x ,

$$\kappa_{abs}(f) = \lim_{\Delta x \rightarrow 0} \left| \frac{f(x + \Delta x) - f(x)}{\Delta x} \right| = \left| \frac{df}{dx}(x) \right|$$

When considering a space of inputs \mathcal{X} it is $\kappa_{abs} = \max_{x \in \mathcal{X}} \left| \frac{df}{dx}(x) \right|$

- ▶ The relative condition number considers relative perturbations in input and output, so that

$$\kappa(f) = \kappa_{rel}(f) = \max_{x \in \mathcal{X}} \lim_{\Delta x \rightarrow 0} \left| \frac{(f(x + \Delta x) - f(x))/f(x)}{\Delta x/x} \right| = \frac{\kappa_{abs}(f)|x|}{|f(x)|}$$

Rounding Error in Floating Point Operations

▶ Addition and Subtraction

- ▶ *Catastrophic cancellation* occurs when the magnitude of the result is much smaller than the magnitude of both operands
- ▶ Cancellation corresponds to losing significant digits, e.g.

$$3.1423 \times 10^5 - 3.1403 \times 10^5 = 2.0 \times 10^2$$

- ▶ Generally, we can bound the error incurred during addition of two real numbers x, y in floating point (ignoring final rounding, which has relative error ϵ) as

$$\frac{|(x + y) - (fl(x) + fl(y))|}{|x + y|} \leq \frac{\epsilon(|x| + |y|)}{|x + y|}$$

by this we can also observe that the condition number of addition of x, y i.e. $f(x, y) = x + y$, is $\kappa(f(x, y)) = (|x| + |y|)/|x + y|$

- ▶ Consequently, when $x + y = 0$ and $x, y \neq 0$ addition is *ill-posed* (has infinite condition number) unless we restrict the space of possible inputs x, y

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 max distance to center is given by the vector maximizing $\max_{\|\mathbf{x}\|=1} \|\mathbf{A}\mathbf{x}\|_2$.*
 - ▶ *The min distance to center is given by the vector minimizing $\min_{\|\mathbf{x}\|=1} \|\mathbf{A}\mathbf{x}\|_2 = 1/(\max_{\|\mathbf{x}\|=1} \|\mathbf{A}^{-1}\mathbf{x}\|_2)$.*
 - ▶ *Thus, we have that $\kappa(\mathbf{A}) = \|\mathbf{A}\|_2 \|\mathbf{A}^{-1}\|_2$*
- ▶ The matrix condition number bounds the worst-case amplification of error in a matrix-vector product: *Consider $\mathbf{y} + \delta\mathbf{y} = \mathbf{A}(\mathbf{x} + \delta\mathbf{x})$, assume $\|\mathbf{x}\|_2 = 1$*
 - ▶ *In the worst case, $\|\mathbf{y}\|_2$ is minimized, that is $\|\mathbf{y}\|_2 = 1/\|\mathbf{A}^{-1}\|_2$*
 - ▶ *In the worst case, $\|\delta\mathbf{y}\|_2$ is maximized, that is $\|\delta\mathbf{y}\|_2 = \|\mathbf{A}\|_2 \|\delta\mathbf{y}\|_2$*
 - ▶ *So $\|\delta\mathbf{y}\|_2/\|\mathbf{y}\|_2$ is at most $\kappa(\mathbf{A})\|\delta\mathbf{x}\|_2/\|\mathbf{x}\|_2$*

Singular Value Decomposition

- ▶ The singular value decomposition (SVD)

We can express any matrix \mathbf{A} as

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$$

where \mathbf{U} and \mathbf{V} are orthogonal, and $\mathbf{\Sigma}$ is square nonnegative and diagonal,

$$\mathbf{\Sigma} = \begin{bmatrix} \sigma_{max} & & & \\ & \ddots & & \\ & & & \sigma_{min} \end{bmatrix}$$

Any matrix is diagonal when expressed as an operator mapping vectors from a coordinate system given by \mathbf{V} to a coordinate system given by \mathbf{U}^T .

- ▶ Condition number in terms of singular values
 - ▶ *We have that $\|\mathbf{A}\|_2 = \sigma_{max}$ and if \mathbf{A}^{-1} exists, $\|\mathbf{A}^{-1}\|_2 = 1/\sigma_{min}$*
 - ▶ *Consequently, $\kappa(\mathbf{A}) = \sigma_{max}/\sigma_{min}$*

Linear Least Squares

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

Since $m \geq n$, the minimizer generally does not attain a zero residual $\mathbf{Ax} - \mathbf{b}$. We can rewrite the optimization problem constraint via

$$\mathbf{x}^* = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{Ax} - \mathbf{b}\|_2^2 = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \left[(\mathbf{Ax} - \mathbf{b})^T (\mathbf{Ax} - \mathbf{b}) \right]$$

- ▶ 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:

- ▶ *The minimizer satisfies $\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T \mathbf{x}^* \cong \mathbf{b}$ and consequently also satisfies*

$$\mathbf{\Sigma}\mathbf{y}^* \cong \mathbf{d} \quad \text{where } \mathbf{y}^* = \mathbf{V}^T \mathbf{x}^* \text{ and } \mathbf{d} = \mathbf{U}^T \mathbf{b}.$$

- ▶ *The minimizer of the reduced problem is $\mathbf{y}^* = \mathbf{\Sigma}^\dagger \mathbf{d}$, so $y_i = d_i / \sigma_i$ for $i \in \{1, \dots, n\}$ and $y_i = 0$ for $i \in \{n+1, \dots, m\}$.*

Normal Equations

Demo: Normal equations vs Pseudoinverse

Demo: Issues with the normal equations

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

If $\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{b}$ then

$$(\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T)^T \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T \mathbf{x} = (\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T)^T \mathbf{b}$$

$$\mathbf{\Sigma}^T \mathbf{\Sigma} \mathbf{V}^T \mathbf{x} = \mathbf{\Sigma}^T \mathbf{U}^T \mathbf{b}$$

$$\mathbf{V}^T \mathbf{x} = (\mathbf{\Sigma}^T \mathbf{\Sigma})^{-1} \mathbf{\Sigma}^T \mathbf{U}^T \mathbf{b} = \mathbf{\Sigma}^\dagger \mathbf{U}^T \mathbf{b}$$

$$\mathbf{x} = \mathbf{V} \mathbf{\Sigma}^\dagger \mathbf{U}^T \mathbf{b} = \mathbf{x}^*$$

- ▶ However, solving the normal equations is a more ill-conditioned problem than the original least squares algorithm

Generally we have $\kappa(\mathbf{A}^T \mathbf{A}) = \kappa(\mathbf{A})^2$ (the singular values of $\mathbf{A}^T \mathbf{A}$ are the squares of those in \mathbf{A}). Consequently, solving the least squares problem via the normal equations may be unstable because it involves solving a problem that has worse conditioning than the initial least squares problem.

Solving the Normal Equations

- ▶ If \mathbf{A} is full-rank, then $\mathbf{A}^T \mathbf{A}$ is symmetric positive definite (SPD):
 - ▶ *Symmetry is easy to check* $(\mathbf{A}^T \mathbf{A})^T = \mathbf{A}^T \mathbf{A}$.
 - ▶ *\mathbf{A} being full-rank implies $\sigma_{\min} > 0$ and further if $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ we have*

$$\mathbf{A}^T \mathbf{A} = \mathbf{V}^T \mathbf{\Sigma}^2 \mathbf{V}$$

which implies that rows of \mathbf{V} are the eigenvectors of $\mathbf{A}^T \mathbf{A}$ with eigenvalues $\mathbf{\Sigma}^2$ since $\mathbf{A}^T \mathbf{A} \mathbf{V}^T = \mathbf{V}^T \mathbf{\Sigma}^2$.

- ▶ Since $\mathbf{A}^T \mathbf{A}$ is SPD we can use Cholesky factorization, to factorize it and solve linear systems:

$$\mathbf{A}^T \mathbf{A} = \mathbf{L}\mathbf{L}^T$$

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 $A = QR$
 - ▶ Given $A^T A = LL^T$, we can take $R = L^T$ and obtain $Q = AL^{-T}$, since $\underbrace{L^{-1}A^T}_{Q^T} \underbrace{AL^{-T}}_Q = I$ implies that Q has orthonormal columns.
- ▶ 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. A full QR factorization gives $Q \in \mathbb{R}^{m \times m}$ and $R \in \mathbb{R}^{m \times n}$, but since R is upper triangular, the latter $m - n$ columns of Q are only constrained so as to keep Q orthogonal. The **reduced QR** factorization is given by taking the first n columns Q and \hat{Q} the upper-triangular block of R , \hat{R} giving $A = \hat{Q}\hat{R}$.
- ▶ We can solve the normal equations (and consequently the linear least squares problem) via reduced QR as follows

$$A^T A x = A^T b \quad \Rightarrow \quad \hat{R}^T \underbrace{\hat{Q}^T \hat{Q}}_I \hat{R} x = \hat{R}^T \hat{Q}^T b \quad \Rightarrow \quad \hat{R} x = \hat{Q}^T b$$

Eigenvalue Decomposition

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

$$A = XDX^{-1}$$

where X are the right eigenvectors, X^{-1} are the left eigenvectors and D are eigenvalues

$$AX = [Ax_1 \quad \cdots \quad Ax_n] = XD = [d_{11}x_1 \quad \cdots \quad d_{nn}x_n].$$

- ▶ If A is symmetric, its right and left singular vectors are the same, and consequently are its eigenvectors.
- ▶ More generally, any *normal* matrix, $A^H A = A A^H$, has unitary eigenvectors.
- ▶ A and B are *similar*, if there exist Z such that $A = ZBZ^{-1}$
 - ▶ Normal matrices are *unitarily similar* ($Z^{-1} = Z^H$) to diagonal matrices
 - ▶ Symmetric real matrices are *orthogonally similar* ($Z^{-1} = Z^T$) to real diagonal matrices
 - ▶ Hermitian matrices are unitarily similar to real diagonal matrices

Similarity of Matrices

<i>matrix</i>	<i>similarity</i>	<i>reduced form</i>
SPD	orthogonal	real positive diagonal
real symmetric	orthogonal	real tridiagonal real diagonal
Hermitian	unitary	real diagonal
normal	unitary	diagonal
real	orthogonal	real Hessenberg
diagonalizable	invertible	diagonal
arbitrary	unitary invertible	triangular bidiagonal

Rayleigh Quotient

- ▶ For any vector \mathbf{x} that is close to an eigenvector, the *Rayleigh quotient* provides an estimate of the associated eigenvalue of \mathbf{A} :

$$\rho_{\mathbf{A}}(\mathbf{x}) = \frac{\mathbf{x}^H \mathbf{A} \mathbf{x}}{\mathbf{x}^H \mathbf{x}}.$$

- ▶ *If \mathbf{x} is an eigenvector of \mathbf{A} , then $\rho_{\mathbf{A}}(\mathbf{x})$ is the associated eigenvalue.*
- ▶ *Moreover, for $\mathbf{y} = \mathbf{A}\mathbf{x}$, the Rayleigh quotient is the best possible eigenvalue estimate given \mathbf{x} and \mathbf{y} , as it is the solution α to $\mathbf{x}\alpha \cong \mathbf{y}$.*
- ▶ *The normal equations for this scalar-output least squares problem are (assuming \mathbf{A} is real),*

$$\mathbf{x}^T \mathbf{x} \alpha = \mathbf{x}^T \mathbf{y} \quad \Rightarrow \quad \alpha = \frac{\mathbf{x}^T \mathbf{y}}{\mathbf{x}^T \mathbf{x}} = \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{x}}.$$

Introduction to Krylov Subspace Methods

- ▶ *Krylov subspace methods* work with information contained in the $n \times k$ matrix

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

We seek to best use the information from the matrix vector product results (columns of \mathbf{K}_k) to solve eigenvalue problems.

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

Letting $\mathbf{k}_n^{(i)} = \mathbf{A}^{i-1}\mathbf{x}$, we observe that

$$\mathbf{A}\mathbf{K}_n = \begin{bmatrix} \mathbf{A}\mathbf{k}_n^{(1)} & \cdots & \mathbf{A}\mathbf{k}_n^{(n-1)} & \mathbf{A}\mathbf{k}_n^{(n)} \end{bmatrix} = \begin{bmatrix} \mathbf{k}_n^{(2)} & \cdots & \mathbf{k}_n^{(n)} & \mathbf{A}\mathbf{k}_n^{(n)} \end{bmatrix},$$

therefore premultiplying by \mathbf{K}_n^{-1} transforms the first $n - 1$ columns of $\mathbf{A}\mathbf{K}_n$ into the last $n - 1$ columns of \mathbf{I} ,

$$\begin{aligned} \mathbf{K}_n^{-1}\mathbf{A}\mathbf{K}_n &= \begin{bmatrix} \mathbf{K}_n^{-1}\mathbf{k}_n^{(2)} & \cdots & \mathbf{K}_n^{-1}\mathbf{k}_n^{(n)} & \mathbf{K}_n^{-1}\mathbf{A}\mathbf{k}_n^{(n)} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{e}_2 & \cdots & \mathbf{e}_n & \mathbf{K}_n^{-1}\mathbf{A}\mathbf{k}_n^{(n)} \end{bmatrix} \end{aligned}$$

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:
 - ▶ *The approximation obtained from $k - 1$ steps of power iteration starting from \mathbf{x}_0 is given by the Rayleigh-quotient of $\mathbf{y} = \mathbf{A}^k \mathbf{x}_0$.*
 - ▶ *This vector is within the Krylov subspace, $\mathbf{y} \in \mathcal{K}_k(\mathbf{A}, \mathbf{x}_0)$.*
 - ▶ *Consequently, Krylov subspace methods will generally obtain strictly better approximations of the dominant eigenpair than power iteration.*

Krylov Subspace Methods

- ▶ The $k \times k$ matrix $\mathbf{H}_k = \mathbf{Q}_k^T \mathbf{A} \mathbf{Q}_k$ minimizes $\|\mathbf{A} \mathbf{Q}_k - \mathbf{Q}_k \mathbf{H}_k\|_2$:
The minimizer \mathbf{X} for the linear least squares problem $\mathbf{Q}_k \mathbf{X} \cong \mathbf{A} \mathbf{Q}_k$ is (via the normal equations) $\mathbf{X} = \mathbf{Q}_k^T \mathbf{A} \mathbf{Q}_k = \mathbf{H}_k$.
- ▶ \mathbf{H}_k is Hessenberg, because the companion matrix \mathbf{C}_k is Hessenberg:

$$\mathbf{H}_k = \mathbf{Q}_k^T \mathbf{A} \mathbf{Q}_k = \mathbf{R}_k \mathbf{K}_k^{-1} \mathbf{A} \mathbf{K}_k \mathbf{R}_k^{-1} = \mathbf{R}_k \mathbf{C}_k \mathbf{R}_k^{-1}$$

is a product of three matrices: upper-triangular \mathbf{R}_k , upper-Hessenberg \mathbf{C}_k , and upper-triangular \mathbf{R}_k^{-1} , which results in upper-Hessenberg \mathbf{H}_k .

Rayleigh-Ritz Procedure

- ▶ The eigenvalues/eigenvectors of \mathbf{H}_k are the *Ritz values/vectors*:

$$\mathbf{H}_k = \mathbf{X} \mathbf{D} \mathbf{X}^{-1}$$

eigenvalue approximations based on Ritz vectors \mathbf{X} are given by $\mathbf{Q}_k \mathbf{X}$.

- ▶ 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 :

Assuming \mathbf{A} is a symmetric matrix with positive eigenvalues, the largest Ritz value $\lambda_{\max}(\mathbf{H}_k)$ will be the maximum Rayleigh quotient of any vector in $\mathcal{K}_k = \text{span}(\mathbf{Q}_k)$,

$$\max_{\mathbf{x} \in \text{span}(\mathbf{Q}_k)} \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{x}} = \max_{\mathbf{y} \neq 0} \frac{\mathbf{y}^T \mathbf{Q}_k^T \mathbf{A} \mathbf{Q}_k \mathbf{y}}{\mathbf{y}^T \mathbf{y}} = \max_{\mathbf{y} \neq 0} \frac{\mathbf{y}^T \mathbf{H}_k \mathbf{y}}{\mathbf{y}^T \mathbf{y}} = \lambda_{\max}(\mathbf{H}_k),$$

which is the best approximation to $\lambda_{\max}(\mathbf{A}) = \max_{\mathbf{x} \neq 0} \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{x}}$ available in \mathcal{K}_k . The quality of the approximation can also be shown to be optimal for other eigenvalues/eigenvectors.

General Multidimensional Optimization

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

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)$$

such that $f(\mathbf{x}_{k+1}) = \min_{\alpha_k} f(\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k))$, i.e. perform a line search (solve 1D optimization problem) 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

$$\|\mathbf{e}_{k+1}\|_{\mathbf{A}} = \mathbf{e}_{k+1}^T \mathbf{A} \mathbf{e}_{k+1} = \frac{\sigma_{\max}(\mathbf{A}) - \sigma_{\min}(\mathbf{A})}{\sigma_{\max}(\mathbf{A}) + \sigma_{\min}(\mathbf{A})} \|\mathbf{e}_k\|_{\mathbf{A}}$$

- ▶ 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

$$\frac{\sigma_{\max}(\mathbf{A}) - \sigma_{\min}(\mathbf{A})}{\sigma_{\max}(\mathbf{A}) + \sigma_{\min}(\mathbf{A})} = \frac{\kappa(\mathbf{A}) - 1}{\kappa(\mathbf{A}) + 1}.$$

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

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k) + \beta_k (\mathbf{x}_k - \mathbf{x}_{k-1})$$

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

$$\|\mathbf{e}_{k+1}\|_{\mathbf{A}} = \frac{\sqrt{\kappa(\mathbf{A})} - 1}{\sqrt{\kappa(\mathbf{A})} + 1} \|\mathbf{e}_k\|_{\mathbf{A}}$$

Nesterov's gradient optimization method is another instance of an extrapolation method that provides further improved optimality guarantees.

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:

$$(\alpha_k, \beta_k) = \operatorname{argmin}_{\alpha_k, \beta_k} \left[f\left(\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k) + \beta_k (\mathbf{x}_k - \mathbf{x}_{k-1})\right) \right]$$

- ▶ *For SPD quadratic programming problems, conjugate gradient is an optimal first order method, converging in n iterations.*
- ▶ *It implicitly computes Lanczos iteration, searching along A -orthogonal directions at each step.*
- ▶ *Parallel tangents* implementation of the method proceeds as follows
 1. *Perform a step of steepest descent to generate $\hat{\mathbf{x}}_k$ from \mathbf{x}_k .*
 2. *Generate \mathbf{x}_{k+1} by minimizing over the line passing through \mathbf{x}_{k-1} and $\hat{\mathbf{x}}_k$.*

The method is equivalent to CG for a quadratic objective function.

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} :
 - ▶ It constructs Krylov subspace $\mathcal{K}_k(\mathbf{A}, \mathbf{c}) = \text{span}(\mathbf{c}, \mathbf{A}\mathbf{c}, \dots, \mathbf{A}^{k-1}\mathbf{c})$.
 - ▶ At the k th step conjugate gradient yields iterate

$$\mathbf{x}_k = -\|\mathbf{c}\|_2 \mathbf{Q}_k \mathbf{T}_k^{-1} \mathbf{e}_1,$$

where \mathbf{Q}_k is an orthogonal basis for Krylov subspace $\mathcal{K}_k(\mathbf{A}, \mathbf{c})$ and $\mathbf{T}_k = \mathbf{Q}_k^T \mathbf{A} \mathbf{Q}_k$.

- ▶ This choice of \mathbf{x}_k minimizes $f(\mathbf{x})$ since

$$\begin{aligned} \min_{\mathbf{x} \in \mathcal{K}_k(\mathbf{A}, \mathbf{c})} f(\mathbf{x}) &= \min_{\mathbf{y} \in \mathbb{R}^k} f(\mathbf{Q}_k \mathbf{y}) \\ &= \min_{\mathbf{y} \in \mathbb{R}^k} \mathbf{y}^T \mathbf{Q}_k^T \mathbf{A} \mathbf{Q}_k \mathbf{y} + \mathbf{c}^T \mathbf{Q}_k \mathbf{y} \\ &= \min_{\mathbf{y} \in \mathbb{R}^k} \mathbf{y}^T \mathbf{T}_k \mathbf{y} + \|\mathbf{c}\|_2 \mathbf{e}_1^T \mathbf{y} \end{aligned}$$

is minimized by $\mathbf{y} = -\|\mathbf{c}\|_2 \mathbf{T}_k^{-1} \mathbf{e}_1$.

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 :

$$f(\mathbf{x}_k + \mathbf{s}) \approx \hat{f}(\mathbf{s}) = f(\mathbf{x}_k) + \mathbf{s}^T \nabla f(\mathbf{x}_k) + \frac{1}{2} \mathbf{s}^T \mathbf{H}_f(\mathbf{x}_k) \mathbf{s}.$$

The minima of this function can be determined by identifying critical points

$$\mathbf{0} = \nabla \hat{f}(\mathbf{s}) = \nabla f(\mathbf{x}_k) + \mathbf{H}_f(\mathbf{x}_k) \mathbf{s},$$

thus to determine \mathbf{s} we solve the linear system,

$$\mathbf{H}_f(\mathbf{x}_k) \mathbf{s} = -\nabla f(\mathbf{x}_k).$$

Assuming invertibility of the Hessian, we can write the Newton's method iteration as

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \underbrace{\mathbf{H}_f(\mathbf{x}_k)^{-1} \nabla f(\mathbf{x}_k)}_{\mathbf{s}}.$$

Quadratic convergence follows by equivalence to Newton's method for solving nonlinear system of optimality equations $\nabla f(\mathbf{x}) = \mathbf{0}$.

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 $f_{\mathbf{x}}(t_i) \approx y_i$:
For example, consider fitting $f_{[x_1, x_2]}(t) = x_1 \sin(x_2 t)$ so that

$$\begin{bmatrix} f_{[x_1, x_2]}(1.5) \\ f_{[x_1, x_2]}(1.9) \\ f_{[x_1, x_2]}(3.2) \end{bmatrix} \approx \begin{bmatrix} -1.2 \\ 4.5 \\ 7.3 \end{bmatrix}.$$

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

Define residual vector function $\mathbf{r}(\mathbf{x})$ so that $r_i(\mathbf{x}) = y_i - f_{\mathbf{x}}(t_i)$ and minimize

$$\phi(\mathbf{x}) = \frac{1}{2} \|\mathbf{r}(\mathbf{x})\|_2^2 = \frac{1}{2} \mathbf{r}(\mathbf{x})^T \mathbf{r}(\mathbf{x}).$$

Now the gradient is $\nabla \phi(\mathbf{x}) = \mathbf{J}_{\mathbf{r}}^T(\mathbf{x}) \mathbf{r}(\mathbf{x})$ and the Hessian is

$$\mathbf{H}_{\phi}(\mathbf{x}) = \mathbf{J}_{\mathbf{r}}^T(\mathbf{x}) \mathbf{J}_{\mathbf{r}}(\mathbf{x}) + \sum_{i=1}^m r_i(\mathbf{x}) \mathbf{H}_{r_i}(\mathbf{x}).$$

Gauss-Newton Method

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

$$\mathbf{H}_\phi(\mathbf{x}) = \mathbf{J}_r^T(\mathbf{x})\mathbf{J}_r(\mathbf{x}) + \sum_{i=1}^m r_i(\mathbf{x})\mathbf{H}_{r_i}(\mathbf{x}).$$

The second term is small when the residual function $r(\mathbf{x})$ is small, so approximate

$$\mathbf{H}_\phi(\mathbf{x}) \approx \hat{\mathbf{H}}_\phi(\mathbf{x}) = \mathbf{J}_r^T(\mathbf{x})\mathbf{J}_r(\mathbf{x}).$$

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

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \hat{\mathbf{H}}_\phi(\mathbf{x}_k)^{-1} \nabla f(\mathbf{x}_k) = \mathbf{x}_k - (\mathbf{J}_r^T(\mathbf{x}_k)\mathbf{J}_r(\mathbf{x}_k))^{-1} \mathbf{J}_r^T(\mathbf{x}_k) \mathbf{r}(\mathbf{x}_k).$$

Recognizing the normal equations, we interpret the Gauss-Newton method as solving linear least squares problems $\mathbf{J}_r(\mathbf{x}_k)\mathbf{s}_k \cong \mathbf{r}(\mathbf{x}_k)$, $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{s}_k$.

Tensors

- ▶ A *tensor* $\mathcal{T} \in \mathbb{R}^{n_1 \times \dots \times n_d}$ has
 - ▶ *Order* d (i.e. d *modes* / *indices*)
 - ▶ *Dimensions* n_1 -by- \dots -by- n_d
 - ▶ *Elements* $t_{i_1 \dots i_d} = t_{\mathbf{i}}$ where $\mathbf{i} \in \bigotimes_{i=1}^d \{1, \dots, n_i\}$
- ▶ Order d tensors represent d -dimensional arrays
 - ▶ $(d \geq 3)$ -dimensional arrays are prevalent in scientific computing
 - ▶ Regular grids, collections of matrices, multilinear operators
 - ▶ Experimental data, visual/graphic data
 - ▶ Higher-order derivatives and correlation

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

$$\mathbf{v} = \text{vec}(\mathbf{A}) \Rightarrow \mathbf{v} \in \mathbb{R}^{mn}, v_{i+jm} = a_{ij}$$

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

$$\mathbf{v} = \text{vec}(\mathcal{T}) \Rightarrow \mathbf{v} \in \mathbb{R}^{n_1 \dots n_d}, v_{i_1+i_2n_1+i_3n_1n_2+\dots} = t_{i_1i_2i_3\dots}$$

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

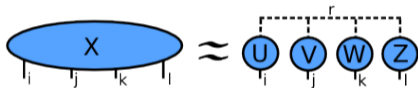
$$\mathbf{U} = o_{n_1 \times \dots \times n_d}(\mathcal{V}) \Rightarrow \mathbf{U} \in \mathbb{R}^{n_1 \times \dots \times n_d}, \quad \text{vec}(\mathbf{U}) = \text{vec}(\mathcal{V})$$

Canonical Polyadic (CP) Decomposition

- ▶ A rank R *CP decomposition* of an $s \times s \times s \times s$ tensor is

$$x_{ijkl} = \sum_{r=1}^R u_{ir} v_{jr} w_{kr} z_{lr}$$

- ▶ We can represent the CP using the following *tensor diagram*:



- ▶ Finding an approximate tensor decomposition corresponds to a nonlinear least squares problem:

$$f(\mathbf{U}, \mathbf{V}, \mathbf{W}, \mathbf{Z}) = \sum_{i,j,k,l} \left(x_{ijkl} - \underbrace{\sum_{r=1}^R u_{ir} v_{jr} w_{kr} z_{lr}}_{\hat{x}_{ijkl}} \right)^2$$

which is the squared Frobenius norm error $\|\mathbf{x} - \hat{\mathbf{x}}\|_F^2$.