

CS 598 EVS: Tensor Computations

Matrix Computations Background

Edgar Solomonik

University of Illinois at Urbana-Champaign

Conditioning

- ▶ **Absolute Condition Number:**

The *absolute condition number* is a property of the problem, which measures its sensitivity to perturbations in input

$$\kappa_{abs}(f) = \lim_{\text{size of input perturbation} \rightarrow 0} \max_{\text{inputs}} \max_{\text{perturbations in input}} \left| \frac{\text{perturbation in output}}{\text{perturbation in input}} \right|$$

For problem f at 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 $\kappa_{abs} = \max_{x \in \mathcal{X}} \left| \frac{df}{dx}(x) \right|$

- ▶ **(Relative) Condition Number:**

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

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

Posedness and Conditioning

- ▶ **What is the condition number of an ill-posed problem?**
 - ▶ *If the condition number is bounded and the solution is unique, the problem is **well-posed***
 - ▶ *An **ill-posed** problem f either has no unique solution or has a (relative) condition number of $\kappa(f) = \infty$*
 - ▶ *This condition implies that the solutions to problem f are continuous and differentiable in the given space of possible inputs to f*
 - ▶ *Sometimes well-posedness is defined to only require continuity*
 - ▶ *Generally, $\kappa(f)$ can be thought of as the reciprocal of the **distance** (in an appropriate geometric embedding of problem configurations) from f **to the nearest 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 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 $\|\delta\mathbf{y}\|_2 = \|\mathbf{A}\|_2 \|\delta\mathbf{x}\|_2$*
 - ▶ *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 A as

$$A = U\Sigma V^T$$

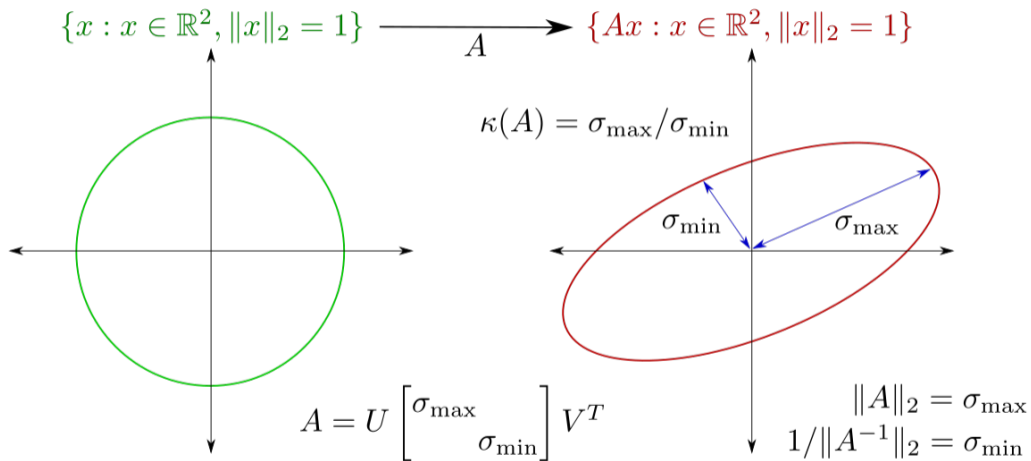
where U and V are orthogonal, and Σ is square nonnegative and diagonal,

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

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

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

Visualization of Matrix Conditioning



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

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

Computing the QR Factorization

- ▶ The Cholesky-QR algorithm uses the normal equations to obtain the QR factorization
 - ▶ Compute $A^T A = LL^T$, take $R = L^T$, and solve for Q triangular linear systems $LQ^T = A^T$
 - ▶ If A is $m \times n$, forming $A^T A$ has cost mn^2 , computing Cholesky factorization has cost $(2/3)n^3$, and solving the triangular systems (if Q is needed) costs mn^2 , yielding total cost $2mn^2 + (2/3)n^3$
 - ▶ However, this method is unstable since $A^T A$ is ill-conditioned. This is addressible by iterating on the computed (nearly-orthogonal) Q factor (CholeskyQR2).
- ▶ Orthogonalization-based methods are most efficient and stable for QR factorization of dense matrices
 - ▶ Apply a sequene of orthogonal transformations Q_1, \dots, Q_k to reduce A to triangulr form $(Q_1 \cdots Q_k)^T A = R$
 - ▶ Householder QR uses rank-1 perturbations of the identity matrix (reflectors) $Q_i = I - 2u_i u_i^T$ to zero-out each sub-column of A
 - ▶ Givens rotations zero-out a single entry at a time
 - ▶ Both approaches have cost $O(mn^2)$ with similar constant to Cholesky-QR

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 $Q_k R_k = K_k$, 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 upper-Hessenberg, because the companion matrix \mathbf{C}_n is upper-Hessenberg:

Note that \mathbf{H}_k is the leading k -by- k minor of \mathbf{H}_n and

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

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

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.

Low Rank Matrix Approximation

- ▶ Given a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ seek rank $r < m, n$ approximation
 - ▶ Given by matrices $\mathbf{U} \in \mathbb{R}^{m \times r}$ and $\mathbf{V} \in \mathbb{R}^{n \times r}$ so

$$\mathbf{A} \approx \mathbf{UV}^T$$

- ▶ Reduces memory footprint and cost of applying \mathbf{A} from mn to $mr + nr$
 - ▶ This factorization is nonunique, $\mathbf{UV}^T = (\mathbf{UM})(\mathbf{VM}^{-T})^T$
- ▶ Eckart-Young (optimal low-rank approximation by SVD) theorem
 - ▶ Truncated SVD approximates \mathbf{A} as

$$\mathbf{A} \approx \tilde{\mathbf{A}} = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T$$

where $\sigma_1, \dots, \sigma_r$ are the largest r singular values, while \mathbf{u}_i and \mathbf{v}_i are the associated left and right singular vectors

- ▶ Eckart-Young theorem demonstrates that the truncated SVD minimizes

$$\underbrace{\|\mathbf{A} - \tilde{\mathbf{A}}\|_2}_{\sigma_{r+1}} \quad \text{and} \quad \underbrace{\|\mathbf{A} - \tilde{\mathbf{A}}\|_F}_{\sum_{i=r+1}^{\min(m,n)} \sigma_i}$$

Rank Revealing Matrix Factorizations

- ▶ Computing the SVD
 - ▶ *Can compute full SVD with $O(mn \min(m, n))$ cost via bidiagonalization*
 - ▶ *unconditionally stable and accurate*
 - ▶ *inefficient for low r or if A is sparse*
 - ▶ *Given any low-rank approximation composed of U and V , compute QR of each and SVD of product of R factors to obtain SVD with total cost $O((m+n)r^2)$*
- ▶ QR with column pivoting
 - ▶ *By selecting columns of largest norm in the trailing matrix during QR factorization, we obtain a pivoted factorization with permutation matrix P*

$$AP = QR$$

- ▶ *Truncating this factorization can be done after applying r Householder reflectors (or another QR algorithm on r columns), with cost $O((m+n)r)$*
- ▶ *Approximation is somewhat suboptimal in theory, but in practice almost always as accurate as truncated SVD*

Orthogonal Iteration

- ▶ For sparse matrices, QR factorization creates fill, so must revert to iterative methods
 - ▶ *Can find SVD of A by implicit products with $A^T A$ or AA^T , since left singular vectors of A are eigenvectors of $A^T A$*
 - ▶ *Krylov subspace methods are effective for computing the largest eigenvector*
 - ▶ *Deflation, e.g., $A \rightarrow (A - \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T)$ can be used to compute other eigenvectors*
- ▶ Orthogonal iteration interleaves deflation and power iteration
 - ▶ *Given starting eigenvector guess $U^{(0)} \in \mathbb{R}^{n \times r}$, compute $V^{(i+1)} = AU^{(i)}$ and obtain $U^{(i+1)}$ as the Q factor of the QR of $V^{(i+1)}$*
 - ▶ *Converges to r largest eigenvectors, for SVD can compute $V^{(i+1)} = A^T(AU^{(i)})$ at each iteration*
 - ▶ *QR factorization serves to orthogonalize each column w.r.t. eigenvectors being converged to by previous columns*

Randomized SVD

- ▶ Orthogonal iteration for SVD can also be viewed as a randomized algorithm
 - ▶ Suppose that we have an exact low-rank factorization $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ with $\mathbf{\Sigma} \in \mathbb{R}^{r \times r}$
 - ▶ If $\mathbf{U}^{(0)}$ is a random orthogonal matrix, so is $\mathbf{V}^T \mathbf{U}^{(0)}$
 - ▶ Consequently, $\mathbf{A}\mathbf{U}^{(0)}$ is a set of r random linear combinations of columns of $\mathbf{U}\mathbf{\Sigma}$
 - ▶ Further, $\mathbf{U} = \mathbf{U}^{(1)}\mathbf{U}^{(1)T}\mathbf{U}$ since

$$\text{span}(\mathbf{U}^{(1)}) = \text{span}(\mathbf{V}^{(1)}) = \text{span}(\mathbf{U}),$$

the latter equality holds with probability 1

- ▶ Consequently, we can compute SVD of $\mathbf{U}^{(1)T}\mathbf{A}$ (with cost $O(nr^2)$) and recover \mathbf{U} by premultiplying the computed left singular vectors by $\mathbf{U}^{(1)}$
- ▶ When \mathbf{A} is not exactly low-rank, span of leading singular vectors can be captured by oversampling (e.g., selecting each $\mathbf{U}^{(i)}$ to have $r + 10$ columns)
- ▶ Initial guess $\mathbf{U}^{(0)}$ need not be orthogonal (Gaussian random performs well, structured pseudo-random enables $O(mn \log n)$ complexity for one-shot randomized SVD), but better accuracy is obtained with orthogonality

Generalized Nyström Algorithm

- ▶ The generalized Nyström algorithm provides an efficient way of computing a sketched low-rank factorization

- ▶ *the rank k factorization of a matrix A is obtained via*

$$\hat{A}_k = \mathbf{A}\mathbf{S}_1^T (\mathbf{S}_2\mathbf{A}\mathbf{S}_1^T)^\dagger \mathbf{S}_2\mathbf{A}$$

- ▶ *where S_1 and S_2 are sketch matrices*
 - ▶ *no need to apply A to a general matrix, can define S_1 and S_2 as sparse or structured (e.g., diagonal matrix times Fourier transform)*
 - ▶ *Sketch matrices can be constructed to be Gaussian random, with awareness of A (e.g., via leverage score sampling) or to be sparse (CountSketch)*

Multidimensional Optimization

- ▶ Minimize $f(\mathbf{x})$
 - ▶ *In the context of constrained optimization, also have equality and or inequality constraints, e.g., $\mathbf{Ax} = \mathbf{b}$ or $\mathbf{x} > 0$*
 - ▶ *Unconstrained local optimality holds if $\nabla f(\mathbf{x}^*) = 0$ and $H_f(\mathbf{x}^*)$ is positive semi-definite*
 - ▶ *Reduces to solving nonlinear equations via optimality condition*
 - ▶ *Unconstrained local optimality conditions are looser, need the gradient to be zero or positive in all unconstrained directions at \mathbf{x}^**
 - ▶ *The condition $\nabla f(\mathbf{x}^*) = 0$ implies poor conditioning, perturbations that change the function value in the k th digit can change the solution in the $(k/2)$ th digit*
- ▶ Quadratic optimization $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \mathbf{Ax} - \mathbf{b}^T \mathbf{x}$
 - ▶ *Quadratic optimization problems can provide local approximations to general nonlinear optimization problems via Newton's method (where \mathbf{A} is the Hessian and \mathbf{b}^T is the gradient)*
 - ▶ *Equivalent to solving linear system $\mathbf{Ax} = \mathbf{b}$ by optimality condition*
 - ▶ *Accordingly, conditioning relative to perturbation in \mathbf{b} is $\kappa(\mathbf{A})$*

Basic Multidimensional Optimization Methods

- ▶ 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 $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{x}$ where \mathbf{A} is symmetric positive definite, the error $\mathbf{e}_k = \mathbf{x}_k - \mathbf{x}^*$ satisfies

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

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

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

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

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

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

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

CG and Krylov Optimization

The solution at the k th step, $\mathbf{y}_k = \|\mathbf{b}\|_2 \mathbf{T}_k^{-1} \mathbf{e}_1$ 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

- ▶ *The Lanczos method constructs \mathbf{T}_{k+1} from \mathbf{T}_k using a matrix-vector product with \mathbf{A}*
- ▶ *The change, $\mathbf{T}_{k+1} - \begin{bmatrix} \mathbf{T}_k & \\ & \mathbf{T}_{k+1}(k+1, k+1) \end{bmatrix}$ is of rank 2*
- ▶ *Consequently, the Sherman-Morrison-Woodbury formula (or an updated factorization), which is for general \mathbf{M} ,*

$$(\mathbf{M} - \mathbf{u}\mathbf{v}^T)^{-1} = \mathbf{M}^{-1} + \frac{\mathbf{M}^{-1}\mathbf{u}\mathbf{v}^T\mathbf{M}^{-1}}{1 - \mathbf{v}^T\mathbf{M}^{-1}\mathbf{u}}$$

may be used to apply \mathbf{T}_{k+1}^{-1} with $O(k)$ cost

- ▶ *CG does this implicitly at each step*
- ▶ *Other Krylov iterative methods are available for solution to general (non-SPD) linear systems, such as the generalized minimum residual method (GMRES) and bi-conjugate gradient, which construct a basis for $\mathbf{A}\mathbf{A}^T$ and $\mathbf{A}^T\mathbf{A}$*

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

- ▶ *need not form $\mathbf{M}^{-1}\mathbf{A}$ but only compute matrix-vector products $\mathbf{M}^{-1}(\mathbf{Ax})$*
 - ▶ *want $\mathbf{M}^{-1}\mathbf{x}$ to be easy to compute (easier than $\mathbf{A}^{-1}\mathbf{x}$)*
 - ▶ *so generally one extracts some $\mathbf{M} \approx \mathbf{A}$ that is easy to solve linear systems with*
 - ▶ *however, $\mathbf{M} \approx \mathbf{A}$ may be insufficient/unnecessary, primary goal is to improve conditioning to accelerate iterative methods, i.e., want $\kappa(\mathbf{M}^{-1}\mathbf{A}) \ll \kappa(\mathbf{A})$*
- ▶ Common preconditioners select parts of \mathbf{A} or perform inexact factorization
 - ▶ *(block-)Jacobi preconditioner takes \mathbf{M} to be (block-)diagonal of \mathbf{A}*
 - ▶ *incomplete LU (ILU) preconditioners compute $\mathbf{M} = \mathbf{LU} \approx \mathbf{A}$ (+pivoting)*
 - ▶ *ILU variants constraint sparsity of \mathbf{L} and \mathbf{U} factors during factorization to be the same or not much more than that of \mathbf{A}*
 - ▶ *good problem-specific preconditioners are often available in practice and theory, applying also to problems beyond linear systems (eigenvalue problems, optimization, approximate graph algorithms)*

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
 - ▶ *To prove this, we can analyze the ‘minimizing’ polynomials in the Krylov subspace in terms of the (real and positive) eigenvalues of \mathbf{A}*
 - ▶ *The approximate solution \mathbf{x}_k obtained by CG after $k - 1$ iterations is given by minimizing $\mathbf{z} \in \mathcal{K}_k(\mathbf{A}, \mathbf{b})$, which means $\mathbf{z} = \rho_{k-1}(\mathbf{A})\mathbf{b} = \rho_{k-1}(\mathbf{A})\mathbf{A}\mathbf{x}$ for some polynomial ρ_{k-1} of degree $k - 1$*
 - ▶ *Now, we observe that minimizing the objective $f(\mathbf{z})$ is equivalent to minimizing*

$$\|\mathbf{b} - \mathbf{A}\mathbf{z}\|_{\mathbf{A}^{-1}}^2 = \phi(\mathbf{z}) = (\mathbf{b} - \mathbf{A}\mathbf{z})^T \mathbf{A}^{-1}(\mathbf{b} - \mathbf{A}\mathbf{z}) = (\mathbf{x} - \mathbf{z})^T \mathbf{A}(\mathbf{x} - \mathbf{z})$$

- ▶ *Observe that*

$$\phi(\mathbf{z}) = (\mathbf{x} - \mathbf{z})^T \mathbf{A}(\mathbf{x} - \mathbf{z}) = \underbrace{\mathbf{z}^T \mathbf{A}\mathbf{z} - 2\mathbf{z}^T \mathbf{b}}_{2f(\mathbf{z})} - \underbrace{\mathbf{x}^T \mathbf{b}}_{\text{constant}}$$

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

Conjugate Gradient Convergence Analysis (II)

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

$$\phi(z) = \left((\mathbf{I} - \rho_{k-1}(\mathbf{A})\mathbf{A})x \right)^T \mathbf{A} \left((\mathbf{I} - \rho_{k-1}(\mathbf{A})\mathbf{A})x \right) = x^T \mathbf{q}_k(\mathbf{A})\mathbf{A}\mathbf{q}_k(\mathbf{A})x$$

where $\mathcal{Q}_k \ni q_k(\xi) = 1 - \rho_{k-1}(\xi) \cdot \xi$ can be any degree k polynomial with $q_k(0) = 1$ (or in matrix form, $\mathbf{q}_k(\mathbf{S}) = \mathbf{I} - \rho_{k-1}(\mathbf{S})\mathbf{S}$ with $\mathbf{q}_k(\mathbf{O}) = \mathbf{I}$), so

$$\phi(\mathbf{x}_k) = \min_{z \in \mathcal{K}_k(\mathbf{A}, \mathbf{b})} \phi(z) = \min_{q_k \in \mathcal{Q}_k} x^T \mathbf{q}_k(\mathbf{A})\mathbf{A}\mathbf{q}_k(\mathbf{A})x$$

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

$$\begin{aligned} \phi(z) &= x^T \mathbf{Q}\mathbf{q}_k(\mathbf{\Lambda})\mathbf{\Lambda}\mathbf{q}_k(\mathbf{\Lambda})\mathbf{Q}^T x \\ &\leq \max_{\lambda_i \in \lambda(\mathbf{A})} (q_k(\lambda_i)^2) \underbrace{x^T \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T x}_{\phi(\mathbf{x}_0)} \end{aligned}$$

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

$$\phi(\mathbf{x}_k) = \min_{q_k \in \mathcal{Q}_k} \phi(\mathbf{z}) \leq \min_{q_k \in \mathcal{Q}_k} \max_{\lambda_i \in \lambda(\mathbf{A})} (q_k(\lambda_i)^2) \phi(\mathbf{x}_0)$$

consequently, the residual norm $\|\mathbf{r}_k\|_{\mathbf{A}^{-1}} = \sqrt{\phi(\mathbf{x}_k)}$ decreases as

$$\frac{\|\mathbf{r}_k\|_{\mathbf{A}^{-1}}}{\|\mathbf{r}_0\|_{\mathbf{A}^{-1}}} \leq \min_{q_k \in \mathcal{Q}_k} \max_{\lambda_i \in \lambda(\mathbf{A})} |q_k(\lambda_i)|$$

- ▶ 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 $q_m(0) = 1$)
 - ▶ Specifically, we select q_m to be zero at each distinct eigenvalue $\lambda_1, \dots, \lambda_m$ of \mathbf{A}

$$q_m(\xi) = \frac{\prod_{j=1}^m (\lambda_j - \xi)}{\prod_{j=1}^m \lambda_j}$$

while also satisfying $q_m(0) = 1$

- ▶ This polynomial implies that $\|\mathbf{r}_m\| = \phi(\mathbf{x}_m) = 0$ since $\max_{\lambda_i \in \lambda(\mathbf{A})} q_m(\lambda_i)^2 = 0$

Round-off Error in Conjugate Gradient

- ▶ CG provides strong convergence guarantees for SPD matrices in exact arithmetic
 - ▶ *Classically, CG was viewed as a direct method, since its guaranteed to convergence in n iterations*
 - ▶ *In practice, round-off error prevents CG from achieving this for realistic matrices, so CG was actually abandoned for a while due to being viewed as unstable*
 - ▶ *Later, it was realized that CG is highly competitive as an iterative method*
- ▶ Due to round-off CG may stagnate / have plateaus in convergence
 - ▶ *A formal analysis of round-off error² reveals that CG with round-off is equivalent to exact CG on a matrix of larger dimension, whose eigenvalues are clustered around those of A*
 - ▶ *Using this view, CG convergence plateaus may be explained by the polynomial q_k developing more and more zeros near the same eigenvalue of A*

²A. Greenbaum and Z. Strakos, SIMAX 1992

Graph and Matrix Duality

- ▶ graphs have have a natural correspondence with sparse matrices
 - ▶ *consider an unweighted undirected graph $G = (V, E)$ with n vertices and m edges*
 - ▶ *the adjacency matrix A of G has n rows/columns for each edge, and $a_{ij} = 1$ if $(i, j) \in E$*
 - ▶ *A is symmetric because G is undirected, weighted and directed graphs can be expressed similarly with A*
- ▶ matrix-based representations of graphs can be used to devise algorithms
 - ▶ *combinatorial algorithms (e.g., breadth-first search or bellman-ford for shortest paths) may be expressed by linear algebra operations on a different semiring*
 - ▶ *for example, for shortest paths, the $(\min, +)$ semiring is used in place of the standard $(+, \times)$*
 - ▶ *writing matrix operations on this semiring as (\otimes, \oplus) , the distance matrix is*

$$I \oplus A \oplus (A \otimes A) \oplus \dots \oplus A^n$$

since $d_{ij} = \min_{k \in \{1, \dots, n\}} \min_{u_1, \dots, u_k} a_{iu_1} + a_{u_1 u_2} + \dots + a_{u_{k-1} u_k}$

- ▶ *approximations to graph problems may also be obtained via numerical optimization*

Graph Partitioning from Eigenvectors

- ▶ The Laplacian matrix provides a model of interactions on a graph that is useful in many contexts
 - ▶ *the Laplacian matrix of an unweighted graph is $D - A$ where D is a diagonal matrix containing vertex degrees and A is the adjacency matrix*
 - ▶ *common 2D/3D grid discretization of numerical partial differential equations yield a Laplacian matrix*
- ▶ The second-smallest-eigenvalue eigenvector of the Laplacian (the Fiedler vector), gives a good partitioning of the graph
 - ▶ *One of the eigenvectors of the Laplacian has eigenvalue 0 and is simply the 1s vector*
 - ▶ *If the graph is disconnected, the null-space has dimension at least 2*

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 $\mathbf{J}_r(\mathbf{x}_k)\mathbf{s}_k \cong \mathbf{r}(\mathbf{x}_k)$, $\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{s}_k$.

Constrained Optimization Problems

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

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

*When f is quadratic, while h, g is linear, this is a *quadratic optimization problem*.*

- ▶ Generally, we will seek to reduce constrained optimization problems to a series of simpler optimization problems:
 - ▶ *sequential quadratic programming*: solve a series of constrained quadratic optimization problems
 - ▶ *interior point methods*: solve a series of more complicated (more ill-conditioned) unconstrained optimization problems

Lagrangian Duality

- ▶ The Lagrangian function with constraints $\mathbf{g}(\mathbf{x}) = \mathbf{0}$ and $\mathbf{h}(\mathbf{x}) \leq \mathbf{0}$ is

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

The constrained minima of $f(\mathbf{x})$ must be saddle points of the Lagrangian function

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

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

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

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

when $\lambda_i^ = 0$, we say the i th constraint is inactive at the minimum point.*

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, $f(\mathbf{x}^*) = q(\boldsymbol{\lambda}^*) = \mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*)$
 - ▶ *The complementarity slackness condition $\max(\boldsymbol{\lambda}^*, \mathbf{h}(\mathbf{x})) = \mathbf{0}$ follows since*

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

so, since \mathbf{x}^ is feasible, we have $\mathbf{h}(\mathbf{x}^*) \geq \mathbf{0}$ and consequently $\lambda_i^* h_i(\mathbf{x}^*) = 0$*

- ▶ *Complementarity slackness must be satisfied along with other KKT conditions by any optimal primal-dual solution if f is differentiable and strong duality holds*
- ▶ *If f is convex, then strong duality holds and further the KKT conditions are not only necessary but sufficient*

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 $\mathcal{L}_f(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \boldsymbol{\lambda}^T \mathbf{g}(\mathbf{x})$:

$$q(\mathbf{x}_k + \mathbf{s}, \boldsymbol{\lambda}_k + \boldsymbol{\delta}) = \mathcal{L}_f(\mathbf{x}_k, \boldsymbol{\lambda}_k) + \mathbf{s}^T (\nabla f(\mathbf{x}_k) + \mathbf{J}_g^T(\mathbf{x}_k) \boldsymbol{\lambda}_k) + \frac{1}{2} \mathbf{s}^T \mathbf{B}(\mathbf{x}_k, \boldsymbol{\lambda}_k) \mathbf{s} \\ + \boldsymbol{\delta}^T (\mathbf{J}_g(\mathbf{x}_k) \mathbf{s} + \mathbf{g}(\mathbf{x}_k))$$

where $\mathbf{B}(\mathbf{x}, \boldsymbol{\lambda}) = \mathbf{H}_f(\mathbf{x}) + \sum_{i=1}^m \lambda_i \mathbf{H}_{g_i}(\mathbf{x})$

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

The above unconstrained quadratic program corresponds to the Lagrangian form of the constrained quadratic program

$$\max_{\mathbf{s}} \mathbf{s}^T (\nabla f(\mathbf{x}_k) + \mathbf{J}_g^T(\mathbf{x}_k) \boldsymbol{\lambda}_k) + \frac{1}{2} \mathbf{s}^T \mathbf{B}(\mathbf{x}_k, \boldsymbol{\lambda}_k) \mathbf{s}$$

with constraint $\mathbf{J}_g(\mathbf{x}_k) \mathbf{s} = -\mathbf{g}(\mathbf{x}_k)$.

Interior Point Methods

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

Inverse barrier function:

$$\phi_{\mu}(\mathbf{x}) = f(\mathbf{x}) - \mu \sum_{i=1}^m \frac{1}{h_i(\mathbf{x})}$$

Logarithmic barrier function:

$$\phi_{\mu}(\mathbf{x}) = f(\mathbf{x}) - \mu \sum_{i=1}^m \log(-h_i(\mathbf{x}))$$

in theory with sufficiently small steps we have $\mathbf{x}_{\mu}^ \rightarrow \mathbf{x}^*$ as $\mu \rightarrow 0$*

- ▶ Interior point methods additionally incorporate Lagrangian optimization
 - ▶ *can be combined with SQP or alternating minimization*
 - ▶ *slack variables with nonnegativity constraints reduce general inequality constraints to nonnegativity and equality constraints*
 - ▶ *optimality conditions for augmented Lagrangian conditions yield linear system*
 - ▶ *conditioning of interior point linear systems suffers as μ decreases*

Karush-Kuhn-Tucker (KKT) conditions

Consider the linear-constrained Quadratic program (QP):

$$\begin{aligned} \min_{x \in \mathbb{R}^n} \quad & \frac{1}{2} x^T H x + x^T c \\ \text{s.t.} \quad & Ax = b, Cx \geq d \end{aligned}$$

Its Lagrangian function may be used to derive an interior point method

$$L(x, \lambda, \nu) = \frac{1}{2} x^T H x + x^T c - \lambda^T (Ax - b) - \nu^T (Cx - d)$$

The first-order optimality (KKT) conditions are

$$\begin{aligned} \nabla_x L(x, \lambda, \nu) &= 0 \\ Ax - b &= 0 \\ Cx - d &\geq 0 \\ \nu^T (Cx - d) &= 0 \\ \nu &\geq 0 \end{aligned}$$

Primal-dual Interior Point Method (IPM)

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

$$Hx + c - A^T \lambda - C^T \nu = 0$$

$$Ax - b = 0$$

$$Cx - d - s = 0$$

$$SVe = \sigma \mu e$$

$$s, \nu > 0$$

where

$$V = \text{diag}(\nu_1, \dots, \nu_{m_2}), \quad S = \text{diag}(s_1, \dots, s_{m_2}), \quad e = [1, \dots, 1]^T \in \mathbb{R}^{m_2}$$

$$\mu = \frac{s^T \nu}{m_2}, \quad \sigma \in [0, 1]$$

Interior Point Method (IPM): KKT system

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

$$\begin{bmatrix} -H & A^T & C^T \\ A & 0 & 0 \\ C & 0 & D^{(k)} \end{bmatrix} \begin{pmatrix} \Delta x^{(k)} \\ \Delta \lambda^{(k)} \\ \Delta \nu^{(k)} \end{pmatrix} = - \begin{pmatrix} r_g^{(k)} \\ r_e^{(k)} \\ r_a^{(k)} \end{pmatrix}$$

where $D^{(k)} = (V^{(k)})^{-1} S^{(k)}$ is diagonal and changing with iteration k .

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

- ▶ the values of $D^{(k)} = (V^{(k)})^{-1} S^{(k)}$ vary greatly in magnitude
- ▶ the values of $S^{(k)}$ go to zero if inequality constraint is active at the local minima
- ▶ at the same time, $\nu_i^{(k)} s_i^{(k)}$ multiply to a fixed value that scales with σ