

CS 598 EVS: Tensor Computations

Matrix Computations Background

Edgar Solomonik

University of Illinois at Urbana-Champaign

Conditioning

- ▶ **Absolute Condition Number:**

- ▶ **(Relative) Condition Number:**

Posedness and Conditioning

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

Matrix Condition Number

- ▶ The matrix condition number $\kappa(\mathbf{A})$ is the ratio between the max and min distance from the surface to the center of the unit ball (norm-1 vectors) transformed by \mathbf{A} :

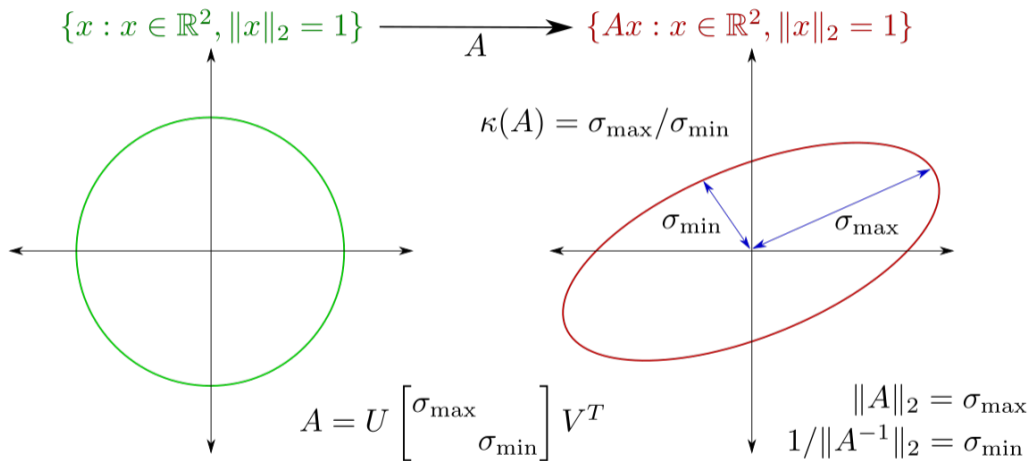
- ▶ The matrix condition number bounds the worst-case amplification of error in a matrix-vector product:

Singular Value Decomposition

- ▶ The singular value decomposition (SVD)

- ▶ Condition number in terms of singular values

Visualization of Matrix Conditioning



Linear Least Squares

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

- ▶ Given the SVD $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ we have $\mathbf{x}^* = \underbrace{\mathbf{V}\mathbf{\Sigma}^\dagger\mathbf{U}^T}_{\mathbf{A}^\dagger} \mathbf{b}$, where $\mathbf{\Sigma}^\dagger$ contains the reciprocal of all nonzeros in $\mathbf{\Sigma}$, and more generally \dagger denotes pseudoinverse:

Normal Equations

Demo: Normal equations vs Pseudoinverse

Demo: Issues with the normal equations

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

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

Solving the Normal Equations

- ▶ If \mathbf{A} is full-rank, then $\mathbf{A}^T \mathbf{A}$ is symmetric positive definite (SPD):

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

QR Factorization

- ▶ If A is full-rank there exists an orthogonal matrix Q and a unique upper-triangular matrix R with a positive diagonal such that $A = QR$

- ▶ A reduced QR factorization (unique part of general QR) is defined so that $Q \in \mathbb{R}^{m \times n}$ has orthonormal columns and R is square and upper-triangular

- ▶ We can solve the normal equations (and consequently the linear least squares problem) via reduced QR as follows

Computing the QR Factorization

- ▶ The Cholesky-QR algorithm uses the normal equations to obtain the QR factorization

- ▶ Orthogonalization-based methods are most efficient and stable for QR factorization of dense matrices

Eigenvalue Decomposition

- ▶ If a matrix A is diagonalizable, it has an *eigenvalue decomposition*
- ▶ A and B are *similar*, if there exist Z such that $A = ZBZ^{-1}$

Similarity of Matrices

<i>matrix</i>	<i>similarity</i>	<i>reduced form</i>
SPD		
real symmetric		
Hermitian		
normal		
real		
diagonalizable		
arbitrary		

Rayleigh Quotient

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

Introduction to Krylov Subspace Methods

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

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

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

Krylov Subspaces

- ▶ Given $\mathbf{Q}_k \mathbf{R}_k = \mathbf{K}_k$, we obtain an orthonormal basis for the Krylov subspace,

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

where p is any polynomial of degree less than k .

- ▶ The Krylov subspace includes the $k - 1$ approximate dominant eigenvectors generated by $k - 1$ steps of power iteration:

Krylov Subspace Methods

- ▶ The $k \times k$ matrix $\mathbf{H}_k = \mathbf{Q}_k^T \mathbf{A} \mathbf{Q}_k$ minimizes $\|\mathbf{A} \mathbf{Q}_k - \mathbf{Q}_k \mathbf{H}_k\|_2$:

- ▶ \mathbf{H}_k is upper-Hessenberg, because the companion matrix \mathbf{C}_n is upper-Hessenberg:

Rayleigh-Ritz Procedure

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

- ▶ The Ritz vectors and values are the *ideal approximations* of the actual eigenvalues and eigenvectors based on only \mathbf{H}_k and \mathbf{Q}_k :

Low Rank Matrix Approximation

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

- ▶ Eckart-Young (optimal low-rank approximation by SVD) theorem

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:

Conjugate Gradient Method

- ▶ The *conjugate gradient method* is capable of making the optimal (for a quadratic objective) choice of α_k and β_k at each iteration of an extrapolation method:

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

Krylov Optimization

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

CG and Krylov Optimization

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

CG for QPs / lin-sys

k_n

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

$$\nabla f(x) = 0 \Rightarrow \underline{A} x = b$$

x_k

$$r = Ax - b$$

$$\text{span} \{ \underbrace{b, Ab, A^2 b, \dots, A^{k-1} b}_{\text{Krylov subspace}} \}$$

Preconditioning

- ▶ Convergence of iterative methods for $Ax = b$ depends on $\kappa(A)$, the goal of a preconditioner M is to obtain x by solving

$$\underline{M^{-1}Ax} = \underline{M^{-1}b}$$

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

$\|e_{k+1}\|_A = \|e_k\|_A \frac{\sqrt{\kappa(A)-1}}{\sqrt{\kappa(A)+1}}$ goal: reduce # of CG iterations
to achieve accuracy

- ▶ Common preconditioners select parts of A or perform inexact factorization

$$A = \overset{M}{\left[\begin{array}{c|c} \square & \\ \hline & \square \end{array} \right]} + \left[\begin{array}{c|c} \square & \\ \hline & \square \end{array} \right]$$

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

$$\underline{\underline{AA}} = XDX^{-1}XDX^{-1} = XD^2X^{-1}$$

$$A^k = XD^kX^{-1}$$

$$p(A) = I + 2A - A^4 = X \underline{\underline{p(D)}} X^{-1}$$

$$p(x) = 1 + 2x - x^4$$

$$\kappa_{k-1}(A, b) = \text{span} \{ b, Ab, \dots, A^{k-1}b \} = \frac{p_{k-1}(A)b}{p_{k-1} \in \prod_{k-1}}$$

Conjugate Gradient Convergence Analysis

$$f(z) = \frac{1}{2} z^T A z - z^T b$$

- ▶ In previous discussion, we assumed K_n is invertible, which may not be the case if A has $m < n$ distinct eigenvalues, however, in exact arithmetic CG converges in $m - 1$ iterations¹

approximate solution x_k

$$\lambda(A) = \begin{matrix} (2, 2) \\ (1, 1) \\ (.001, .001) \end{matrix}$$

$$z \in \mathcal{K}_k(A, b), \quad z = \beta_{\text{opt}}(A) \underbrace{b}_{Ax}$$

$$\|b - Az\|_{A^{-1}}^2 = \underbrace{\varphi(z)} = \underbrace{(b - Az)^T}_{Ax} A^{-1} \underbrace{(b - Az)}_{Ax}$$

$$= \underbrace{(x - z)^T}_{\text{constant}} A (x - z)$$

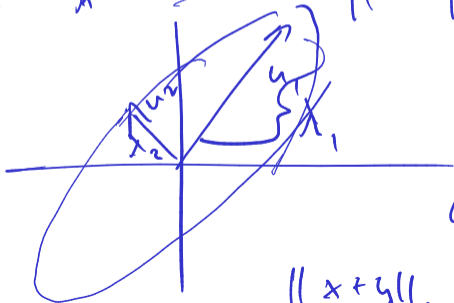
$$= \underbrace{x^T A x}_b - 2 \underbrace{x^T A z}_b + z^T A z - 2 b^T z + z^T A z$$

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

$$\|x\|_A^2 = x^T A x$$

↑
SPD

$$\|x\|_A^2 = x^T \sqrt{A} \sqrt{A} x$$

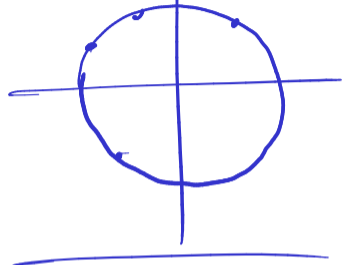


$$V Q Q^T V = V^T V$$

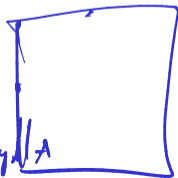
$$\|v\| = \|v\|_2$$

$$Q Q^T = I$$

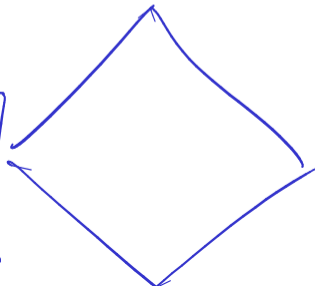
A is SPD then $\sqrt{A} = A^{1/2} = X D X^{-1}$
2-norm



1-norm



∞-norm



$$\|x+y\|_A \leq \|x\|_A + \|y\|_A$$

Conjugate Gradient Convergence Analysis (II)

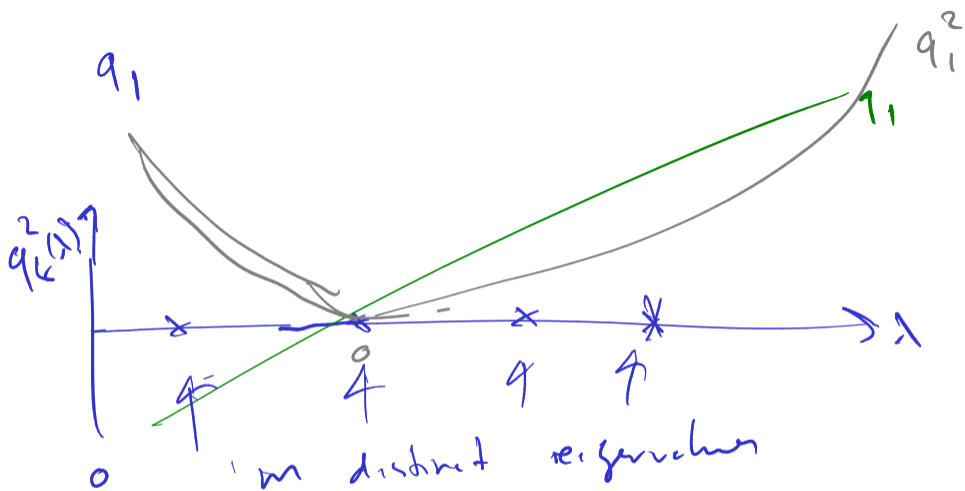
- Using $z = \rho_{k-1}(A)Ax$, we can simplify $\phi(z) = (x - z)^T A(x - z)$ as

$$\begin{aligned} \phi(z) &= (x - \rho_{k-1}(A)Ax)^T A(x - \rho_{k-1}(A)Ax) \\ &= \underbrace{(\underbrace{I - \rho_{k-1}(A)A}_{q_k(A)} x)^T}_{q_k(\alpha)} A(x - \rho_{k-1}(A)Ax) \end{aligned}$$

$q_k(\alpha) = 1 - \rho_{k-1}(\alpha)\alpha$
 $q_k(0) = 1$

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

$$\begin{aligned} \phi(z) &= x^T Q \underbrace{q_k(\lambda)}_{q_k(\lambda)} \lambda \underbrace{q_k(\lambda)}_{q_k(\lambda)} Q^T x \\ &\leq \max_{\lambda_i \in \lambda(A)} (q_k(\lambda_i)^2) x^T \underbrace{Q \Lambda Q^T}_A x \\ &= \max_{\lambda_i \in \lambda(A)} (q_k(\lambda_i)^2) \phi(x_0) \end{aligned}$$



then $\exists q_m \in \Pi_m$ with 0_s or
 $q_m(x) = \frac{\prod_{i=1}^m (\lambda_i - x)}{\prod_{i=1}^m \lambda_i} \quad \lambda_1, \dots, \lambda_m$

Conjugate Gradient Convergence Analysis (III)

Q_m are all degree ≤ m poly with q(0)=1

- Using our bound on the square of the residual norm $\phi(z)$, we can see why CG converges after $m - 1$ iterations if there are only $m < n$ distinct eigenvalues

Ax=b

$$e(x_k) \leq \min_{q_k \in \mathcal{Q}_k} \max_{\lambda_i \in \lambda(A)} (q_k(\lambda_i))^2 \phi(x_0)$$

$$\|r_k\|_{A^{-1}} = \sqrt{e(x_k)} \quad \left| \quad \frac{\|r_k\|_{A^{-1}}}{\|r_0\|_{A^{-1}}} \leq \min_{q_k \in \mathcal{Q}_k} \max_{\lambda_i \in \lambda(A)} |q_k(\lambda_i)| \right|$$

- 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 $q_m(0) = 1$)

$$\|r_m\| = 0 \Rightarrow f(y) = 0$$

$$Ax_m = b$$

Round-off Error in Conjugate Gradient

- ▶ CG provides strong convergence guarantees for SPD matrices in exact arithmetic

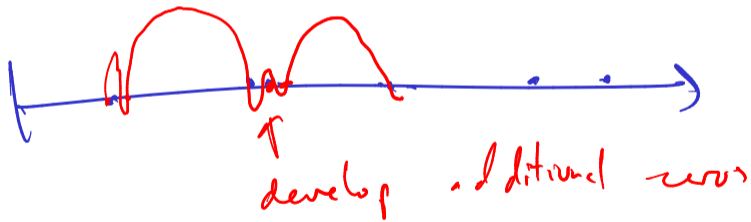
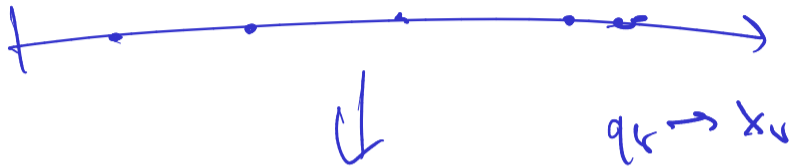
• classic view of CG (as direct) direct

• due to round-off, not so exact

• robust as an iterative method

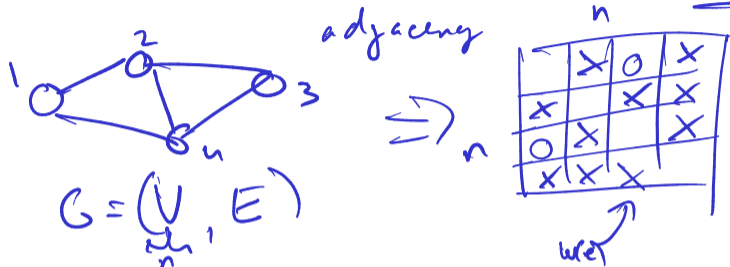
- ▶ Due to round-off CG may stagnate / have plateaus in convergence

• as CG on a larger matrix, whose eigenvalues are clustered around those of A



Graph and Matrix Duality

- ▶ graphs have a natural correspondence with sparse matrices



$w(u, v)$

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

- combinatorial algorithms (BFS, SSSP, BF)

- semirings $(\min, +)$ $(+, \cdot)$

$$w = A \otimes v \Rightarrow \underline{w}_i = \min_{j \in N(i)} a_{ij} + v_j$$

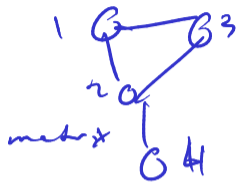
- $D = I \oplus (A \oplus (A \otimes A) \oplus \dots \oplus (A^{\otimes n})) \Rightarrow d_{ij} = \min_{k \in \{1, \dots, n\}} \min_{u_1, \dots, u_{k-1}} w(i, u_1) + w(u_1, u_2) + \dots + w(u_{k-1}, j)$

Graph Partitioning from Eigenvectors

- ▶ The Laplacian matrix provides a model of interactions on a graph that is useful in many contexts

$$L = \text{Laplacian} = D - A$$

\uparrow \uparrow
 diagonal adjacency
 degree matrix
 matrix



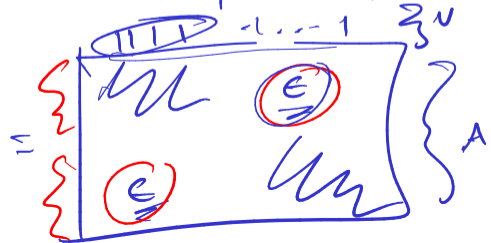
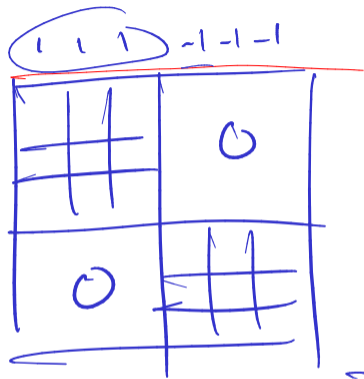
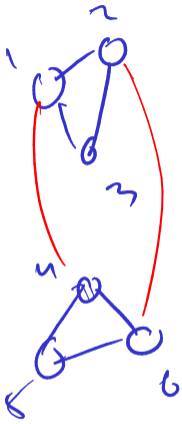
$$\begin{bmatrix} 2 & -1 & -1 & 0 \\ -1 & 3 & -1 & -1 \\ -1 & -1 & 2 & 0 \\ 0 & -1 & 0 & 1 \end{bmatrix}$$

- ▶ The second-smallest-eigenvalue eigenvector of the Laplacian (the Fiedler vector), gives a good partitioning of the graph

$$v = L \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix} = D \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix} - A \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

$$v_i = d_i - d_i = 0$$

L has at least 1
0 eigenvalue



$$P_A(U) = \frac{U^T A U}{U^T U} \approx E$$

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:

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

- ▶ Generally, we will seek to reduce constrained optimization problems to a series of simpler optimization problems:

Lagrangian Duality

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

- ▶ 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

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

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

