

CS257 Linear and Convex Optimization

Lecture 12

Bo Jiang

John Hopcroft Center for Computer Science
Shanghai Jiao Tong University

November 23, 2020

Recap: Damped Newton's Method

- 1: initialization $\mathbf{x} \leftarrow \mathbf{x}_0 \in \mathbb{R}^n$
- 2: **while** $\|\nabla f(\mathbf{x})\| > \delta$ **do**
- 3: $\mathbf{d} \leftarrow -[\nabla^2 f(\mathbf{x})]^{-1} \nabla f(\mathbf{x})$
- 4: find t with backtracking line search
- 5: $\mathbf{x} \leftarrow \mathbf{x} + t\mathbf{d}$
- 6: **end while**
- 7: **return** \mathbf{x}

Pure vs damped Newton's method (under appropriate conditions).

- Pure Newton's method has fast local convergence with no global guarantee

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq \frac{M}{2m} \|\mathbf{x}_k - \mathbf{x}^*\|^2$$

- Damped Newton's method guarantees global convergence with a slow damped phase and a fast pure phase.

Recap: Equality Constrained Optimization Problem

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & \mathbf{Ax} = \mathbf{b} \end{aligned}$$

where $\mathbf{A}^T = (\mathbf{a}_1, \dots, \mathbf{a}_k) \in \mathbb{R}^{n \times k}$, $\mathbf{b} \in \mathbb{R}^k$, f is differentiable and convex.

Assume the feasible set $X = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{Ax} = \mathbf{b}\} \neq \emptyset$, and the constraints are independent, i.e. $\text{rank} \mathbf{A} = k$.

First-order optimality condition. $\mathbf{x}^* \in X$ is optimal iff

$$\nabla f(\mathbf{x}^*) \perp \text{Null}(\mathbf{A})$$

or equivalently,

$$\nabla f(\mathbf{x}^*) \in \text{Range}(\mathbf{A}^T)$$

i.e.

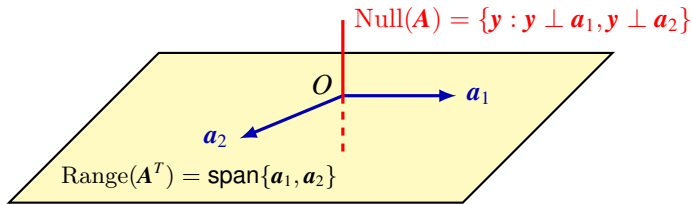
$$\nabla f(\mathbf{x}^*) + \mathbf{A}^T \boldsymbol{\lambda}^* = \nabla f(\mathbf{x}^*) + \sum_{i=1}^k \lambda_i^* \mathbf{a}_i = \mathbf{0}, \quad \text{for some } \boldsymbol{\lambda}^* \in \mathbb{R}^k$$

The constants $\lambda_1^*, \dots, \lambda_k^*$ are called **Lagrange multipliers**.

Appendix

Lemma. $\text{Null}(\mathbf{A})^\perp = \text{Range}(\mathbf{A}^T)$, where $\text{Range}(\mathbf{A}^T) = \{\mathbf{A}^T \mathbf{v} : \mathbf{v} \in \mathbb{R}^k\}$ and $\text{Null}(\mathbf{A})^\perp$ is the **orthogonal complement** of $\text{Null}(\mathbf{A})$, i.e.

$$\mathbf{x} \in \text{Null}(\mathbf{A})^\perp \iff \mathbf{x} \perp \mathbf{y}, \quad \forall \mathbf{y} \in \text{Null}(\mathbf{A})$$



Proof. Show $\text{Range}(\mathbf{A}^T) \subset \text{Null}(\mathbf{A})^\perp$ is a subspace with the same dimension, so $\text{Range}(\mathbf{A}^T) = \text{Null}(\mathbf{A})^\perp$.

- $\mathbf{x} \in \text{Range}(\mathbf{A}^T) \implies \mathbf{x} = \mathbf{A}^T \mathbf{z}$ for some \mathbf{z}
- $\forall \mathbf{y} \in \text{Null}(\mathbf{A}), \mathbf{x}^T \mathbf{y} = \mathbf{z}^T \mathbf{A} \mathbf{y} = \mathbf{z}^T \mathbf{0} = 0$, i.e. $\mathbf{x} \perp \mathbf{y}$, so $\mathbf{x} \in \text{Null}(\mathbf{A})^\perp$.
- $\dim \text{Range}(\mathbf{A}^T) = \text{rank } \mathbf{A} = n - \dim \text{Null}(\mathbf{A}) = \dim \text{Null}(\mathbf{A})^\perp$

Lagrange Condition

Define **Lagrangian** (or **Lagrange function**) by

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \boldsymbol{\lambda}^T (\mathbf{A}\mathbf{x} - \mathbf{b}) = f(\mathbf{x}) + \sum_{i=1}^k \lambda_i (\mathbf{a}_i^T \mathbf{x} - b_i)$$

The optimality condition becomes the following **KKT equations**¹

$$\begin{cases} \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*) = \nabla f(\mathbf{x}^*) + \mathbf{A}^T \boldsymbol{\lambda}^* = \mathbf{0} \\ \nabla_{\boldsymbol{\lambda}} \mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*) = \mathbf{A}\mathbf{x}^* - \mathbf{b} = \mathbf{0} \end{cases}$$

where $\nabla_{\mathbf{x}}$ and $\nabla_{\boldsymbol{\lambda}}$ are partial gradient² w.r.t. \mathbf{x} and $\boldsymbol{\lambda}$. or

$$\nabla \mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*) = \mathbf{0}$$

i.e. $(\mathbf{x}^*, \boldsymbol{\lambda}^*)$ is a stationary point of \mathcal{L} .

¹KKT stands for Karush-Kuhn-Tucker. We'll see later why it is called such.

²We use a similar notation $\nabla_{\mathbf{d}} f(\mathbf{x}) = \nabla f(\mathbf{x})^T \mathbf{d}$ to denote the directional derivative of f along the direction \mathbf{d} . The context should make it clear which is which.

Example

Consider

$$\begin{aligned} \min_{x_1, x_2} \quad & f(x_1, x_2) = \frac{1}{2}x_1^2 + \frac{1}{2}x_2^2 \\ \text{s.t.} \quad & x_1 + 2x_2 = 1 \end{aligned}$$

Method 1. Reduction to an equivalent unconstrained problem.

$$g(x_2) \triangleq f(1 - 2x_2, x_2) = \frac{1}{2}(1 - 2x_2)^2 + \frac{1}{2}x_2^2$$

$$\min_{x_2} g(x_2) \implies g'(x_2^*) = 0 \implies x_2^* = \frac{2}{5} \implies x_1^* = 1 - 2x_2^* = \frac{1}{5}$$

Method 2. Lagrangian multipliers method. The Lagrangian is

$$\mathcal{L}(x_1, x_2, \lambda) = \frac{1}{2}x_1^2 + \frac{1}{2}x_2^2 + \lambda(x_1 + 2x_2 - 1)$$

By the Lagrange condition,

$$\begin{cases} \frac{\partial \mathcal{L}}{\partial x_1} = x_1 + \lambda = 0 \\ \frac{\partial \mathcal{L}}{\partial x_2} = x_2 + 2\lambda = 0 \\ \frac{\partial \mathcal{L}}{\partial \lambda} = x_1 + 2x_2 - 1 = 0 \end{cases} \implies \begin{cases} x_1^* = \frac{1}{5} \\ x_2^* = \frac{2}{5} \\ \lambda^* = -\frac{1}{5} \end{cases}$$

Example (cont'd)

$$\begin{aligned} \min_{x_1, x_2} \quad & f(x_1, x_2) = \frac{1}{2}x_1^2 + \frac{1}{2}x_2^2 \\ \text{s.t.} \quad & x_1 + 2x_2 = 1 \end{aligned}$$

- normal vector to the feasible set X

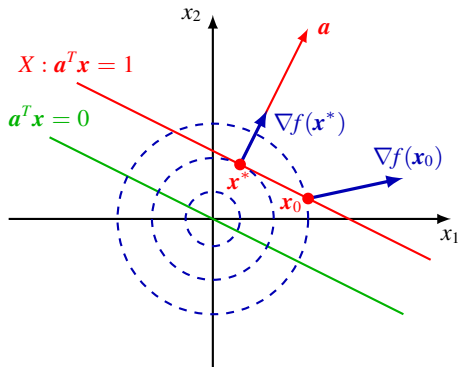
$$\mathbf{a} = (1, 2)^T$$

- gradient

$$\nabla f(\mathbf{x}) = \mathbf{x}$$

- at \mathbf{x}^* ,

$$\nabla f(\mathbf{x}^*) = -\lambda^* \mathbf{a} \perp X$$



Example

$$\begin{array}{ll} \min_{\mathbf{x}} & f(\mathbf{x}) = \frac{1}{2} \|\mathbf{x}\|^2, \\ \text{s.t.} & \mathbf{Ax} = \mathbf{b} \end{array}, \quad \text{where } \mathbf{A} = \begin{bmatrix} 1 & 2 & 0 \\ 2 & 2 & 1 \end{bmatrix}, \mathbf{b} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

Method 1. Reduction to an equivalent unconstrained problem.

- $\text{rank } \mathbf{A} = 2$. Find two independent columns of \mathbf{A} , e.g. the first and third columns, and solve for the corresponding x_i 's in terms of the others. Let $\mathbf{A}_1 = \begin{bmatrix} 1 & 0 \\ 2 & 1 \end{bmatrix}$, $\mathbf{A}_2 = \begin{bmatrix} 2 \\ 2 \end{bmatrix}$. The constraints become

$$\mathbf{A}_1 \begin{bmatrix} x_1 \\ x_3 \end{bmatrix} + \mathbf{A}_2 x_2 = \mathbf{b} \implies \begin{bmatrix} x_1 \\ x_3 \end{bmatrix} = \mathbf{A}_1^{-1} \mathbf{b} - \mathbf{A}_1^{-1} \mathbf{A}_2 x_2 = \begin{bmatrix} 1 - 2x_2 \\ 2x_2 - 1 \end{bmatrix}$$

- Substitution into f yields

$$g(x_2) = f(1 - 2x_2, x_2, 2x_2 - 1) = (2x_2 - 1)^2 + \frac{1}{2}x_2^2 \implies x_2^* = \frac{4}{9}$$

- $x_1^* = 1 - 2x_2^* = \frac{1}{9}$, $x_3^* = 2x_2^* - 1 = -\frac{1}{9}$

Example (cont'd)

$$\begin{array}{ll} \min_{\mathbf{x}} & f(\mathbf{x}) = \frac{1}{2} \|\mathbf{x}\|^2, \\ \text{s.t.} & \mathbf{Ax} = \mathbf{b} \end{array}, \quad \text{where } \mathbf{A} = \begin{bmatrix} 1 & 2 & 0 \\ 2 & 2 & 1 \end{bmatrix}, \mathbf{b} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

Method 2. Lagrange multipliers method.

- The Lagrangian is

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = \frac{1}{2} \|\mathbf{x}\|^2 + \boldsymbol{\lambda}^T (\mathbf{Ax} - \mathbf{b})$$

- Lagrange condition

$$\begin{cases} \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = \mathbf{x} + \mathbf{A}^T \boldsymbol{\lambda} = \mathbf{0} \\ \nabla_{\boldsymbol{\lambda}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = \mathbf{Ax} - \mathbf{b} = \mathbf{0} \end{cases} \quad \text{or} \quad \begin{bmatrix} \mathbf{I} & \mathbf{A}^T \\ \mathbf{A} & \mathbf{O} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \boldsymbol{\lambda} \end{bmatrix} = \begin{bmatrix} \mathbf{0} \\ \mathbf{b} \end{bmatrix}$$

- Solve for $\mathbf{x}, \boldsymbol{\lambda}$ e.g. by block Gaussian elimination,

$$\begin{cases} \mathbf{x}^* = -\mathbf{A}^T \boldsymbol{\lambda}^* = \mathbf{A}^T (\mathbf{AA}^T)^{-1} \mathbf{b} \\ \boldsymbol{\lambda}^* = -(\mathbf{AA}^T)^{-1} \mathbf{b} \end{cases} \quad \implies \quad \begin{cases} \mathbf{x}^* = (\frac{1}{9}, \frac{4}{9}, -\frac{1}{9})^T \\ \boldsymbol{\lambda}^* = (-\frac{1}{3}, \frac{1}{9})^T \end{cases}$$

Example (cont'd)

Block Gaussian elimination.

- The augmented matrix is

$$\begin{bmatrix} I & A^T & \mathbf{0} \\ A & O & b \end{bmatrix}$$

- Multiply the first “row” by $-A$ and add to the second “row”,

$$\begin{bmatrix} I & A^T & \mathbf{0} \\ O & -AA^T & b \end{bmatrix}$$

- Multiply the second “row” by $-(AA^T)^{-1}$ (why invertible?),

$$\begin{bmatrix} I & A^T & \mathbf{0} \\ O & I & -(AA^T)^{-1}b \end{bmatrix}$$

- Multiply the second “row” by $-A^T$ and add to the first “row”,

$$\begin{bmatrix} I & O & A^T(AA^T)^{-1}b \\ O & I & -(AA^T)^{-1}b \end{bmatrix}$$

Example (cont'd)

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) = \frac{1}{2} \|\mathbf{x}\|^2, \\ \text{s.t.} \quad & \mathbf{A}\mathbf{x} = \mathbf{b} \end{aligned}, \quad \text{where } \mathbf{A} = \begin{bmatrix} 1 & 2 & 0 \\ 2 & 2 & 1 \end{bmatrix}, \mathbf{b} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

- vectors normal to the feasible set X

$$\text{span}\{\mathbf{a}_1, \mathbf{a}_2\}$$

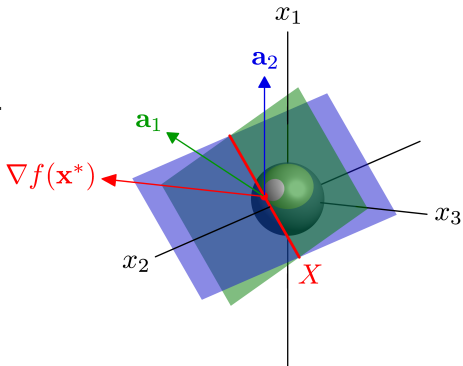
with $\mathbf{a}_1 = (1, 2, 0)^T$, $\mathbf{a}_2 = (2, 2, 1)^T$.

- gradient

$$\nabla f(\mathbf{x}) = \mathbf{x}$$

- at \mathbf{x}^* ,

$$\nabla f(\mathbf{x}^*) = -\lambda_1^* \mathbf{a}_1 - \lambda_2^* \mathbf{a}_2 \perp X$$



Equality Constrained Convex QP

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{g}^T \mathbf{x} + c \\ \text{s.t.} \quad & \mathbf{A} \mathbf{x} = \mathbf{b} \end{aligned} \quad (\star)$$

where $\mathbf{Q} \in \mathbb{R}^n$, $\mathbf{Q} \succeq \mathbf{O}$, $\mathbf{A} \in \mathbb{R}^{k \times n}$, $\text{rank} \mathbf{A} = k$.

Note. This is the basis for an extension of Newton's method to equality constrained problems.

- The Lagrangian is

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{g}^T \mathbf{x} + c + \boldsymbol{\lambda}^T (\mathbf{A} \mathbf{x} - \mathbf{b})$$

- The Lagrange condition is

$$\begin{cases} \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = \mathbf{Q} \mathbf{x} + \mathbf{g} + \mathbf{A}^T \boldsymbol{\lambda} = \mathbf{0} \\ \nabla_{\boldsymbol{\lambda}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = \mathbf{A} \mathbf{x} - \mathbf{b} = \mathbf{0} \end{cases} \quad \text{or} \quad \begin{bmatrix} \mathbf{Q} & \mathbf{A}^T \\ \mathbf{A} & \mathbf{O} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \boldsymbol{\lambda} \end{bmatrix} = \begin{bmatrix} -\mathbf{g} \\ \mathbf{b} \end{bmatrix}$$

This is the **KKT system** of the problem (\star) . The coefficient matrix

$\mathbf{K} = \begin{bmatrix} \mathbf{Q} & \mathbf{A}^T \\ \mathbf{A} & \mathbf{O} \end{bmatrix}$ is called the **KKT matrix**.

Solving KKT System When $Q \succ O$

$$\begin{cases} Qx + g + A^T \lambda = 0 \\ Ax - b = 0 \end{cases} \quad \text{or} \quad \begin{bmatrix} Q & A^T \\ A & O \end{bmatrix} \begin{bmatrix} x \\ \lambda \end{bmatrix} = \begin{bmatrix} -g \\ b \end{bmatrix}$$

1. Solving for x in term of λ from $Qx + g + A^T \lambda = 0$,

$$x = -Q^{-1}g - Q^{-1}A^T \lambda$$

2. Substituting into $Ax - b = 0$,

$$-AQ^{-1}g - AQ^{-1}A^T \lambda = b$$

3. Since $AQ^{-1}A^T \succ O$ (why?), solving for λ ,

$$\lambda = -[AQ^{-1}A^T]^{-1}[AQ^{-1}g + b]$$

4. Plugging into step 1,

$$x = -Q^{-1}g + Q^{-1}A^T[AQ^{-1}A^T]^{-1}[AQ^{-1}g + b]$$

Note. We can also use block Gaussian elimination (cf. slide 9).

Unsolvable KKT System

Example.

$$\begin{aligned} \min_{x_1, x_2} \quad & f(x_1, x_2) = \frac{1}{2}x_2^2 + x_1 \\ \text{s.t.} \quad & x_2 = 0 \end{aligned}$$

This is a convex QP with

$$\mathbf{Q} = \text{diag}\{0, 1\}, \quad \mathbf{g} = (1, 0)^T, \quad \mathbf{A} = (0, 1), \quad \mathbf{b} = 0$$

The KKT system is

$$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 1 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \lambda \end{bmatrix} = \begin{bmatrix} -1 \\ 0 \\ 0 \end{bmatrix}$$

which has no solution, since $0 \cdot x_1 + 0 \cdot x_2 + 0 \cdot \lambda \neq -1$.

Note $f^* = -\infty$ for this problem.

Unsolvable KKT System (cont'd)

If the KKT system has **no** solution, then the problem (\star) is either **infeasible** or **unbounded below**.

- KKT system has no solution iff

$$\begin{bmatrix} -\mathbf{g} \\ \mathbf{b} \end{bmatrix} \notin \text{Range}(\mathbf{K}) = \text{Range}(\mathbf{K}^T) = \text{Null}(\mathbf{K})^\perp$$

- There exists $\begin{bmatrix} \mathbf{v} \\ \mathbf{w} \end{bmatrix} \in \text{Null}(\mathbf{K})$ s.t. $\begin{bmatrix} -\mathbf{g} \\ \mathbf{b} \end{bmatrix}^T \begin{bmatrix} \mathbf{v} \\ \mathbf{w} \end{bmatrix} \neq 0$, i.e.

$$\mathbf{Q}\mathbf{v} + \mathbf{A}^T\mathbf{w} = \mathbf{0}, \quad \mathbf{A}\mathbf{v} = \mathbf{0}, \quad -\mathbf{g}^T\mathbf{v} + \mathbf{b}^T\mathbf{w} \neq 0$$

- If \mathbf{x}_0 is feasible, then $\mathbf{x}_0 + t\mathbf{v}$ is feasible for any $t \in \mathbb{R}$,

$$\begin{aligned} f(\mathbf{x}_0 + t\mathbf{v}) &= f(\mathbf{x}_0) + t(\mathbf{x}_0^T\mathbf{Q}\mathbf{v} + \mathbf{g}^T\mathbf{v}) + \frac{1}{2}t^2\mathbf{v}^T\mathbf{Q}\mathbf{v} \\ &= f(\mathbf{x}_0) + t(-\mathbf{x}_0^T\mathbf{A}^T\mathbf{w} + \mathbf{g}^T\mathbf{v}) - \frac{1}{2}t^2\mathbf{w}^T\mathbf{A}\mathbf{v} \quad (\text{use } \mathbf{Q}\mathbf{v} = -\mathbf{A}^T\mathbf{w}) \\ &= f(\mathbf{x}_0) - t(\mathbf{b}^T\mathbf{w} - \mathbf{g}^T\mathbf{v}) \quad (\text{use } \mathbf{A}\mathbf{v} = \mathbf{0} \text{ and } \mathbf{A}\mathbf{x}_0 = \mathbf{b}) \end{aligned}$$

which goes to $-\infty$, as $t \rightarrow \text{sign}(\mathbf{b}^T\mathbf{w} - \mathbf{g}^T\mathbf{v}) \cdot \infty$.

Nonsingularity of KKT Matrix

If the KKT matrix \mathbf{K} is nonsingular, then the KKT system has a unique solution, which is optimal.

Recall $\mathbf{Q} \succeq \mathbf{O}$ and $\text{rank } \mathbf{A} = k$. The following conditions are equivalent

1. \mathbf{K} is nonsingular
2. $\text{Null}(\mathbf{Q}) \cap \text{Null}(\mathbf{A}) = \{\mathbf{0}\}$, i.e. \mathbf{Q} and \mathbf{A} have no nontrivial common nullspace, i.e. $\mathbf{Ax} = \mathbf{0}$, $\mathbf{Qx} = \mathbf{0}$ only have the trivial solution $\mathbf{x} = \mathbf{0}$.
3. $\mathbf{Ax} = \mathbf{0}, \mathbf{x} \neq \mathbf{0} \implies \mathbf{x}^T \mathbf{Qx} > 0$, i.e. \mathbf{Q} is positive definite on the nullspace of \mathbf{A} .
4. $\mathbf{F}^T \mathbf{QF} \succ \mathbf{O}$ for any $\mathbf{F} \in \mathbb{R}^{n \times (n-k)}$ s.t. $\text{Range}(\mathbf{F}) = \text{Null}(\mathbf{A})$, i.e. the columns of \mathbf{F} are linearly independent solutions of $\mathbf{Ax} = \mathbf{0}$.

In particular, if $\mathbf{Q} \succ \mathbf{O}$, then \mathbf{K} is nonsingular (by 3).

Proof

We show $1 \Rightarrow 2 \Rightarrow 3 \Rightarrow 4 \Rightarrow 1$.

- (1 \Rightarrow 2). If $\mathbf{0} \neq \mathbf{x} \in \text{Null}(\mathbf{Q}) \cap \text{Null}(\mathbf{A})$, then

$$\begin{bmatrix} \mathbf{Q} & \mathbf{A}^T \\ \mathbf{A} & \mathbf{O} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{0} \end{bmatrix} = \mathbf{0}$$

contradicting the nonsingularity of \mathbf{K} .

- (2 \Rightarrow 3.) Assume $\mathbf{Ax} = \mathbf{0}$ and $\mathbf{x}^T \mathbf{Qx} = 0$. We show $\mathbf{x} = \mathbf{0}$. Since $\mathbf{Q} \succeq \mathbf{O}$, $\mathbf{x}^T \mathbf{Qx} = 0$ iff $\mathbf{Qx} = \mathbf{0}$ ³. By 2, $\mathbf{Ax} = \mathbf{0}$ and $\mathbf{Qx} = \mathbf{0}$ implies $\mathbf{x} = \mathbf{0}$.

³Proof of necessity: Let $\mathbf{x}_1, \dots, \mathbf{x}_n$ be an orthonormal eigenbasis of \mathbf{Q} and $\mathbf{Qx}_i = \nu_i \mathbf{x}_i$. Then $\mathbf{Q} = \sum_{i=1}^n \nu_i \mathbf{x}_i \mathbf{x}_i^T$. Note $\nu_i \geq 0$, since $\mathbf{Q} \succeq \mathbf{O}$. Then

$$0 = \mathbf{x}^T \mathbf{Qx} = \mathbf{x}^T \left(\sum_{i=1}^n \nu_i \mathbf{x}_i \mathbf{x}_i^T \right) \mathbf{x} = \sum_{i=1}^n \nu_i \|\mathbf{x}_i^T \mathbf{x}\|^2 \implies \mathbf{x}_i^T \mathbf{x} = 0 \text{ if } \nu_i > 0$$

i.e. either $\nu_i = 0$ or $\mathbf{x}_i^T \mathbf{x} = 0$. Thus $\mathbf{Qx} = \sum_{i=1}^n \nu_i \mathbf{x}_i \mathbf{x}_i^T \mathbf{x} = \mathbf{0}$.

Proof (cont'd)

- (3 \Rightarrow 4.) $\text{rank } \mathbf{F} = \dim \text{Null}(\mathbf{A}) = n - \text{rank } \mathbf{A} = n - k$, so \mathbf{F} has full column rank. If $\mathbf{z} \neq \mathbf{0}$, then $\mathbf{x} = \mathbf{Fz} \neq \mathbf{0}$ and $\mathbf{x} \in \text{Range}(\mathbf{F}) = \text{Null}(\mathbf{A})$. By 3, $\mathbf{z}^T (\mathbf{F}^T \mathbf{Q} \mathbf{F}) \mathbf{z} = \mathbf{x}^T \mathbf{Q} \mathbf{x} > 0$.
- (4 \Rightarrow 1.) To show \mathbf{K} is nonsingular, we assume

$$\begin{bmatrix} \mathbf{Q} & \mathbf{A}^T \\ \mathbf{A} & \mathbf{O} \end{bmatrix} \begin{bmatrix} \mathbf{v} \\ \mathbf{w} \end{bmatrix} = \mathbf{0}$$

and show $\mathbf{v} = \mathbf{0}$ and $\mathbf{w} = \mathbf{0}$. Note $\mathbf{v} \in \text{Null}(\mathbf{A})$ and

$$\mathbf{Q} \mathbf{v} = -\mathbf{A}^T \mathbf{w} \implies \mathbf{v}^T \mathbf{Q} \mathbf{v} = -\mathbf{v}^T \mathbf{A}^T \mathbf{w} = -(\mathbf{A} \mathbf{v})^T \mathbf{w} = 0$$

Let $\mathbf{F} \in \mathbb{R}^{n \times (n-k)}$ be a matrix whose columns consist of a basis of $\text{Null}(\mathbf{A})$. Then $\text{Range}(\mathbf{F}) = \text{Null}(\mathbf{A})$ and $\mathbf{v} = \mathbf{Fz}$ for some \mathbf{z} . Now

$$0 = \mathbf{v}^T \mathbf{Q} \mathbf{v} = \mathbf{z}^T \mathbf{F}^T \mathbf{Q} \mathbf{F} \mathbf{z}$$

By 4, $\mathbf{z} = \mathbf{0}$, so $\mathbf{v} = \mathbf{Fz} = \mathbf{0}$. Then $\mathbf{A}^T \mathbf{w} = -\mathbf{Q} \mathbf{v} = \mathbf{0}$. Since \mathbf{A}^T has full column rank, $\mathbf{w} = \mathbf{0}$.

Example

$$\begin{aligned} \min_{x_1, x_2} \quad & f(x_1, x_2) = \frac{1}{2}x_2^2 \\ \text{s.t.} \quad & x_1 + 2x_2 = b \end{aligned}$$

Trivial with solution $x_1^* = b, x_2^* = 0$.

But let's check the condition on slide 15.

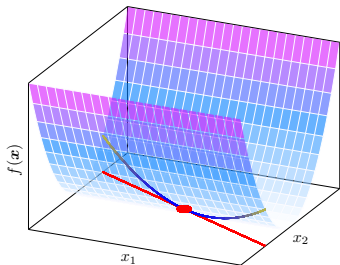
$$\mathbf{Q} = \text{diag}\{0, 1\}, \quad \mathbf{A} = (1, 2)$$

Let $\mathbf{F} = (2, -1)^T$. Then $\text{Range}(\mathbf{F}) = \text{Null}(\mathbf{A})$, and

$$\mathbf{F}^T \mathbf{Q} \mathbf{F} = [1] \succ \mathbf{0}$$

By 4 of slide 15, the KKT matrix is nonsingular, so \exists a unique solution.

Note. The unconstrained problem $\min_{\mathbf{x}} f(\mathbf{x})$ has infinitely many solutions. But this does not prevent the constrained problem from having a unique solution, as $\mathbf{Q} \succ \mathbf{0}$ on $\text{Null}(\mathbf{A})$ (see 3 on slide 15).



Newton Direction for Equality Constrained Problem

Consider the second-order Taylor approximation for f at a **feasible** \mathbf{x}_k ,

$$\begin{aligned} \min_{\mathbf{d}} \quad & h(\mathbf{d}) \triangleq \hat{f}(\mathbf{x}_k + \mathbf{d}) = f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T \mathbf{d} + \frac{1}{2} \mathbf{d}^T \nabla^2 f(\mathbf{x}_k) \mathbf{d} \\ \text{s.t.} \quad & \mathbf{A}(\mathbf{x}_k + \mathbf{d}) = \mathbf{b} \end{aligned}$$

Using $\mathbf{A}\mathbf{x}_k = \mathbf{b}$,

$$\begin{aligned} \min_{\mathbf{d}} \quad & h(\mathbf{d}) = f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T \mathbf{d} + \frac{1}{2} \mathbf{d}^T \nabla^2 f(\mathbf{x}_k) \mathbf{d} \\ \text{s.t.} \quad & \mathbf{A}\mathbf{d} = \mathbf{0} \end{aligned}$$

KKT system for this quadratic problem is

$$\begin{bmatrix} \nabla^2 f(\mathbf{x}_k) & \mathbf{A}^T \\ \mathbf{A} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{d} \\ \boldsymbol{\lambda} \end{bmatrix} = \begin{bmatrix} -\nabla f(\mathbf{x}_k) \\ \mathbf{0} \end{bmatrix}$$

The **Newton direction** \mathbf{d}_k is given by the solution to the KKT system. We will assume the KKT matrix is nonsingular (cf. slide 15)

Newton's Method for Equality Constrained Problem

1: initialization $\mathbf{x} \leftarrow \mathbf{x}_0 \in X$ ▷ \mathbf{x}_0 is feasible, i.e. $A\mathbf{x}_0 = \mathbf{b}$

2: **repeat**

3: **Compute Newton's direction \mathbf{d} by solving**

$$\begin{bmatrix} \nabla^2 f(\mathbf{x}) & A^T \\ A & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{d} \\ \lambda \end{bmatrix} = \begin{bmatrix} -\nabla f(\mathbf{x}) \\ \mathbf{0} \end{bmatrix}$$

4: $t \leftarrow 1$ ▷ backtracking line search on lines 4-7

5: **while** $f(\mathbf{x} + t\mathbf{d}) > f(\mathbf{x}) + \alpha t \nabla f(\mathbf{x})^T \mathbf{d}$ **do**

6: $t \leftarrow \beta t$

7: **end while**

8: $\mathbf{x} \leftarrow \mathbf{x} + t\mathbf{d}$

9: **until** $\|\mathbf{d}\| \leq \delta$

10: **return** \mathbf{x}

Note. We **cannot** use $\|\nabla f(\mathbf{x})\| \leq \delta$ as stopping criterion now, as $\nabla f(\mathbf{x}^*) = \mathbf{0}$ no longer holds in general. [BV] uses $\sqrt{\mathbf{d}^T \nabla^2 f(\mathbf{x}) \mathbf{d}} \leq \delta$.

Note. This is called a **feasible descent method**, since all \mathbf{x}_k are feasible and $f(\mathbf{x}_{k+1}) < f(\mathbf{x}_k)$ unless \mathbf{x}_k is optimal.

Newton's Method and Constraint Elimination

Let $F \in \mathbb{R}^{n \times (n-k)}$ be a matrix whose columns are linearly independent solutions to $Ax = \mathbf{0}$. For a fixed feasible $\tilde{x} \in X$,

$$X = \{x : Ax = b\} = \{\tilde{x} + Fz : z \in \mathbb{R}^{n-k}\}$$

Constrained problem reduces to unconstrained problem by $x = Fz + \tilde{x}$,

$$\begin{cases} \min_x & f(x) \\ \text{s.t.} & Ax = b \end{cases} \iff \min_z g(z) = f(Fz + \tilde{x})$$

Note (slides 8 and 17 of Lecture 2),

$$\nabla g(z) = F^T \nabla f(Fz + \tilde{x}), \quad \nabla^2 g(z) = F^T \nabla^2 f(Fz + \tilde{x}) F \quad (\dagger)$$

Apply Newton's method to both problems with initial points x_0 and z_0 . If $x_0 = Fz_0 + \tilde{x}$, we show by induction that $x_k = Fz_k + \tilde{x}$, so Newton's method converges for the constrained problem if it does for the unconstrained problem.

Proof

The Newton direction $\Delta \mathbf{x}_k$ for the constrained problem satisfies

$$\begin{bmatrix} \nabla^2 f(\mathbf{x}_k) & \mathbf{A}^T \\ \mathbf{A} & \mathbf{O} \end{bmatrix} \begin{bmatrix} \Delta \mathbf{x}_k \\ \boldsymbol{\lambda}_k \end{bmatrix} = \begin{bmatrix} -\nabla f(\mathbf{x}_k) \\ \mathbf{0} \end{bmatrix}$$

1. By the induction hypothesis $\mathbf{x}_k = \mathbf{F}\mathbf{z}_k + \tilde{\mathbf{x}}$ and (†) on the previous slide

$$\mathbf{F}^T \nabla^2 f(\mathbf{x}_k) \mathbf{F} = \nabla^2 g(\mathbf{z}_k), \quad \mathbf{F}^T \nabla f(\mathbf{x}_k) = \nabla g(\mathbf{z}_k)$$

2. By 1 above and 3 on slide 15, the KKT matrix in 1 is nonsingular iff $\nabla^2 g(\mathbf{z}_k) \succ \mathbf{O}$, so $\Delta \mathbf{x}_k$ is well-defined iff the Newton direction $\Delta \mathbf{z}_k = -[\nabla^2 g(\mathbf{z}_k)]^{-1} \nabla g(\mathbf{z}_k)$ is well-defined.
3. Since $\text{Null}(\mathbf{A}) = \text{Range}(\mathbf{F})$, $\mathbf{A} \Delta \mathbf{x}_k = \mathbf{0} \iff \Delta \mathbf{x}_k = \mathbf{F}\mathbf{u}$ for some \mathbf{u} .
4. Plugging $\Delta \mathbf{x}_k = \mathbf{F}\mathbf{u}$ into the first KKT equation,

$$\nabla^2 f(\mathbf{x}_k) \mathbf{F}\mathbf{u} + \mathbf{A}^T \boldsymbol{\lambda}_k = -\nabla f(\mathbf{x}_k)$$

5. Pre-multiplying by \mathbf{F}^T ,

$$\mathbf{F}^T \nabla^2 f(\mathbf{x}_k) \mathbf{F}\mathbf{u} + (\mathbf{A}\mathbf{F})^T \boldsymbol{\lambda}_k = -\mathbf{F}^T \nabla f(\mathbf{x}_k)$$

Proof (cont'd)

6. Since the columns of F are solutions to $Ax = \mathbf{0}$, $AF = \mathbf{O}$

7. By 1, 5 and 6,

$$\nabla^2 g(\mathbf{z}_k) \mathbf{u} = -\nabla g(\mathbf{z}_k) \implies \mathbf{u} = -[\nabla^2 g(\mathbf{z}_k)]^{-1} \nabla g(\mathbf{z}_k) = \Delta \mathbf{z}_k$$

so

$$\Delta \mathbf{x}_k = F \mathbf{u} = F \Delta \mathbf{z}_k$$

8. By 7, backtracking line search gives the same step size t_k , since

$$f(\mathbf{x}_k + t \Delta \mathbf{x}_k) = f(F(\mathbf{z}_k + t \Delta \mathbf{z}_k) + \tilde{\mathbf{x}}) = g(\mathbf{z}_k + t \Delta \mathbf{z}_k)$$

9. By 7, 8, and the induction hypothesis $\mathbf{x}_k = F \mathbf{z}_k + \tilde{\mathbf{x}}$,

$$\mathbf{x}_{k+1} = \mathbf{x}_k + t_k \Delta \mathbf{x}_k = F \mathbf{z}_k + \tilde{\mathbf{x}} + t_k F \Delta \mathbf{z}_k = F(\mathbf{z}_k + t_k \Delta \mathbf{z}_k) + \tilde{\mathbf{x}} = F \mathbf{z}_{k+1} + \tilde{\mathbf{x}}$$

completing the induction.