

CS257 Linear and Convex Optimization

Lecture 11

Bo Jiang

John Hopcroft Center for Computer Science
Shanghai Jiao Tong University

November 16, 2020

Recap: Line Search

Exact line search.

$$t_k = \arg \min_s f(\mathbf{x}_k - s \nabla f(\mathbf{x}_k))$$

Backtracking line search (Armijo's rule).

$$f(\mathbf{x}_k) - f(\mathbf{x}_k - t_k \nabla f(\mathbf{x}_k)) \geq \alpha t_k \|\nabla f(\mathbf{x}_k)\|_2^2$$

- 1: initialization $\mathbf{x} \leftarrow \mathbf{x}_0 \in \mathbb{R}^n$
- 2: **while** $\|\nabla f(\mathbf{x})\| > \delta$ **do**
- 3: choose direction \mathbf{d} $\triangleright \mathbf{d} = -\nabla f(\mathbf{x})$ for gradient descent
- 4: $t \leftarrow t_0$
- 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: **end while**
- 10: **return** \mathbf{x}

Recap: Convergence of Gradient Descent

For m -strongly convex and L -smooth f with minimum \mathbf{x}^*

- gradient descent with constant step size $t \in (0, \frac{1}{L}]$ satisfies

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \frac{L(1 - mt)^k}{m} [f(\mathbf{x}_0) - f(\mathbf{x}^*)]$$

- gradient descent with exact line search satisfies

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \left(1 - \frac{m}{L}\right)^k [f(\mathbf{x}_0) - f(\mathbf{x}^*)]$$

- gradient descent with backtracking line search satisfies

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq c^k [f(\mathbf{x}_0) - f(\mathbf{x}^*)]$$

where

$$c = 1 - \min \left\{ 2m\alpha t_0, \frac{4m\beta\alpha(1 - \alpha)}{L} \right\}$$

Recap: Newton's Method

Newton's method for solving optimization problem $\min_{\mathbf{x}} f(\mathbf{x})$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k)$$

Newton's method for solving $\mathbf{g}(\mathbf{x}) = \mathbf{0}$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - [D\mathbf{g}(\mathbf{x}_k)]^{-1} \mathbf{g}(\mathbf{x}_k)$$

Connection. First-order optimality condition

$$\nabla f(\mathbf{x}^*) = \mathbf{0}$$

Today

- Analysis of Newton's method
- Damped Newton's method
- Equality Constrained Optimization

Contents

1. Analysis of Newton's method

2. Damped Newton's Method

3. Equality Constrained Optimization

Convergence of Newton's Method

Example. Consider the minimization of $f(x) = \sqrt{1+x^2}$.

$$f'(x) = \frac{x}{\sqrt{1+x^2}}, \quad f''(x) = \frac{1}{(1+x^2)^{3/2}}$$

The Newton direction is

$$d_k = -f'(x_k)/f''(x_k) = -x_k - x_k^3$$

The Newton step is

$$x_{k+1} = x_k + d_k = -x_k^3$$

Note $x_k \rightarrow x^* = 0$ iff $|x_0| < 1$. When $|x_0| > 1$, x_k diverges, and

$$f(x_{k+1}) > f(x_k)$$

In general, Newton's method does **not** guarantee global convergence. When it does converge, the convergence is usually very fast.

Convergence Analysis: 1D Case

Theorem. If f is m -strongly convex, f'' is M -Lipschitz continuous, and x^* is a minimum of f , then the sequence $\{x_k\}$ produced by Newton's method satisfies

$$|x_{k+1} - x^*| \leq \frac{M}{2m} |x_k - x^*|^2$$

Notes. Let $\xi_k = \frac{M}{2m} |x_k - x^*|$. The above inequality becomes $\xi_{k+1} \leq \xi_k^2$.

- If $\xi_k = 10^{-p}$, then $\xi_{k+1} \leq 10^{-2p}$, the number of significant digits doubles in each iteration!
- If $\xi_0 < 1$ i.e. $|x_0 - x^*| < \frac{2m}{M}$, then $\xi_k \leq \xi_0^{2^k}$ converges to 0 extremely fast. The number of iterations to ensure $\xi_k \leq \epsilon$ is $k \geq \log_2 \log_{\frac{1}{\xi_0}} \frac{1}{\epsilon}$.
For $\epsilon = 10^{-p}$, $k \geq \log_2 p + \log_2 \log_{\frac{1}{\xi_0}} 10$, only logarithmic in the number of digits. Very few iterations are required!
- This theorem is a **local** convergence result. Fast convergence if x_0 is close enough to x^* , i.e. $|x_0 - x^*| < \frac{2m}{M}$. No guarantee if $|x_0 - x^*|$ is large.

Proof: 1D Case

$$\begin{aligned} & |x_{k+1} - x^*| \\ &= |x_k - x^* - [f''(x_k)]^{-1} f'(x_k)| \\ &= |f''(x_k)|^{-1} \cdot |f'(x^*) - f'(x_k) - f''(x_k)(x^* - x_k)| \\ &= \frac{|x_k - x^*|}{|f''(x_k)|} \cdot \left| \int_0^1 [f''(x_k + t(x^* - x_k)) - f''(x_k)] dt \right| \\ &\leq \frac{|x_k - x^*|}{|f''(x_k)|} \cdot \int_0^1 |f''(x_k + t(x^* - x_k)) - f''(x_k)| dt \\ &\leq \frac{|x_k - x^*|}{|f''(x_k)|} \cdot \int_0^1 Mt|x_k - x^*| dt \\ &= \frac{M}{2|f''(x_k)|} |x_k - x^*|^2 \\ &\leq \frac{M}{2m} |x_k - x^*|^2 \end{aligned}$$

Newton step

$$f'(x^*) = 0$$

Newton-Leibniz

$$\left| \int f \right| \leq \int |f|$$

M -Lipschitz of f''

m -strong convexity

Matrix Norm

The set of $m \times n$ matrices $\mathbb{R}^{m \times n}$ is a mn -dimensional vector space

A **matrix norm** on $\mathbb{R}^{m \times n}$ is a function $\|\cdot\| : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$ s.t.

1. $\|\mathbf{A}\| \geq 0, \forall \mathbf{A} \in \mathbb{R}^{m \times n}$
2. $\|\mathbf{A}\| = 0$ iff $\mathbf{A} = \mathbf{O}$
3. $\|c\mathbf{A}\| = |c| \cdot \|\mathbf{A}\|, \forall c \in \mathbb{R}, \mathbf{A} \in \mathbb{R}^{m \times n}$ (**positive homogeneity**)
4. $\|\mathbf{A} + \mathbf{B}\| \leq \|\mathbf{A}\| + \|\mathbf{B}\|, \forall \mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times n}$ (**triangle inequality**)

Example. The **Frobenius norm** on $\mathbb{R}^{m \times n}$ is the 2-norm on \mathbb{R}^{mn} .

$$\|\mathbf{A}\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n a_{ij}^2} \quad \text{for } \mathbf{A} = (a_{ij}) \in \mathbb{R}^{m \times n}$$

Operator Norm

A matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ defines a linear transformation from \mathbb{R}^n to \mathbb{R}^m

$$\begin{aligned}\mathbf{A} &: \mathbb{R}^n \rightarrow \mathbb{R}^m \\ \mathbf{x} &\mapsto \mathbf{Ax}\end{aligned}$$

Given two vector norms $\|\cdot\|_a$ and $\|\cdot\|_b$ on \mathbb{R}^n and \mathbb{R}^m , respectively, the **operator norm** or **induced norm** of \mathbf{A} is defined by

$$\|\mathbf{A}\|_{a,b} = \max_{\mathbf{x}:\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{Ax}\|_b}{\|\mathbf{x}\|_a} = \max_{\mathbf{x}:\|\mathbf{x}\|_a=1} \|\mathbf{Ax}\|_b = \max_{\mathbf{x}:\|\mathbf{x}\|_a \leq 1} \|\mathbf{Ax}\|_b$$

Exercise. Show the three definitions are equivalent.

The induced norm has the following important property.

Proposition (compatibility of norms).

$$\|\mathbf{Ax}\|_b \leq \|\mathbf{A}\|_{a,b} \|\mathbf{x}\|_a$$

Spectral Norm

When the norms on \mathbb{R}^n and \mathbb{R}^m are both 2-norms, the induced norm on $\mathbb{R}^{n \times m}$ is simply called the **2-norm** or **spectral norm**, denoted by $\|\cdot\|_2$.

Proposition.

$$\|\mathbf{A}\|_2 = \sqrt{\lambda_{\max}(\mathbf{A}^T \mathbf{A})},$$

where $\lambda_{\max}(\mathbf{A}^T \mathbf{A})$ is the maximum eigenvalue of $\mathbf{A}^T \mathbf{A}$.

Proof. Let $\|\mathbf{x}\|_2 = 1$. By slide 15 of Lecture 8,

$$\|\mathbf{A}\mathbf{x}\|_2^2 = \mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x} \leq \lambda_{\max}(\mathbf{A}^T \mathbf{A}) \|\mathbf{x}\|_2^2 = \lambda_{\max}(\mathbf{A}^T \mathbf{A}), \quad \forall \mathbf{x} \in \mathbb{R}^n$$

with equality iff \mathbf{x} is an eigenvector of $\mathbf{A}^T \mathbf{A}$ associated with $\lambda_{\max}(\mathbf{A}^T \mathbf{A})$.

Corollary. If \mathbf{A} is symmetric,

$$\|\mathbf{A}\|_2 = \max\{|\lambda_{\max}(\mathbf{A})|, |\lambda_{\min}(\mathbf{A})|\}$$

If $\mathbf{A} \succeq \mathbf{O}$, then $\|\mathbf{A}\|_2 = \lambda_{\max}(\mathbf{A})$.

Examples

Example.

$$\mathbf{A} = \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix}$$

To find the 2-norm,

$$\mathbf{A}^T \mathbf{A} = \begin{pmatrix} 1 & 3 \\ 2 & 4 \end{pmatrix} \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} = \begin{pmatrix} 10 & 14 \\ 14 & 20 \end{pmatrix}$$

$$\|\mathbf{A}\|_2 = \sqrt{\lambda_{\max}(\mathbf{A}^T \mathbf{A})} = \sqrt{15 + \sqrt{221}} \approx 5.465$$

Example.

$$\mathbf{A} = \begin{pmatrix} 1 & 2 \\ 2 & 4 \end{pmatrix} \succeq \mathbf{O}$$

$$\|\mathbf{A}\|_2 = \sqrt{\lambda_{\max}(\mathbf{A}^T \mathbf{A})} = \sqrt{\lambda_{\max}(\mathbf{A}^2)} = \sqrt{\lambda_{\max}^2(\mathbf{A})} = \lambda_{\max}(\mathbf{A}) = 5$$

Convergence Analysis

$\nabla^2 f$ is M -Lipschitz continuous if

$$\|\nabla^2 f(\mathbf{x}) - \nabla^2 f(\mathbf{y})\|_2 \leq M\|\mathbf{x} - \mathbf{y}\|_2, \quad \forall \mathbf{x}, \mathbf{y}$$

Theorem. If f is m -strongly convex, $\nabla^2 f$ is M -Lipschitz continuous, and \mathbf{x}^* is a minimum of f , then the sequence $\{\mathbf{x}_k\}$ produced by Newton's method satisfies

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

Note. The same remarks on slide 7 apply here with $|x_k - x^*|$ replaced by $\|\mathbf{x}_k - \mathbf{x}^*\|$. In particular, if $\|\mathbf{x}_0 - \mathbf{x}^*\| < \frac{2m}{M}$, then

$$\|\mathbf{x}_k - \mathbf{x}^*\| \leq \frac{2m}{M} \left(\frac{M}{2m} \|\mathbf{x}_0 - \mathbf{x}^*\| \right)^{2^k}$$

The proof is also very similar with only minor modifications.

Proof

$$\begin{aligned} & \|\mathbf{x}_{k+1} - \mathbf{x}^*\| \\ &= \|\mathbf{x}_k - \mathbf{x}^* - [\nabla^2 f(\mathbf{x}_k)]^{-1} [\nabla f(\mathbf{x}_k) - \nabla f(\mathbf{x}^*)]\| \end{aligned} \quad (1)$$

$$\leq \|[\nabla^2 f(\mathbf{x}_k)]^{-1}\| \cdot \|\nabla f(\mathbf{x}^*) - \nabla f(\mathbf{x}_k) - \nabla^2 f(\mathbf{x}_k)(\mathbf{x}^* - \mathbf{x}_k)\| \quad (2)$$

$$= \|[\nabla^2 f(\mathbf{x}_k)]^{-1}\| \cdot \left\| \int_0^1 [\nabla^2 f(\mathbf{x}_k + t(\mathbf{x}^* - \mathbf{x}_k)) - \nabla^2 f(\mathbf{x}_k)](\mathbf{x}^* - \mathbf{x}_k) dt \right\| \quad (3)$$

$$\leq \|[\nabla^2 f(\mathbf{x}_k)]^{-1}\| \int_0^1 \|[\nabla^2 f(\mathbf{x}_k + t(\mathbf{x}^* - \mathbf{x}_k)) - \nabla^2 f(\mathbf{x}_k)](\mathbf{x}^* - \mathbf{x}_k)\| dt \quad (4)$$

$$\leq \|[\nabla^2 f(\mathbf{x}_k)]^{-1}\| \int_0^1 \|\nabla^2 f(\mathbf{x}_k + t(\mathbf{x}^* - \mathbf{x}_k)) - \nabla^2 f(\mathbf{x}_k)\| \cdot \|\mathbf{x}^* - \mathbf{x}_k\| dt \quad (5)$$

$$\leq \|[\nabla^2 f(\mathbf{x}_k)]^{-1}\| \int_0^1 Mt \|\mathbf{x}^* - \mathbf{x}_k\|^2 dt \quad (6)$$

$$= \|[\nabla^2 f(\mathbf{x}_k)]^{-1}\| \cdot \frac{M}{2} \|\mathbf{x}^* - \mathbf{x}_k\|^2 \quad (7)$$

$$\leq \frac{M}{2m} \|\mathbf{x}_k - \mathbf{x}^*\|^2 \quad (8)$$

Proof (cont'd)

1. Step (1) uses the Newton updating rule

$$\mathbf{x}_{k+1} = \mathbf{x}_k - [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k)$$

and the optimality condition $\nabla f(\mathbf{x}^*) = \mathbf{0}$.

2. Step (2) applies the compatibility of norms on slide 10 to

$$[\nabla^2 f(\mathbf{x}_k)]^{-1} [\nabla f(\mathbf{x}^*) - \nabla f(\mathbf{x}_k) - \nabla^2 f(\mathbf{x}_k)(\mathbf{x}^* - \mathbf{x}_k)]$$

3. Step (3) applies the Newton-Leibniz formula to the function $\mathbf{h}(t) = \nabla f(\mathbf{x}_k + t(\mathbf{x}^* - \mathbf{x}_k))$,

$$\nabla f(\mathbf{x}^*) - \nabla f(\mathbf{x}_k) = \mathbf{h}(1) - \mathbf{h}(0) = \int_0^1 \mathbf{h}'(t) dt$$

where $\mathbf{h}'(t)$ is given by the chain rule,

$$\mathbf{h}'(t) = \nabla^2 f(\mathbf{x}_k + t(\mathbf{x}^* - \mathbf{x}_k))(\mathbf{x}^* - \mathbf{x}_k)$$

Proof (cont'd)

4. Step (4) uses the following inequality

$$\left\| \int \mathbf{f}(t) dt \right\| \leq \int \|\mathbf{f}(t)\| dt$$

Proof. Let $\mathbf{z} = \int \mathbf{f}(t) dt$.

$$\|\mathbf{z}\|^2 = \mathbf{z}^T \int \mathbf{f}(t) dt \stackrel{(a)}{=} \int \mathbf{z}^T \mathbf{f}(t) dt \stackrel{(b)}{\leq} \int \|\mathbf{z}\| \cdot \|\mathbf{f}(t)\| dt = \|\mathbf{z}\| \int \|\mathbf{f}(t)\| dt,$$

where (a) uses linearity of integration and (b) Cauchy-Schwarz.

5. Step (5) again applies the compatibility of norms on slide 10
6. Step (6) uses the Lipschitz continuity of $\nabla^2 f$
7. Step (7) performs the integration over t
8. Step (8) uses the m -strong convexity of f

$$\|[\nabla^2 f(\mathbf{x}_k)]^{-1}\| = \lambda_{\max}([\nabla^2 f(\mathbf{x}_k)]^{-1}) = \frac{1}{\lambda_{\min}(\nabla^2 f(\mathbf{x}_k))} \leq \frac{1}{m}$$

Contents

1. Analysis of Newton's method

2. Damped Newton's Method

3. Equality Constrained Optimization

Damped Newton's Method

The Newton direction $-\left[\nabla^2 f(\mathbf{x})\right]^{-1} \nabla f(\mathbf{x})$ is a descent direction, but with step size 1, Newton's method does not guarantee $f(\mathbf{x}_{k+1}) < f(\mathbf{x}_k)$.

To ensure $f(\mathbf{x}_{k+1}) < f(\mathbf{x}_k)$, **damped Newton's method** does backtracking line search along the Newton direction.

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: $\mathbf{d} \leftarrow -\left[\nabla^2 f(\mathbf{x})\right]^{-1} \nabla f(\mathbf{x})$
- 4: $t \leftarrow 1$
- 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: **end while**
- 10: **return** \mathbf{x}

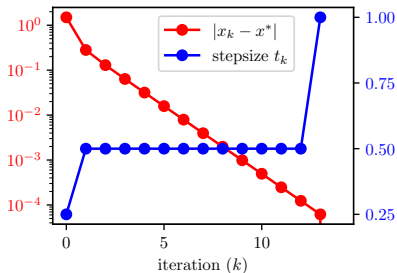
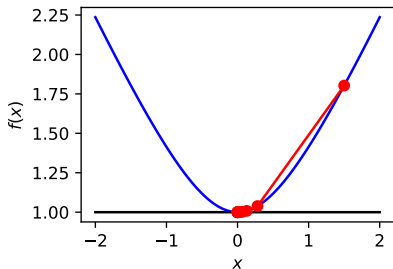
where $\alpha, \beta \in (0, 1)$

Example

$$f(x) = \sqrt{1 + x^2}$$

Recall pure Newton's method converges iff $|x_0| < 1$.

Damped Newton's method converges globally, e.g. for $x_0 = 1.5$.



Convergence Analysis

Theorem. Assume f is m -strongly convex and L -smooth, $\nabla^2 f$ is M -Lipschitz, and \mathbf{x}^* is a minimum of f . Damped Newton's method satisfies the following error bounds

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \begin{cases} f(\mathbf{x}_0) - f(\mathbf{x}^*) - \gamma k, & \text{if } k \leq k_0 \\ \frac{2m^3}{M^2} \left(\frac{1}{2}\right)^{2^{k-k_0+1}}, & \text{if } k > k_0 \end{cases}$$

where $\gamma = 2\alpha\bar{\alpha}\beta\eta^2m/L^2$, $\eta = \min\{1, 3(1 - 2\alpha)\}m^2/M$, and k_0 is the number of steps until $\|\nabla f(\mathbf{x}_{k_0+1})\| \leq \eta$.

Notes.

- Damped Newton's method guarantees **global** convergence.
- To get $f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \epsilon$, we need at most

$$\frac{f(\mathbf{x}_0) - f(\mathbf{x}^*)}{\gamma} + \log_2 \log_2 \frac{\epsilon_0}{\epsilon}$$

where $\epsilon_0 = \frac{2m^3}{M^2}$. It can be slow if γ is small.

Convergence Analysis (cont'd)

Detailed analysis shows that the convergence follows two stages

- **Damped Newton phase.** When $\|\nabla f(\mathbf{x}_k)\| > \eta$, backtracking selects a step size $t_k \leq 1$, and

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k) \leq -\gamma$$

Summing over k from 0 to $k_0 - 1$,

$$f(\mathbf{x}^*) - f(\mathbf{x}_0) \leq f(\mathbf{x}_{k_0}) - f(\mathbf{x}_0) \leq -k_0\gamma \implies k_0 \leq \frac{f(\mathbf{x}_0) - f(\mathbf{x}^*)}{\gamma}$$

- **Pure Newton phase.** When $\|\nabla f(\mathbf{x}_k)\| \leq \eta$, backtracking always selects step size $t_k = 1$, and

$$\|\nabla f(\mathbf{x}_{k+1})\| \leq \frac{M}{2m^2} \|\nabla f(\mathbf{x}_k)\|^2 \leq \frac{1}{2} \|\nabla f(\mathbf{x}_k)\|$$

Once we are in the pure Newton phase, we will remain so.

Contents

1. Analysis of Newton's method
2. Damped Newton's Method
3. Equality Constrained Optimization

Equality Constrained Optimization Problems

Consider the equality constrained convex optimization problem

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & \mathbf{a}_i^T \mathbf{x} = b_i, \quad i = 1, 2, \dots, k \end{aligned}$$

where f is convex with $\text{dom} f = \mathbb{R}^n$. In a more compact form,

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

where $\mathbf{A}^T = (\mathbf{a}_1, \dots, \mathbf{a}_k) \in \mathbb{R}^{n \times k}$, $\mathbf{b} = (b_1, \dots, b_k)^T \in \mathbb{R}^k$.

The feasible set is

$$X = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{A}\mathbf{x} = \mathbf{b}\}$$

We assume $X \neq \emptyset$. We also assume the constraints are independent, i.e. $\text{rank} \mathbf{A} = k$ (What if $\text{rank} \mathbf{A} < k$?)

Optimality Condition

Lemma. Assume f is differentiable. $\mathbf{x}^* \in X$ is optimal iff

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

where $\text{Null}(\mathbf{A}) = \{\mathbf{x} : \mathbf{A}\mathbf{x} = \mathbf{0}\}$ is the **null space** of \mathbf{A} .

Proof. Recall (slide 20 of Lecture 6) $\mathbf{x}^* \in X$ is optimal iff

$$\nabla f(\mathbf{x}^*)^T (\mathbf{x} - \mathbf{x}^*) \geq 0, \quad \forall \mathbf{x} \in X$$

Note $\mathbf{x} \in X$ i.e. $\mathbf{A}\mathbf{x} = \mathbf{b}$ iff $\mathbf{x} - \mathbf{x}^* \in \text{Null}(\mathbf{A})$. The above condition becomes

$$\nabla f(\mathbf{x}^*)^T \mathbf{y} \geq 0, \quad \forall \mathbf{y} \in \text{Null}(\mathbf{A})$$

Note $\mathbf{y} \in \text{Null}(\mathbf{A}) \iff -\mathbf{y} \in \text{Null}(\mathbf{A})$. The condition then reduces to

$$\nabla f(\mathbf{x}^*)^T \mathbf{y} = 0, \quad \forall \mathbf{y} \in \text{Null}(\mathbf{A})$$

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

Optimality Condition (cont'd)

Second Proof. Let $\mathbf{y}_1, \dots, \mathbf{y}_{n-k}$ be a basis of $\text{Null}(\mathbf{A})$. Then $\mathbf{x} \in X$ iff

$$\mathbf{x} = \mathbf{x}^* + \sum_{i=1}^{n-k} z_i \mathbf{y}_i = \mathbf{x}^* + \mathbf{Fz}$$

where $\mathbf{F} = (\mathbf{y}_1, \dots, \mathbf{y}_{n-k})$. Let $g(\mathbf{z}) = f(\mathbf{x}^* + \mathbf{Fz})$. Note \mathbf{x}^* is optimal for the constrained problem (EC) iff $\mathbf{0}$ is an unconstrained minimum of g . By the chain rule, the optimality condition is

$$\nabla g(\mathbf{0}) = \mathbf{F}^T \nabla f(\mathbf{x}^*) = \mathbf{0}$$

or

$$\frac{\partial g(\mathbf{0})}{\partial z_i} = \mathbf{y}_i^T \nabla f(\mathbf{x}^*) = 0, \quad i = 1, \dots, n-k$$

Since $\mathbf{y}_1, \dots, \mathbf{y}_{n-k}$ is a basis of $\text{Null}(\mathbf{A})$,

$$\mathbf{y}^T \nabla f(\mathbf{x}^*) = 0, \quad \forall \mathbf{y} \in \text{Null}(\mathbf{A})$$

Optimality Condition (cont'd)

Theorem. Assume f is differentiable. $\mathbf{x}^* \in X$ is optimal iff there exists $\boldsymbol{\lambda}^* = (\lambda_1^*, \dots, \lambda_k^*)^T \in \mathbb{R}^k$ s.t.

$$\nabla f(\mathbf{x}^*) + \mathbf{A}^T \boldsymbol{\lambda}^* = \mathbf{0},$$

or written out,

$$\nabla f(\mathbf{x}^*) + \sum_{i=1}^k \lambda_i^* \mathbf{a}_i = \mathbf{0}.$$

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

Proof. By the previous lemma, $\mathbf{x}^* \in X$ is optimal iff $\nabla f(\mathbf{x}^*) \perp \text{Null}(\mathbf{A})$.
Since

$$\text{Null}(\mathbf{A})^\perp = \text{Range}(\mathbf{A}^T) \triangleq \{\mathbf{A}^T \mathbf{v} : \mathbf{v} \in \mathbb{R}^k\},$$

\mathbf{x}^* is optimal iff

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

i.e. there exists \mathbf{v}^* s.t. $\nabla f(\mathbf{x}^*) = \mathbf{A}^T \mathbf{v}^* = -\mathbf{A}^T \boldsymbol{\lambda}^*$ with $\boldsymbol{\lambda}^* = -\mathbf{v}^*$.