

# CS257 Linear and Convex Optimization

## Lecture 13

Bo Jiang

John Hopcroft Center for Computer Science  
Shanghai Jiao Tong University

November 30, 2020

## Recap: Equality Constrained Convex Problem

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

where  $\mathbf{A} \in \mathbb{R}^{k \times n}$ ,  $\mathbf{b} \in \mathbb{R}^k$ ,  $f$  is differentiable and convex.

Lagrangian.

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

Lagrange condition.  $\mathbf{x}^*$  is optimal iff  $\exists$  Lagrange multiplier  $\boldsymbol{\lambda}^* \in \mathbb{R}^k$  s.t.

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

Convex QP.  $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{Qx} + \mathbf{g}^T \mathbf{x} + c$ , where  $\mathbf{Q} \succeq \mathbf{O}$ . KKT system

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

# Recap: Newton's Method

Solve an approximate quadratic problem in each iteration.

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

1: initialization  $\mathbf{x} \leftarrow \mathbf{x}_0$  s.t.  $\mathbf{A} \mathbf{x}_0 = \mathbf{b}$

2: **repeat**

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

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

4:     find stepsize  $t$  by backtracking line search

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

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

7: **return**  $\mathbf{x}$

# Contents

1. General Equality Constrained Problems

2. Inequality Constrained Problem

## Optimization on 2D Circle

Let  $\mathbf{x} = (x, y)^T \in \mathbb{R}^2$ . Consider the following nonconvex (why?) problem

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & h(\mathbf{x}) = \|\mathbf{x}\|^2 - 1 = 0 \end{aligned}$$

Parameterize the feasible set by  $\mathbf{x}(t) = (\cos t, \sin t)^T$  and reduce the above constrained problem to the following unconstrained problem

$$\min_t g(t) \triangleq f(\mathbf{x}(t)) = f(\cos t, \sin t)$$

If  $\mathbf{x}^* = \mathbf{x}(t^*)$  is a **local** minimum of the constrained problem, then  $t^*$  is a local minimum of  $g$ , so

$$g'(t^*) = \frac{\partial f(\mathbf{x}^*)}{\partial x} x'(t^*) + \frac{\partial f(\mathbf{x}^*)}{\partial y} y'(t^*) = 0$$

On the other hand,  $h(\mathbf{x}(t)) = 0$ . Differentiating w.r.t.  $t$  at  $t^*$ ,

$$\frac{\partial h(\mathbf{x}^*)}{\partial x} x'(t^*) + \frac{\partial h(\mathbf{x}^*)}{\partial y} y'(t^*) = 0$$

## Optimization on 2D Circle (cont'd)

Combining the previous two equations,

$$\begin{bmatrix} \frac{\partial f(\mathbf{x}^*)}{\partial x} & \frac{\partial f(\mathbf{x}^*)}{\partial y} \\ \frac{\partial h(\mathbf{x}^*)}{\partial x} & \frac{\partial h(\mathbf{x}^*)}{\partial y} \end{bmatrix} \begin{bmatrix} x'(t^*) \\ y'(t^*) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \text{or} \quad \begin{bmatrix} \nabla f(\mathbf{x}^*)^T \\ \nabla h(\mathbf{x}^*)^T \end{bmatrix} \mathbf{x}'(t^*) = \mathbf{0}$$

- The linear system has a solution  $\mathbf{x}'(t^*) = (-\sin t^*, \cos t^*)^T \neq \mathbf{0}$ , so  $\nabla f(\mathbf{x}^*)$  and  $\nabla h(\mathbf{x}^*)$  must be linearly dependent.
- Note  $\nabla h(\mathbf{x}^*) = \mathbf{x}^* \neq \mathbf{0}$  (why?), so there exists  $\lambda^*$  s.t.

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

Define the **Lagrangian** by

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

**Lagrange condition.**  $\mathbf{x}^*$  is a local optimum **only if** there exists  $\lambda^*$  s.t.

$$\begin{cases} \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}^*, \lambda^*) = \nabla f(\mathbf{x}^*) + \lambda^* \nabla h(\mathbf{x}^*) \\ \nabla_{\lambda} \mathcal{L}(\mathbf{x}^*, \lambda^*) = h(\mathbf{x}^*) = 0 \end{cases}$$

**Note.** This is only a necessary condition for nonconvex problems.

## Example

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) = x + 2y \\ \text{s.t.} \quad & h(\mathbf{x}) = \|\mathbf{x}\|^2 - 1 = 0 \end{aligned}$$

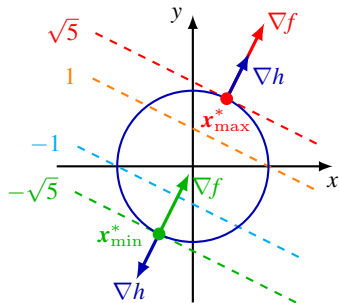
- Lagrange condition

$$\begin{cases} \frac{\partial f(\mathbf{x})}{\partial x} + \lambda \frac{\partial h(\mathbf{x})}{\partial x} = 1 + 2\lambda x = 0 \implies x = -\frac{1}{2\lambda} \\ \frac{\partial f(\mathbf{x})}{\partial y} + \lambda \frac{\partial h(\mathbf{x})}{\partial y} = 2 + 2\lambda y = 0 \implies y = -\frac{1}{\lambda} \\ h(\mathbf{x}^*) = x^2 + y^2 - 1 = 0 \end{cases}$$

- solutions to the above equations

$$(1) \begin{cases} x = -\frac{\sqrt{5}}{5} \\ y = -\frac{2\sqrt{5}}{5} \\ \lambda = \frac{\sqrt{5}}{2} \end{cases} \quad (2) \begin{cases} x = \frac{\sqrt{5}}{5} \\ y = \frac{2\sqrt{5}}{5} \\ \lambda = -\frac{\sqrt{5}}{2} \end{cases}$$

- (1) global minimum, (2) global maximum
- at all extrema,  $\nabla f \parallel \nabla h \perp X$



## Example

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) = x^2 - y \\ \text{s.t.} \quad & h(\mathbf{x}) = \|\mathbf{x}\|^2 - 1 = 0 \end{aligned}$$

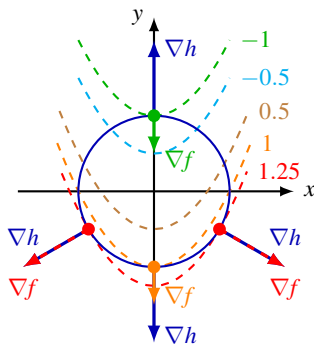
- Lagrange condition

$$\begin{cases} \frac{\partial f(\mathbf{x})}{\partial x} + \lambda \frac{\partial h(\mathbf{x})}{\partial x} = 2x + 2\lambda x = 0 \\ \frac{\partial f(\mathbf{x})}{\partial y} + \lambda \frac{\partial h(\mathbf{x})}{\partial y} = -1 + 2\lambda y = 0 \\ h(\mathbf{x}^*) = x^2 + y^2 - 1 = 0 \end{cases}$$

- solutions to above equations

$$(1) \begin{cases} x = 0 \\ y = 1 \\ \lambda = \frac{1}{2} \end{cases} \quad (2) \begin{cases} x = 0 \\ y = -1 \\ \lambda = -\frac{1}{2} \end{cases} \quad (3) \begin{cases} x = \frac{\sqrt{3}}{2} \\ y = -\frac{1}{2} \\ \lambda = -1 \end{cases} \quad (4) \begin{cases} x = -\frac{\sqrt{3}}{2} \\ y = -\frac{1}{2} \\ \lambda = -1 \end{cases}$$

- (1) global minimum, (2) local minimum, (3)(4) global maxima
- at all extrema,  $\nabla f \parallel \nabla h \perp X$



**Exercise.** Solve equivalent problem  $g(y) = 1 - y^2 - y$  s.t.  $|y| \leq 1$ .



## Implicit Function Theorem in 2D

The derivation on slides 4-5 works for general  $h$  of two variables, as long as we can parameterize the feasible set in a neighborhood of  $\mathbf{x}^*$  by  $\mathbf{x}(t)$ , i.e.  $h(\mathbf{x}(t)) = 0$ , s.t.  $\mathbf{x}'(t^*) \neq \mathbf{0}$  and  $\nabla h(\mathbf{x}^*) \neq \mathbf{0}$ . The Implicit Function Theorem guarantees this is possible if  $\nabla h(\mathbf{x}^*) \neq \mathbf{0}$ .

**Implicit Function Theorem.** If  $F(x, y)$  is continuously differentiable in a neighborhood of  $(x_0, y_0)$ , and satisfies

$$F(x_0, y_0) = 0, \quad \frac{\partial F(x_0, y_0)}{\partial y} \neq 0$$

then there exists a continuously differentiable function  $y = \phi(x)$  defined in a neighborhood of  $x_0$  s.t.

$$F(x, \phi(x)) = 0, \quad \phi'(x) = - \left[ \frac{\partial F(x, \phi(x))}{\partial y} \right]^{-1} \frac{\partial F(x, \phi(x))}{\partial x}$$

# Implicit Function Theorem and Parameterization

If  $\nabla h(x_0, y_0) \neq \mathbf{0}$ , then either  $\frac{\partial h(x_0, y_0)}{\partial x} \neq 0$  or  $\frac{\partial h(x_0, y_0)}{\partial y} \neq 0$ .

- If  $\frac{\partial h(x_0, y_0)}{\partial y} \neq 0$ , we can parameterize the feasible set by  $t = x$ ,

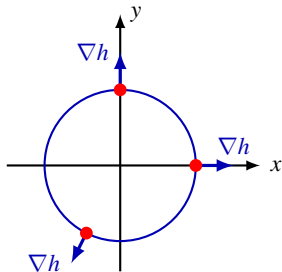
$$\mathbf{x}(t) = (t, \phi(t))^T \quad \text{with} \quad \mathbf{x}'(t) = (1, \phi'(t))^T \neq \mathbf{0}$$

- If  $\frac{\partial h(x_0, y_0)}{\partial x} \neq 0$ , we can parameterize the feasible set by  $t = y$ ,

$$\mathbf{x}(t) = (\psi(t), t)^T \quad \text{with} \quad \mathbf{x}'(t) = (\psi'(t), 1)^T \neq \mathbf{0}$$

**Example.** For  $h(\mathbf{x}) = \|\mathbf{x}\|^2 - 1$ .

- at  $\mathbf{x}_0 = (1, 0)^T$ , use  $\mathbf{x}(t) = (\sqrt{1-t^2}, t)^T$
- at  $\mathbf{x}_0 = (0, 1)^T$ , use  $\mathbf{x}(t) = (t, \sqrt{1-t^2})^T$
- at  $\mathbf{x}_0$  in the 3rd quadrant, we can use  $\mathbf{x}(t) = (t, -\sqrt{1-t^2})^T$  or  $\mathbf{x}(t) = (-\sqrt{1-t^2}, t)^T$



## First-order Necessary Condition in 2D

A point  $\mathbf{x}$  is called a **regular point** of a function  $h$  if  $\nabla h(\mathbf{x}) \neq \mathbf{0}$ ; otherwise it is called a **critical point**.

**Theorem.** If  $\mathbf{x}^*$  is a local extremum (maximum or minimum) of  $f$  s.t.  $h(\mathbf{x}) = 0$ , and  $\mathbf{x}^*$  is a **regular** point of  $h$ , then there exists  $\lambda^*$  s.t.

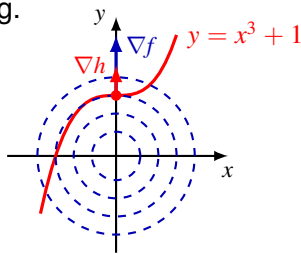
$$\nabla f(\mathbf{x}^*) + \lambda^* \nabla h(\mathbf{x}^*) = \mathbf{0}$$

**Note.**  $\mathbf{x}^*$  satisfying the above Lagrange condition may be neither a maximum nor a minimum. E.g.

$$\begin{aligned} f(\mathbf{x}) &= \|\mathbf{x}\|^2 \\ h(\mathbf{x}) &= y - x^3 - 1 \end{aligned}$$

At  $\mathbf{x}^* = (0, 1)^T$ ,

$$\nabla f(\mathbf{x}^*) = (0, 2)^T, \quad \nabla h(\mathbf{x}^*) = (0, 1)^T$$



Second-order conditions can help distinguish different cases ([CZ, LY])

## Critical Points

The Lagrange condition may fail at critical points.

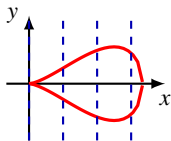
**Example.**

$$\begin{aligned} \min_{x,y} \quad & f(x, y) = x + y \\ \text{s. t.} \quad & h(x, y) = x^2 + y^2 = 0 \end{aligned}$$

The feasible set is  $X = \{\mathbf{0}\}$ , so  $\mathbf{x}^* = \mathbf{0}$  is the global minimum. There is no  $\lambda^* \in \mathbb{R}$  satisfying the Lagrange condition  $\nabla f(\mathbf{x}^*) + \lambda^* \nabla h(\mathbf{x}^*) = \mathbf{0}$ , as  $\nabla f(\mathbf{x}^*) = (1, 1)^T$  and  $\nabla h(\mathbf{x}^*) = \mathbf{0}$ .

**Example.**

$$\begin{aligned} \min_{x,y} \quad & f(x, y) = x \\ \text{s. t.} \quad & h(x, y) = y^2 + x^4 - x^3 = 0 \end{aligned}$$



Note  $x^3 - x^4 = y^2 \geq 0$  implies  $x \in [0, 1]$ , so  $\mathbf{x}^* = \mathbf{0}$  is the global minimum. Lagrange condition fails as  $\nabla f(\mathbf{x}^*) = (1, 0)^T$ ,  $\nabla h(\mathbf{x}^*) = \mathbf{0}$ .

**Note.** To find the minimum, we need to check both regular points satisfying the Lagrange condition and feasible critical points.

# First-order Necessary Condition

Let  $\mathbf{x} \in \mathbb{R}^n$  and  $n > k$ . Consider the equality constrained problem

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & h_i(\mathbf{x}) = 0, \quad i = 1, 2, \dots, k \end{aligned} \tag{ECP}$$

A point  $\mathbf{x}$  is a **regular point** of  $\mathbf{h} = (h_1, \dots, h_k)^T$  if  $\nabla h_1(\mathbf{x}), \dots, \nabla h_k(\mathbf{x})$  are linearly independent; otherwise it is a **critical point** of  $\mathbf{h}$ .

**Theorem.** If  $\mathbf{x}^*$  is a local extremum of  $f$  s.t.  $\mathbf{h}(\mathbf{x}) = \mathbf{0}$ , and  $\mathbf{x}^*$  is a regular point of  $\mathbf{h}$ , then there exist **Lagrange multipliers**  $\lambda_1^*, \dots, \lambda_k^* \in \mathbb{R}$  s.t.

$$\nabla f(\mathbf{x}^*) + (\boldsymbol{\lambda}^*)^T \mathbf{h}(\mathbf{x}^*) = \nabla f(\mathbf{x}^*) + \sum_{i=1}^k \lambda_i^* \nabla h_i(\mathbf{x}^*) = \mathbf{0}$$

Define the **Lagrangian** of (ECP) by

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \boldsymbol{\lambda}^T \mathbf{h}(\mathbf{x}) = f(\mathbf{x}) + \sum_{i=1}^k \lambda_i h_i(\mathbf{x})$$

Then the **Lagrange condition** is  $\nabla \mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*) = \mathbf{0}$ .

## Appendix: Implicit Function Theorem

Write  $F : \mathbb{R}^{n+k} \rightarrow \mathbb{R}^k$  as  $F(\mathbf{x}, \mathbf{y})$  with  $\mathbf{x} \in \mathbb{R}^n$  and  $\mathbf{y} \in \mathbb{R}^k$ . Let  $F = (F_1, \dots, F_k)^T$ , and

$$\frac{\partial F}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial F_1}{\partial x_1} & \cdots & \frac{\partial F_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial F_k}{\partial x_1} & \cdots & \frac{\partial F_k}{\partial x_n} \end{bmatrix}, \quad \frac{\partial F}{\partial \mathbf{y}} = \begin{bmatrix} \frac{\partial F_1}{\partial y_1} & \cdots & \frac{\partial F_1}{\partial y_k} \\ \vdots & \ddots & \vdots \\ \frac{\partial F_k}{\partial y_1} & \cdots & \frac{\partial F_k}{\partial y_k} \end{bmatrix}$$

**Implicit Function Theorem.** If  $F : \mathbb{R}^{n+k} \rightarrow \mathbb{R}^k$  is continuously differentiable in a neighborhood  $(\mathbf{x}_0, \mathbf{y}_0)$ , and satisfies

$$F(\mathbf{x}_0, \mathbf{y}_0) = \mathbf{0}, \quad \det \frac{\partial F(\mathbf{x}_0, \mathbf{y}_0)}{\partial \mathbf{y}} \neq 0,$$

then there exists continuously differentiable function  $\mathbf{y} = \phi(\mathbf{x})$  defined in a neighborhood of  $\mathbf{x}_0$  s.t.

$$F(\mathbf{x}, \phi(\mathbf{x})) = 0, \quad \frac{\partial \phi(\mathbf{x})}{\partial \mathbf{x}} = - \left[ \frac{\partial F(\mathbf{x}, \phi(\mathbf{x}))}{\partial \mathbf{y}} \right]^{-1} \frac{\partial F(\mathbf{x}, \phi(\mathbf{x}))}{\partial \mathbf{x}}$$

## Appendix: Proof of Lagrange Condition

Let  $\mathbf{h} = (h_1, \dots, h_k)^T$ . Note that the Jacobian matrix of  $\mathbf{h}$  is

$$\frac{\partial \mathbf{h}(\mathbf{x}^*)}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial h_1(\mathbf{x}^*)}{\partial x_1} & \cdots & \frac{\partial h_1(\mathbf{x}^*)}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial h_k(\mathbf{x}^*)}{\partial x_1} & \cdots & \frac{\partial h_k(\mathbf{x}^*)}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \nabla h_1(\mathbf{x}^*)^T \\ \vdots \\ \nabla h_k(\mathbf{x}^*)^T \end{bmatrix}$$

Since  $\mathbf{x}^*$  is regular,  $\text{rank} \frac{\partial \mathbf{h}(\mathbf{x}^*)}{\partial \mathbf{x}} = k$ . By re-indexing  $x_1, \dots, x_n$  if necessary, we assume the last  $m$  columns are linearly independent.

Let  $\mathbf{y} = (x_1, \dots, x_{n-k})^T$ ,  $\mathbf{z} = (x_{n-k+1}, \dots, x_n)^T$ .

By the Implicit Function Theorem, there is a continuously differentiable function  $\mathbf{z} = \phi(\mathbf{y})$  s.t.  $\mathbf{h}(\mathbf{y}, \phi(\mathbf{y})) = \mathbf{0}$ , i.e. we can parameterize the feasible set  $X$  by<sup>1</sup>

$$\mathbf{x}(\mathbf{y}) = (\mathbf{y}, \phi(\mathbf{y}))$$

(ECP) reduces to

$$\min_{\mathbf{y}} g(\mathbf{y}) = f(\mathbf{y}, \phi(\mathbf{y}))$$

---

<sup>1</sup>we are sloppy about the shape here, but it should not cause any confusion.

## Proof (cont'd)

Since  $\mathbf{x}^* = \mathbf{x}(\mathbf{y}^*) = (\mathbf{y}^*, \phi(\mathbf{y}^*))$  is a local extremum of (ECP),  $\mathbf{y}^*$  is a local extreme of  $g$ . Recalling  $y_\ell = x_\ell$  for  $\ell = 1, \dots, n - k$ ,

$$\frac{\partial g(\mathbf{y}^*)}{\partial y_\ell} = \frac{\partial f(\mathbf{x}^*)}{\partial x_\ell} + \sum_{j=n-k+1}^n \frac{\partial f(\mathbf{x}^*)}{\partial x_j} \frac{\partial \phi_{j-(n-k)}(\mathbf{y}^*)}{\partial y_\ell} = 0, \quad \ell = 1, \dots, n - k$$

Differentiating  $h_i(\mathbf{x}(\mathbf{y})) = h_i(\mathbf{y}, \phi(\mathbf{y})) = 0$  at  $\mathbf{y}^*$ ,

$$\frac{\partial h_i(\mathbf{x}^*)}{\partial y_\ell} = \frac{\partial h_i(\mathbf{x}^*)}{\partial y_\ell} + \sum_{j=n-k+1}^n \frac{\partial h_j(\mathbf{x}^*)}{\partial x_j} \frac{\partial \phi_{j-(n-k)}(\mathbf{y}^*)}{\partial y_\ell} = 0, \quad \ell = 1, \dots, n - k$$

In matrix form,

$$\begin{bmatrix} 1 & \cdots & 0 & \frac{\partial \phi_1(\mathbf{x}^*)}{\partial y_1} & \cdots & \frac{\partial \phi_k(\mathbf{x}^*)}{\partial y_1} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & 1 & \frac{\partial \phi_1(\mathbf{x}^*)}{\partial y_{n-k}} & \cdots & \frac{\partial \phi_k(\mathbf{x}^*)}{\partial y_{n-k}} \end{bmatrix} \begin{bmatrix} \frac{\partial f(\mathbf{x}^*)}{\partial x_1} & \frac{\partial h_1(\mathbf{x}^*)}{\partial x_1} & \cdots & \frac{\partial h_k(\mathbf{x}^*)}{\partial x_1} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f(\mathbf{x}^*)}{\partial x_n} & \frac{\partial h_1(\mathbf{x}^*)}{\partial x_n} & \cdots & \frac{\partial h_k(\mathbf{x}^*)}{\partial x_n} \end{bmatrix} = \mathbf{0}$$



## Proof (cont'd)

The matrix equation takes the form

$$\left[ \mathbf{I}_{n-k} \quad \frac{\partial \phi(\mathbf{x}^*)}{\partial \mathbf{y}}^T \right] \left[ \nabla f(\mathbf{x}^*) \quad \nabla h_1(\mathbf{x}^*) \quad \dots \quad \nabla h_k(\mathbf{x}^*) \right] = \mathbf{0}_{(n-k) \times (k+1)}$$

meaning  $\nabla f(\mathbf{x}^*), \nabla h_1(\mathbf{x}^*), \dots, \nabla h_k(\mathbf{x}^*)$  are all in  $\text{Null}(\mathbf{A})$ , where

$$\mathbf{A} = \left[ \mathbf{I}_{n-k} \quad \frac{\partial \phi(\mathbf{x}^*)}{\partial \mathbf{y}}^T \right] \in \mathbb{R}^{(n-k) \times n}.$$

Note  $\dim \text{Null}(\mathbf{A}) = k$ , so  $\nabla f(\mathbf{x}^*), \nabla h_1(\mathbf{x}^*), \dots, \nabla h_k(\mathbf{x}^*)$  are linearly dependent. But  $\nabla h_1(\mathbf{x}^*), \dots, \nabla h_k(\mathbf{x}^*)$  are linearly independent, since  $\mathbf{x}^*$  is a regular point. Thus there exist  $\lambda_1^*, \dots, \lambda_k^* \in \mathbb{R}$  s.t.

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

## Example

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^3} \quad & f(\mathbf{x}) = x_1 + 2x_2 + x_3 \\ \text{s.t.} \quad & h_1(\mathbf{x}) = x_1 + x_2 + 2x_3 = 0 \\ & h_2(\mathbf{x}) = \|\mathbf{x}\|^2 - 1 = 0 \end{aligned}$$

A critical point  $\mathbf{x}$  satisfies  $\nabla h_2(\mathbf{x}) \parallel \nabla h_1(\mathbf{x})$ , so  $\mathbf{x} \propto (1, 1, 2)^T$ , infeasible.

The Lagrangian is

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = x_1 + 2x_2 + x_3 + \lambda_1(x_1 + x_2 + 2x_3) + \lambda_2(x_1^2 + x_2^2 + x_3^2 - 1)$$

The Lagrange condition is

$$\begin{cases} \partial_{x_1} \mathcal{L} = 1 + \lambda_1 + 2\lambda_2 x_1 = 0 & (1) \\ \partial_{x_2} \mathcal{L} = 2 + \lambda_1 + 2\lambda_2 x_2 = 0 & (2) \\ \partial_{x_3} \mathcal{L} = 1 + 2\lambda_1 + 2\lambda_2 x_3 = 0 & (3) \\ \partial_{\lambda_1} \mathcal{L} = x_1 + x_2 + 2x_3 = 0 & (4) \\ \partial_{\lambda_2} \mathcal{L} = x_1^2 + x_2^2 + x_3^2 - 1 = 0 & (5) \end{cases}$$

## Example (cont'd)

- $(1)+(2)+(3) \times 2$ ,

$$5 + 6\lambda_1 + 2\lambda_2(x_1 + x_2 + 2x_3) = 0 \quad (6)$$

- Plugging (4) into (6) yields  $\lambda_1 = -\frac{5}{6}$ .
- Plugging  $\lambda_1$  into (1)(2)(3), and noting that  $\lambda_2 \neq 0$ ,

$$x_1 = -\frac{1}{12\lambda_2}, \quad x_2 = -\frac{7}{12\lambda_2}, \quad x_3 = \frac{1}{3\lambda_2} \quad (7)$$

- Plugging (8) into (5) yields  $\lambda_2 = \pm\sqrt{\frac{33}{72}}$ , so

$$(1) \begin{cases} x_1 = -\frac{1}{\sqrt{66}} \\ x_2 = -\frac{7}{\sqrt{66}} \\ x_3 = \frac{4}{\sqrt{66}} \\ \lambda_1 = -\frac{5}{6} \\ \lambda_2 = \sqrt{\frac{33}{72}} \end{cases} \quad \text{or} \quad (2) \begin{cases} x_1 = \frac{1}{\sqrt{66}} \\ x_2 = \frac{7}{\sqrt{66}} \\ x_3 = -\frac{4}{\sqrt{66}} \\ \lambda_1 = -\frac{5}{6} \\ \lambda_2 = -\sqrt{\frac{33}{72}} \end{cases}$$

- (1) global minimum, (2) global maximum

# Contents

1. General Equality Constrained Problems

2. Inequality Constrained Problem

## Active and Inactive Constraints

Let  $\mathbf{x} \in \mathbb{R}^n$  and  $n > k$ . Consider

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & h_i(\mathbf{x}) = 0, \quad i = 1, 2, \dots, k \\ & g_j(\mathbf{x}) \leq 0, \quad j = 1, 2, \dots, m \end{aligned} \tag{ICP}$$

We do **not** assume it is a convex problem. Assume the domain is  $\mathbb{R}^n$ .  
The feasible set is

$$X = \{\mathbf{x} : h_i(\mathbf{x}) = 0, 1 \leq i \leq k; g_j(\mathbf{x}) \leq 0, 1 \leq j \leq m\}$$

Let  $\mathbf{x}_0 \in X$ . The  $j$ -th inequality constraint  $g_j(\mathbf{x}) \leq 0$  is called **active** at  $\mathbf{x}_0$  if  $g_j(\mathbf{x}_0) = 0$ , and **inactive** at  $\mathbf{x}_0$  if  $g_j(\mathbf{x}_0) < 0$ . Denote by  $J(\mathbf{x}_0)$  the set of indices of the active inequality constraints at  $\mathbf{x}_0$ ,

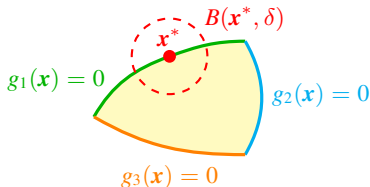
$$J(\mathbf{x}_0) = \{j : g_j(\mathbf{x}_0) = 0\}$$

By convention, equality constraints are considered active at all  $\mathbf{x} \in X$ .

# Reduction to Equality Constrained Problem

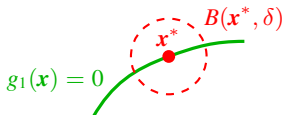
Suppose  $\mathbf{x}^*$  is a local minimum of (ICP). It is the solution to

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & h_i(\mathbf{x}) = 0, \quad i = 1, 2, \dots, k \\ & g_j(\mathbf{x}) \leq 0, \quad j = 1, 2, \dots, m \\ & \mathbf{x} \in B(\mathbf{x}^*, \delta) \end{aligned}$$



for some small enough  $\delta$ . It is equivalent to

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & h_i(\mathbf{x}) = 0, \quad i = 1, 2, \dots, k \\ & g_j(\mathbf{x}) = 0, \quad j \in J(\mathbf{x}^*) \\ & \mathbf{x} \in B(\mathbf{x}^*, \delta) \end{aligned}$$



If it is known a priori which constraints are active at  $\mathbf{x}^*$ , we can find  $\mathbf{x}^*$  by solving the above equality constrained problem.

## Reduction to Equality Constrained Problem (cont'd)

A local minimum  $\mathbf{x}^*$  of (ICP) is also a local minimum of the following

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & h_i(\mathbf{x}) = 0, \quad i = 1, 2, \dots, k \\ & g_j(\mathbf{x}) = 0, \quad j \in J(\mathbf{x}^*) \end{aligned}$$

$\mathbf{x}^* \in X$  is a **regular point** if  $\nabla h_i(\mathbf{x}^*)$ ,  $1 \leq i \leq k$  and  $\nabla g_j(\mathbf{x}^*)$ ,  $j \in J(\mathbf{x}^*)$  are linearly independent.

At a regular local minimum, Lagrange condition yields

$$\nabla f(\mathbf{x}^*) + \sum_{i=1}^k \lambda_i^* \nabla h_i(\mathbf{x}^*) + \sum_{j \in J(\mathbf{x}^*)} \mu_j^* \nabla g_j(\mathbf{x}^*) = \mathbf{0}$$

Setting  $\mu_j^* = 0$  for inactive constraints, i.e.  $j \notin J(\mathbf{x}^*)$ ,

$$\nabla f(\mathbf{x}^*) + \sum_{i=1}^k \lambda_i^* \nabla h_i(\mathbf{x}^*) + \sum_{j=1}^m \mu_j^* \nabla g_j(\mathbf{x}^*) = \mathbf{0}$$

# Karush-Kuhn-Tucker (KKT) Conditions

**Theorem.** If  $\mathbf{x}^*$  is a local minimum of (ICP) and also a regular point, then there exist **Lagrange multipliers**<sup>2</sup>  $\lambda_1^*, \dots, \lambda_k^*, \mu_1^*, \dots, \mu_m^* \in \mathbb{R}$  s.t. the following KKT conditions hold,

1.  $\mu_j^* \geq 0, j = 1, 2, \dots, m$
2.  $\nabla f(\mathbf{x}^*) + \sum_{i=1}^k \lambda_i^* \nabla h_i(\mathbf{x}^*) + \sum_{j=1}^m \mu_j^* \nabla g_j(\mathbf{x}^*) = \mathbf{0}$
3.  $\mu_j^* g_j(\mathbf{x}^*) = 0, j = 1, 2, \dots, m$

**Note.** Condition 2 says  $\nabla_x \mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) = \mathbf{0}$  for the **Lagrangian**

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{x}) + \boldsymbol{\lambda}^T \mathbf{h}(\mathbf{x}) + \boldsymbol{\mu}^T \mathbf{g}(\mathbf{x}) = f(\mathbf{x}) + \sum_{i=1}^k \lambda_i h_i(\mathbf{x}) + \sum_{j=1}^m \mu_j g_j(\mathbf{x})$$

**Note.** Condition 3 is called **complementary slackness condition**, as it, together with 1 and  $g_j(\mathbf{x}) \leq 0$ , implies either  $\mu_j^* = 0$  or  $g_j(\mathbf{x}^*) = 0$ .

---

<sup>2</sup>Sometimes also called KKT multipliers. Sometimes  $\lambda_i$  are called Lagrange multipliers while  $\mu_j$  are called KKT multipliers.



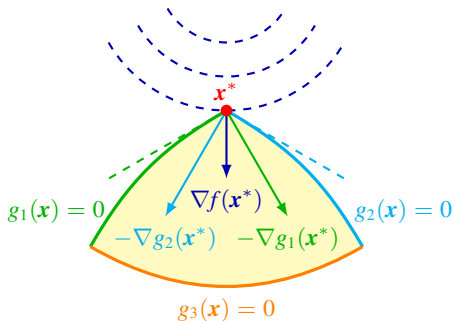
# Geometric Interpretation

Let  $\mathbf{x} \in \mathbb{R}^2$ . Consider

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & g_j(\mathbf{x}) \leq 0, \quad j = 1, 2, 3 \end{aligned}$$

Suppose  $\mathbf{x}^*$  is a local minimum and only  $g_1$  and  $g_2$  are active at  $\mathbf{x}^*$ . The KKT condition says  $\mu_1^* \geq 0$ ,  $\mu_2^* \geq 0$ ,  $\mu_3^* = 0$  and

$$\nabla f(\mathbf{x}^*) = -\mu_1^* \nabla g_1(\mathbf{x}^*) - \mu_2^* \nabla g_2(\mathbf{x}^*)$$



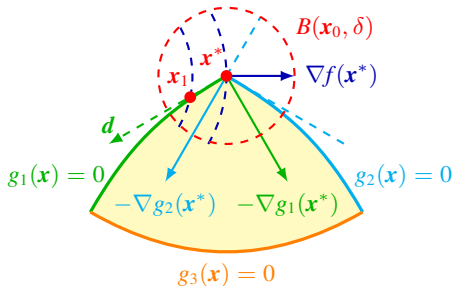
## Geometric Interpretation (cont'd)

Why  $\mu_j^* \geq 0$ ? Assume  $\mu_2^* < 0$  and we show a contradiction.

- Let  $\mathbf{d}$  be a tangent vector of  $g_1(\mathbf{x}) = 0$  at  $\mathbf{x}^*$ , so  $\mathbf{d} \perp \nabla g_1(\mathbf{x}^*)$ .
- $\mathbf{d}^T \nabla g_2(\mathbf{x}^*) \neq 0$ ; otherwise,  $\mathbf{d} \perp \nabla g_2(\mathbf{x}^*)$ , so  $\nabla g_1(\mathbf{x}^*) \parallel \nabla g_2(\mathbf{x}^*)$ , contradicting the regularity of  $\mathbf{x}^*$ .
- Replacing  $\mathbf{d}$  by  $-\mathbf{d}$  if necessary, we can assume  $\mathbf{d}^T \nabla g_2(\mathbf{x}^*) < 0$ .
- Move along the curve  $g_1(\mathbf{x}) = 0$  in the direction of  $\mathbf{d}$  from  $\mathbf{x}^*$  to  $\mathbf{x}_1$ .

$$\mathbf{d}^T \nabla f(\mathbf{x}^*) = \mathbf{d}^T [-\mu_1 \nabla g_1(\mathbf{x}^*) - \mu_2 \nabla g_2(\mathbf{x}^*)] = -\mu_2^* \mathbf{d}^T \nabla g_2(\mathbf{x}^*) < 0.$$

For a small move,  $f(\mathbf{x}_1) < f(\mathbf{x}^*)$ , contradicting minimality of  $f(\mathbf{x}^*)$ .



## Appendix: Proof for $\mu \geq 0$

Suppose  $\mu_{j_0}^* < 0$  for some  $j_0$ . Let  $J'(\mathbf{x}^*) = J(\mathbf{x}^*) \setminus \{j_0\}$ , and  $S$  the set determined by all active constraints other than  $g_{j_0}$ ,

$$S = \{\mathbf{x} : h_i(\mathbf{x}) = 0, i = 1, 2, \dots, k; g_j(\mathbf{x}) = 0, j \in J'(\mathbf{x}^*)\}$$

We will show we can move away from  $\mathbf{x}^*$  on  $S$  so that feasibility is maintained but  $f$  decreases, contradicting the minimality of  $\mathbf{x}^*$ .

1. There exists a direction  $\mathbf{d}_0$  tangent to  $S$  s.t.  $\nabla g_{j_0}(\mathbf{x}^*)^T \mathbf{d}_0 < 0$
2. KKT then implies  $\nabla f(\mathbf{x}^*)^T \mathbf{d}_0 < 0$
3. By Implicit Function Theorem, there exists a curve  $\mathbf{x}(t) \subset S$  s.t.  $\mathbf{x}(0) = \mathbf{x}^*$ ,  $\mathbf{x}'(0) = \mathbf{d}_0$ . Thus  $g_j(\mathbf{x}(t)) = 0$  for  $j \in J'(\mathbf{x}^*)$ .
4. By continuity,  $g_j(\mathbf{x}(t)) < 0$  for small  $t$  and  $j \notin J(\mathbf{x}^*)$
5. By the chain rule,

$$\left. \frac{d}{dt} g_{j_0}(\mathbf{x}(t)) \right|_{t=0} = \nabla g_{j_0}(\mathbf{x}^*)^T \mathbf{x}'(0) = \nabla g_{j_0}(\mathbf{x}^*)^T \mathbf{d}_0 < 0$$

For small  $t > 0$ ,  $g_{j_0}(\mathbf{x}(t)) < g_{j_0}(\mathbf{x}^*) = 0$ . Similarly,  $f(\mathbf{x}(t)) < f(\mathbf{x}^*)$ .

## Proof for $\mu \geq 0$ (cont'd)

1. There exists a direction  $\mathbf{d}_0$  tangent to  $S$  s.t.  $\nabla g_{j_0}(\mathbf{x}^*)^T \mathbf{d}_0 < 0$

Proof.

- Let  $\mathbf{A}$  be a matrix whose columns are  $\nabla h_i(\mathbf{x}^*)$  and  $\nabla g_j(\mathbf{x}^*)$ , i.e.

$$\mathbf{A} = [\nabla h_i(\mathbf{x}^*), i = 1, \dots, k; \nabla g_j(\mathbf{x}^*), j \in J'(\mathbf{x}^*)]$$

- The tangent “plane” (or more precisely, tangent space) of  $S$  at  $\mathbf{x}^*$  is

$$T(\mathbf{x}^*) = \text{Null}(\mathbf{A}^T) = \{\mathbf{d} : \nabla h_i(\mathbf{x}^*)^T \mathbf{d} = 0, \forall i; \nabla g_j(\mathbf{x}^*)^T \mathbf{d} = 0, j \in J'(\mathbf{x}^*)\}$$

- By regularity of  $\mathbf{x}^*$ ,

$$\nabla g_{j_0}(\mathbf{x}^*) \notin \text{span}\{\nabla h_i(\mathbf{x}^*), \forall i; \nabla g_j(\mathbf{x}^*), j \in J'(\mathbf{x}^*)\} = \text{Range}(\mathbf{A}) = \text{Null}(\mathbf{A}^T)^\perp$$

so there exists  $\mathbf{d}_0 \in T(\mathbf{x}^*)$  s.t.  $\nabla g_{j_0}(\mathbf{x}^*)^T \mathbf{d}_0 \neq 0$ .

- Replacing  $\mathbf{d}_0$  by  $-\mathbf{d}_0 \in T(\mathbf{x}^*)$  if necessary, we have

$$\nabla g_{j_0}(\mathbf{x}^*)^T \mathbf{d}_0 < 0$$

## Proof for $\mu \geq \mathbf{0}$ (cont'd)

2. For  $\mathbf{d}_0$  given by step 1,  $\nabla f(\mathbf{x}^*)^T \mathbf{d}_0 < 0$

Proof.

- By KKT Conditions 2 and 3

$$\nabla f(\mathbf{x}^*) = - \sum_{i=1}^k \lambda_i^* \nabla h_i(\mathbf{x}^*) - \sum_{j \in J(\mathbf{x}^*)} \mu_j^* \nabla g_j(\mathbf{x}^*) - \underbrace{\sum_{j \notin J(\mathbf{x}^*)} \mu_j^* \nabla g_j(\mathbf{x}^*)}_{=0}$$

- Since  $\mathbf{d}_0 \in \text{Null}(\mathbf{A}^T)$ ,  $\mu_{j_0}^* < 0$ ,

$$\begin{aligned} \nabla f(\mathbf{x}^*)^T \mathbf{d}_0 &= - \sum_{i=1}^k \lambda_i^* \underbrace{\nabla h_i(\mathbf{x}^*)^T \mathbf{d}_0}_{=0} - \sum_{j \in J'(\mathbf{x}^*)} \mu_j^* \underbrace{\nabla g_j(\mathbf{x}^*)^T \mathbf{d}_0}_{=0} - \mu_{j_0}^* \underbrace{\nabla g_{j_0}(\mathbf{x}^*)^T \mathbf{d}_0}_{<0} \\ &< 0 \end{aligned}$$

## Proof for $\mu \geq \mathbf{0}$ (cont'd)

3. There exists a curve  $\mathbf{x}(t) \subset S$  s.t.  $\mathbf{x}(0) = \mathbf{x}^*$ ,  $\mathbf{x}'(0) = \mathbf{d}_0$

**Proof.** For notational simplicity, denote  $g_j, j \in J'(\mathbf{x}^*)$  by  $h_{k+1}, \dots, h_K$ .

- Define  $\tilde{\mathbf{x}}(t, \boldsymbol{\alpha}) = \mathbf{x}^* + t\mathbf{d}_0 + \sum_{i=1}^K \alpha_i \nabla h_i(\mathbf{x}^*)$  and

$$\tilde{h}_p(t, \boldsymbol{\alpha}) = h_p(\tilde{\mathbf{x}}(t, \boldsymbol{\alpha})), \quad 1 \leq p \leq K$$

- Note  $\tilde{h}_p(0, \mathbf{0}) = h_p(\mathbf{x}^*) = 0$ , and

$$\frac{\partial \tilde{h}_p(0, \mathbf{0})}{\partial \alpha_q} = \sum_{\ell=1}^n \frac{\partial h_p(\mathbf{x}^*)}{\partial x_\ell} \frac{\partial \tilde{x}_\ell(0, \mathbf{0})}{\partial \alpha_q} = \nabla h_p(\mathbf{x}^*)^T \nabla h_q(\mathbf{x}^*)$$

- By regularity of  $\mathbf{x}^*$ ,  $\mathbf{A} = [\nabla h_1(\mathbf{x}^*), \dots, \nabla \tilde{h}_K(\mathbf{x}^*)]$  has rank  $K$ , so

$$\frac{\partial \tilde{\mathbf{h}}(0, \mathbf{0})}{\partial \boldsymbol{\alpha}} = \mathbf{A}^T \mathbf{A} \succ \mathbf{0}$$

and hence nonsingular.

## Proof for $\mu \geq \mathbf{0}$ (cont'd)

### Proof (cont'd).

- By Implicit Function Theorem, there exists  $\alpha(t)$  for small  $|t|$  s.t.  $\alpha(0) = \mathbf{0}$  and  $\tilde{\mathbf{h}}(t, \alpha(t)) = \mathbf{h}(\tilde{\mathbf{x}}(t, \alpha(t))) = \mathbf{0}$ . Furthermore,

$$\alpha'(0) = \left[ \frac{\partial \tilde{\mathbf{h}}(0, \mathbf{0})}{\partial \alpha} \right]^{-1} \frac{\partial \tilde{\mathbf{h}}(0, \mathbf{0})}{\partial t} = \left[ \frac{\partial \tilde{\mathbf{h}}(0, \mathbf{0})}{\partial \alpha} \right]^{-1} \begin{bmatrix} \nabla h_1(\mathbf{x}^*)^T \mathbf{d}_0 \\ \vdots \\ \nabla h_K(\mathbf{x}^*)^T \mathbf{d}_0 \end{bmatrix} = \mathbf{0}$$

since  $\mathbf{d}_0 \in T(\mathbf{x}^*) = \text{Null}(\mathbf{A}^T)$ .

- Let  $\mathbf{x}(t) = \tilde{\mathbf{x}}(t, \alpha(t))$ . Then  $\mathbf{h}(\mathbf{x}(t)) = \mathbf{0}$ , so  $\mathbf{x}(t) \subset S$ .

$$\mathbf{x}'(0) = \mathbf{d}_0 + \sum_{i=1}^K \alpha'_i(0) \nabla h_i(\mathbf{x}^*) = \mathbf{d}_0$$

# Sufficiency of KKT Conditions for Convex Problems

**Theorem.** For convex (ICP), i.e.  $f$  and  $g_j$  are convex, and  $h_i$  are affine, if there exist  $\lambda_1^*, \dots, \lambda_k^*$  and  $\mu_1^*, \dots, \mu_m^*$  s.t. the KKT conditions are satisfied at a feasible  $\mathbf{x}^* \in X$ , then  $\mathbf{x}^*$  is a global minimum of (ICP).

**Note.** The previous necessary conditions assume  $\mathbf{x}^*$  is regular point. The sufficient conditions here assume convexity but **not** regularity.

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

1. By the KKT conditions,

$$\nabla f(\mathbf{x}^*)^T(\mathbf{x} - \mathbf{x}^*) = - \sum_i \lambda_i^* \nabla h_i(\mathbf{x}^*)^T(\mathbf{x} - \mathbf{x}^*) - \sum_{j \in J(\mathbf{x}^*)} \mu_j^* \nabla g_j(\mathbf{x}^*)^T(\mathbf{x} - \mathbf{x}^*)$$

It suffices to show  $\nabla h_i(\mathbf{x}^*)^T(\mathbf{x} - \mathbf{x}^*) = 0$  and  $\nabla g_j(\mathbf{x}^*)^T(\mathbf{x} - \mathbf{x}^*) \leq 0$ .

2. Since  $h_i(\mathbf{x}) = \mathbf{a}_i^T \mathbf{x} - b_i$  is affine, and  $h_i(\mathbf{x}^*) = h(\mathbf{x}^*) = 0$  by feasibility,

$$\nabla h_i(\mathbf{x}^*)^T(\mathbf{x} - \mathbf{x}^*) = \mathbf{a}_i^T(\mathbf{x} - \mathbf{x}^*) = h_i(\mathbf{x}) - h_i(\mathbf{x}^*) = 0$$

3. For  $j \in J(\mathbf{x}^*)$ ,  $g_j(\mathbf{x}^*) = 0$  and  $g_j(\mathbf{x}) \leq 0$ . By the convexity of  $g_j$ ,

$$\nabla g_j(\mathbf{x}^*)^T(\mathbf{x} - \mathbf{x}^*) \leq g_j(\mathbf{x}) - g_j(\mathbf{x}^*) \leq 0$$