

CS 2601 Linear and Convex Optimization

9. Lagrange condition

Bo Jiang

John Hopcroft Center for Computer Science
Shanghai Jiao Tong University

Fall 2022

Outline

- Convex problems with equality constraints
- General equality constrained problems

Equality constrained convex problems

Consider the equality constrained convex optimization problem

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

where f is convex with $\text{dom } f = \mathbb{R}^n$.

In a more compact form,

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

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

We assume f is differentiable and the problem is feasible.

Feasible set

The feasible set is

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

Given any $\mathbf{x}_0 \in X$,

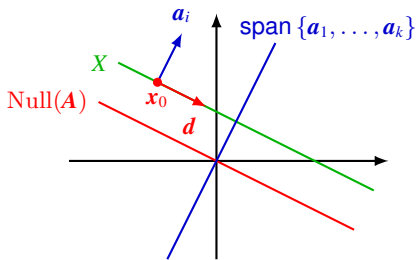
$$X = \mathbf{x}_0 + \text{Null}(\mathbf{A})$$

where $\text{Null}(\mathbf{A}) = \{\mathbf{x} : \mathbf{A}\mathbf{x} = \mathbf{0}\} = \{\mathbf{x} : \mathbf{a}_i^T \mathbf{x} = 0, i = 1, \dots, k\}$ is the **null space** of \mathbf{A} .

$\text{Null}(\mathbf{A})$ is precisely the set of feasible directions (at any $\mathbf{x}_0 \in X$)

$$\mathbf{x}_0 + \mathbf{d} \in X \iff \mathbf{a}_i^T \mathbf{d} = 0, \forall i$$

- \mathbf{a}_i is a normal vector to X
- $\mathbf{d} \in \text{Null}(\mathbf{A})$ is a tangent vector to X , the velocity $\mathbf{x}'(0)$ of a path $\mathbf{x}(t) = \mathbf{x}_0 + t\mathbf{d} \subset X$

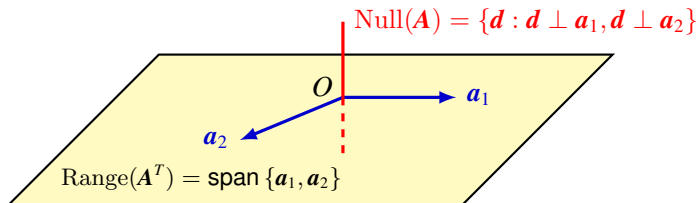


$$\text{Range}(\mathbf{A}^T) = \text{span}\{\mathbf{a}_1, \dots, \mathbf{a}_k\}$$

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.

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



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

Optimality condition

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

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

Note. Geometrically, $\nabla f(\mathbf{x}^*) \perp \text{Null}(\mathbf{A})$ means $\nabla f(\mathbf{x}^*)$ is perpendicular to all feasible directions, which are also tangent vectors at \mathbf{x}^* .

Proof. Recall (slide 7 of §5 part 1) $\mathbf{x}^* \in X$ is optimal iff

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

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

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

Since $\mathbf{d} \in \text{Null}(\mathbf{A}) \iff -\mathbf{d} \in \text{Null}(\mathbf{A})$, the condition then reduces to

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

Note. If f is nonconvex and \mathbf{x}^* a local minimum, then $\nabla f(\mathbf{x}^*) \perp \text{Null}(\mathbf{A})$ is a necessary condition. For a proof, note $t = 0$ is a local minimum of $g(t) = f(\mathbf{x}^* + t\mathbf{d})$, so $g'(0) = \nabla f(\mathbf{x}^*)^T \mathbf{d} = 0$.

Lagrange condition

Theorem. $\mathbf{x}^* \in X$ 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

$$\text{Null}(\mathbf{A})^\perp = \text{Range}(\mathbf{A}^T)$$

\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 $\boldsymbol{\lambda}^* = -\mathbf{v}^*$.

Lagrange condition (cont'd)

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 **Lagrange condition**, aka **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 gradient 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 as such.

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

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

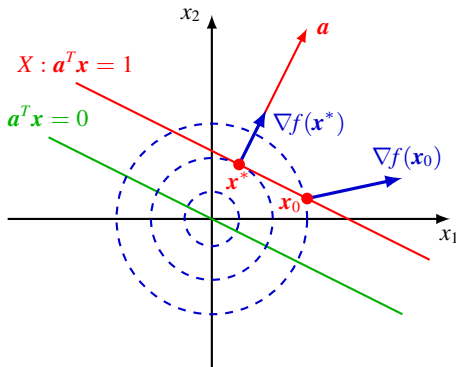
- gradient

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

- at \mathbf{x}^* ,

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

Note X is parallel to $\text{Null}(\mathbf{a}^T)$.



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)

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

Method 2. Lagrange multipliers method.

- The Lagrangian is

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

- Lagrange condition

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

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

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

Example (cont'd)

Block Gaussian elimination.

- The augmented matrix is

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

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

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

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

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

- Left 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)

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

- vectors normal to the feasible set X

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

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

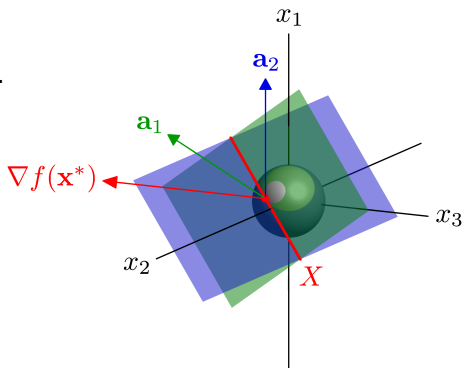
- gradient

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

- at \mathbf{x}^* ,

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

Note X is parallel to $\text{Null}(\mathbf{A}^T)$.



Outline

- Convex problems with equality constraints
- General equality constrained problems

Optimization on 2D circle

Consider the constraint in \mathbb{R}^2 ,

$$h(\mathbf{x}) = \|\mathbf{x}\|^2 - 1 = 0$$

Feasible set $X = \{\mathbf{x} : \|\mathbf{x}\| = 1\}$. At $\mathbf{x}_0 \in X$,

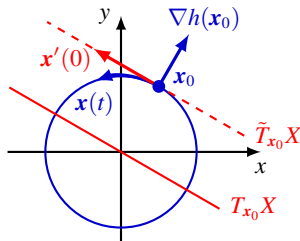
- A **tangent vector** is the initial velocity $\mathbf{x}'(0)$ of a feasible local path $\mathbf{x}(t)$ starting at \mathbf{x}_0 , i.e. $\mathbf{x}(0) = \mathbf{x}_0$, $h(\mathbf{x}(t)) = 0$ for small t . Note

$$h'(\mathbf{x}_0)\mathbf{x}'(0) = \nabla h(\mathbf{x}_0)^T \mathbf{x}'(0) = 0 \quad \text{i.e.} \quad \mathbf{x}'(0) \in \text{Null}(h'(\mathbf{x}_0))$$

- A tangent vector \mathbf{d} is a feasible direction in the sense that there is a feasible path $\mathbf{x}(t)$ in that direction, i.e. $\mathbf{d} = \mathbf{x}'(0)$.
- The **tangent space** $T_{\mathbf{x}_0}X$ is the set of tangent vectors. It turns out

$$T_{\mathbf{x}_0}X = \text{Null}(h'(\mathbf{x}_0)) = \{\mathbf{d} : \nabla h(\mathbf{x}_0)^T \mathbf{d} = 0\}$$

- Think of the tangent line $\tilde{T}_{\mathbf{x}_0}X$ as $T_{\mathbf{x}_0}X$ attached at \mathbf{x}_0



Optimization on 2D circle

Consider the smooth nonconvex (why?) problem

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

Let $\mathbf{x}^* \in X$ be a **local** minimum. Given $\mathbf{d} \in \text{Null}(h'(\mathbf{x}^*))$, let $\mathbf{x}(t)$ be a feasible local path² with $\mathbf{x}(0) = \mathbf{x}^*$, $\mathbf{x}'(0) = \mathbf{d}$ and $h(\mathbf{x}(t)) = 0$ for small t .

Since $\mathbf{x}^* = \mathbf{x}(0)$ is a local minimum of the constrained problem, $t = 0$ is a local minimum of $g(t) = f(\mathbf{x}(t))$, so

$$0 = g'(0) = \nabla f(\mathbf{x}^*)^T \mathbf{x}'(0) = \nabla f(\mathbf{x}^*)^T \mathbf{d}$$

Since $\mathbf{d} \in \text{Null}(h'(\mathbf{x}^*))$ is arbitrary,

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

²For example, if $\mathbf{x}^* = (\cos \phi_0, \sin \phi_0)$, then $\mathbf{d} = (-a \sin \phi_0, a \sin \phi_0)$ for some $a \in \mathbb{R}$. Then $\mathbf{x}(t) = (\cos(at + \phi_0), \sin(at + \phi_0))$ satisfies the requirement.

Optimization on 2D circle (cont'd)

By $\text{Null}(\mathbf{A})^\perp = \text{Range}(\mathbf{A}^T)$,

$$\nabla f(\mathbf{x}^*) \in \text{Range}(h'(\mathbf{x}^*)^T) = \text{span}\{\nabla h(\mathbf{x}^*)\}$$

so there exists a λ^* 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 λ^* s.t.

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

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

Example

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

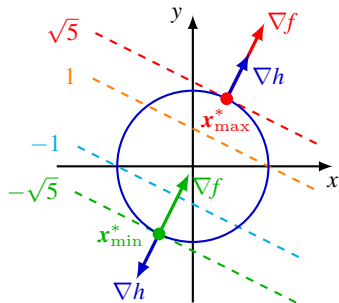
- Lagrange condition

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

- 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, $\nabla f \parallel \nabla h$ and $\nabla f \perp X$



Example

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

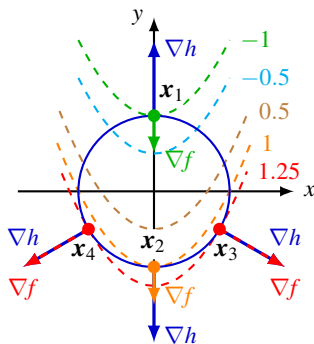
- Lagrange condition

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

- 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 (and certain other points), $\nabla f \parallel \nabla h$ and $\nabla f \perp X$



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

General equality constraints

Consider a general equality constraint function \mathbf{h} , where $\mathbf{h} : \mathbb{R}^n \rightarrow \mathbb{R}^k$ has smooth components h_1, \dots, h_k . The feasible set is

$$X = \{\mathbf{x} : \mathbf{h}(\mathbf{x}) = \mathbf{0}\}$$

A point \mathbf{x}_0 is a **regular point** of \mathbf{h} if

$$\mathbf{h}'(\mathbf{x}_0) = \begin{bmatrix} \nabla h_1(\mathbf{x}_0)^T \\ \vdots \\ \nabla h_k(\mathbf{x}_0)^T \end{bmatrix}$$

has full (row) rank k , or equivalently, $\nabla h_1(\mathbf{x}_0), \dots, \nabla h_k(\mathbf{x}_0)$ are linearly independent; otherwise it is a **critical point** of \mathbf{h} .

At a regular point \mathbf{x}_0 , the local geometry of X can be well characterized by the first order information $\mathbf{h}'(\mathbf{x}_0)$, or $\nabla h_1(\mathbf{x}_0), \dots, \nabla h_k(\mathbf{x}_0)$, and the derivation on slides 16-17 carries over.

Tangent space and normal space

A **tangent vector** of X at $\mathbf{x}_0 \in X$ is the initial velocity $\mathbf{x}'(0)$ of a feasible local path $\mathbf{x}(t)$ starting at \mathbf{x}_0 , i.e. $\mathbf{x}(0) = \mathbf{x}_0$, $h(\mathbf{x}(t)) = 0$ for small t . Note

$$\left. \frac{d}{dt} \mathbf{h}(\mathbf{x}(t)) \right|_{t=0} = \mathbf{h}'(\mathbf{x}_0) \mathbf{x}'(0) = \mathbf{0} \quad \text{i.e.} \quad \mathbf{x}'(0) \in \text{Null}(\mathbf{h}'(\mathbf{x}_0))$$

The **tangent space** $T_{\mathbf{x}_0}X$ of X at \mathbf{x}_0 is the set of all tangent vectors at \mathbf{x}_0 .

The **normal space** $N_{\mathbf{x}_0}X$ of X at \mathbf{x}_0 is the orthogonal complement of $T_{\mathbf{x}_0}X$,

$$N_{\mathbf{x}_0}X = [T_{\mathbf{x}_0}X]^\perp$$

Theorem. At a regular point $\mathbf{x}_0 \in X$,

$$T_{\mathbf{x}_0}X = \text{Null}(\mathbf{h}'(\mathbf{x}_0)) = \{\mathbf{d} : \nabla h_i(\mathbf{x}_0)^T \mathbf{d} = 0, \quad i = 1, 2, \dots, k\}$$

and

$$N_{\mathbf{x}_0}X = \text{span} \{ \nabla h_1(\mathbf{x}_0), \dots, \nabla h_k(\mathbf{x}_0) \}$$

Tangent space and normal space

Proof

We already know

$$T_{x_0}X \subset \text{Null}(\mathbf{h}'(\mathbf{x}_0))$$

For $\text{Null}(\mathbf{h}'(\mathbf{x}_0)) \subset T_{x_0}X$, we have the following

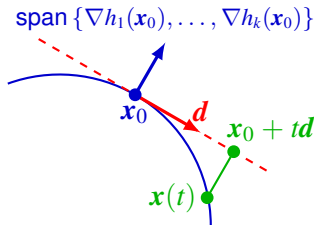
Lemma. If \mathbf{x}_0 is a regular point, then for any \mathbf{d} s.t. $\mathbf{h}'(\mathbf{x}_0)\mathbf{d} = \mathbf{0}$, there exists a local path $\mathbf{x}(t)$ s.t. $\mathbf{h}(\mathbf{x}(t)) = \mathbf{0}$, $\mathbf{x}(0) = \mathbf{x}_0$ and $\mathbf{x}'(0) = \mathbf{d}$.

Proof. Let

$$\begin{aligned}\tilde{\mathbf{x}}(t, \alpha) &= \mathbf{x}_0 + t\mathbf{d} + \mathbf{h}'(\mathbf{x}_0)^T \alpha, \\ &= \mathbf{x}_0 + t\mathbf{d} + \sum_{i=1}^k \alpha_i \nabla h_i(\mathbf{x}_0)\end{aligned}$$

and

$$\mathbf{F}(t, \alpha) = \mathbf{h}(\tilde{\mathbf{x}}(t, \alpha))$$



Proof of lemma (cont'd)

Note

$$\mathbf{F}(0, \mathbf{0}) = \mathbf{h}(\mathbf{x}_0) = \mathbf{0}, \quad \frac{\partial \mathbf{F}(0, \mathbf{0})}{\partial \boldsymbol{\alpha}} = \mathbf{h}'(\mathbf{x}_0) \mathbf{h}'(\mathbf{x}_0)^T \succ \mathbf{0}$$

since $\mathbf{h}'(\mathbf{x}_0)^T$ has full rank k by regularity at \mathbf{x}_0 .

By the Implicit Function Theorem, there exists $\boldsymbol{\alpha} = \boldsymbol{\phi}(t)$ for small t s.t. $\boldsymbol{\phi}(0) = \mathbf{0}$, $\mathbf{F}(t, \boldsymbol{\phi}(t)) = \mathbf{0}$ and

$$\boldsymbol{\phi}'(0) = - \left[\frac{\partial \mathbf{F}(0, \mathbf{0})}{\partial \boldsymbol{\alpha}} \right]^{-1} \frac{\partial \mathbf{F}(0, \mathbf{0})}{\partial t} = - \left[\frac{\partial \mathbf{F}(0, \mathbf{0})}{\partial \boldsymbol{\alpha}} \right]^{-1} \mathbf{h}'(\mathbf{x}_0) \mathbf{d} = \mathbf{0}$$

Then

$$\mathbf{x}(t) = \tilde{\mathbf{x}}(t, \boldsymbol{\phi}(t)) = \mathbf{x}_0 + t\mathbf{d} + \mathbf{h}'(\mathbf{x}_0)^T \boldsymbol{\phi}(t) = \mathbf{x}_0 + t\mathbf{d} + \sum_{i=1}^k \phi_i(t) \nabla h_i(\mathbf{x}_0)$$

satisfies the requirement.

Appendix: Implicit function theorem

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

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

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

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

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

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

First-order necessary condition

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

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

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.

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

Note. This simply says $\nabla f(\mathbf{x}^*) \in N_{\mathbf{x}_0} X = \text{span} \{ \nabla h_1(\mathbf{x}_0), \dots, \nabla h_k(\mathbf{x}_0) \}$.

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

Proof

Let $\mathbf{d} \in T_{\mathbf{x}_0}X$ and $\mathbf{x}(t)$ a feasible local path at \mathbf{x}^* with $\mathbf{x}'(0) = \mathbf{d}$. Then $t = 0$ is a local minimum of $g(t) = f(\mathbf{x}(t))$, so

$$0 = g'(0) = \nabla f(\mathbf{x}^*)^T \mathbf{d}$$

Since \mathbf{d} is arbitrary,

$$\nabla f(\mathbf{x}^*) \perp T_{\mathbf{x}_0}X$$

and hence

$$\nabla f(\mathbf{x}^*) \in [T_{\mathbf{x}_0}X]^\perp = N_{\mathbf{x}_0}X$$

Since

$$N_{\mathbf{x}_0}X = \text{span} \{ \nabla h_1(\mathbf{x}_0), \dots, \nabla h_k(\mathbf{x}_0) \}$$

at a regular point \mathbf{x}_0 , there exist $\lambda_1^*, \dots, \lambda_k^* \in \mathbb{R}$ s.t.

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

Insufficiency of Lagrange condition

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

$$f(\mathbf{x}) = \|\mathbf{x}\|^2$$

$$h(\mathbf{x}) = y - x^3 - 1$$

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

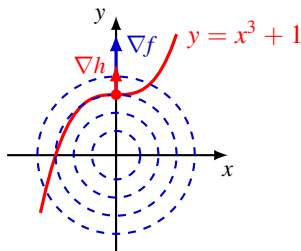
$$\nabla f(\mathbf{x}^*) = (0, 2)^T, \quad \nabla h(\mathbf{x}^*) = (0, 1)^T$$

so

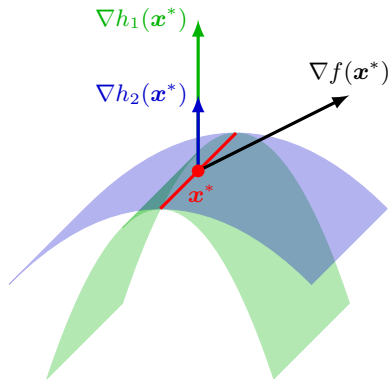
$$\nabla f(\mathbf{x}^*) - 2\nabla h(\mathbf{x}^*) = \mathbf{0}$$

but \mathbf{x}^* is neither a maximum nor a minimum.

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



Critical points



Critical points

The Lagrange condition may fail at critical points.

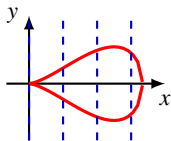
Example.

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

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.

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 λ_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 (7) 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

Example (cont'd)

