

# 535520: Optimization Algorithms

## Lecture 1 – Fundamentals

Ping-Chun Hsieh (謝秉均)

September 2, 2024

# This Lecture

1. Optimization Problems: Formulation and Terminology

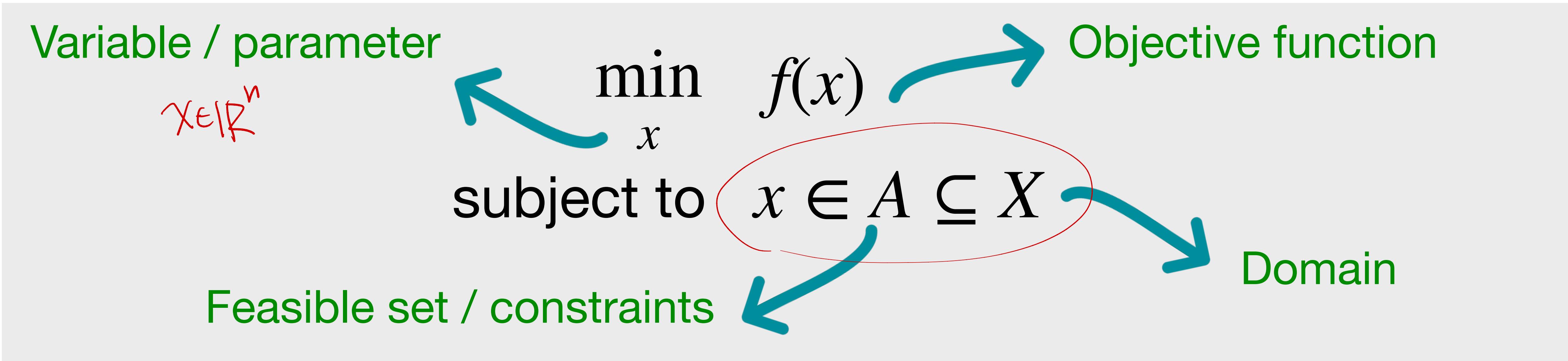
2. Optimality Conditions

3. Subgradients and Subdifferentials

Reading material:

- Chapters 1.1 and 2.1 of Dimitri Bertsekas's textbook “Nonlinear Programming”
- Chapters 2 and 3 of Stephen Boyd's textbook “Convex Optimization”

# Basic Formulation of an Optimization Problem



- Nice properties of an objective function
  - Continuity?
  - Smoothness?
  - Convexity?
  - Differentiability?
  - Separability?
- Nice properties of a feasible set
  - Compactness?
  - Convex sets?
  - Unconstrained?
  - Linear constraints?
  - Discrete?

# A Motivating Example: Portfolio Selection

$$W = 1$$

## ► Portfolio selection (in hindsight)

- Let  $r_t = (r_{t,1}, \dots, r_{t,n})$  denote “price ratio” of the  $n$  assets at each time  $t$
- Suppose initially we have total wealth  $w > 0$
- We want to choose an initial portfolio allocation vector  $a$  in hindsight such that the total wealth after  $T$  iterations is maximized

**Question:** How to write down the optimization problem (e.g., objective function, constraints)?



$$a = (a_1, \dots, a_n)$$

Max

$$\sum_{i=1}^n r_{1,i} r_{2,i} \dots r_{T,i} a_i$$

subject to

$$a_1 + \dots + a_n = w$$

$$a_i \geq 0, \text{ for all } i=1, 2, \dots, n$$

# Optimization in ML and Beyond

## Machine Learning

Empirical Risk  
Minimization

Online Learning

Federated Learning

## RL and Robotics

Policy Optimization  
Adaptive Control  
Learning From  
Demonstrations

## Deep Learning

Seq2Seq  
GANs  
Representation  
Learning

## Information Processing

Image Processing  
Speech Signal  
Processing  
Data compression

## Computational Science

Physics  
Chemistry  
Bioinformatics

## Network Optimization

Network Utility  
Maximization  
Social Networks  
Packet Scheduling

# Optimization: 3 Questions to Answer

1. **Characterization:** Sufficient / necessary conditions of an optimal solution?  
*(Our focus today)*
2. **Algorithms:** Iterative algorithms that find an optimal solution?
3. **Convergence:** Do the iterates converge to an optimum? How fast?

# **Optimality Conditions**

**(Structural information about optimal solutions)**

# Optimality Conditions (Necessary / Sufficient)

## Unconstrained cases:

- C1. FONC: First-order necessary conditions for local optimality
- C2. SONC: Second-order necessary conditions for local optimality
- C3. FOSC: First-order sufficient conditions for global optimality
- C4. SOSC: Second-order sufficient conditions for local optimality

## Constrained cases:

- C5. FONC-C: First-order necessary conditions for constrained local optimality
- C6. FOSC-C: First-order sufficient conditions for constrained global optimality

## Non-differentiable cases:

- C7. Fermat's Rule

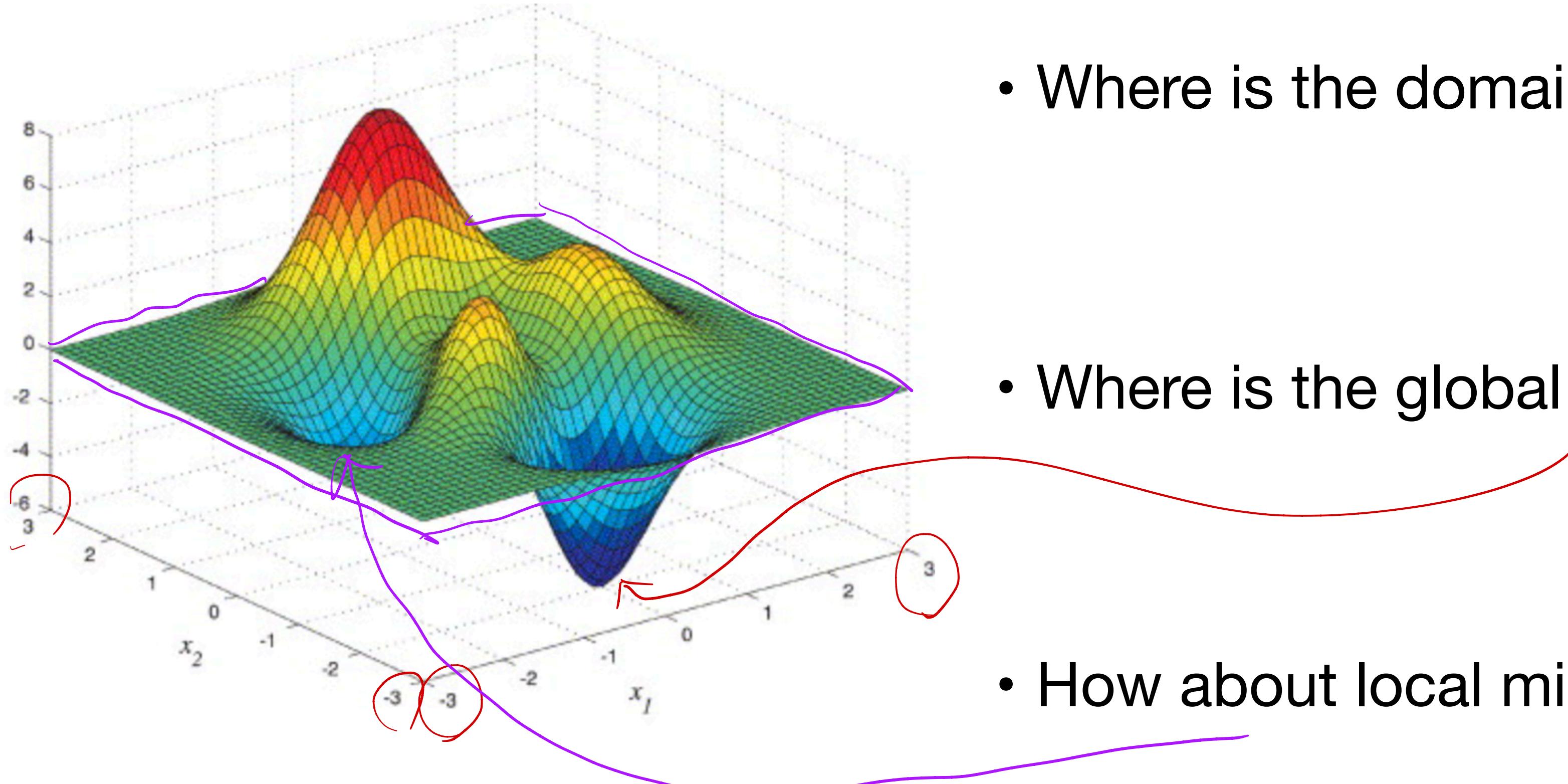
# Notation & Assumptions for This Lecture

Unless stated otherwise:

- $\|\cdot\|_p$  denotes the  $\ell_p$  norm
- $\|\cdot\| \equiv \|\cdot\|_2$  denotes the Euclidean norm
- We focus on **multivariate single-objective minimization** problems (i.e.,  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ )
- The objective function is assumed differentiable

# Local and Global Minima

Intuitively, let's make some observations:



- Where is the domain  $X$ ?
- Where is the global minimizer?
- How about local minimizer(s)?

$$X = \left\{ (x_1, x_2) : -3 \leq x_1 \leq 3 \right. \\ \left. -3 \leq x_2 \leq 3 \right\}$$

# Local and Global Minima (Formally)

**Definition:** Given  $f: X \rightarrow \mathbb{R}$ , a vector  $x^* \in X$  is a *local minimizer* if there exists some  $\epsilon > 0$  such that

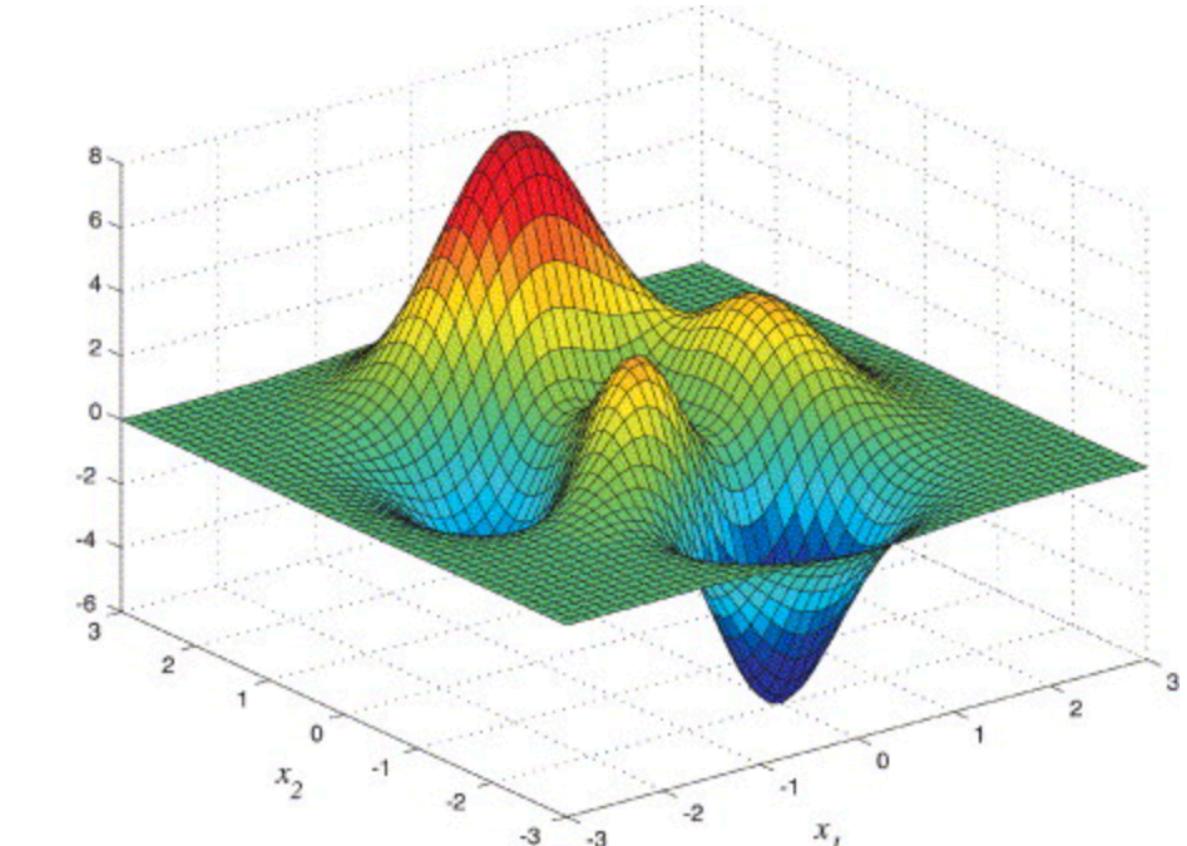
$$f(x^*) \leq f(x), \text{ for all } x \in X \text{ with } \|x - x^*\| < \epsilon$$

$\epsilon$ -neighborhood

**Definition:** Given  $f: X \rightarrow \mathbb{R}$ , a vector  $x^* \in X$  is a *global minimizer* if

$$f(x^*) \leq f(x), \text{ for all } x \in X$$

**Remark:** *Strict* local / global minimizers if “ $\leq$ ” is replaced by “ $<$ ” for all  $x \neq x^*$



# A Quick Recap of Notations, Calculus, and Linear Algebra (1/3)

Given  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $x = (x_1, \dots, x_n)$

- Gradient vector

$$\nabla f(x) = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \vdots \\ \frac{\partial f(x)}{\partial x_n} \end{bmatrix}$$

- # • Hessian matrix

$\nabla^2 f(x)$  =   
 $\frac{\partial^2 f(x)}{\partial x_i \partial x_j}$

# • Useful Properties:

- (1)  $H \succcurlyeq 0$  if and only if all its eigenvalues are non-negative
  - (2)  $H \succ 0$  if and only if all its eigenvalues are positive

(When discussing psd/pd, we can assume  $H$  is symmetric, without loss of generality)

A real square matrix  $H \in \mathbb{R}^{n \times n}$  is said to be

- Symmetric if

$$H = H^T$$

- Positive semidefinite

(psd), or  $H \geq 0$ , if: For all  $x \in \mathbb{R}^n$ ,  $x^T H x \geq 0$ .

- Positive definite (pd)

or  $H \succeq 0$  if:

For all  $x \in \mathbb{R}^n$ ,  $x^T H x \geq 0$

$(x \neq 0)$

$$x^T H x = x^T \left( \underbrace{H + H^T}_{\text{symmetric}} \right) x$$

# A Quick Recap of Notations, Calculus, and Linear Algebra (2/3)

In this course, we will leverage Taylor's Theorem a lot! (for both intuition and analysis)

**Taylor's Theorem (First-Order Version):** Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be continuously differentiable on an open neighborhood  $S$  of a vector  $x$ .

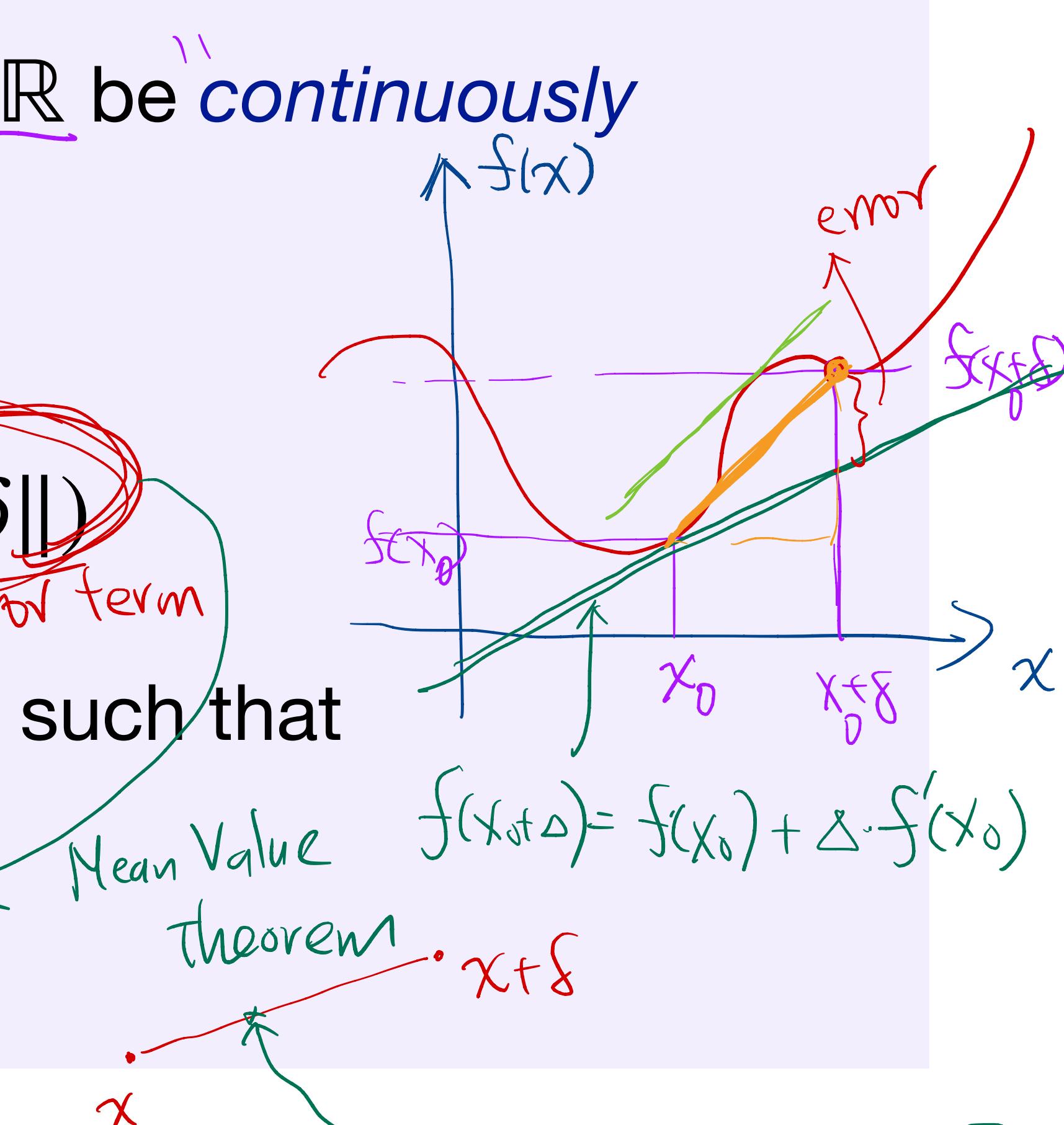
(1) For all  $\delta$  such that  $x + \delta \in S$ , we have

$$f(x + \delta) = f(x) + \delta^\top \nabla f(x) + o(\|\delta\|)$$

*error term*

(2) For all  $\delta$  such that  $x + \delta \in S$ , there exists  $\alpha \in [0, 1]$  such that

$$f(x + \delta) = f(x) + \delta^\top \nabla f(x + \alpha\delta)$$



**Question:** Why do we need “continuous differentiability”?

(for Mean Value Theorem)

# A Quick Recap of Notations, Calculus, and Linear Algebra (3/3)

Higher-order version of Taylor theorem:

Single-variate:  $\frac{1}{2} \cdot \delta^T f''(x) \cdot \delta$

**Taylor's Theorem (Second-Order Version):** Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be *twice continuously differentiable* on an open neighborhood  $S$  of a vector  $x$ .

(1) For all  $\delta$  such that  $x + \delta \in S$ , we have

$$f(x + \delta) = f(x) + \underbrace{\delta^T \nabla f(x)}_{\text{linear term}} + \underbrace{\frac{1}{2} \delta^T \nabla^2 f(x) \delta}_{\text{error term}} + o(\|\delta\|^2)$$

(2) For all  $\delta$  such that  $x + \delta \in S$ , there exists  $\alpha \in [0, 1]$  such that

$$f(x + \delta) = f(x) + \delta^T \nabla f(x) + \frac{1}{2} \delta^T \nabla^2 f(x + \alpha\delta) \delta$$

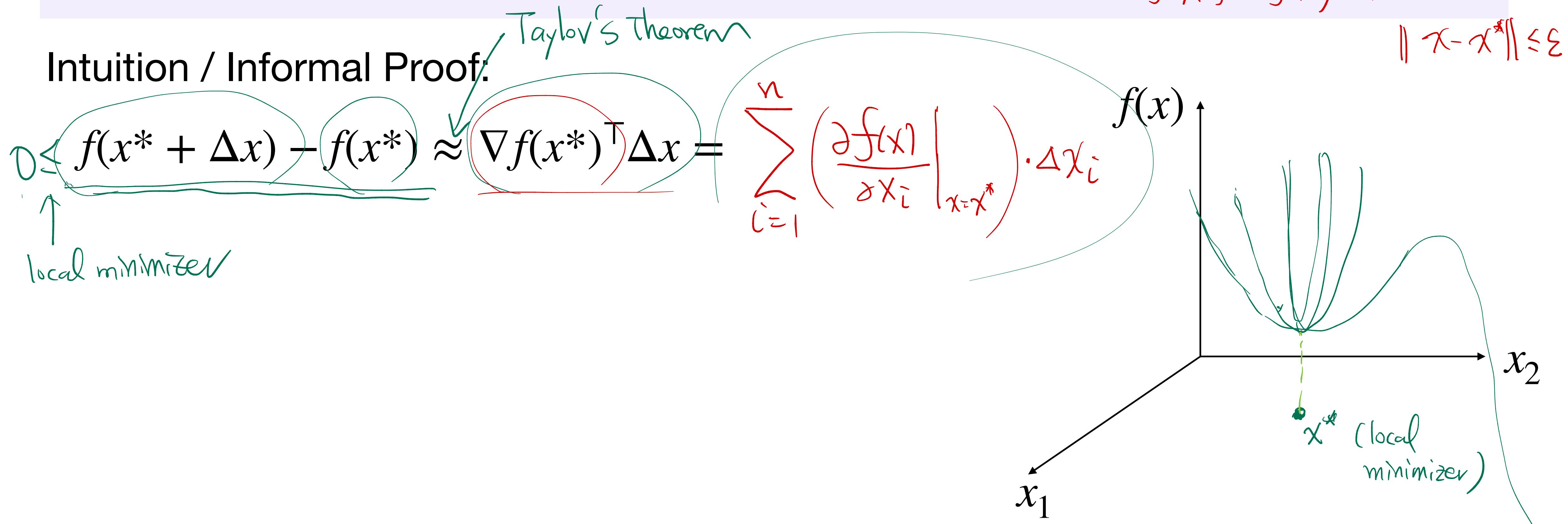
# C1. Necessary Conditions for Local Optimality: Unconstrained

**Theorem (First-Order Necessary Condition, FONC):** Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be continuously differentiable on a neighborhood of  $x^*$ , which is a local minimizer.

Then, we have  $\nabla f(x^*) = 0$ .

There exists  $\varepsilon > 0$ ,  $f(x^*) \leq f(x)$  for all  $x$  with  $\|x - x^*\| \leq \varepsilon$

Intuition / Informal Proof:



# C1. Necessary Conditions for Local Optimality: Unconstrained

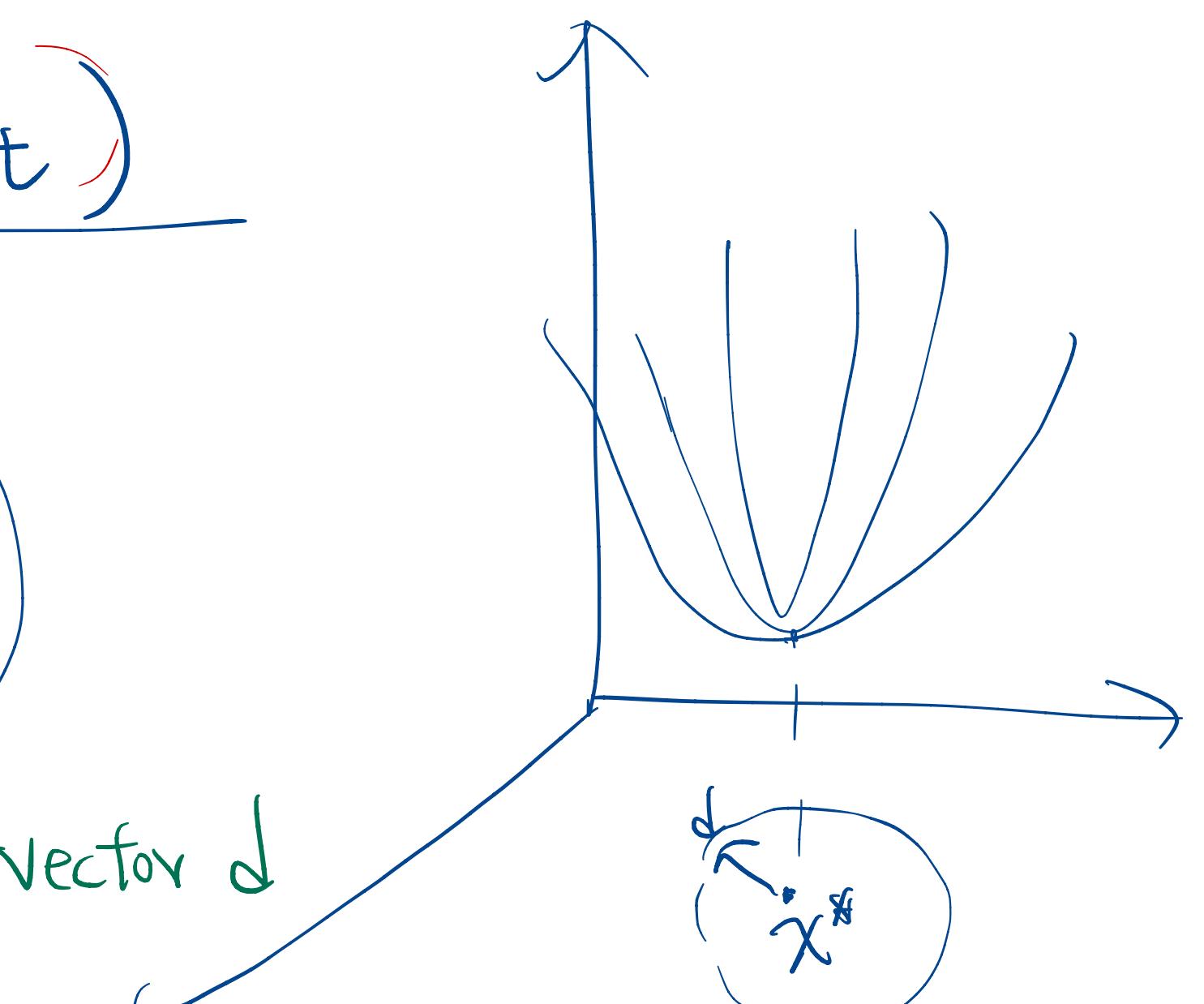
**Theorem (First-Order Necessary Condition, FONC):** Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be continuously differentiable on a neighborhood of  $x^*$ , which is a local minimizer. Then, we have  $\nabla f(x^*) = 0$ .

*Fix  $\|d\| = 1$  for simplicity*

*Proof:* Construct a function  $g(t) = f(x^* + td)$ , where  $d \in \mathbb{R}^n$  and  $t > 0$

$$\begin{aligned} 0 &\leq \lim_{t \downarrow 0} \frac{f(x^* + td) - f(x^*)}{t} = \lim_{t \downarrow 0} \frac{\nabla f(x^*)^T \cdot (td) + o(t)}{t} \\ &= \lim_{t \downarrow 0} \nabla f(x^*)^T \cdot d + \frac{o(t)}{t} \\ &= \nabla f(x^*)^T d \quad \text{for all unit vector } d \end{aligned}$$

Therefore, we must have  $\nabla f(x^*) = 0$ .  $\square$





## METHODUS

Ad disqurendam maximam &amp; minimam.

M N I S de inventione maximæ & minimæ doctrina, duabus positionibus ignotis innititur, & hac unica præceptione statuatur quilibet quæstionis terminus esse A, sive planum, sive solidum, aut longitudo, prout proposito satisfieri par est, & inventa maxima aut minima in terminis sub A, gradu ut libet iuuoluris; Ponatur rursus idem qui prius esse terminus A, + E, iterumque inveniatur maxima aut minima in terminis sub A & E, gradibus ut libet coefficientibus. Adæquentur, ut loquitur Diophantus, duo homogenea maximæ aut minimæ æqualia & demptis communibus ( quo peracto homogenea omnia ex parte alterutra ( ab E, vel ipsius gradibus afficiuntur ) applicentur omnia ad E, vel ad elatiorem ipsius gradum, donec aliquid ex homogeneis, ex parte utravis affectione sub E, omnino liberetur.

Elidantur deinde utrumque homogenea sub E, aut ipsius gradibus quomodolibet involuta & reliqua æquentur. Aut si ex una parte nihil superest æquentur sane, quod codem recidit, negata ad firmatis. Refolutio ultimæ istius æqualitatis dabit ualorem A, quâ cognita, maxima aut minima ex repetitis prioris resolutionis vestigiis innoteat.

Exemplum subijcimus

Sit recta AC, ita dividenda in E, ut rectang. AE C, sit maximum; Recta AC, dividatur B.

A      E      C

ponatur par altera B, esse A, ergo reliqua erit B, — A, & rectang. sub segmentis erit B, in A, — A<sup>2</sup> quod debet inueniri maximum. Ponatur rursus pars altera ipsius B, esse A, + E, ergo reliqua erit B, — A — E, & rectang. Sub. segmentis erit B, in A, — A<sup>2</sup> + B, in E, — E, quod debet adæquati superiori rectang. B, in A, — A<sup>2</sup>, demptis communibus B, in E, adæquabitur A, in E — E<sup>2</sup>, & omnibus per E, divisis B, adæquabitur <sup>2</sup>A + E, elidatur E, B, adæquabitur <sup>2</sup>A, igitur B, bifariam est dividenda, ad solutionem propositi, nec potest generalior dari methodus.

\*\*\*\*\*

## De Tangentibus linearum curvarum.

A D superiore methodum inventionem Tangentium ad data puncta in lineis quibuscumque curvis reducimus.

## An Interesting Fact:

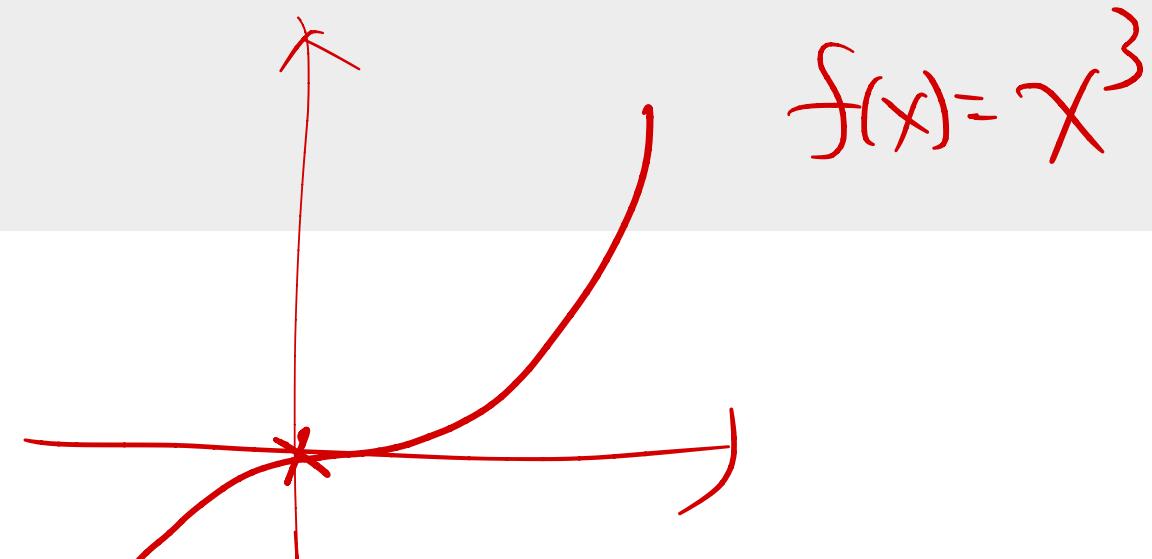
- This necessary condition was originally formulated by Fermat in 1637, *without proof* (as expected!)



Pierre de Fermat  
(1607-1665)

# Remarks on FONC

- The condition  $\nabla f(x^*) = 0$  is necessary but **not sufficient** for local optimality  
(any counterexample?)



- Despite this, the condition  $\nabla f(x^*) = 0$  is still useful as it provides a candidate set of locally optimal solutions.
- **Terminology:** A point  $x$  with  $\nabla f(x) = 0$  is termed as a “stationary point” or a “critical point” in the optimization literature.

# Critical Points and Saddle Points

Given a differentiable  $f: X \rightarrow \mathbb{R}$ :

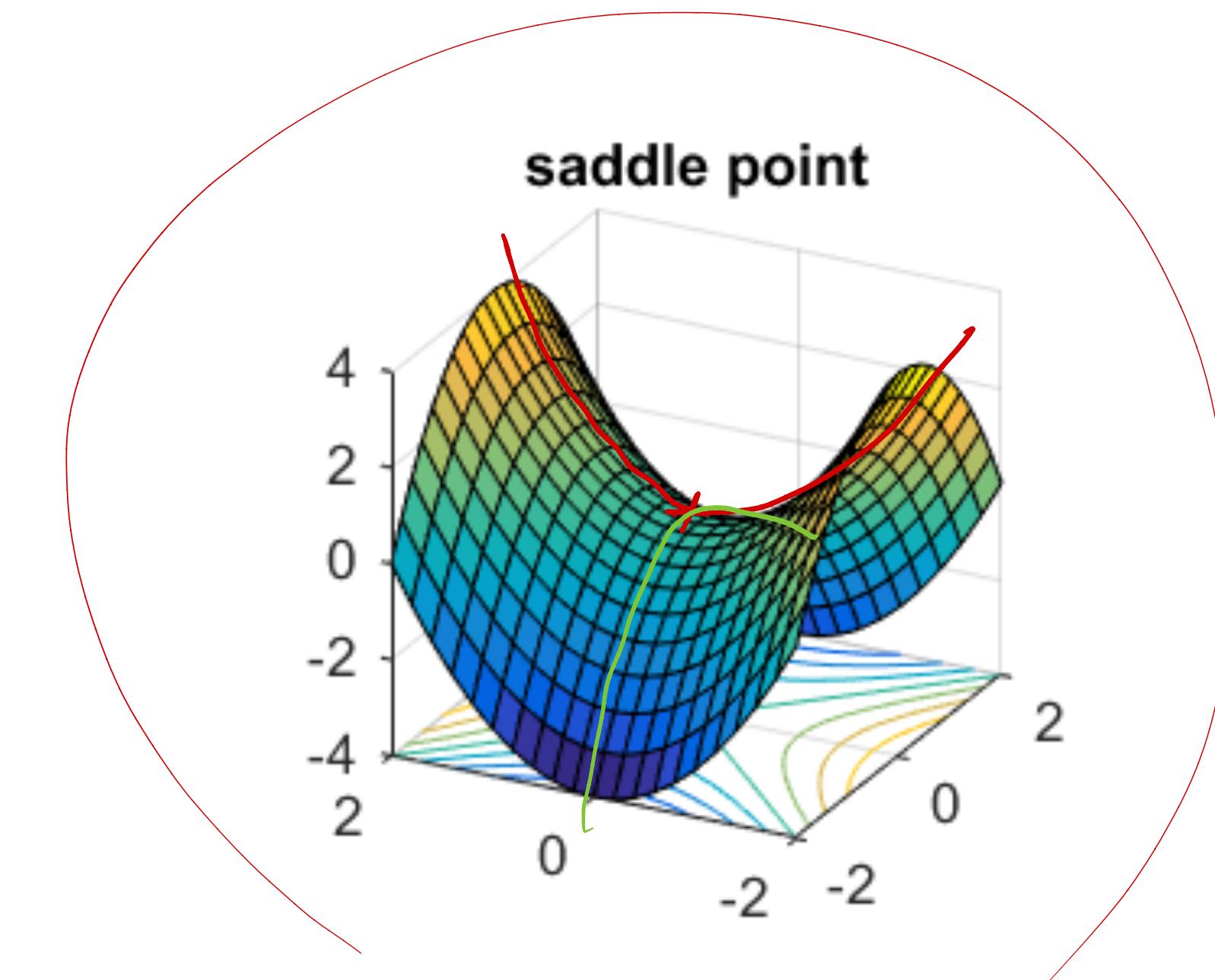
**Definition:** A vector  $x_0 \in X$  is a *critical point*

if  $\nabla f(x)|_{x=x_0} = 0$

**Definition:** A vector  $x_0 \in X$  is a *saddle point* if

$\nabla f(x)|_{x=x_0} = 0$  and  $x_0$  is not a local minimizer

nor a local maximizer



- The existence of saddle point suggests that  $\nabla f(x^*) = 0$  is **not sufficient** for local optimality

## C2. Second-Order Necessary Condition: Unconstrained Cases

**Theorem (Second-Order Necessary Condition, SONC):** Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be a twice continuously differentiable on a neighborhood of  $x^*$ , and  $x^*$  is a local minimizer. Then, in addition to  $\nabla f(x^*) = 0$ , we must also have

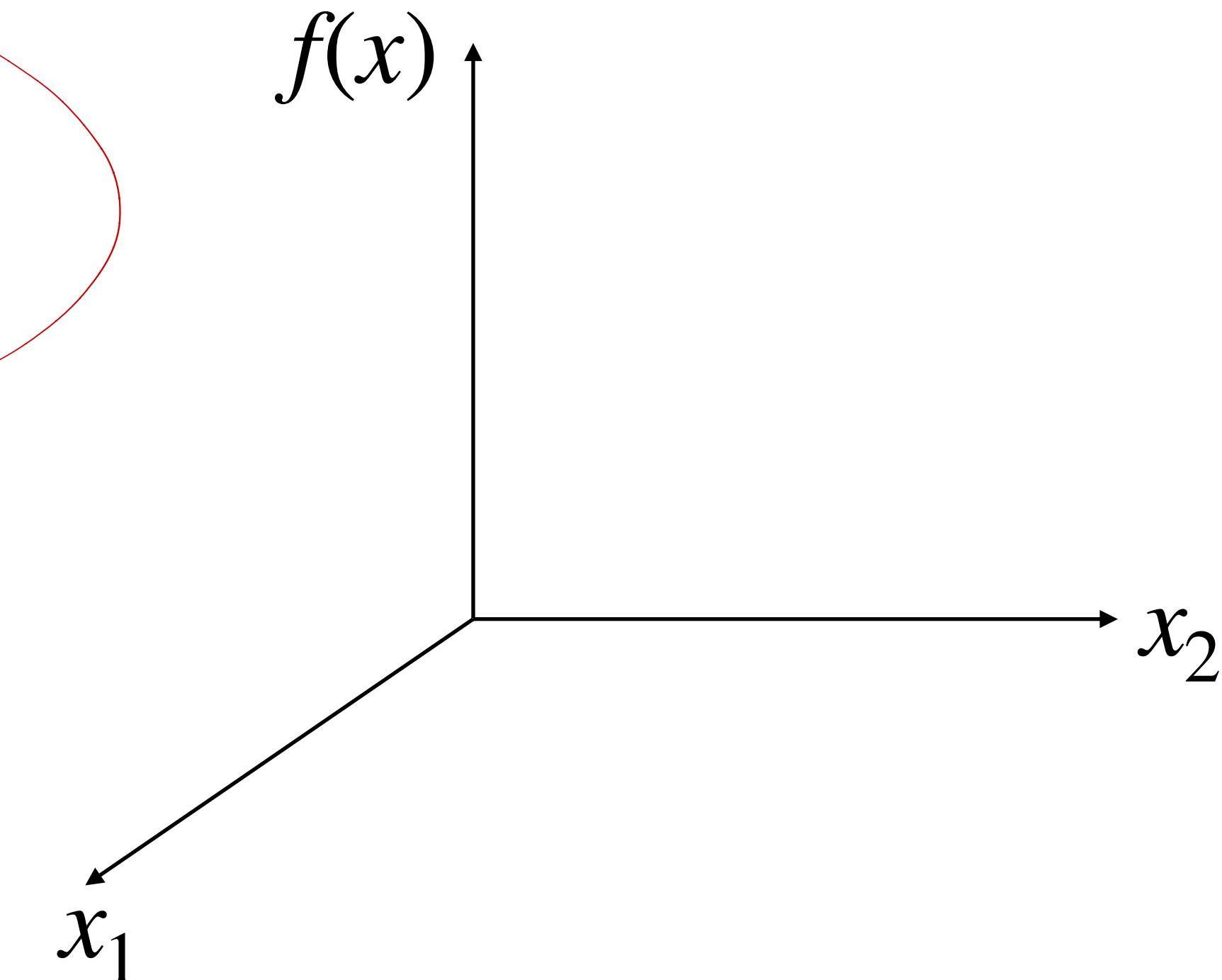
$$\nabla^2 f(x^*) \succeq 0$$

Intuition / Informal Proof:

$$f(x^* + \Delta x) - f(x^*) \approx \nabla f(x^*)^\top \Delta x + \frac{1}{2} \Delta x^\top \nabla^2 f(x^*) \Delta x$$

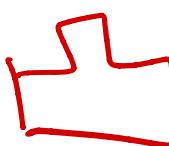
if  $\Delta x$  is small

$\nabla^2 f(x^*)$  is psd



Proof: You will be asked to prove this in HW0 :)

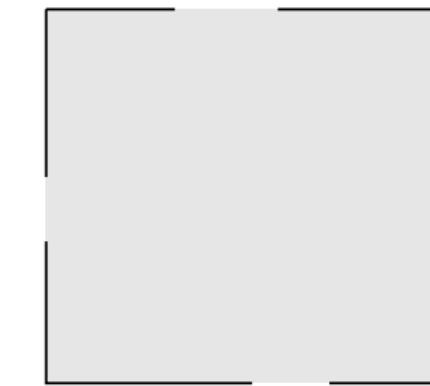
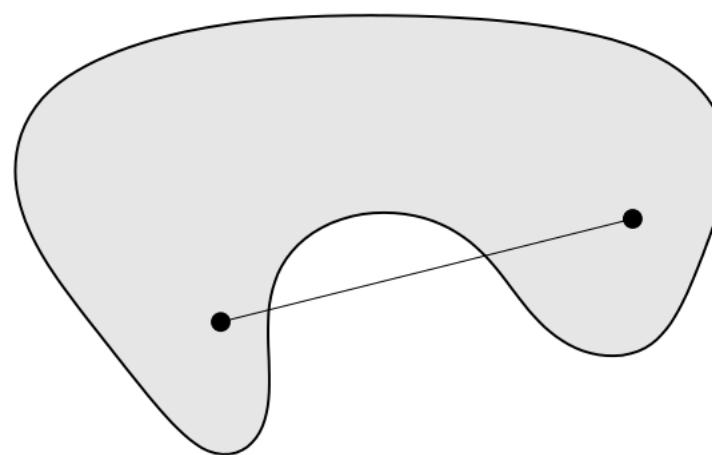
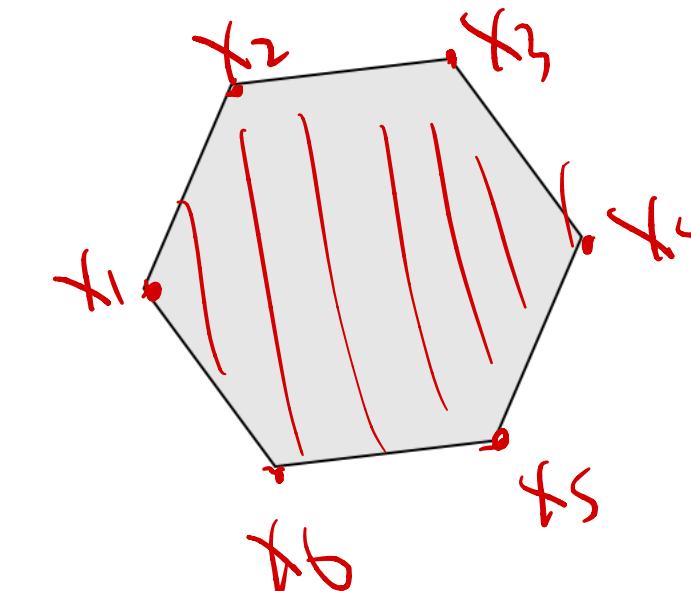
**Next question: Are these conditions *sufficient*?  
(If so, under what conditions?)**



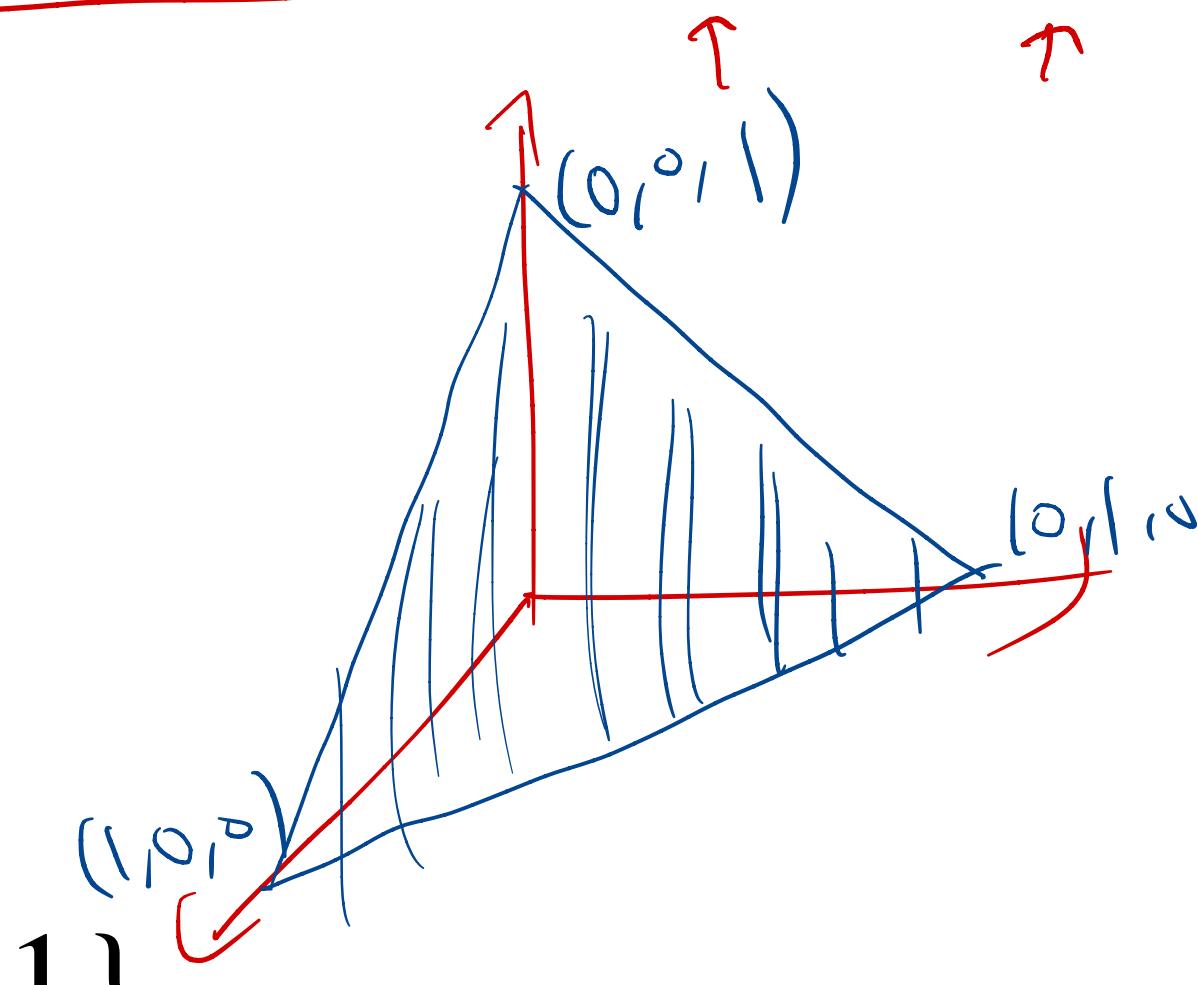
# Convex Sets

**Definition:** A set  $S \subset \mathbb{R}^n$  is called **convex** if for any  $x, y \in S$ , the line segment  $\{\alpha x + (1 - \alpha)y, \alpha \in [0,1]\}$  is also in  $S$ .

**Examples:**



- Convex hull: Let  $x_1, \dots, x_k \in \mathbb{R}^d$ . The convex hull  $\text{CH}(x_1, \dots, x_k) := \{ \sum_i \alpha_i x_i : \alpha_i \geq 0, \sum_i \alpha_i = 1 \}$
- Halfspace:  $\{x : a^\top x \leq b\}$
- Hyperplane:  $\{x : a^\top x = b\}$
- Ellipsoid:  $\{x : (x - a)^\top A(x - a) \leq 1\}$
- Probability simplex:  $\{x : x \geq 0, \sum_i x_i = 1\}$



(For more properties of convex sets, please check Stephen Boyd's lecture slides)

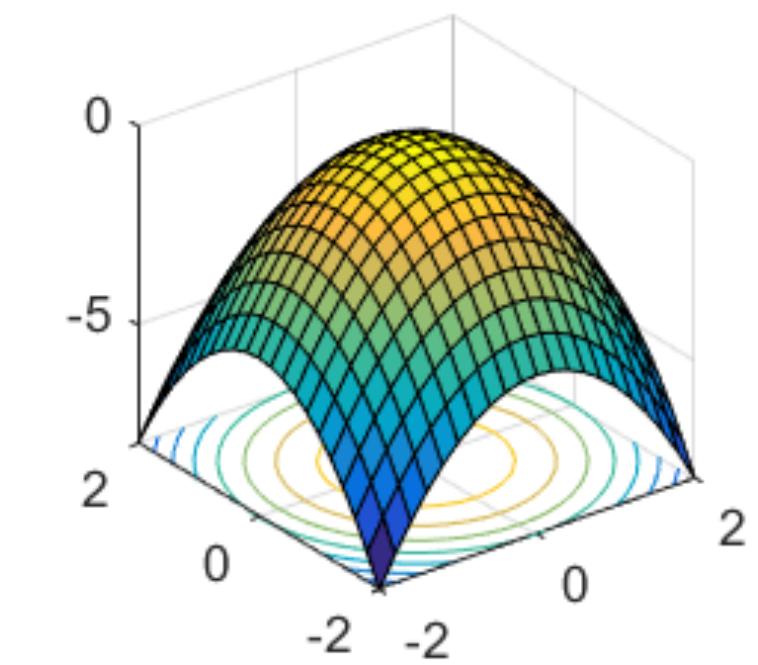
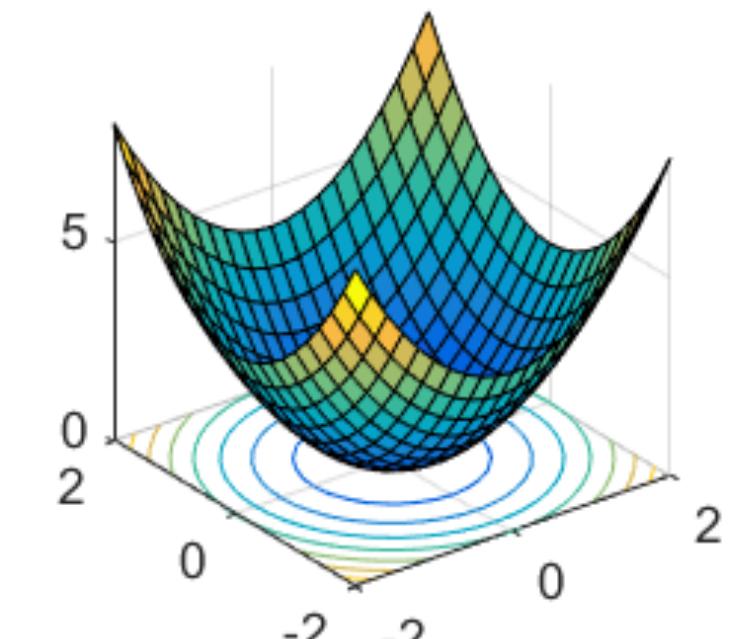
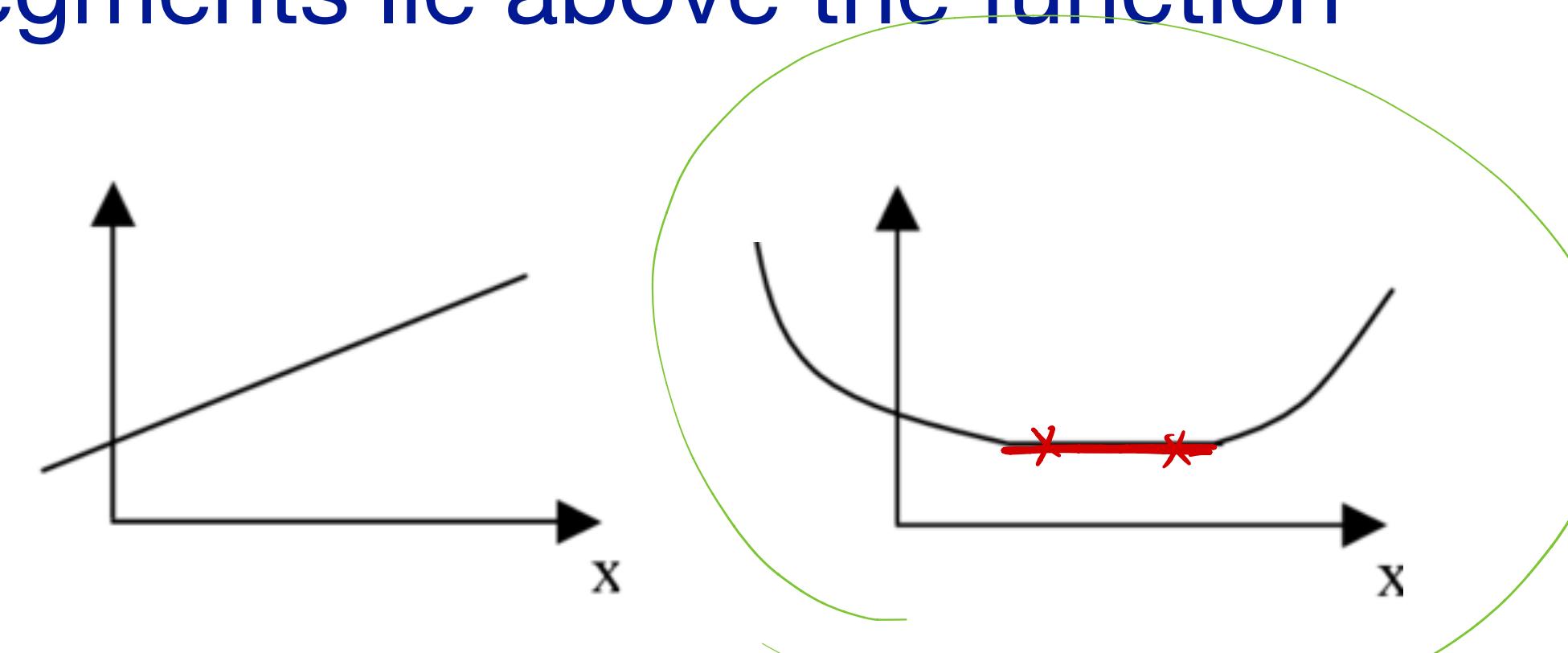
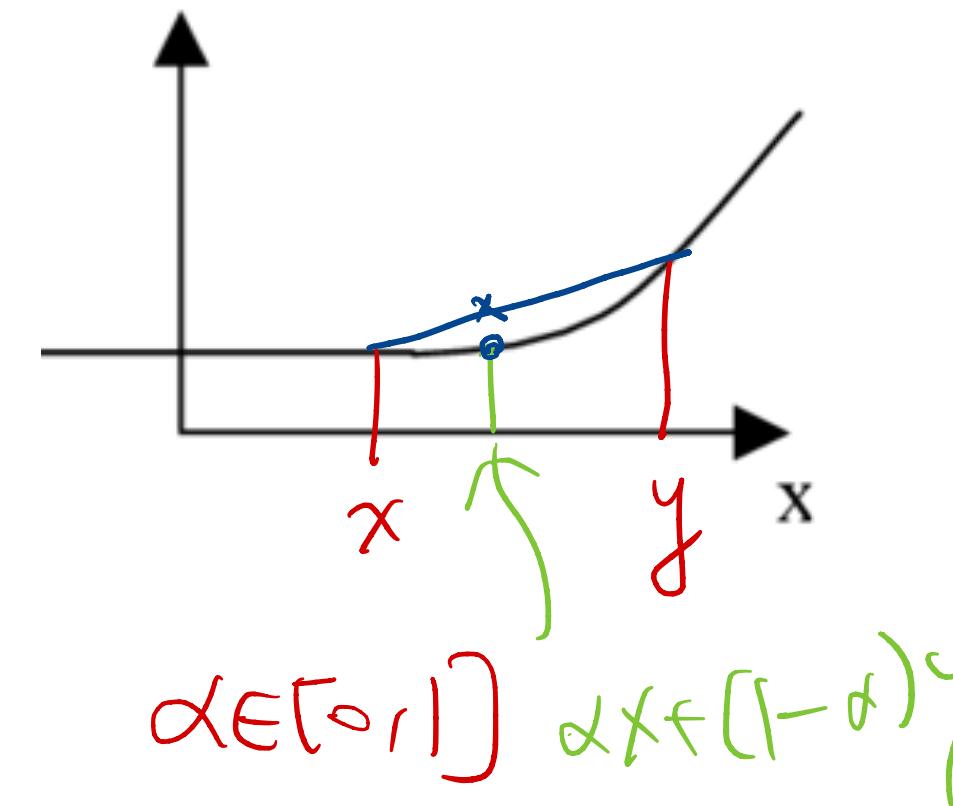
# Convex and Concave Functions

( $f$  may not be differentiable)

**Definition:** A function  $f: X \rightarrow \mathbb{R}$  is called a *convex function* if its domain  $X$  is a convex set and for any  $x, y \in X$  and any  $\alpha \in [0, 1]$ , we have

$$f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y)$$

**Intuition:** “All the segments lie above the function”



**Remark:** A function  $h$  is called concave if  $-h$  is convex

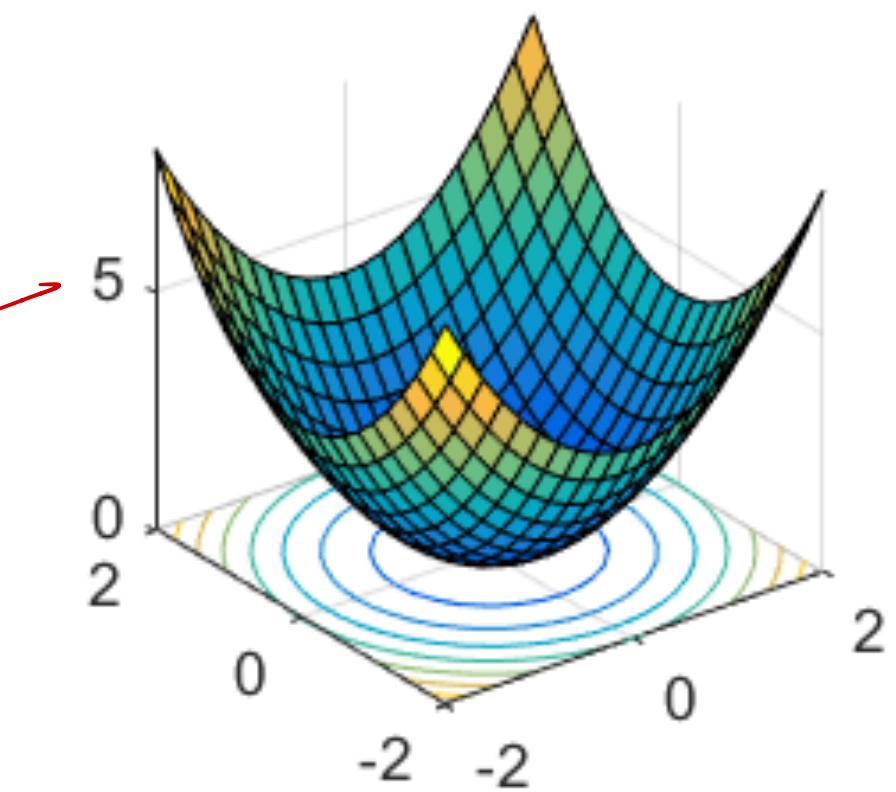
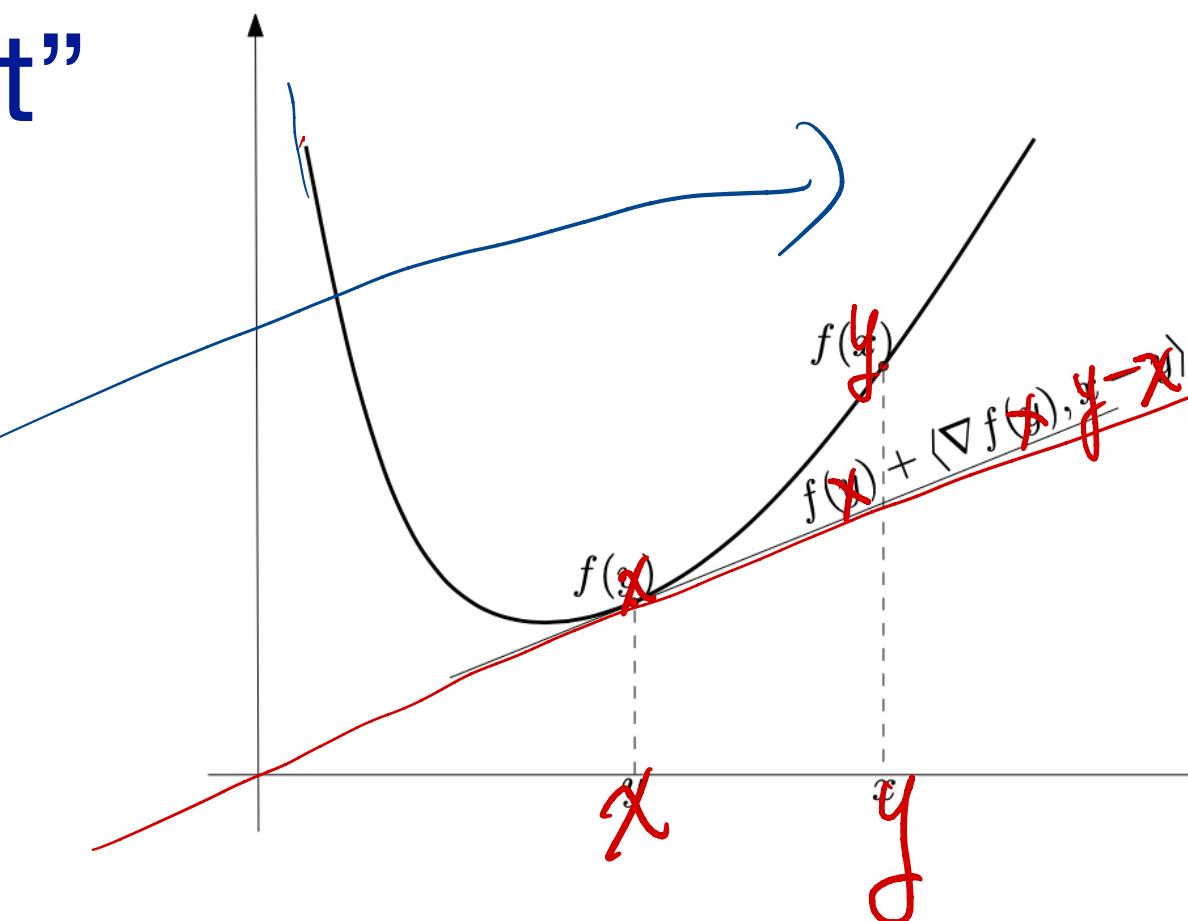
**Question:** Why do we need the domain  $X$  to be convex? To ensure that  $\alpha x + (1-\alpha)y \in X$  for all  $\alpha$ , all  $x, y$

# Characterizing Convex Functions Under Differentiability

- If  $f: X \rightarrow \mathbb{R}$  is differentiable, then  $f$  is convex if and only if  $X$  is a convex set and  $f(y) \geq f(x) + \nabla f(x)^T(y - x)$ , for all  $x, y \in X$ . (a.k.a. first-order condition of convexity)

Intuition: “The function lies above the tangent”

Tangent



- If  $f: X \rightarrow \mathbb{R}$  is twice differentiable, then  $f$  is convex if and only if  $X$  is a convex open set and  $\nabla^2 f(x) \succeq 0$ , for all  $x \in X$ .

$$\leq f(y) = f(x) + \nabla f(x)^T(y - x) + \frac{1}{2}(y - x)^T \nabla^2 f(x)(y - x)$$

Exercise: Prove the above two properties

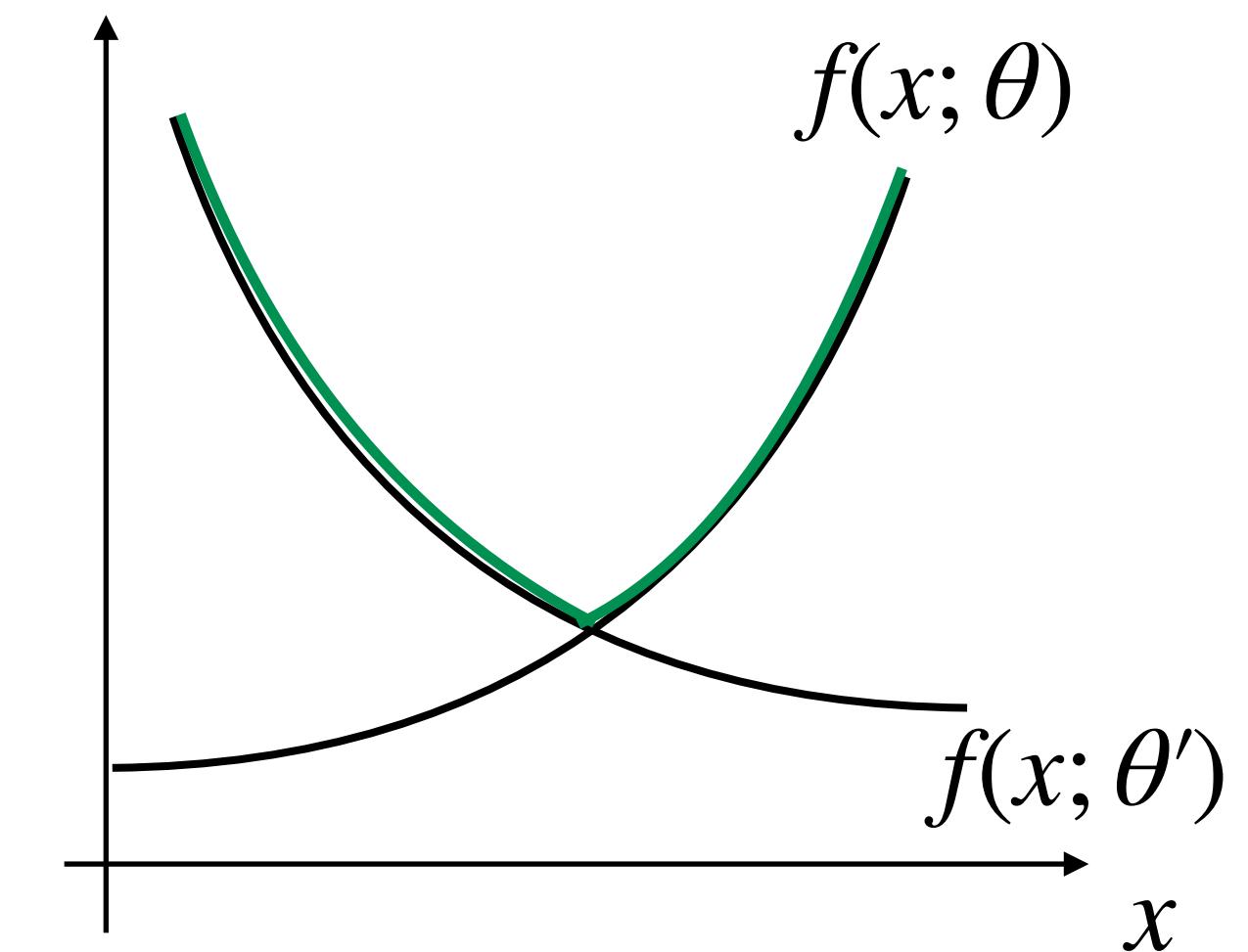
# Example: Pointwise Maximum of Convex Functions

The **pointwise maximum** of a family of convex functions is still convex

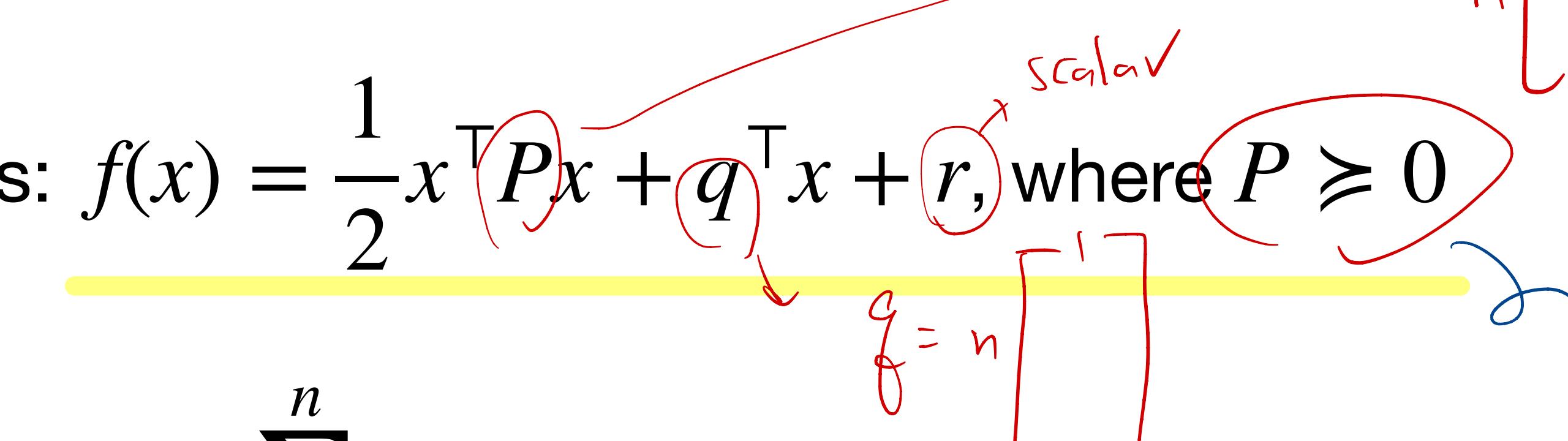
- Let  $f(x; \theta)$  be a convex function of  $x$  for every  $\theta \in \Theta$ , where  $\Theta$  is an arbitrary index set.

Define  $F(x) := \max_{\theta \in \Theta} f(x; \theta)$

- Question:** Is  $F(x)$  a convex function?
- 



# Popular Examples of Convex Functions

- ✓ Quadratic functions:  $f(x) = \frac{1}{2}x^T P x + q^T x + r$ , where  $P \geq 0$   

$$\nabla^2 f(x) = P$$
- Negative entropy:  $f(x) = \sum_{i=1}^n x_i \log x_i$ , where  $x$  is a probability vector
- Log-sum-exp:  $f(x) = \log(\exp(x_1) + \dots + \exp(x_n))$
- Log-determinant of pd matrices:  $f(X) = \log \det(X)$

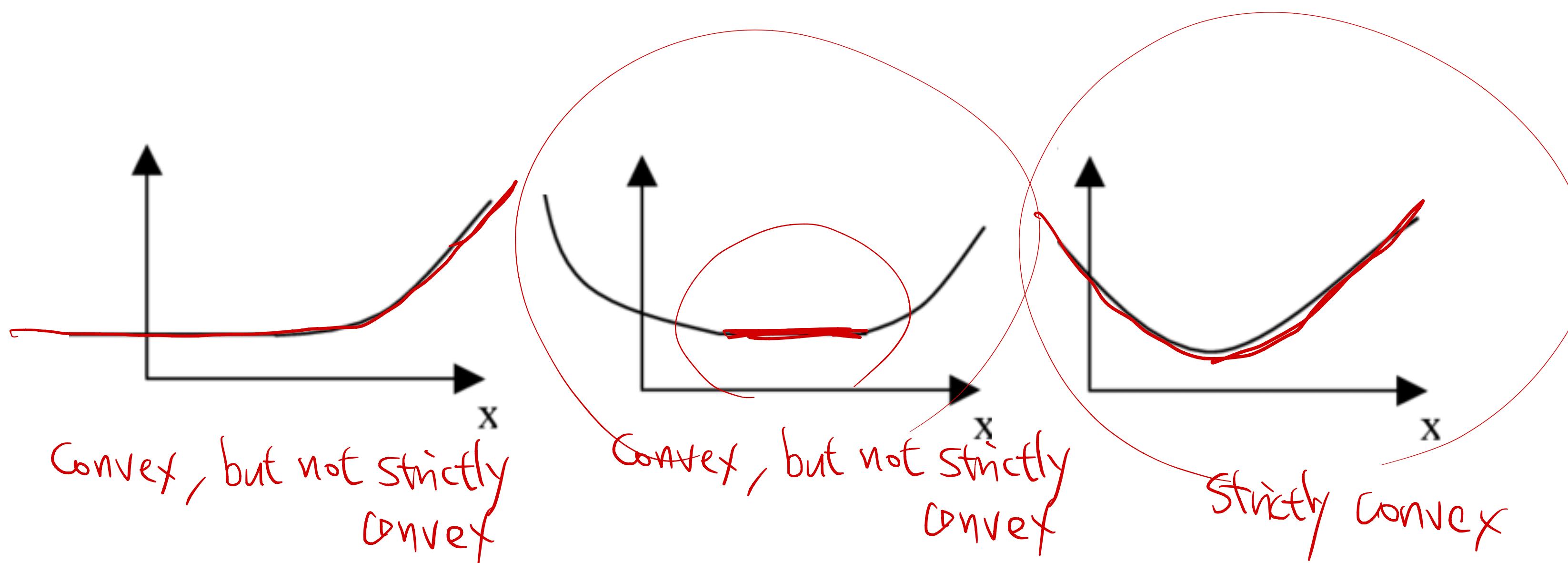
You will be asked to verify this in HW0 :)

# Strictly Convex Functions

**Definition:** A function  $f: X \rightarrow \mathbb{R}$  is called **strictly convex** if its domain  $X$  is a convex set and for any  $x, y \in X$  with  $x \neq y$  and any  $\alpha \in (0,1)$ , we have

$$f(\alpha x + (1 - \alpha)y) < \alpha f(x) + (1 - \alpha)f(y)$$

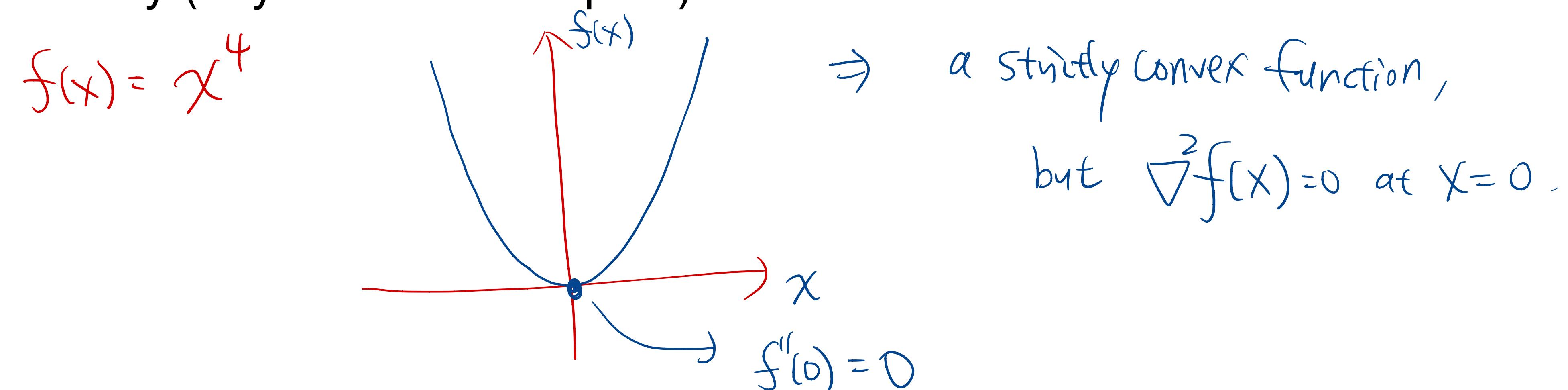
Intuition: “All the line segments lie strictly above the function”



# Characterizing Strictly Convex Functions Under Differentiability

- If  $f: X \rightarrow \mathbb{R}$  is twice differentiable, then  $f$  is strictly convex if  $X$  is a convex set and  $\underline{\nabla^2 f(x) > 0}$ , for all  $x \in X$ .

**Remark:** The condition  $\nabla^2 f(x) > 0$  is only sufficient but **not necessary** for strict convexity (any counterexample?)



**Remark:** We will mention a related concept “strong convexity” when discussing gradient descent in Lecture 4

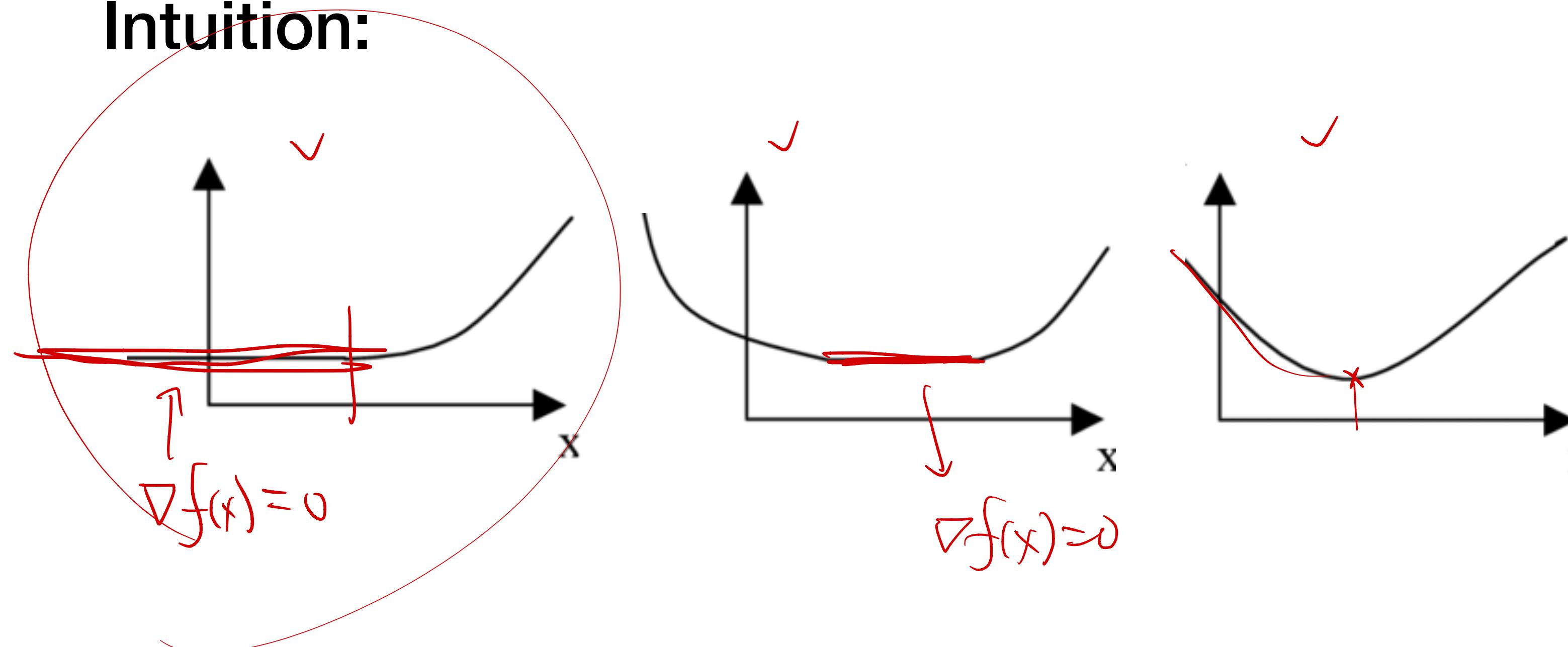
# Nice Properties of Convex Functions

- Fact 1: If  $f$  is convex over  $X$ , then a local minimum is also a global minimum  
(This makes convex optimization special!)
- Fact 2: If  $f$  is strictly convex over  $X$ , then there exists at most one global minimum  
(Could you explain these by the definition of strict convexity?)

# C3. First-Order Sufficient Condition for Global Optimality (Intuition)

**Theorem (FOSC):** If  $f: X \rightarrow \mathbb{R}$  is convex and the set  $X$  is convex, then  $\nabla f(x^*) = 0$  is sufficient for  $x^* \in X$  to be a global minimizer.

Intuition:



### C3. First-Order Sufficient Conditions for Global Optimality (Formally)

**Theorem (FOSC):** If  $f: X \rightarrow \mathbb{R}$  is **convex** and the set  $X$  is convex, then  $\nabla f(x^*) = 0$  is sufficient for  $x^* \in X$  to be a global minimizer.

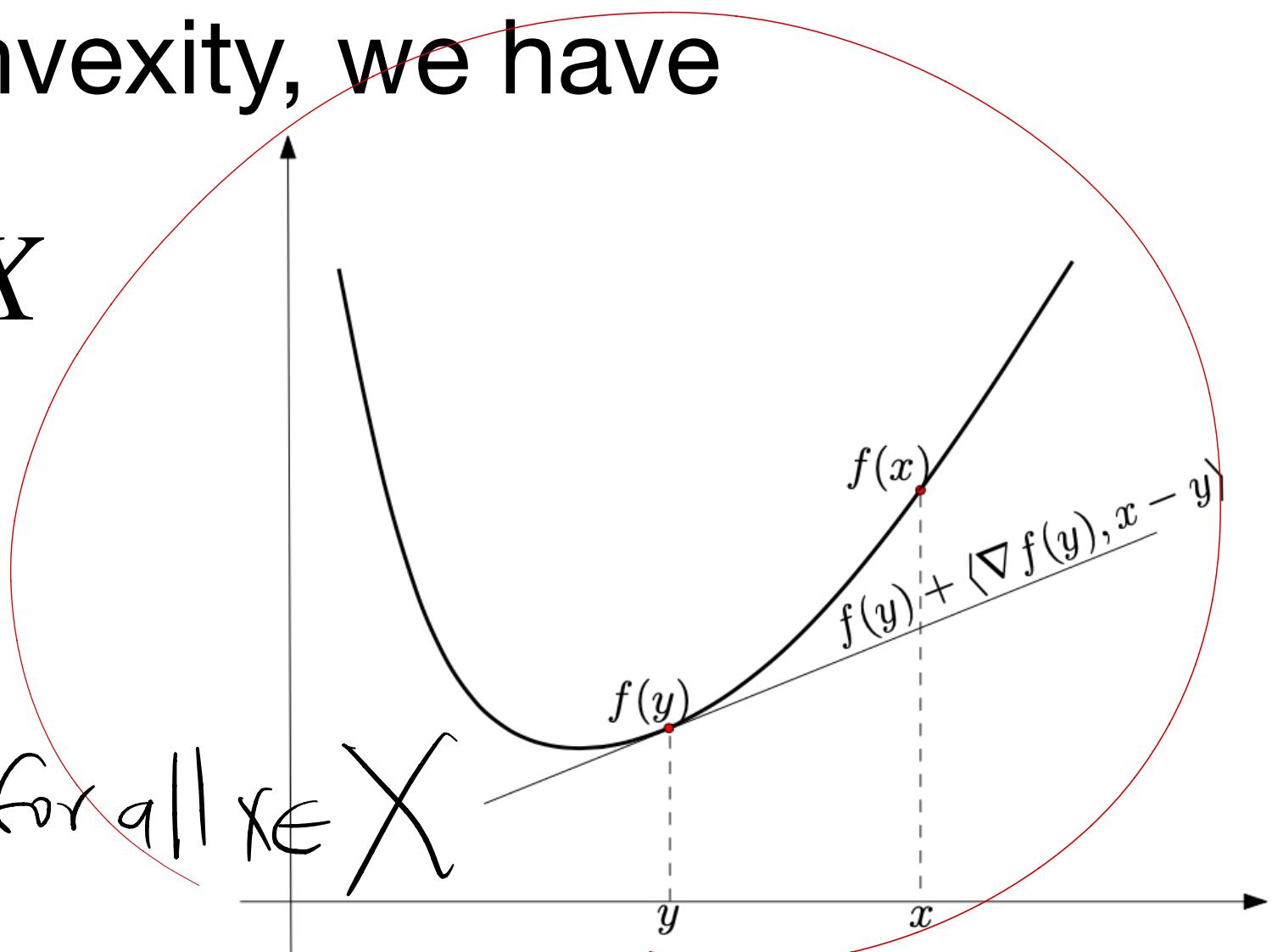
Moreover, if  $X$  is also an open set, then  $\nabla f(x^*) = 0$  is both necessary and sufficient for  $x^* \in X$  to be a global minimizer.

*Proof:* Let  $x^* \in X$  be a global minimizer. Then, by convexity, we have

$$f(x) \geq f(x^*) + \nabla f(x^*)^\top (x - x^*), \quad \text{for all } x \in X$$

"function lies above the tangent"

By the condition that  $\nabla f(x^*) = 0$ , then  $f(x) \geq f(x^*)$  for all  $x \in X$



**Question:** Why is "openness of domain  $X$ " needed?

## C4. Second-Order Sufficient Condition (SOSC) for Local Optimality

**Theorem (SOSC):** Let  $f: X \rightarrow \mathbb{R}$  be twice continuously differentiable. Then,

$x^* \in X$  is a strict local minimizer of  $f$  if  $x^*$  satisfies:

(i)  $\nabla f(x^*) = 0$  and (ii)  $\nabla^2 f(x^*) > 0$ .

$\nabla^2 f(x^*)$  is pd, and all the eigenvalues are positive

*Proof:*  $f(x^* + d) - f(x^*) = \nabla f(x^*)^\top d + \frac{1}{2} d^\top \nabla^2 f(x^*) d + o(\|d\|^2)$ , for all  $d \in \mathbb{R}^n$

$\lambda_{\min}$  denote the smallest eigenvalue

$$\begin{aligned} &= \underbrace{\nabla f(x^*)^\top d}_{0} + \underbrace{\frac{1}{2} d^\top \nabla^2 f(x^*) d}_{\geq 0} + o(\|d\|^2) \\ &\geq \frac{1}{2} \underbrace{\lambda_{\min}}_{>0} \|d\|^2 + o(\|d\|^2) \\ &> 0 \quad (\text{for sufficiently small } d) \end{aligned}$$

Moreover, if your  $f$  is  
a convex function,

then (i)-(ii) are sufficient

for  $x^*$  to be <sup>strict</sup> global  
minimizer

Therefore,  $x^*$  is a strict local minimizer.

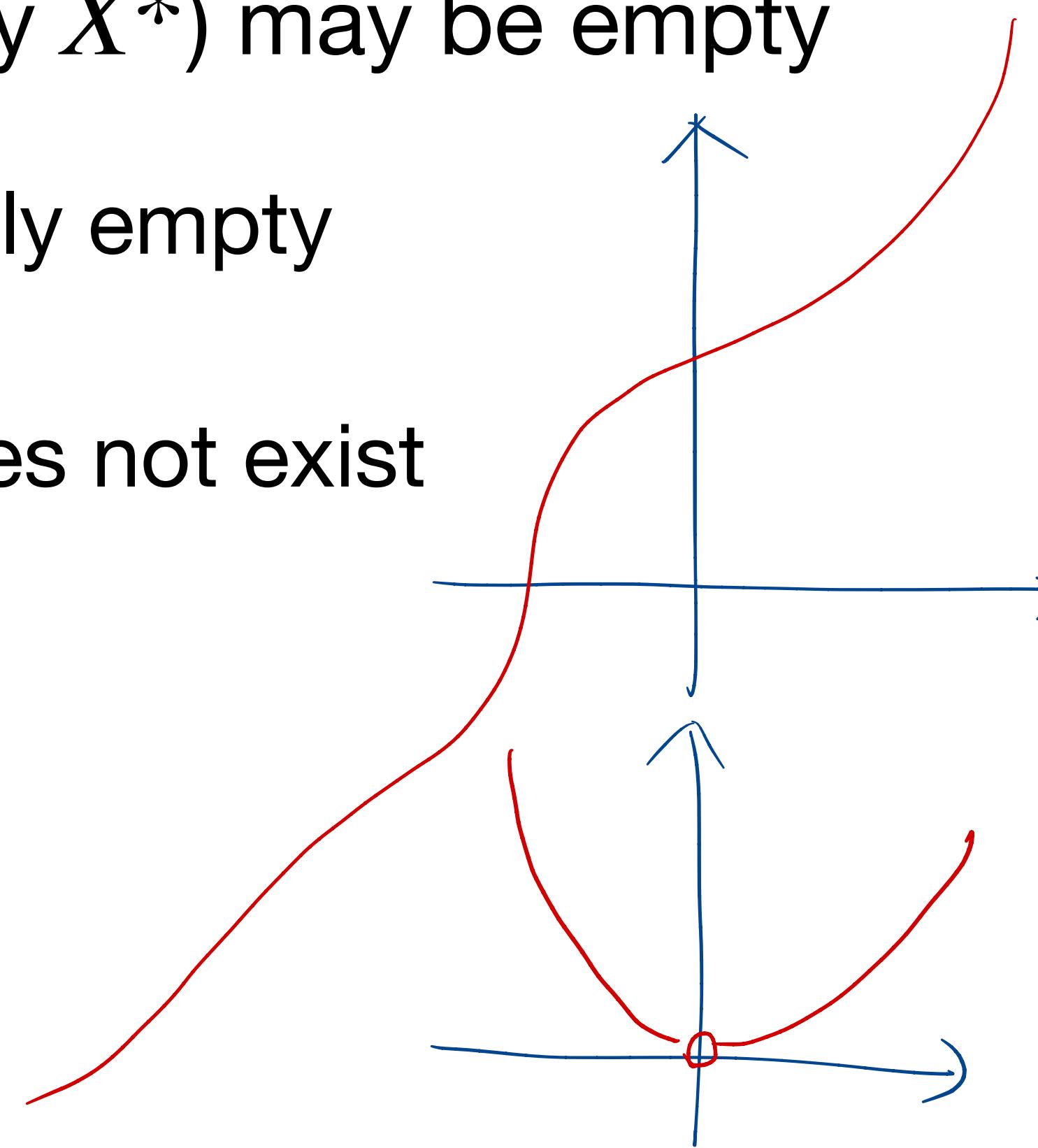
# Set of Optimal Solutions

- Fact 1: The set of optimal solutions (denoted by  $X^*$ ) may be empty

**Example:** If the domain  $X$  is empty, then  $X^*$  is surely empty

**Example:** When only the “inf” exists but “min” does not exist

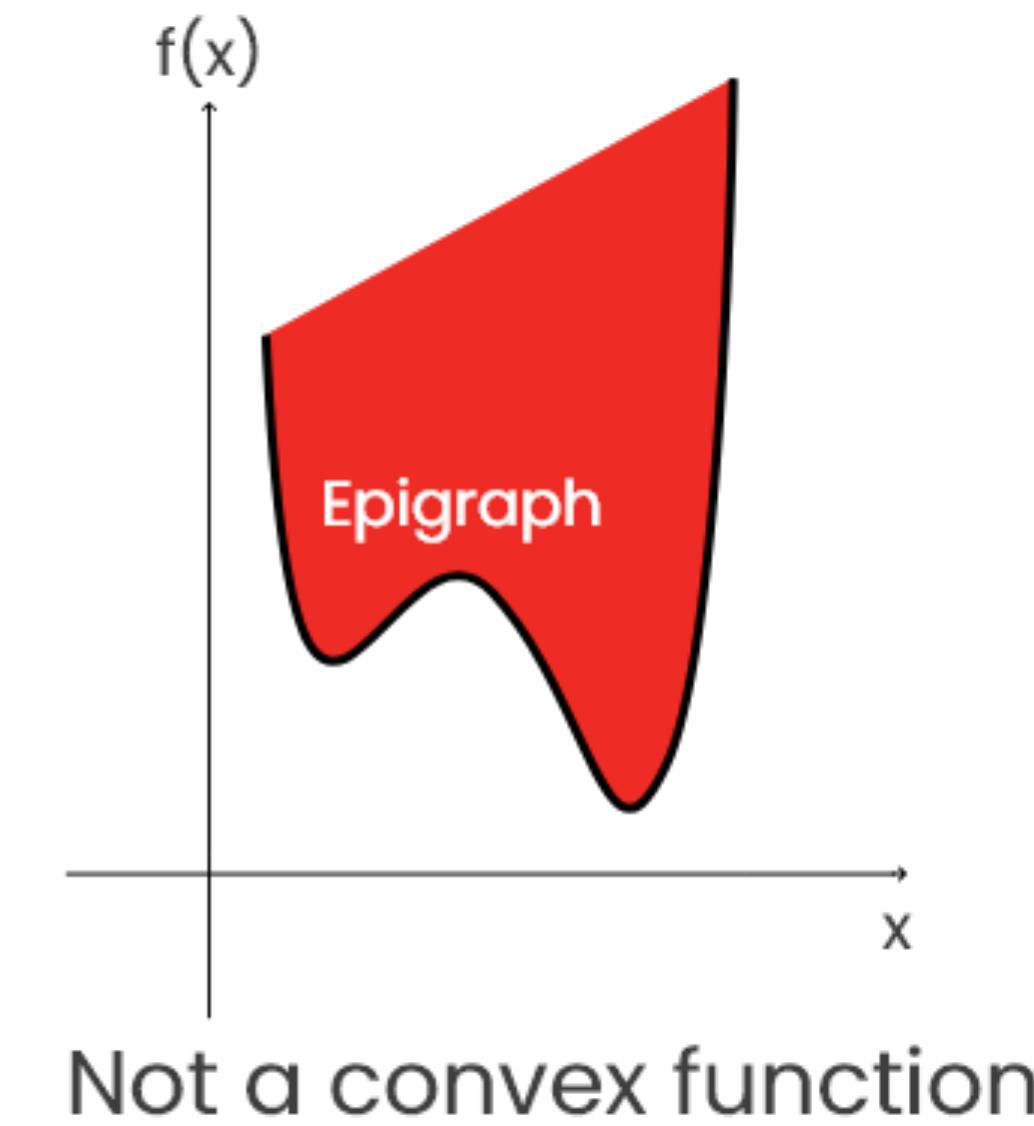
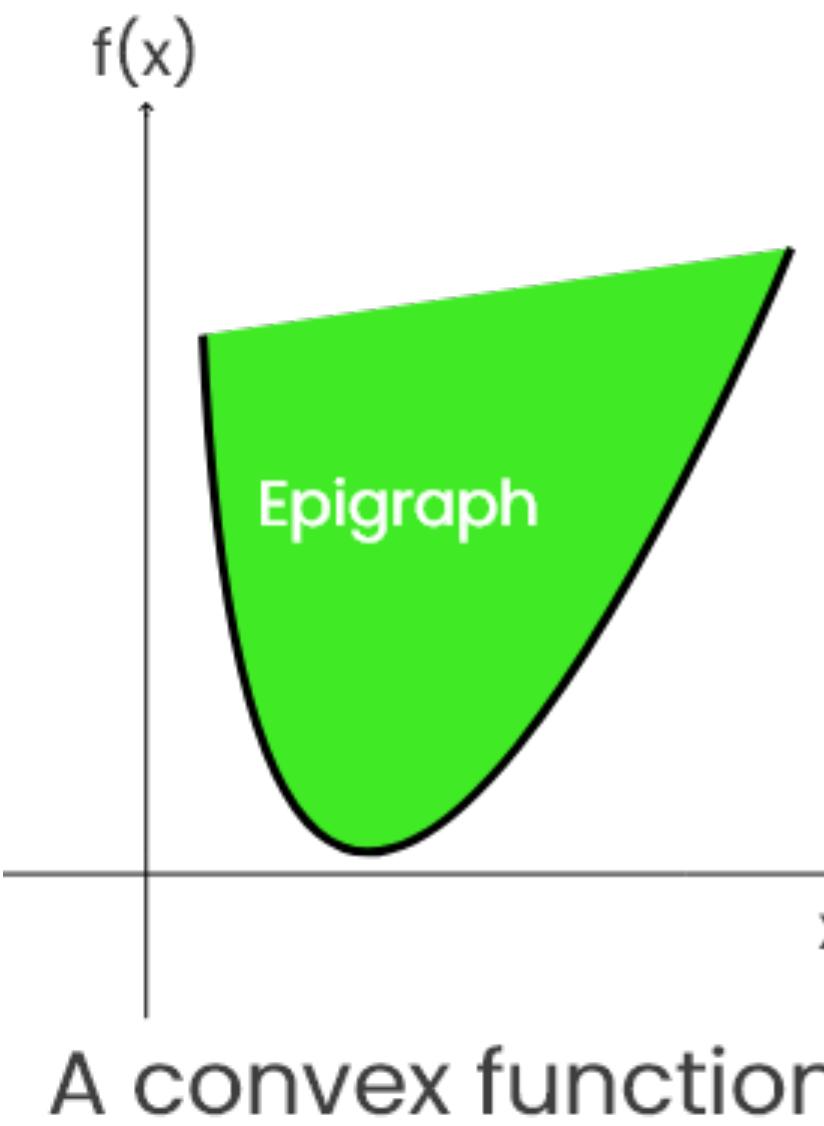
$$f(x) = -\|x\|^2$$



- Fact 2: Suppose the domain  $X$  is a convex set and  $f$  is a convex function. If  $X^*$  is not empty, then  $X^*$  must be a convex set (why?)

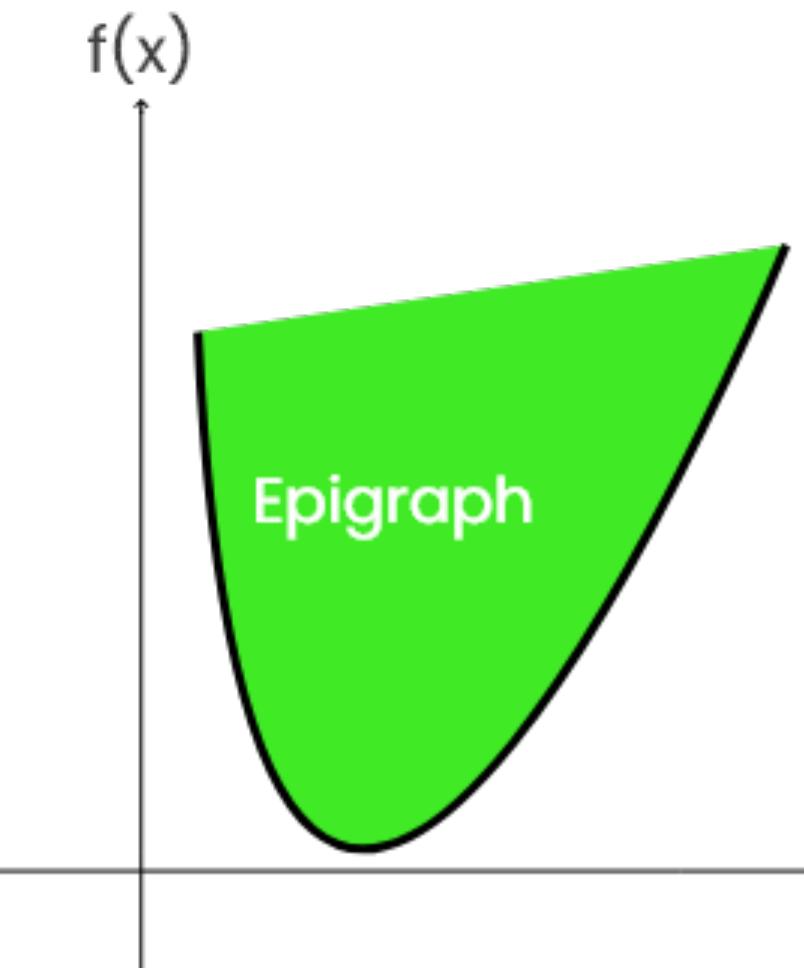
# Remark: Definition of Convex Functions via *Epigraph*

Intuition: Can you find anything special in the regions in green and red?

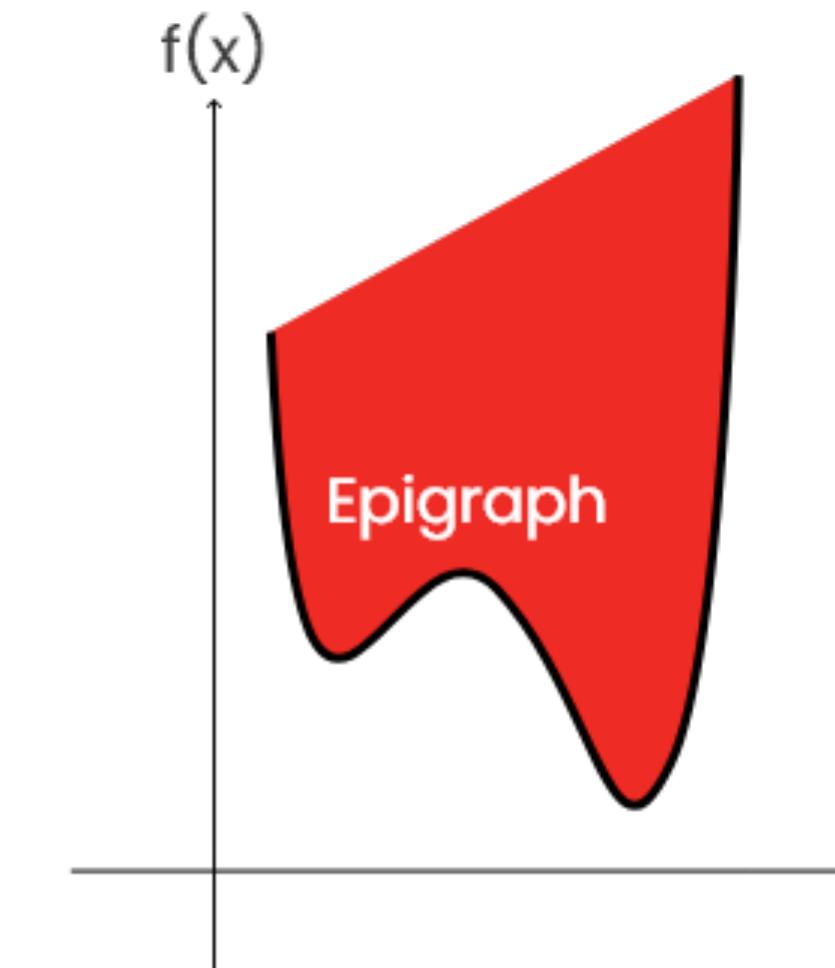


# Remark: Definition of Convex Functions via *Epigraph*

Intuition: Can you find anything special in the regions in green and red?



A convex function



Not a convex function

**Definition:** The **epigraph** of a function

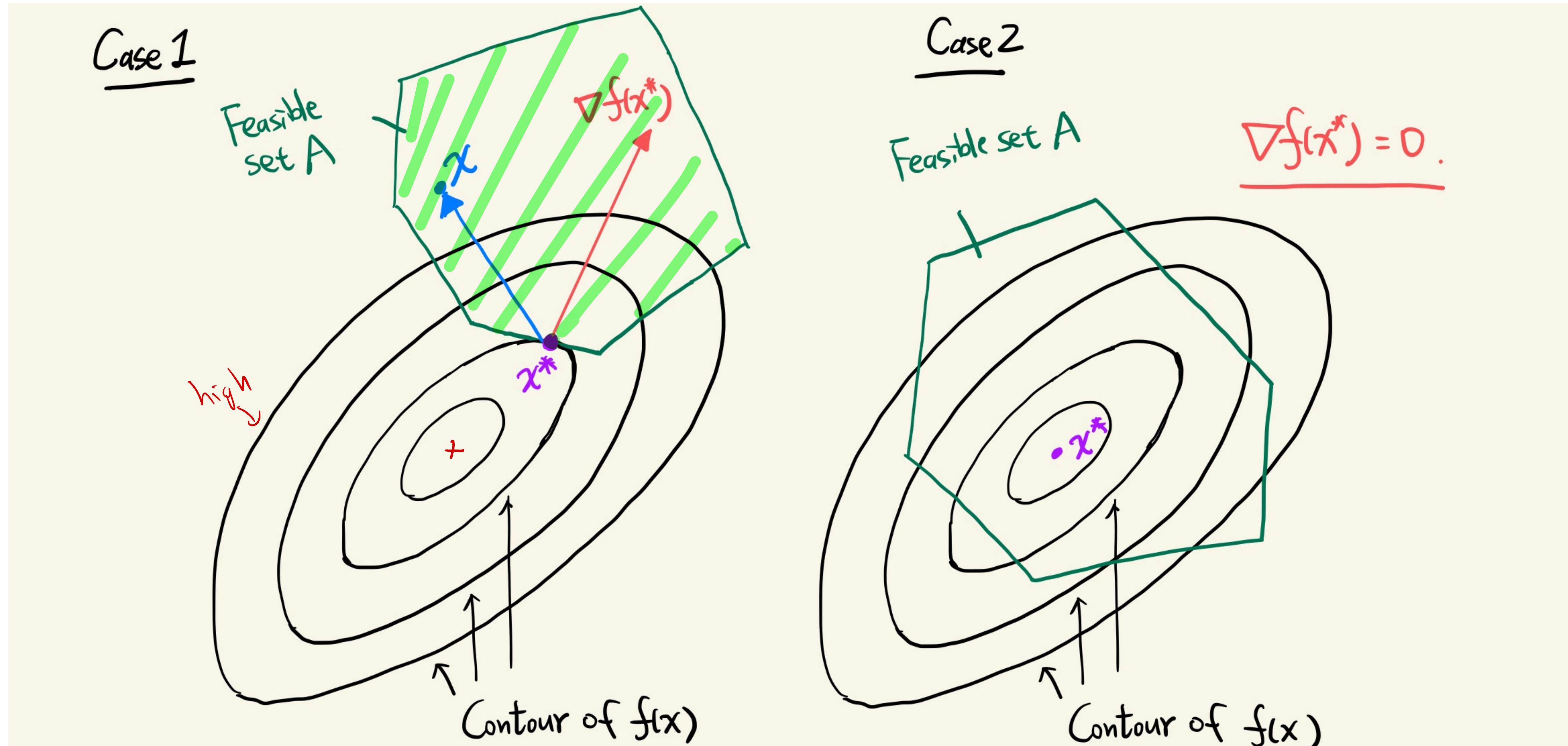
$f : X \rightarrow \mathbb{R}$  is defined as

$$\text{epi}(f) := \{(x, \gamma) : x \in X, f(x) \leq \gamma\}$$

**Property:** A function  $f$  is convex if and only if its epigraph is a convex set.

# Optimality Conditions for Constrained Problems?

Let's start with some intuition! By “constrained”: Feasible set  $A \subseteq X$



# C5 & C6. Optimality Conditions for Constrained Problems (Formally)

**Theorem:** Let  $f: X \rightarrow \mathbb{R}$  be continuously differentiable and let  $A \subseteq X$  be a convex feasible set.

(C5) If  $x^*$  is a local minimizer of  $f$  over  $A$ , then we have (Necessary)

$$\nabla f(x^*)^\top (x - x^*) \geq 0, \quad \forall x \in A \quad \dots\dots (*)$$

(C6) If  $f$  is a convex function over  $A$ , then the condition  $(*)$  is also sufficient for  $x^*$  to be a global minimizer of  $f$  over  $A$  (Sufficient)

**Remark:** If  $A = \mathbb{R}^n$  (i.e., unconstrained), then  $(*)$  reduces to  $\nabla f(x^*) = 0$  (why?)

## C5 & C6. Optimality Conditions for Constrained Problems (Formally)

(C5) If  $x^*$  is a local minimizer of  $f$  over  $A$ , then we have

$$\nabla f(x^*)^\top (x - x^*) \geq 0, \quad \forall x \in A \quad \dots\dots (*)$$

---

**Proof of (C5): Prove this by contradiction**

Step 1: Suppose there exists some  $x \in A$  such that  $\nabla f(x)^\top (x - x^*) < 0$

Step 2: By Taylor's Theorem, for any  $\varepsilon > 0$ , we have

$$f(x^* + \varepsilon(x - x^*)) = f(x^*) + \varepsilon \nabla f(x')^\top (x - x^*)$$

where  $x' = x^* + \alpha\varepsilon(x - x^*), \alpha \in [0,1]$

Step 3: Since  $\nabla f(x)$  is continuous, we have that for all sufficiently small  $\varepsilon$

- (i)  $\nabla f(\textcolor{red}{x}')^\top (x - x^*) < 0$
- (ii)  $x' \in A$  (why?)

These imply that  $f(x^* + \varepsilon(x - x^*)) < f(x^*)$ , for all sufficiently small  $\varepsilon > 0$

This contradicts the fact that  $x^*$  is a local minimizer

## C5 & C6. Optimality Conditions for Constrained Problems (Formally)

(C6) If  $f$  is a **convex** function over  $A$ , then the condition (\*) is also **sufficient** for  $x^*$  to be a global minimizer of  $f$  over  $A$

---

Proof of (C6):

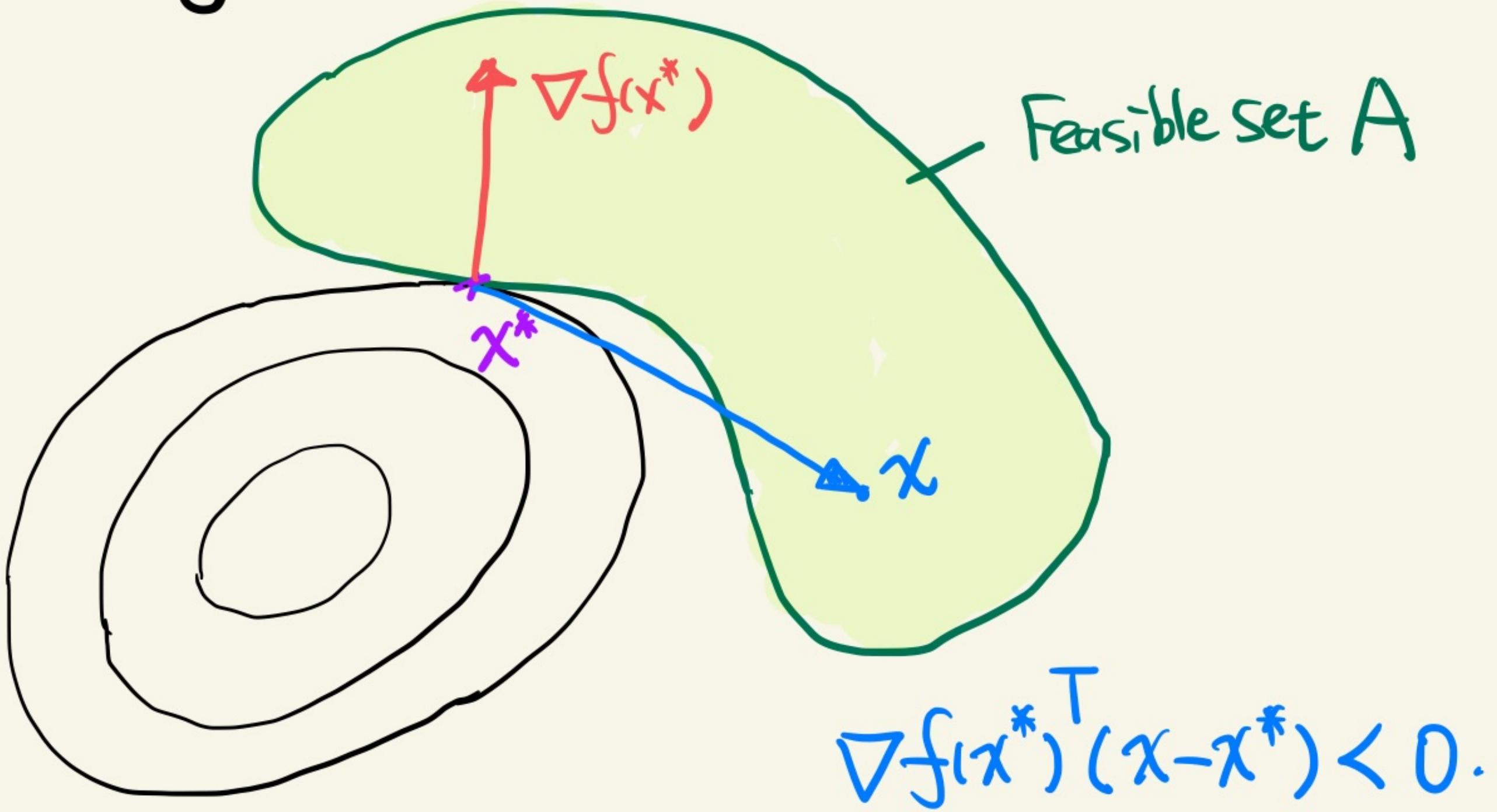
Step 1. By the convexity of  $f$ , we have

$$\begin{aligned} f(x) &\geq f(x^*) + \nabla f(x^*)^\top (x - x^*), \text{ for all } x \in A \\ &\quad \underline{\geq 0, \text{ for all } x \in A} \end{aligned}$$

Step 2. Therefore, we have  $f(x) \geq f(x^*)$ , for all  $x \in A$

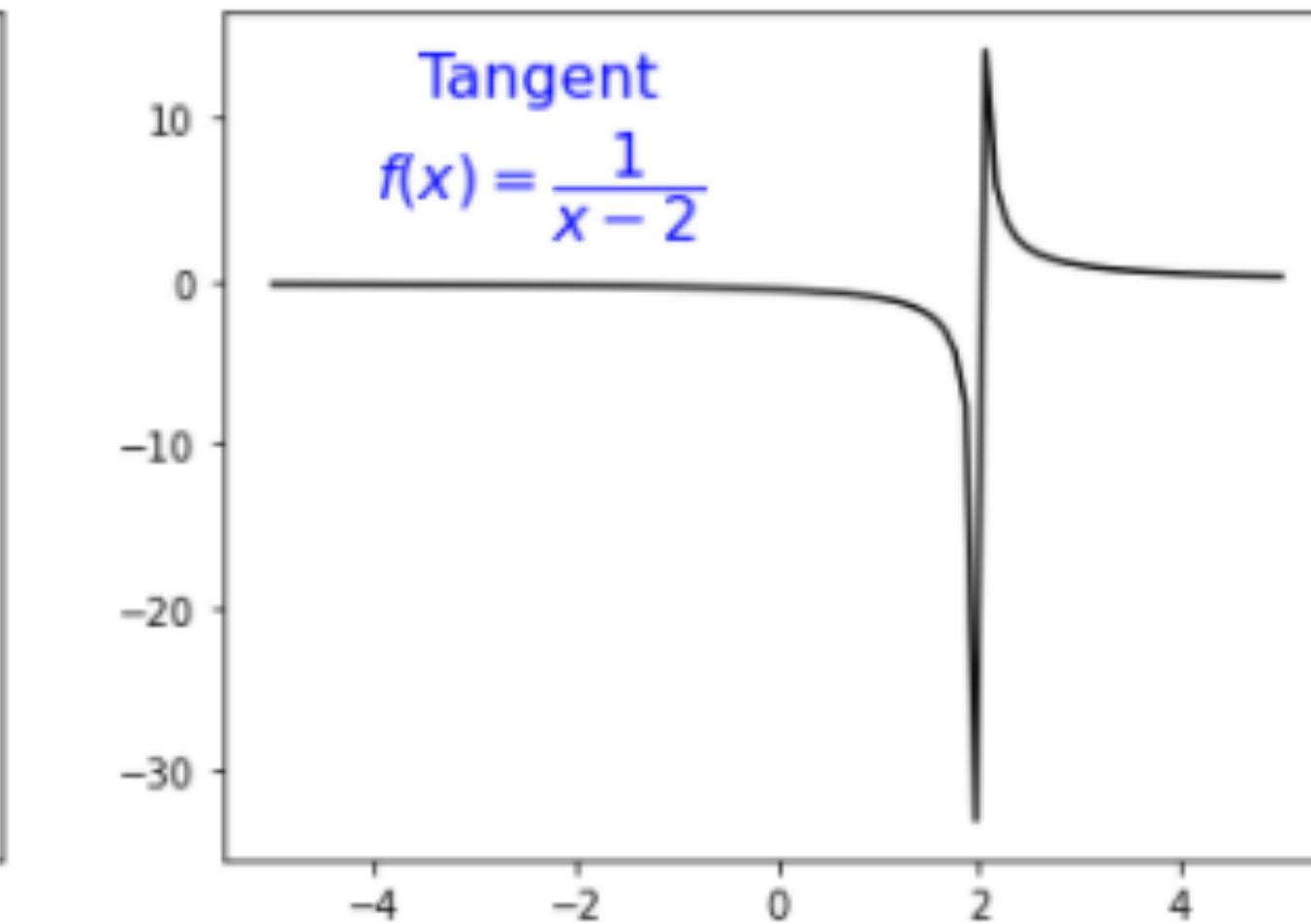
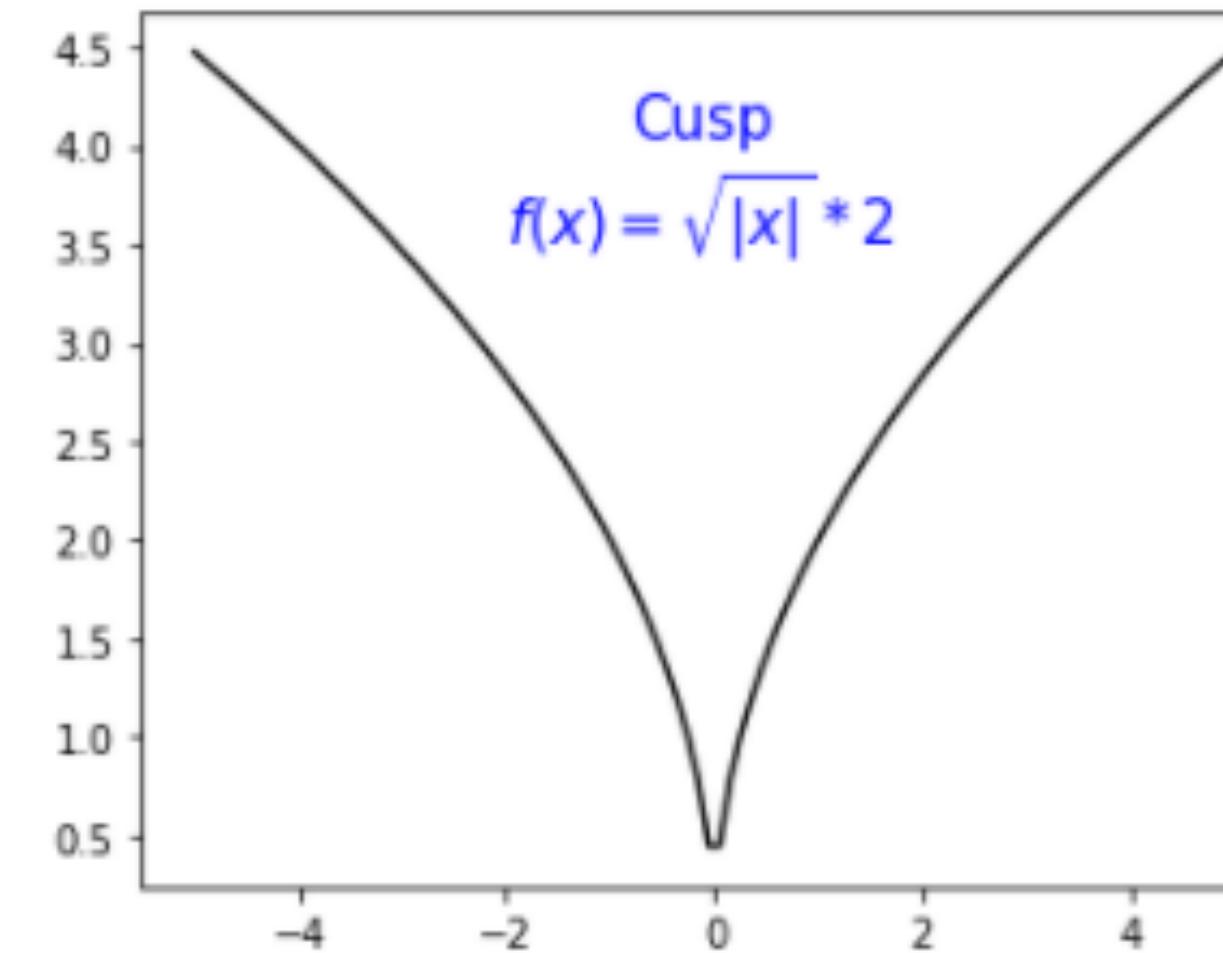
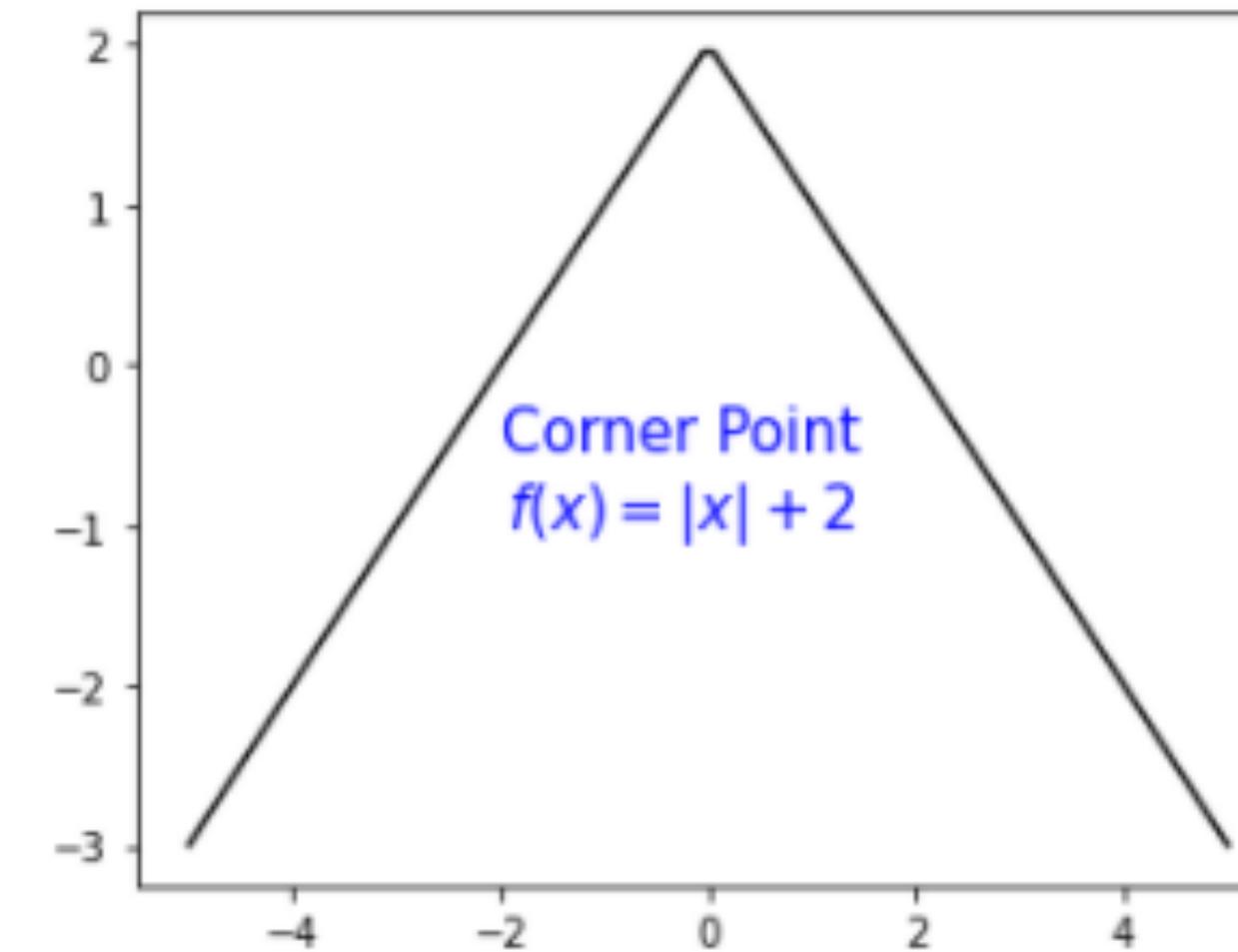
# Why "Convex" Feasible Sets?

The necessary condition  $\nabla f(x^*)^T(x - x^*) \geq 0$  may fail when A is not convex



# Optimality Conditions Beyond Differentiability?

- What if there are some non-differentiable points in  $f$ ?



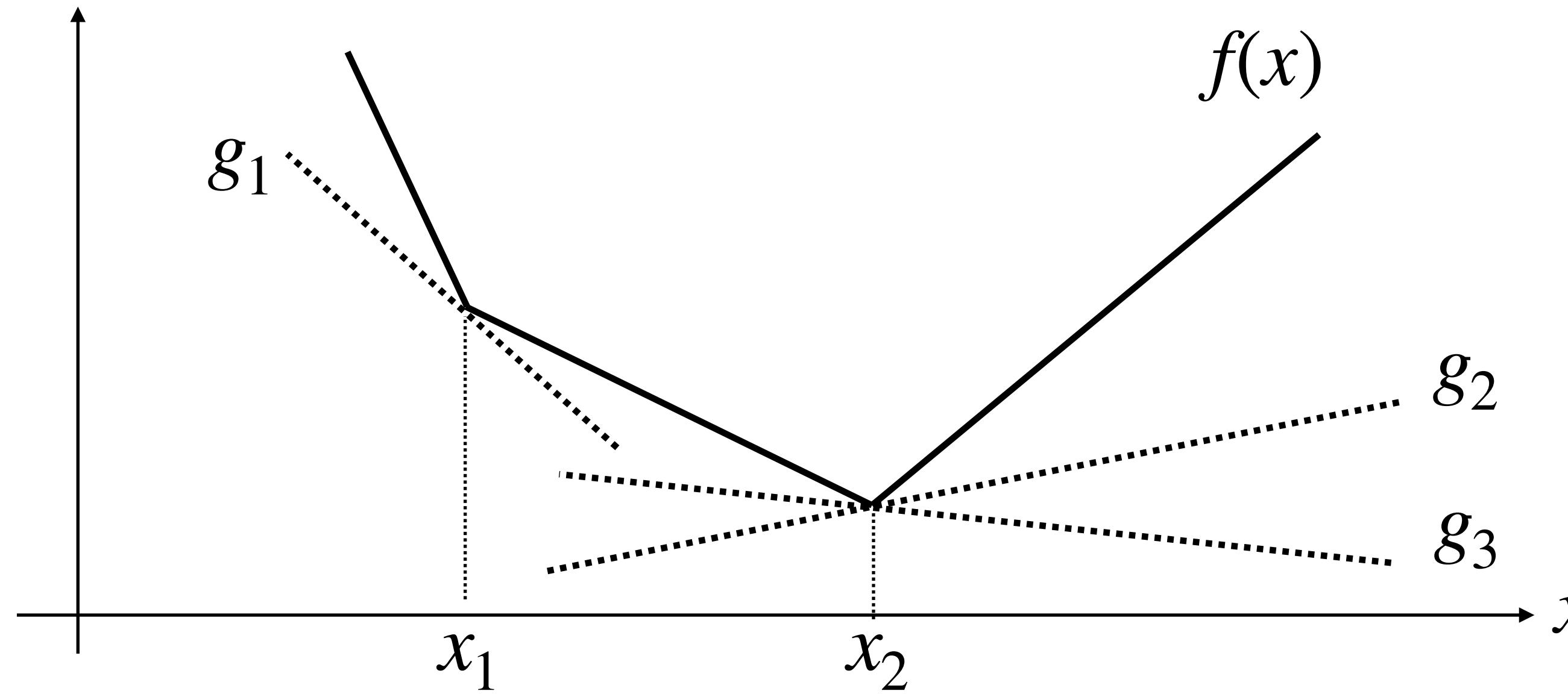
Could we extend the notion of gradients and optimality conditions?

(Figure Credit: <https://python.plainenglish.io/descent-carefully-on-your-gradient-c0f030ddef81>)

# **Subgradients and Subdifferential**

**(The slides are partially adapted from Stephen Boyd's EE364B)**

# Subgradients



(In plain English:  $f(x) + g^\top(z - x)$  is a global underestimate)

**Definition:**  $g \in \mathbb{R}^n$  is a *subgradient* of  $f$  (possibly non-convex) at  $x$  if

$$f(z) \geq f(x) + g^\top(z - x), \text{ for all } z \in X$$

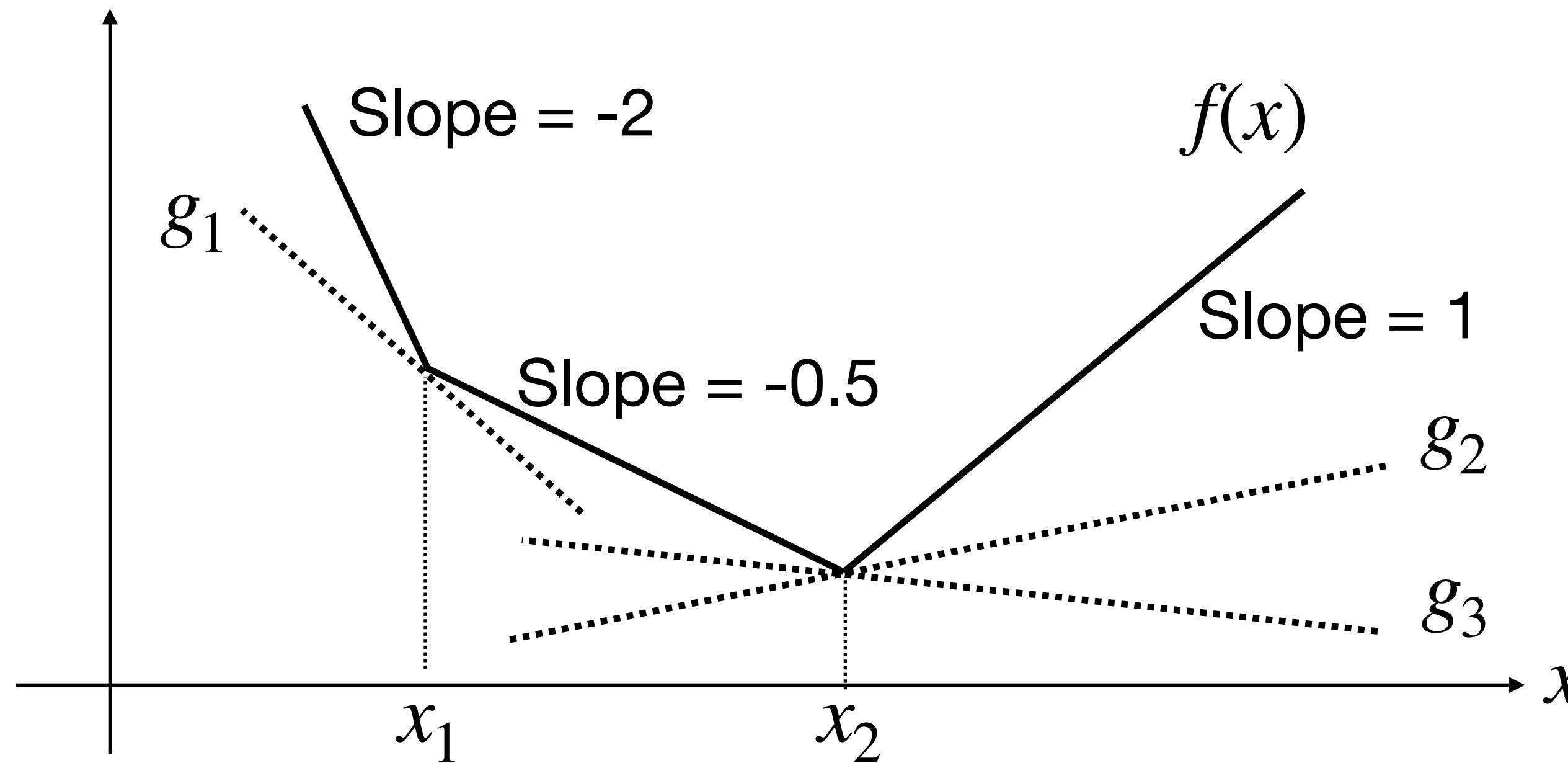
**Question:** If  $f$  is differentiable, then could you find a natural subgradient?

# Subdifferentials

**Definition:** The *subdifferential* of  $f$  at  $x$ , denoted by  $\partial f(x)$ , is defined as the set of all subgradients of  $f$  at  $x$ .

**Example:**

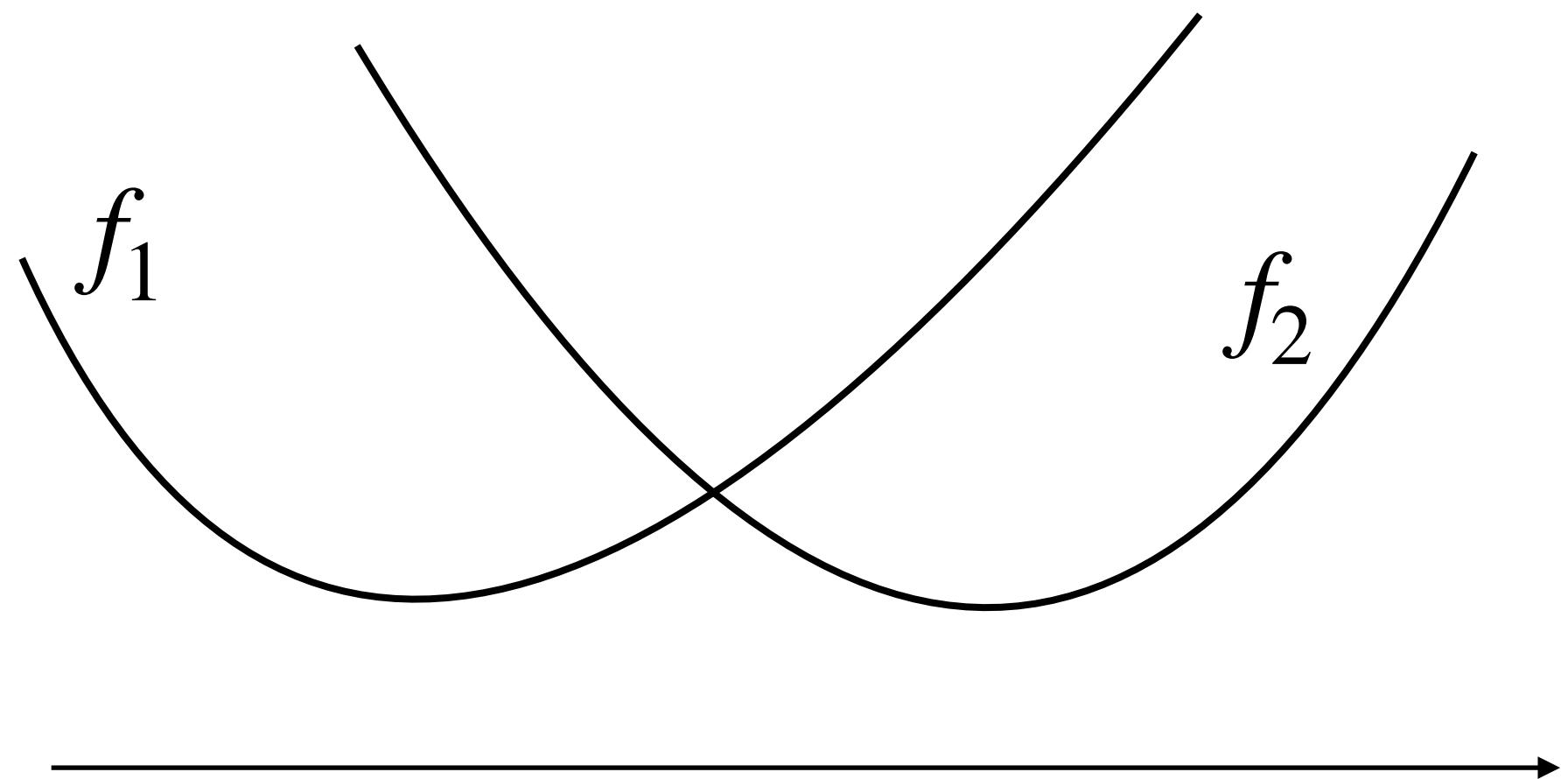
- Subdifferential at  $x_1$ ?



- Subdifferential at  $x_2$ ?

# More Examples of Subdifferentials

Suppose  $f = \max\{f_1, f_2\}$ , where  $f_1 : \mathbb{R} \rightarrow \mathbb{R}$  and  $f_2 : \mathbb{R} \rightarrow \mathbb{R}$  are both convex and differentiable



Subdifferential of  $f$ ?

- For  $x$  with  $f_1(x) > f_2(x)$ :
- For  $x$  with  $f_1(x) < f_2(x)$ :
- For  $x$  with  $f_1(x) = f_2(x)$ :

# Subdifferentials of Convex Functions

If  $f: X \rightarrow \mathbb{R}$  is convex, then  $\partial f(x)$  has some nice properties

- If  $x$  is in the relative interior of  $X$ , then  $\partial f(x) \neq \emptyset$
- If  $f$  is differentiable at  $x$ , then  $\partial f(x) = \{ \nabla f(x) \}$
- If  $\partial f(x) = \{g\}$ , then  $f$  is differentiable and  $g = \nabla f(x)$
- If  $f$  is differentiable at  $x$ , then  $\partial f(x) = \{ \nabla f(x) \}$

# Basic Calculus Rules of Subdifferentials

- **Nonnegative scaling:** For any  $\alpha > 0$ ,  $\partial(\alpha f)(x) = \{\alpha g : g \in \partial f(x)\}$
- **Addition (General):**  $\partial f_1(x) + \partial f_2(x) \subset \partial(f_1 + f_2)(x)$  (Vice versa? See the next page)  
(Set addition / Minkowski addition)
- **Addition (Convex cases):** If  $f_1, f_2$  are convex, then  $\partial f_1(x) + \partial f_2(x) = \partial(f_1 + f_2)(x)$

(For more properties of subdifferentials, please check  
Stephen Boyd's slides for Lecture 1 of EE364B)

# We Do Not Have $\partial(f_1 + f_2)(x) = \partial f_1(x) + \partial f_2(x)$ in General

Let's construct two functions  $f_1 : \mathbb{R} \rightarrow \mathbb{R}, f_2 : \mathbb{R} \rightarrow \mathbb{R}$

$$f_1(x) := \begin{cases} -\sqrt{x}, & \text{if } x \geq 0 \\ -x + \sqrt{-x}, & \text{if } x < 0 \end{cases}$$

$$f_2(x) := \begin{cases} x + \sqrt{x}, & \text{if } x \geq 0 \\ -\sqrt{-x}, & \text{if } x < 0 \end{cases}$$

**Exercise:** Try to verify the following

- $\partial f_1(0) = \emptyset$  and  $\partial f_2(0) = \emptyset$
- $\partial(f_1 + f_2)(0) = [-1, 1]$

# C7. Optimality Conditions Revisited: Without Differentiability

**Theorem (Fermat's Rule):** Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  (not necessarily differentiable).

Then, we have

$$\arg \min f = \{x \in \mathbb{R}^n : 0 \in \partial f(x)\}$$

Proof: (1) RHS  $\subseteq$  LHS

(2) LHS  $\subseteq$  RHS

# Example: Indicator Function

Given any set  $A \subset \mathbb{R}^n$ , let  $\mathbf{1}_A$  be the indicator function for  $A$ :

$$\mathbf{1}_A(x) := \begin{cases} 0, & \text{if } x \in A, \\ \infty, & \text{otherwise.} \end{cases}$$

**Property:**  $A$  is a convex set if and only if  $\mathbf{1}_A(x)$  is convex

---

**Question:** Subdifferential of  $\mathbf{1}_A(x)$  for any  $x \in X$ ?

# Connecting Indicator Functions and Constrained Problems

Given any set  $A \in \mathbb{R}^n$ , let  $\mathbf{1}_A$  be the indicator function for  $A$ :

$$\mathbf{1}_A(x) := \begin{cases} 0, & \text{if } x \in A, \\ \infty, & \text{otherwise.} \end{cases}$$

---

Convert a *constrained* problem into an *unconstrained* one:

$$\begin{array}{l} \min f(x) \\ \text{subject to } x \in A \subseteq X \end{array} \xrightarrow{\hspace{1cm}} \min_{x \in X} f(x) + \mathbf{1}_A(x)$$

An alternative derivation of “optimality condition for convex constrained problems”

- If  $f$  is convex and differentiable, then  $x^*$  is a global minimizer iff  $0 \in \partial(f + \mathbf{1}_A)(x^*)$
- $\partial(f + \mathbf{1}_A)(x^*) = \nabla f(x^*) + \partial\mathbf{1}_A(x^*) = \nabla f(x^*) + \{g : g^\top(y - x^*) \leq 0, \forall y \in X\}$
- Hence,  $\nabla f(x^*)^\top(y - x^*) \geq 0, \forall y \in X$  !

# Remark: Indicator Functions and Constrained Problems

- $0 \in \nabla f(x^*) + \partial \mathbf{1}_A(x^*)$  is an elegant and very general condition for optimality in convex optimization problems
- However, it is not always easy to play with
- We will discuss some simpler conditions (e.g., KKT conditions) later!

# Appendix

# Mean Value Theorem (For Multivariate Functions)

**Mean Value Theorem:** Let  $f: X \rightarrow \mathbb{R}$  be a differentiable function. Let  $a, b$  be points in  $X$  such that the line segment of  $a, b$  lies in  $U$ . Then, there must exist some  $z = \alpha a + (1 - \alpha)b$  with  $\alpha \in [0,1]$  such that

$$f(b) - f(a) = \nabla f(z)^T (b - a)$$

The proof can be found at: <https://links.uwaterloo.ca/amath731docs/meanvalue.pdf>