

# Introduction to Machine Learning CentraleSupélec Paris — Fall 2017

## 2. Elements of convex optimization

**Chloé-Agathe Azencott**

Centre for Computational Biology, Mines ParisTech  
[chloe-agathe.azencott@mines-paristech.fr](mailto:chloe-agathe.azencott@mines-paristech.fr)



# Why talk about optimization?

- Supervised ML: **empirical risk minimization**

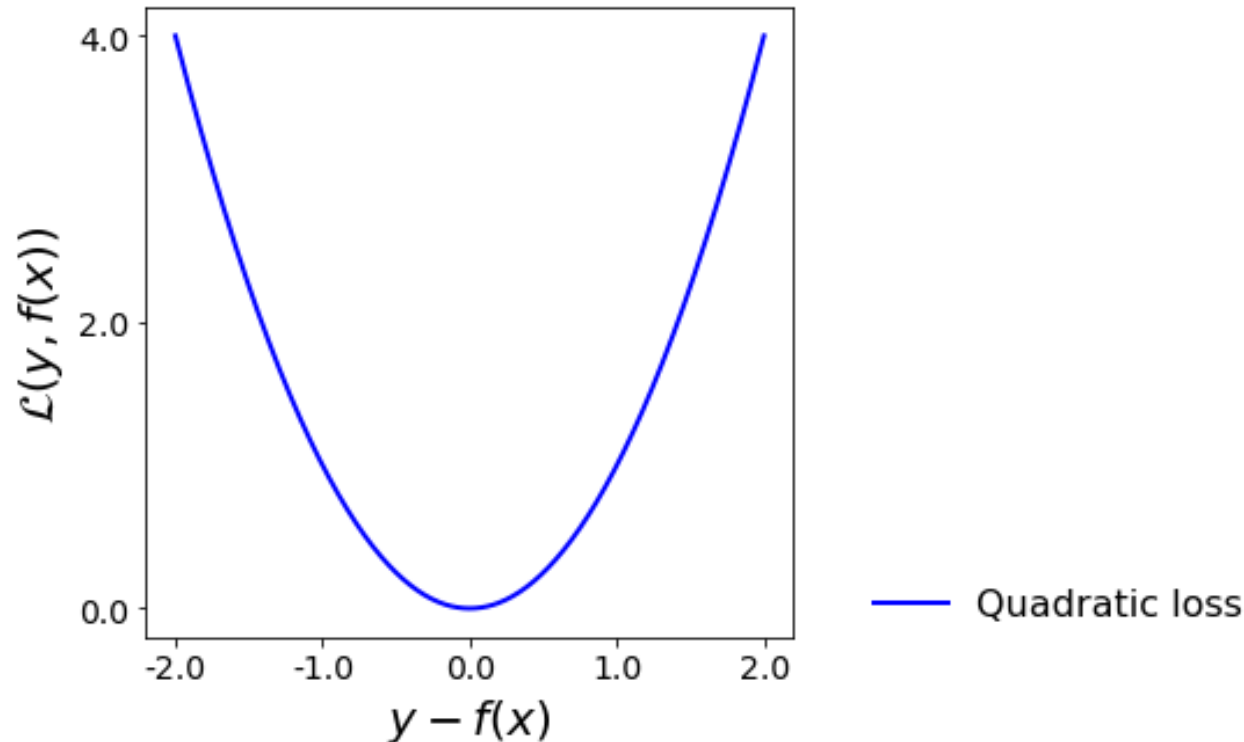
$$f^* = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n \mathcal{L}(y^i, f(x^i))$$

# Why talk about optimization?

- Supervised ML: **empirical risk minimization**

$$f^* = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n \mathcal{L}(y^i, f(x^i))$$

- Quadratic loss**  $\mathcal{L}(y, f(x)) = (y - f(x))^2$

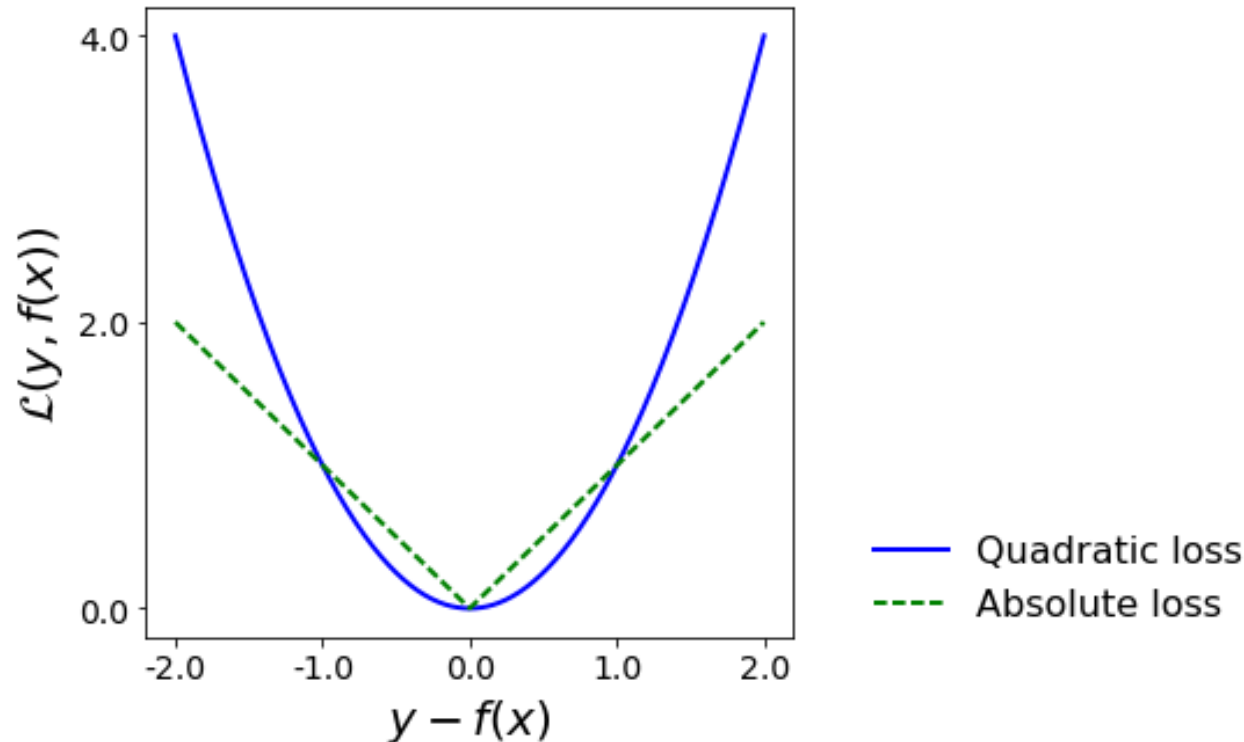


# Why talk about optimization?

- Supervised ML: **empirical risk minimization**

$$f^* = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n \mathcal{L}(y^i, f(\mathbf{x}^i))$$

- Absolute loss**  $\mathcal{L}(y, f(\mathbf{x})) = |y - f(\mathbf{x})|$



# Why talk about optimization?

- Supervised ML: **empirical risk minimization**

$$f^* = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n \mathcal{L}(y^i, f(\mathbf{x}^i))$$

- 0/1 loss**

$$\mathcal{L}(y, f(\mathbf{x})) = \begin{cases} 0 & \text{if } f(\mathbf{x}) = y \\ 1 & \text{otherwise} \end{cases}$$

# Why talk about optimization?

- **Unsupervised machine learning** also involves minimizing functions.

Examples:

- **Dimensionality reduction:** find a set of  $m$  features,  $m < p$ , such that the data projected on these  $m$  features retains maximal information.
- **Clustering:** find  $K$  groups of samples such that the between-groups variance is high and the within-group variance is small.

# Learning objectives

- **Recognize** a convex optimization problem.
- Solve an **unconstrained convex optimization** problem
  - Exactly when possible
  - By **gradient descent** and a number of its variants.
- Solve a **quadratic program**
  - Formulate the **dual problem**
  - Write down **Karush-Kuhn-Tucker conditions**.
- Transform inequality constraints with **slack variables**.

# Convex functions



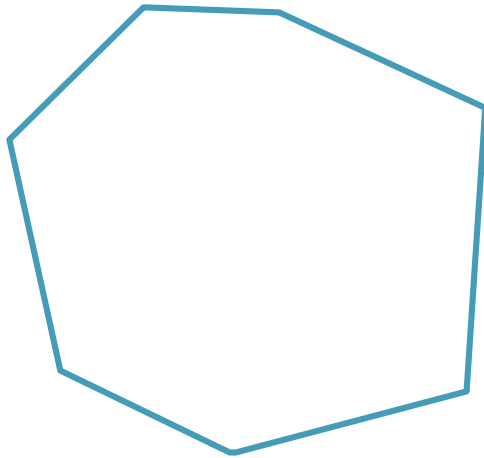
# Convex set

$\mathcal{S} \subseteq \mathbb{R}^n$  is a **convex set** iff:

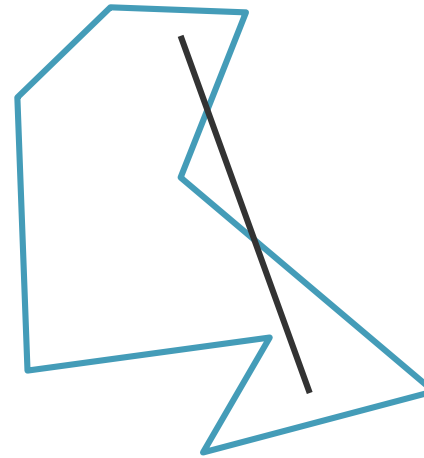
$$tu + (1 - t)v \in \mathcal{S}$$

for all  $u, v \in \mathcal{S}$  and  $0 \leq t \leq 1$

Line segments between 2 points of  $\mathcal{S}$  lie entirely in  $\mathcal{S}$ .



Convex set of  $\mathbb{R}^2$



Non-convex set of  $\mathbb{R}^2$

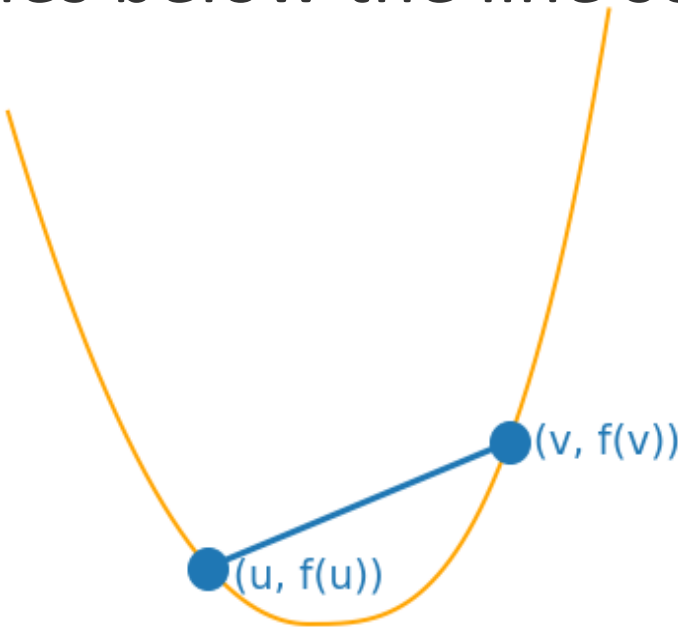
# Convex function

$f : \mathbb{R}^n \rightarrow \mathbb{R}$  is **convex** iff:

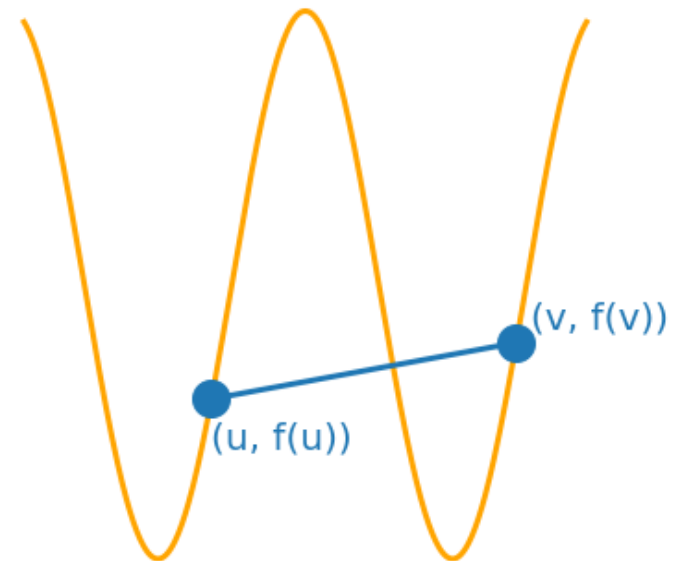
- its domain is a **convex set**
- $f(t\mathbf{u} + (1 - t)\mathbf{v}) \leq tf(\mathbf{u}) + (1 - t)f(\mathbf{v})$

for all  $\mathbf{u}, \mathbf{v} \in \text{dom}(f)$  and  $0 \leq t \leq 1$

$f$  lies below the line segment joining  $f(\mathbf{u})$  and  $f(\mathbf{v})$ .



Convex function of  $\mathbb{R} \rightarrow \mathbb{R}$



Non-convex function of  $\mathbb{R} \rightarrow \mathbb{R}$

# Concave function

$f : \mathbb{R}^n \rightarrow \mathbb{R}$  is **concave** iff:

- its domain is a **convex set**

- $f(t\mathbf{u} + (1 - t)\mathbf{v}) \geq tf(\mathbf{u}) + (1 - t)f(\mathbf{v})$

for all  $\mathbf{u}, \mathbf{v} \in \text{dom}(f)$  and  $0 \leq t \leq 1$

$f$  concave  $\Leftrightarrow -f$  convex

# Are the following univariate functions convex?

- $f : \mathbb{R} \rightarrow \mathbb{R} \quad u \mapsto 2u^2 + 3$
- $f : \mathbb{R} \rightarrow \mathbb{R} \quad u \mapsto 2u^2 - 3$
- $f : [0, 1] \cup [3, +\infty[ \rightarrow \mathbb{R} \quad u \mapsto 2u^2 - 3$
- $f : \mathbb{R} \rightarrow \mathbb{R} \quad u \mapsto 2u^3 - 3$
- $f : \mathbb{R}_+ \rightarrow \mathbb{R} \quad u \mapsto 2u^3 - 3$
- $f : \mathbb{R} \rightarrow \mathbb{R} \quad u \mapsto \max(u, 2)$
- $f : \mathbb{R} \rightarrow \mathbb{R} \quad u \mapsto \sin(u)$

# Are the following univariate functions convex?

- $f : \mathbb{R} \rightarrow \mathbb{R} \quad u \mapsto 2u^2 + 3$
- $f : \mathbb{R} \rightarrow \mathbb{R} \quad u \mapsto 2u^2 - 3$
- $f : [0, 1] \cup [3, +\infty[ \rightarrow \mathbb{R} \quad u \mapsto 2u^2 - 3$
- $f : \mathbb{R} \rightarrow \mathbb{R} \quad u \mapsto 2u^3 - 3$
- $f : \mathbb{R}_+ \rightarrow \mathbb{R} \quad u \mapsto 2u^3 - 3$
- $f : \mathbb{R} \rightarrow \mathbb{R} \quad u \mapsto \max(u, 2)$
- $f : \mathbb{R} \rightarrow \mathbb{R} \quad u \mapsto \sin(u)$



# Are the following univariate functions convex?

- $f : \mathbb{R} \rightarrow \mathbb{R} \quad u \mapsto 2u^2 + 3$
- $f : \mathbb{R} \rightarrow \mathbb{R} \quad u \mapsto 2u^2 - 3$  **Yes!**
- $f : [0, 1] \cup [3, +\infty[ \rightarrow \mathbb{R} \quad u \mapsto 2u^2 - 3$
- $f : \mathbb{R} \rightarrow \mathbb{R} \quad u \mapsto 2u^3 - 3$
- $f : \mathbb{R}_+ \rightarrow \mathbb{R} \quad u \mapsto 2u^3 - 3$
- $f : \mathbb{R} \rightarrow \mathbb{R} \quad u \mapsto \max(u, 2)$
- $f : \mathbb{R} \rightarrow \mathbb{R} \quad u \mapsto \sin(u)$

# Are the following univariate functions convex?

- $f : \mathbb{R} \rightarrow \mathbb{R} \quad u \mapsto 2u^2 + 3$
- $f : \mathbb{R} \rightarrow \mathbb{R} \quad u \mapsto 2u^2 - 3$
- $f : [0, 1] \cup [3, +\infty[ \rightarrow \mathbb{R} \quad u \mapsto 2u^2 - 3$
- $f : \mathbb{R} \rightarrow \mathbb{R} \quad u \mapsto 2u^3 - 3$
- $f : \mathbb{R}_+ \rightarrow \mathbb{R} \quad u \mapsto 2u^3 - 3$
- $f : \mathbb{R} \rightarrow \mathbb{R} \quad u \mapsto \max(u, 2)$
- $f : \mathbb{R} \rightarrow \mathbb{R} \quad u \mapsto \sin(u)$



# Are the following univariate functions convex?

- $f : \mathbb{R} \rightarrow \mathbb{R} \quad u \mapsto 2u^2 + 3$
- $f : \mathbb{R} \rightarrow \mathbb{R} \quad u \mapsto 2u^2 - 3$
- $f : [0, 1] \cup [3, +\infty[ \rightarrow \mathbb{R} \quad u \mapsto 2u^2 - 3$  **No!**
- $f : \mathbb{R} \rightarrow \mathbb{R} \quad u \mapsto 2u^3 - 3$
- $f : \mathbb{R}_+ \rightarrow \mathbb{R} \quad u \mapsto 2u^3 - 3$
- $f : \mathbb{R} \rightarrow \mathbb{R} \quad u \mapsto \max(u, 2)$
- $f : \mathbb{R} \rightarrow \mathbb{R} \quad u \mapsto \sin(u)$



# Are the following univariate functions convex?

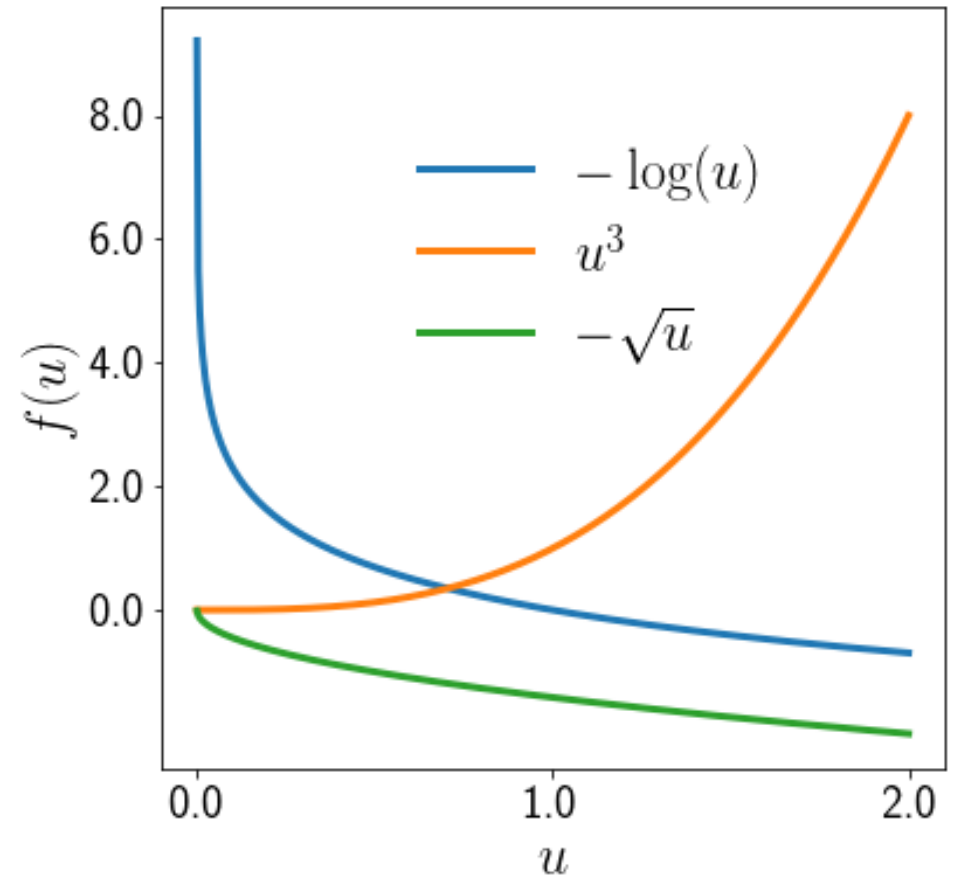
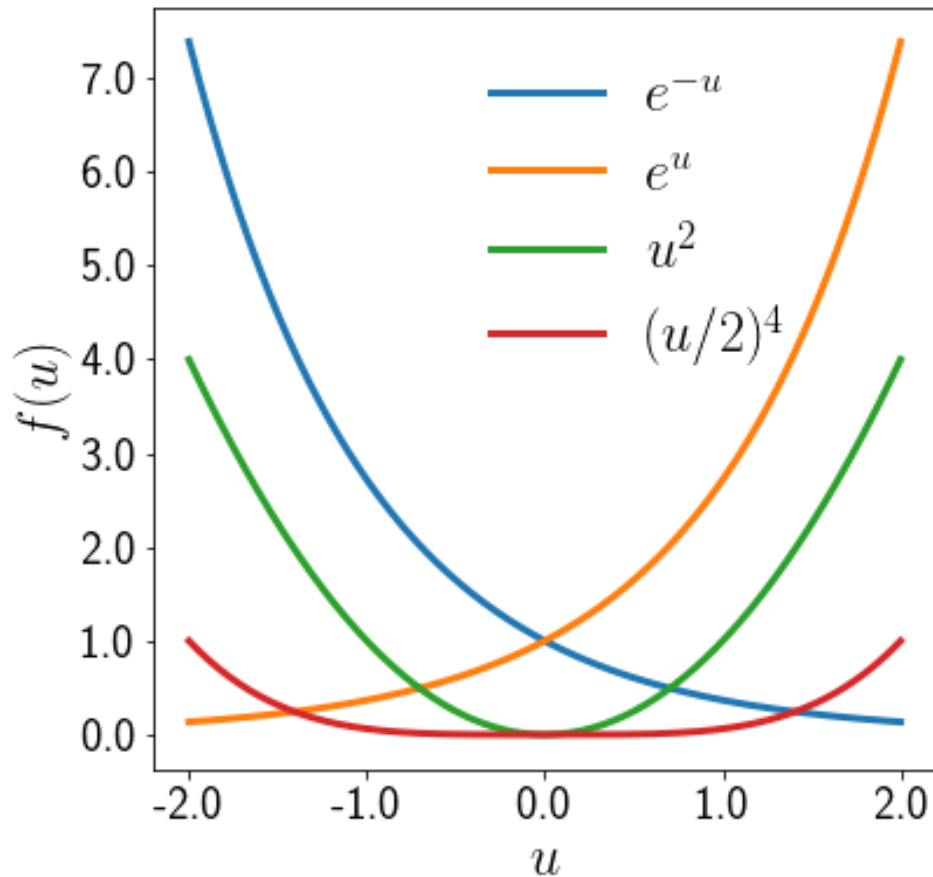
- $f : \mathbb{R} \rightarrow \mathbb{R} \quad u \mapsto 2u^2 + 3$  **Yes**
- $f : \mathbb{R} \rightarrow \mathbb{R} \quad u \mapsto 2u^2 - 3$  **Yes**
- $f : [0, 1] \cup [3, +\infty[ \rightarrow \mathbb{R} \quad u \mapsto 2u^2 - 3$  **No**
- $f : \mathbb{R} \rightarrow \mathbb{R} \quad u \mapsto 2u^3 - 3$  **No**
- $f : \mathbb{R}_+ \rightarrow \mathbb{R} \quad u \mapsto 2u^3 - 3$  **Yes**
- $f : \mathbb{R} \rightarrow \mathbb{R} \quad u \mapsto \max(u, 2)$  **Yes**
- $f : \mathbb{R} \rightarrow \mathbb{R} \quad u \mapsto \sin(u)$  **No**

# Univariate examples

- **Exponential:**  $u \mapsto e^{au} \quad \forall a \in \mathbb{R}$
- **Logarithmic:**  $u \mapsto -\log(au) \quad \forall a > 0 \quad \text{on } \mathbb{R}_+^*$
- **Power functions:**

$$u \mapsto u^a \quad \forall a \geq 1 \text{ or } a \leq 0 \text{ on } \mathbb{R}_+^*$$
$$u \mapsto -u^a \quad \forall 0 \leq a \leq 1 \text{ on } \mathbb{R}_+^*$$

$$u \mapsto u^a \quad \forall a = 2n, n \in \mathbb{N}$$



# More examples

- **Affine functions** are both convex and concave

$$u \mapsto a^\top u + b \quad a \in \mathbb{R}^n$$

# More examples

- **Affine functions** are both convex and concave

$$u \mapsto a^\top u + b \quad a \in \mathbb{R}^n$$

- **Quadratic functions**

$$u \mapsto \frac{1}{2} u^\top Q u + b^\top u + c \quad \boxed{Q \succeq 0} \quad b \in \mathbb{R}^n \quad c \in \mathbb{R}$$

Q **positive semi-definite**



# More examples

- **Affine functions** are both convex and concave

$$u \mapsto a^\top u + b \quad a \in \mathbb{R}^n$$

- **Quadratic functions**

$$u \mapsto \frac{1}{2} u^\top Q u + b^\top u + c \quad \boxed{Q \succeq 0} \quad b \in \mathbb{R}^n \quad c \in \mathbb{R}$$

$Q$  **positive semi-definite**


$$u^\top Q u \geq 0 \quad \forall u \in \mathbb{R}^n$$

- All eigenvalues of  $Q$  are non-negative
- The bilinear form  $u, v \mapsto v^\top Q u$  is an inner product
- $Q$  is a Gram matrix of independent vectors  $Q_{ij} = \langle v_i, v_j \rangle$
- Unique Cholesky decomposition  $Q = LL^\top$

# More examples

- **Affine functions** are both convex and concave

$$u \mapsto a^\top u + b \quad a \in \mathbb{R}^n$$

- **Quadratic functions**

$$u \mapsto \frac{1}{2} u^\top Q u + b^\top u + c \quad \boxed{Q \succeq 0} \quad b \in \mathbb{R}^n \quad c \in \mathbb{R}$$

$Q$  **positive semi-definite**

- **Lp norms**

$$u \mapsto \|u\|_p = \left( \sum_{j=1}^n |u_j|^p \right)^{1/p}$$

- $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is **strictly convex** iff:

$$\forall \mathbf{u} \neq \mathbf{v} \in \text{dom}(f), \quad \forall 0 < t < 1$$

$$f(t\mathbf{u} + (1-t)\mathbf{v}) < tf(\mathbf{u}) + (1-t)f(\mathbf{v}).$$

$f$  is convex and has **greater curvature than a linear function**.

- $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is **strongly convex of parameter  $m > 0$**  iff:

$$f - \frac{m}{2} \|\mathbf{u}\|_2^2 \text{ is convex.}$$

$f$  is convex and has curvature **as least as great as a quadratic function**.

**strongly convex  $\Rightarrow$  strictly convex  $\Rightarrow$  convex**

# First-order characterization

- If  $f$  is **differentiable**, then  $f$  is **convex** if and only if:
  - its domain is a convex set
  - for all  $\mathbf{u}, \mathbf{v} \in \text{dom}(f)$

$$f(\mathbf{v}) \geq f(\mathbf{u}) + \nabla f(\mathbf{u})^\top (\mathbf{v} - \mathbf{u})$$



# First-order characterization

- If  $f$  is **differentiable**, then  $f$  is **convex** if and only if:
  - its domain is a convex set
  - for all  $u, v \in \text{dom}(f)$

$$f(v) \geq f(u) + \nabla f(u)^\top (v - u)$$



**Gradient of  $f$**

$$\nabla f = \begin{pmatrix} \frac{\partial f}{\partial u_1} \\ \frac{\partial f}{\partial u_2} \\ \dots \\ \frac{\partial f}{\partial u_n} \end{pmatrix}$$

# First-order characterization

- If  $f$  is **differentiable**, then  $f$  is **convex** if and only if:
  - its domain is a convex set
  - for all  $\mathbf{u}, \mathbf{v} \in \text{dom}(f)$

$$f(\mathbf{v}) \geq f(\mathbf{u}) + \nabla f(\mathbf{u})^\top (\mathbf{v} - \mathbf{u})$$



# First-order characterization

- If  $f$  is **differentiable**, then  $f$  is **convex** if and only if:
  - its domain is a convex set
  - for all  $\mathbf{u}, \mathbf{v} \in \text{dom}(f)$

$$f(\mathbf{v}) \geq f(\mathbf{u}) + \boxed{\nabla f(\mathbf{u})^\top (\mathbf{v} - \mathbf{u})}$$

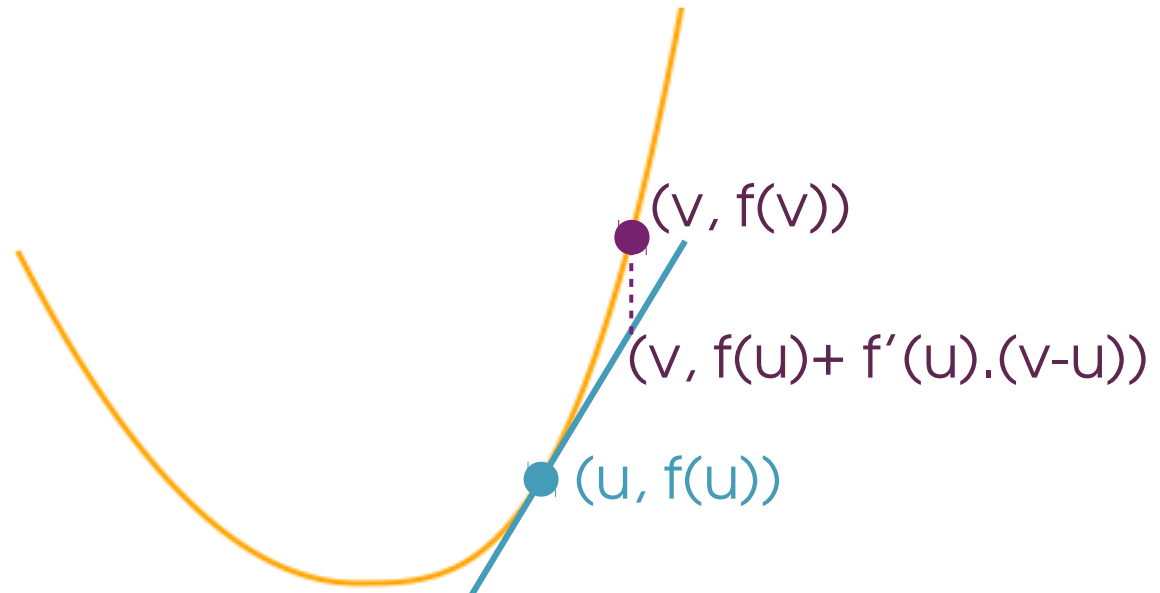
First-order Taylor expansion of  $f$  in  $\mathbf{u}$

# First-order characterization

- If  $f$  is **differentiable**, then  $f$  is **convex** if and only if:
  - its domain is a convex set
  - for all  $u, v \in \text{dom}(f)$

$$f(v) \geq f(u) + \nabla f(u)^\top (v - u)$$

First-order Taylor expansion of  $f$  in  $u$



# First-order characterization

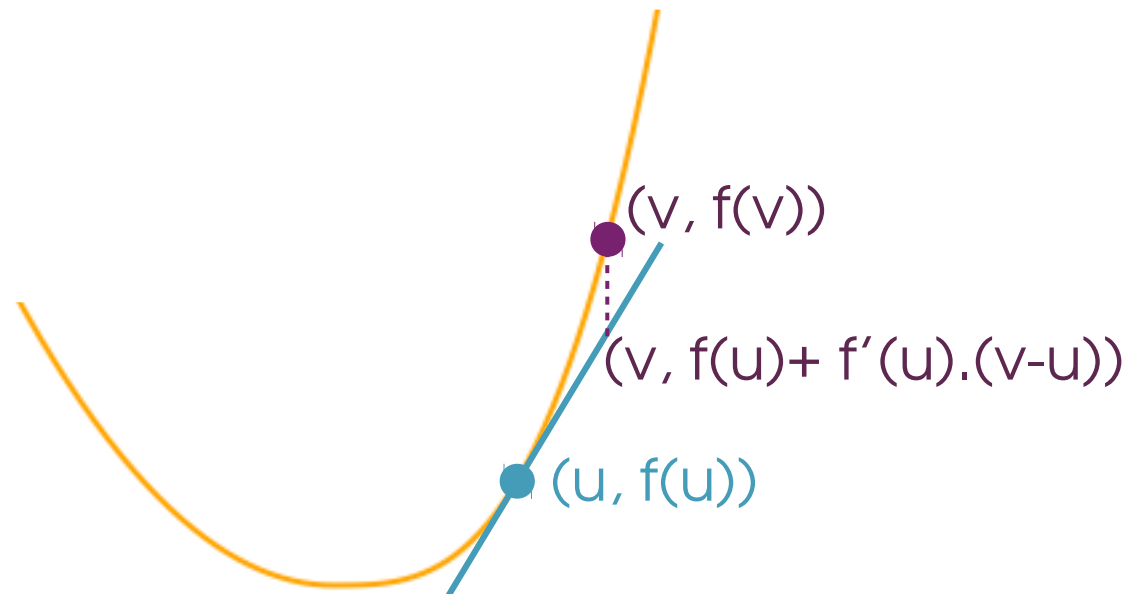
- If  $f$  is **differentiable**, then  $f$  is **convex** if and only if:
  - its domain is a convex set
  - for all  $u, v \in \text{dom}(f)$

$$f(v) \geq f(u) + \nabla f(u)^\top (v - u)$$

First-order Taylor expansion of  $f$  in  $u$

What does it mean if

$$\nabla f(u) = 0 \quad ?$$



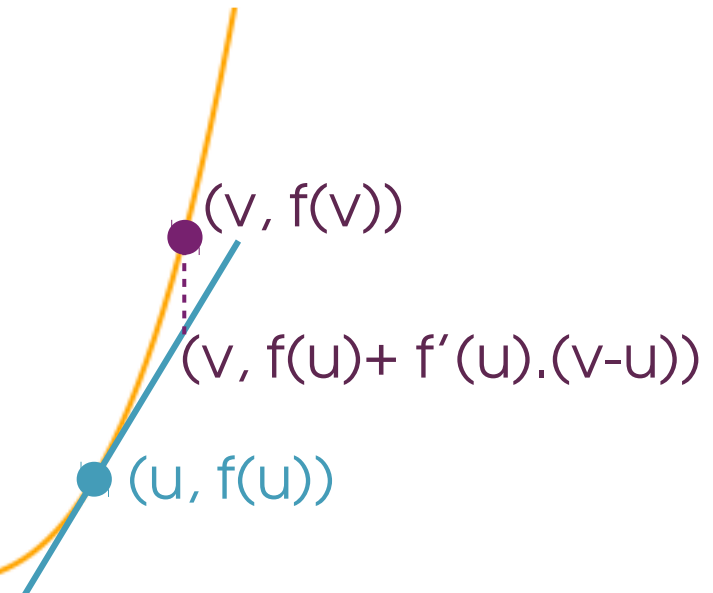
# First-order characterization

- If  $f$  is **differentiable**, then  $f$  is **convex** if and only if:
  - its domain is a convex set
  - for all  $\mathbf{u}, \mathbf{v} \in \text{dom}(f)$


$$f(\mathbf{v}) \geq f(\mathbf{u}) + \nabla f(\mathbf{u})^\top (\mathbf{v} - \mathbf{u})$$

First-order Taylor expansion of  $f$  in  $\mathbf{u}$

$$\nabla f(\mathbf{u}) = 0 \Leftrightarrow \mathbf{u} \text{ minimizes } f$$



# Second-order characterization


- If  $f$  is **twice differentiable**, then  $f$  is **convex** iff:
  - its domain is a convex set
  - for all  $u \in \text{dom}(f)$  
$$\nabla^2 f(u) \succeq 0$$

# Second-order characterization

- If  $f$  is **twice differentiable**, then  $f$  is **convex** iff:
  - its domain is a convex set
  - for all  $u \in \text{dom}(f)$

$$\nabla^2 f(u) \succeq 0$$

**Hessian of  $f$**


$$\nabla^2 f = \begin{pmatrix} \frac{\partial^2 f}{\partial u_1^2} & \frac{\partial^2 f}{\partial u_1 u_2} & \cdots & \frac{\partial^2 f}{\partial u_1 u_n} \\ \frac{\partial^2 f}{\partial u_2 u_1} & \frac{\partial^2 f}{\partial u_2^2} & \cdots & \frac{\partial^2 f}{\partial u_2 u_n} \\ \cdots & \cdots & \cdots & \cdots \\ \frac{\partial^2 f}{\partial u_n u_1} & \frac{\partial^2 f}{\partial u_n u_2} & \cdots & \frac{\partial^2 f}{\partial u_n^2} \end{pmatrix}$$



# Second-order characterization

- If  $f$  is **twice differentiable**, then  $f$  is **convex** iff:
  - its domain is a convex set
  - for all  $u \in \text{dom}(f)$

$$\boxed{\nabla^2 f(u)} \succeq 0$$

Hessian of  $f$

- $f$  has **positive curvature** in any point  $u$ .



# Operations preserving convexity

- **Non-negative linear combination**

If  $f_1, f_2, \dots, f_m$  convex and  $a_1, a_2, \dots, a_m \geq 0$   
then  $a_1 f_1 + a_2 f_2 + \dots + a_m f_m$  is convex.

- **Pointwise maximization**

If  $f_1, f_2, \dots, f_m$  convex, then

$u \mapsto \max_{1, \dots, m} f_k(u)$  is convex (also true for an infinite number of functions  $f_k$ ).

- **Partial minimization**

If  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  convex and  $C$  is a convex set, then

$(u_1, u_2, \dots, u_{n-1}) \mapsto \min_{v \in C} f(u_1, u_2, \dots, u_{n-1}, v)$  is convex.

# Convex optimization

# Unconstrained convex optimization

Unconstrained convex optimization program/problem:

$$\min_{\boldsymbol{u} \in \text{dom}(f)} f(\boldsymbol{u})$$

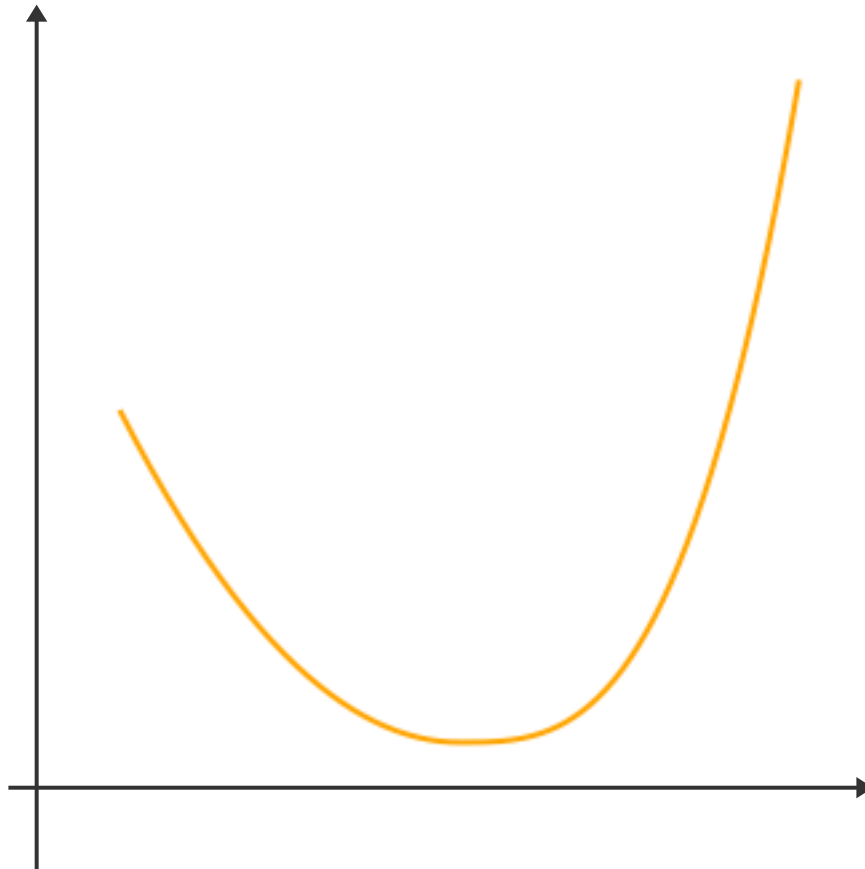
where  $f$  is convex.

# Unconstrained convex optimization

Unconstrained convex optimization program/problem:

$$\min_{u \in \text{dom}(f)} f(u)$$

where  $f$  is convex.

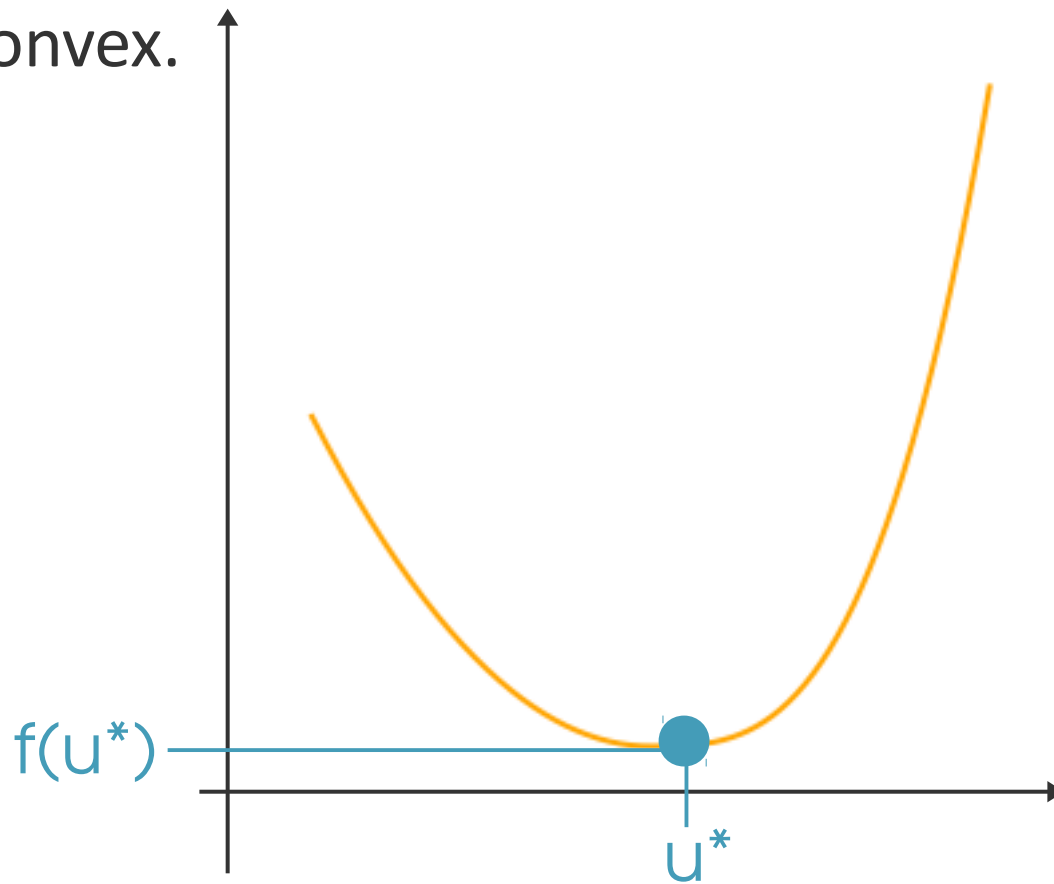


# Unconstrained convex optimization

Unconstrained convex optimization program/problem:

$$\min_{u \in \text{dom}(f)} f(u)$$

where  $f$  is convex.



# Constrained convex optimization

- Convex optimization program/problem:

$$\begin{aligned} \min_{\mathbf{u} \in D} f(\mathbf{u}) \\ \text{subject to } g_i(\mathbf{u}) \leq 0, i = 1, \dots, m \\ h_j(\mathbf{u}) = 0, j = 1, \dots, r \end{aligned}$$

- $f$  is **convex**
- $g_i, i = 1, \dots, m$  are **convex**
- $h_j, j = 1, \dots, r$  are **affine**  $h_j : \mathbf{u} \mapsto \mathbf{a}_j^\top \mathbf{u} + b_j$
- $D$  is the **common domain** of all the functions.

$$D = \text{dom}(f) \cap \bigcap_{i=1}^m \text{dom}(g_i) \bigcap_{j=1}^r \text{dom}(h_j)$$

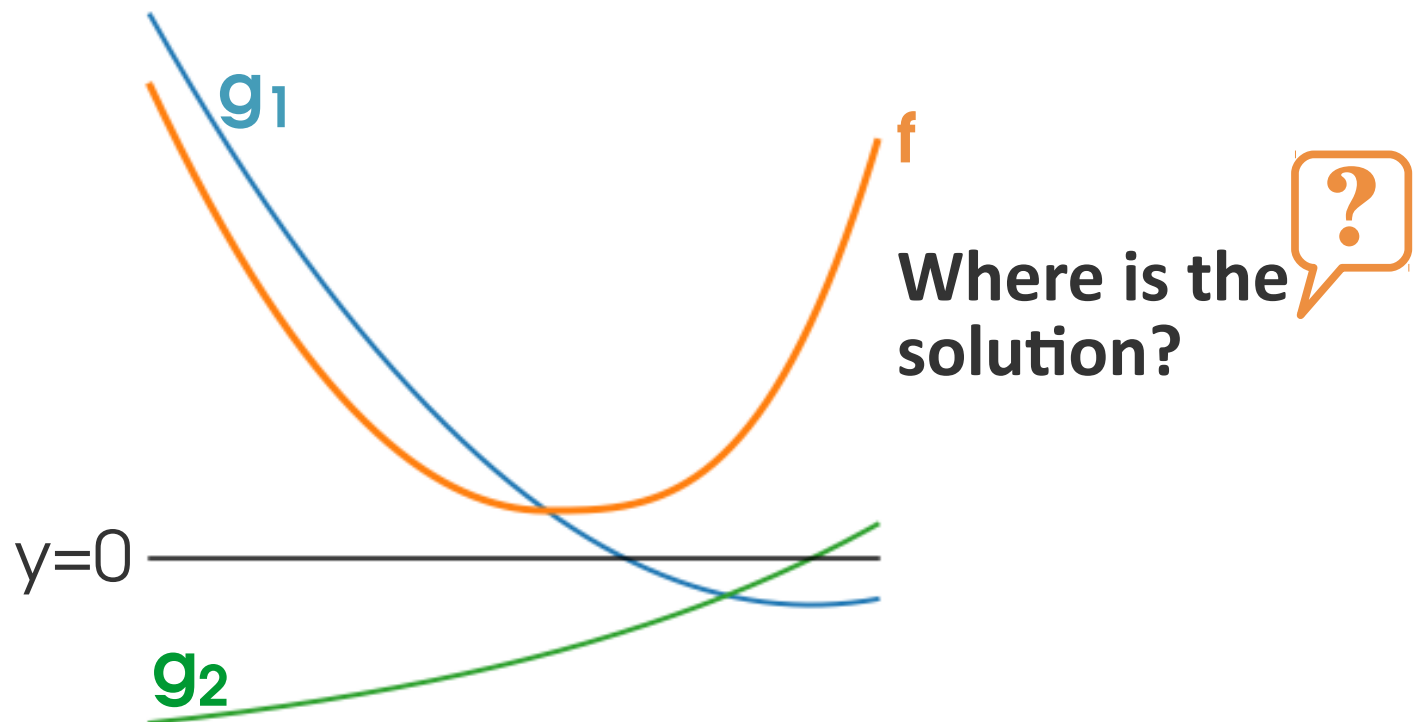
# Constrained convex optimization

- Convex optimization program/problem:

$$\min_{\mathbf{u} \in D} f(\mathbf{u})$$

subject to  $g_i(\mathbf{u}) \leq 0, i = 1, \dots, m$

$$h_j(\mathbf{u}) = 0, j = 1, \dots, r$$





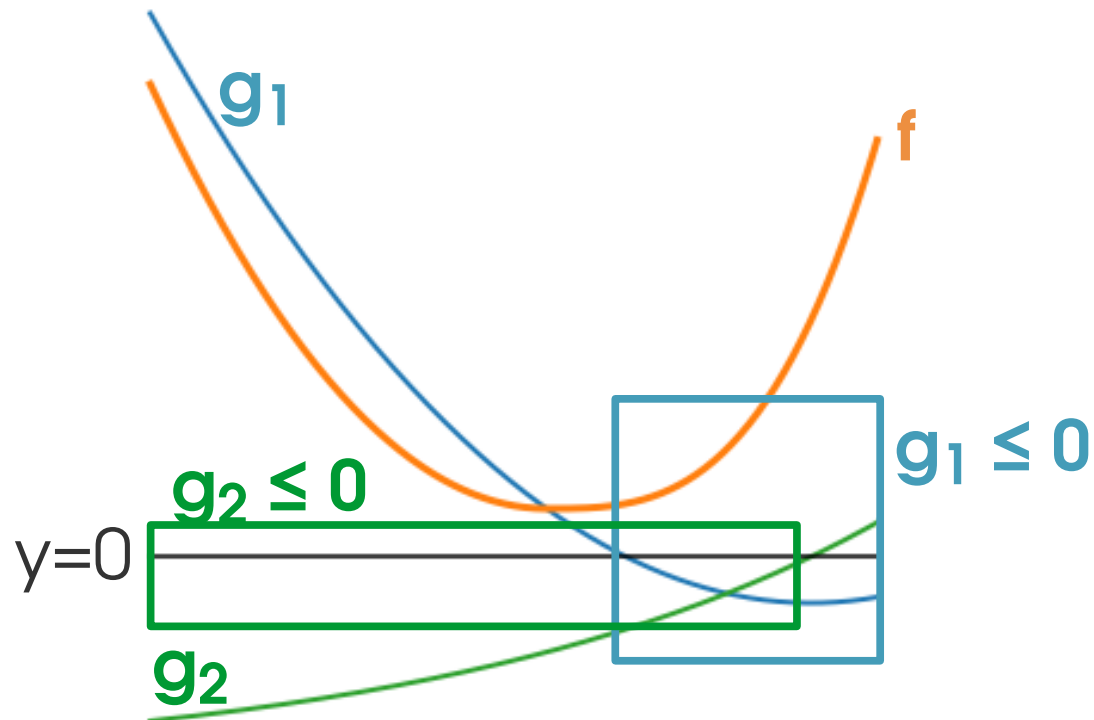
# Constrained convex optimization

- Convex optimization program/problem:

$$\min_{\mathbf{u} \in D} f(\mathbf{u})$$

subject to  $g_i(\mathbf{u}) \leq 0, i = 1, \dots, m$

$$h_j(\mathbf{u}) = 0, j = 1, \dots, r$$



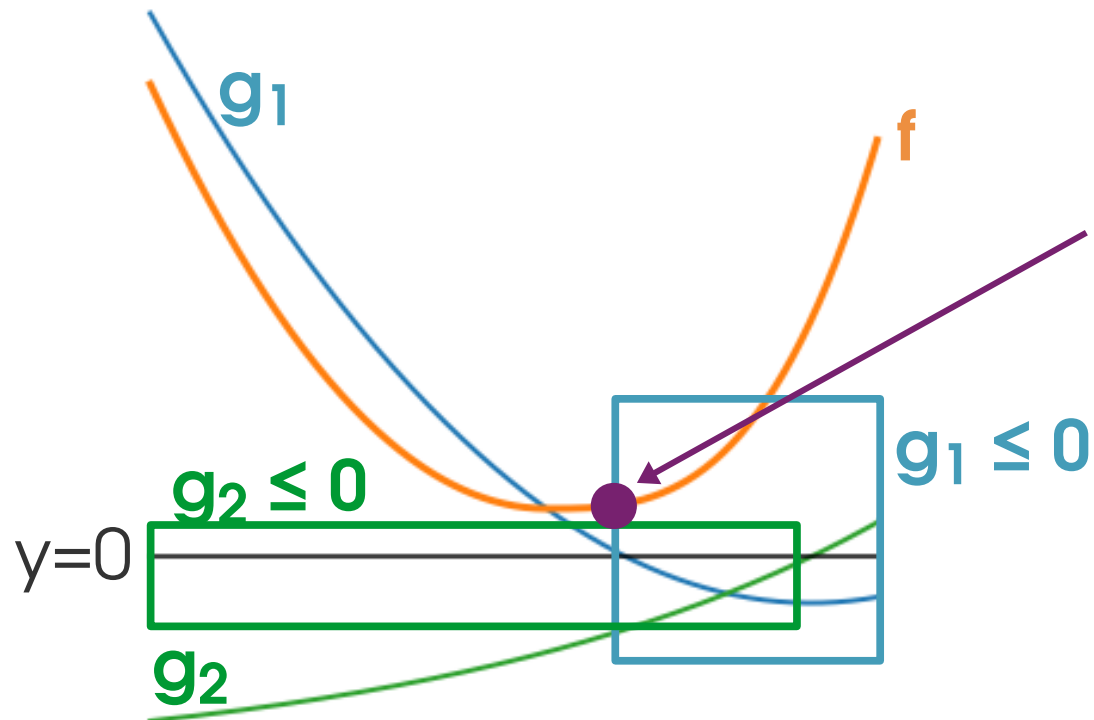
# Constrained convex optimization

- Convex optimization program/problem:

$$\min_{\mathbf{u} \in D} f(\mathbf{u})$$

subject to  $g_i(\mathbf{u}) \leq 0, i = 1, \dots, m$

$$h_j(\mathbf{u}) = 0, j = 1, \dots, r$$



# Constrained convex optimization

- **Convex optimization program/problem:**

$$\min_{\mathbf{u} \in D} f(\mathbf{u})$$

subject to  $g_i(\mathbf{u}) \leq 0, i = 1, \dots, m$

$$h_j(\mathbf{u}) = 0, j = 1, \dots, r$$

- $f$  is the **objective function**
- $g_i, i = 1, \dots, m$  are the **inequality constraints**
- $h_j, j = 1, \dots, r$  are the **equality constraints**
- $\mathbf{v} \in D$  that verifies all constraints is a **feasible point**  
 $g_i(\mathbf{v}) \leq 0, i = 1, \dots, m$  and  $h_j(\mathbf{v}) = 0, j = 1, \dots, r$
- The set of all feasible points is the **feasible region**  
 $\{\mathbf{v} : \mathbf{v} \in D; g_i(\mathbf{v}) \leq 0, i = 1, \dots, m; h_j(\mathbf{v}) = 0, j = 1, \dots, r\}$

# Constrained convex optimization

- **Convex optimization program/problem:**

$$\begin{aligned} & \min_{\mathbf{u} \in D} f(\mathbf{u}) \\ & \text{subject to } g_i(\mathbf{u}) \leq 0, i = 1, \dots, m \\ & \quad h_j(\mathbf{u}) = 0, j = 1, \dots, r \end{aligned}$$

- Assuming it exists, the solution  $f^*$ , that is to say, the minimum value of  $f$  over all feasible points, is the **optimal value (optimum)**
- $\mathbf{u}$  feasible such that  $f(\mathbf{u}) = f^*$  is called **optimal**, or a **minimizer** (it needs not be unique).
- If  $\mathbf{u}$  is feasible and  $g_i(\mathbf{u}) = 0$  then  $g_i$  is **active** at  $\mathbf{u}$ .

# Local & global optima

For convex optimization problems,

**local minima are global minima!**

If  $\mathbf{u}$  is feasible and minimizes  $f$  in a local neighborhood:

$$f(\mathbf{u}) \leq f(\mathbf{v}) \text{ for all feasible } \mathbf{v}, \|\mathbf{u} - \mathbf{v}\|_2^2 \leq \epsilon$$

then  $\mathbf{u}$  minimizes  $f$  globally.

# Local & global optima

For convex optimization problems,

**local minima are global minima!**

If  $\mathbf{u}$  is feasible and minimizes  $f$  in a local neighborhood:

$$f(\mathbf{u}) \leq f(\mathbf{v}) \text{ for all feasible } \mathbf{v}, \|\mathbf{u} - \mathbf{v}\|_2^2 \leq \epsilon$$

then  $\mathbf{u}$  minimizes  $f$  globally.

Suppose  $\bar{\mathbf{u}}$  is feasible and a local optimum of  $f$ :  $\exists \epsilon > 0 : f(\bar{\mathbf{u}}) \leq f(\mathbf{v})$  for all feasible  $\mathbf{v}$  such that  $\|\bar{\mathbf{u}} - \mathbf{v}\| \leq \epsilon$ .  
Suppose  $\mathbf{u}^*$  is a global optimum,  $\mathbf{u}^* \neq \bar{\mathbf{u}}$ . Then

$$f(\mathbf{u}^*) < f(\bar{\mathbf{u}}) \tag{1}$$

and

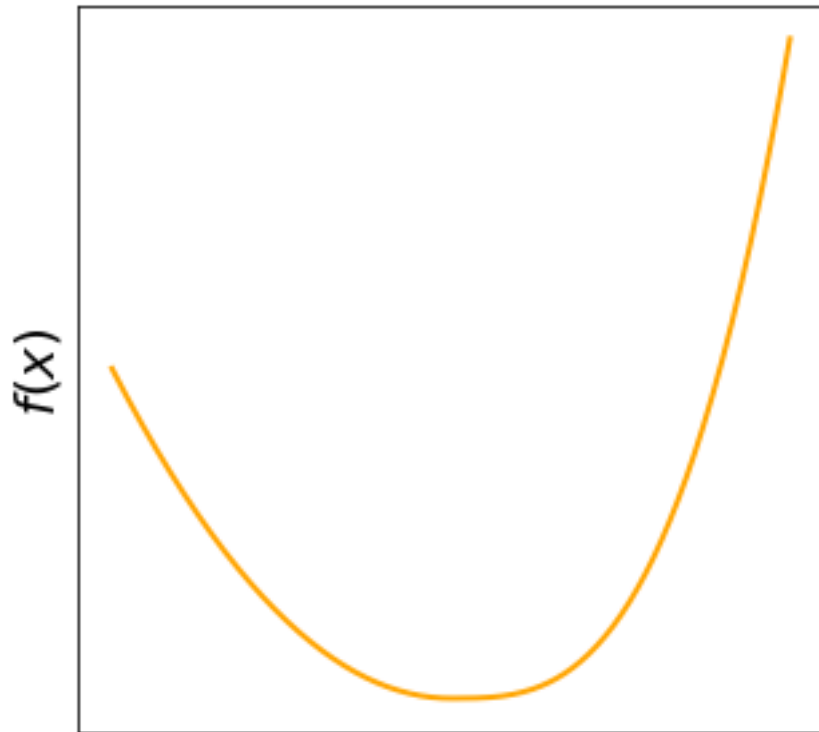
$$\|\bar{\mathbf{u}} - \mathbf{u}^*\| > \epsilon \tag{2}$$

Consider  $\mathbf{u} = (1 - \lambda)\bar{\mathbf{u}} + \lambda\mathbf{u}^*$ , where  $\lambda = \frac{1}{2} \frac{\epsilon}{\|\bar{\mathbf{u}} - \mathbf{u}^*\|}$ . Because of (2),  $0 \leq \lambda < 1$ . Because the feasible set is convex and both  $\bar{\mathbf{u}}$  and  $\mathbf{u}^*$  are feasible,  $\mathbf{u}$  is feasible.

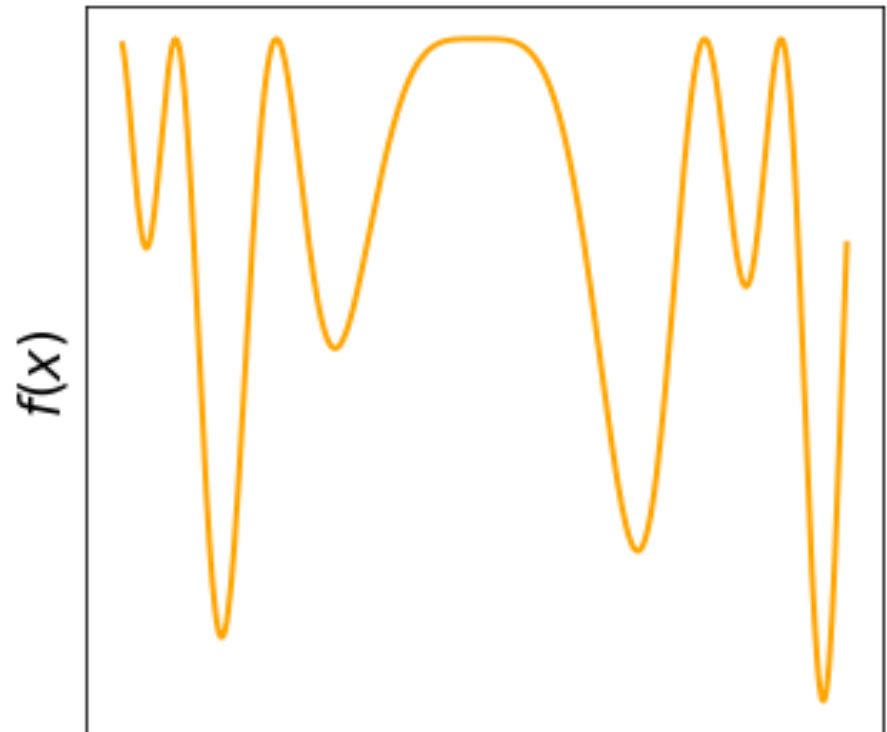
$\|\mathbf{u} - \bar{\mathbf{u}}\| = \|(1 - \lambda)\bar{\mathbf{u}} + \lambda\mathbf{u}^*\| = \lambda\|\mathbf{u}^* - \bar{\mathbf{u}}\| = \epsilon/2 < \epsilon$ , hence  $f(\mathbf{u}) \geq f(\bar{\mathbf{u}})$  (as  $\bar{\mathbf{u}}$  is a local minimum).

But because  $f$  is convex,  $f(\mathbf{u}) \leq (1 - \lambda)f(\bar{\mathbf{u}}) + \lambda f(\mathbf{u}^*) = f(\bar{\mathbf{u}}) + \lambda(f(\mathbf{u}^*) - f(\bar{\mathbf{u}})) < f(\bar{\mathbf{u}})$ . (The last inequality comes from (1).) The contradiction implies that it is impossible that  $\mathbf{u}^* \neq \bar{\mathbf{u}}$ , hence the local optimum is also global.

# Why talk about convex optimization?



**convex**



**non-convex**

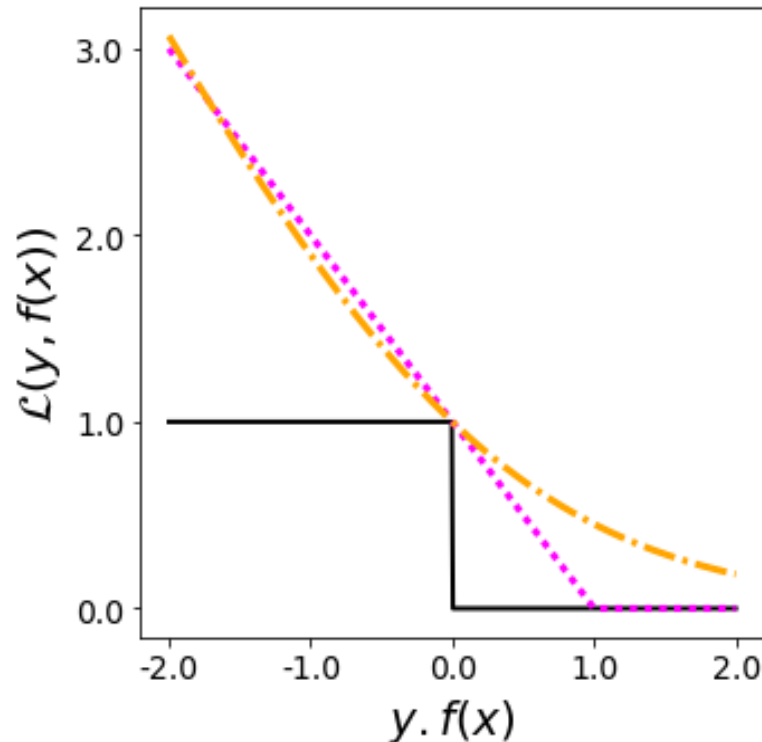
- Convex optimization is “easy”.
- We’ll often try to formulate ML problems as convex optimization problems.

# Why talk about convex optimization?

- Supervised ML: **empirical risk minimization**

$$f^* = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n \mathcal{L}(y^i, f(x^i))$$

- Losses for classification**



The 0/1 loss is non-convex.   
We'll replace it with other losses.

$y \in \{-1, 1\}$

- 0/1 loss
- Hinge loss
- · - · Logistic loss



# Unconstrained convex optimization

# First-order characterization

- Suppose  $f$  differentiable
- **Given the first-order characterization of convex functions, how can we solve  $\min_{u \in \text{dom}(f)} f(u)$  ?**

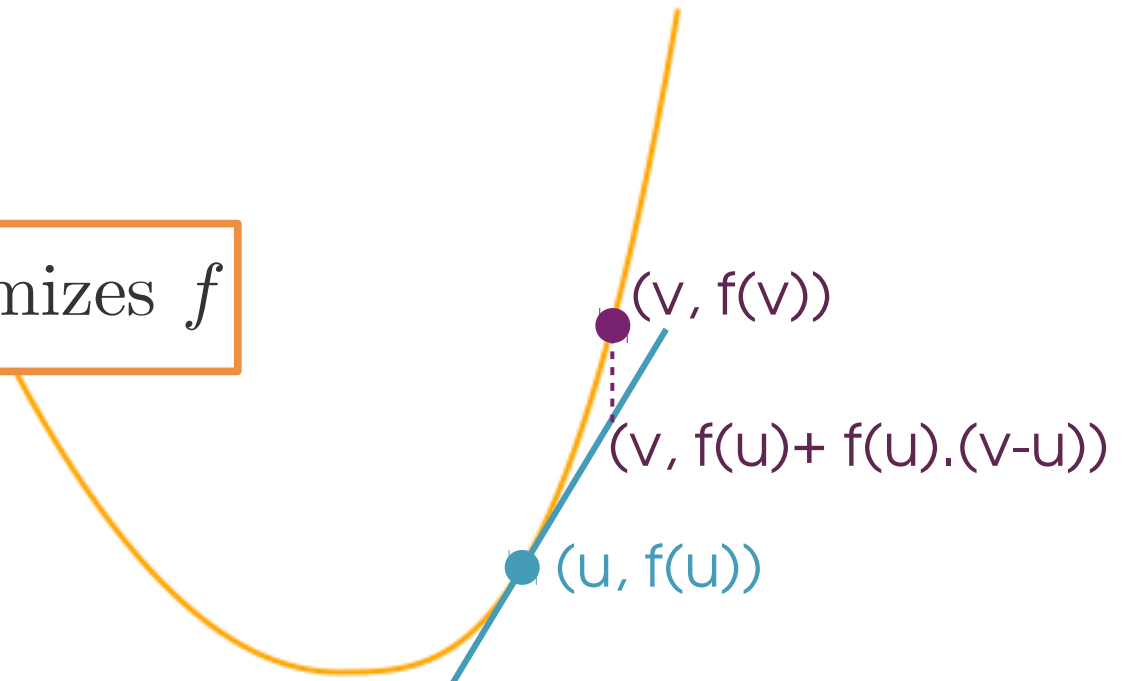


# First-order characterization

- Suppose  $f$  differentiable
- **Given the first-order characterization of convex functions, how can we solve  $\min_{u \in \text{dom}(f)} f(u)$  ?**



$$\nabla f(u) = 0 \Leftrightarrow u \text{ minimizes } f$$

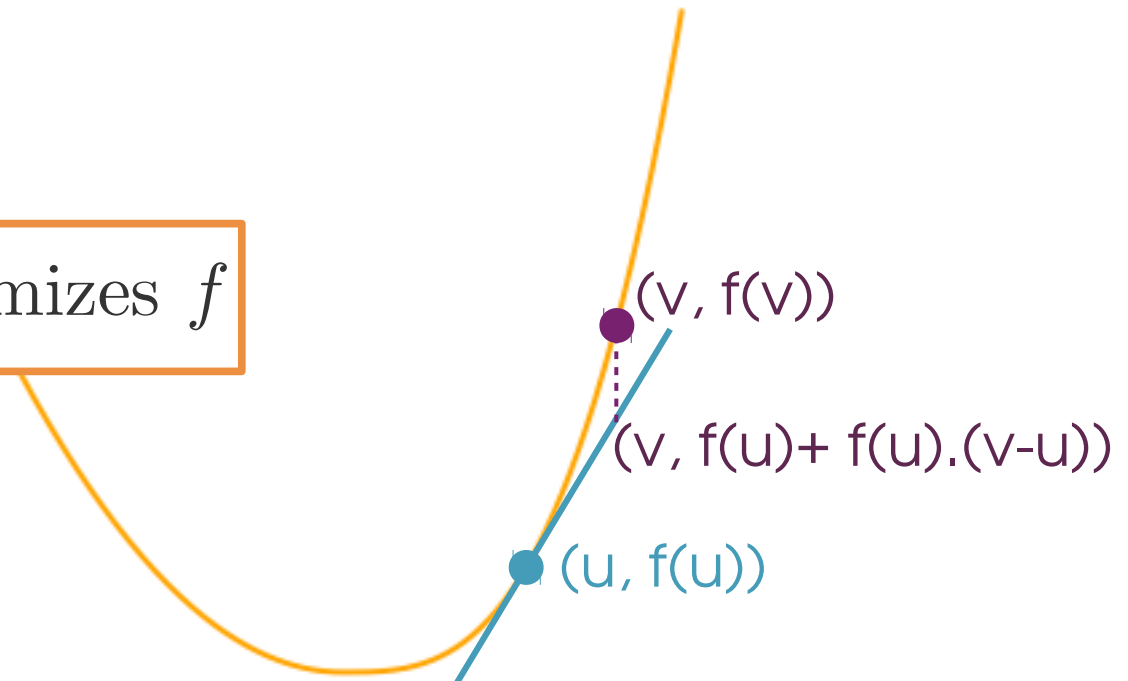


# First-order characterization

- Suppose  $f$  differentiable
- **Given the first-order characterization of convex functions, how can we solve  $\min_{u \in \text{dom}(f)} f(u)$  ?**

Set the gradient of  $f$  to 0

$$\nabla f(u) = 0 \Leftrightarrow u \text{ minimizes } f$$



# First-order characterization

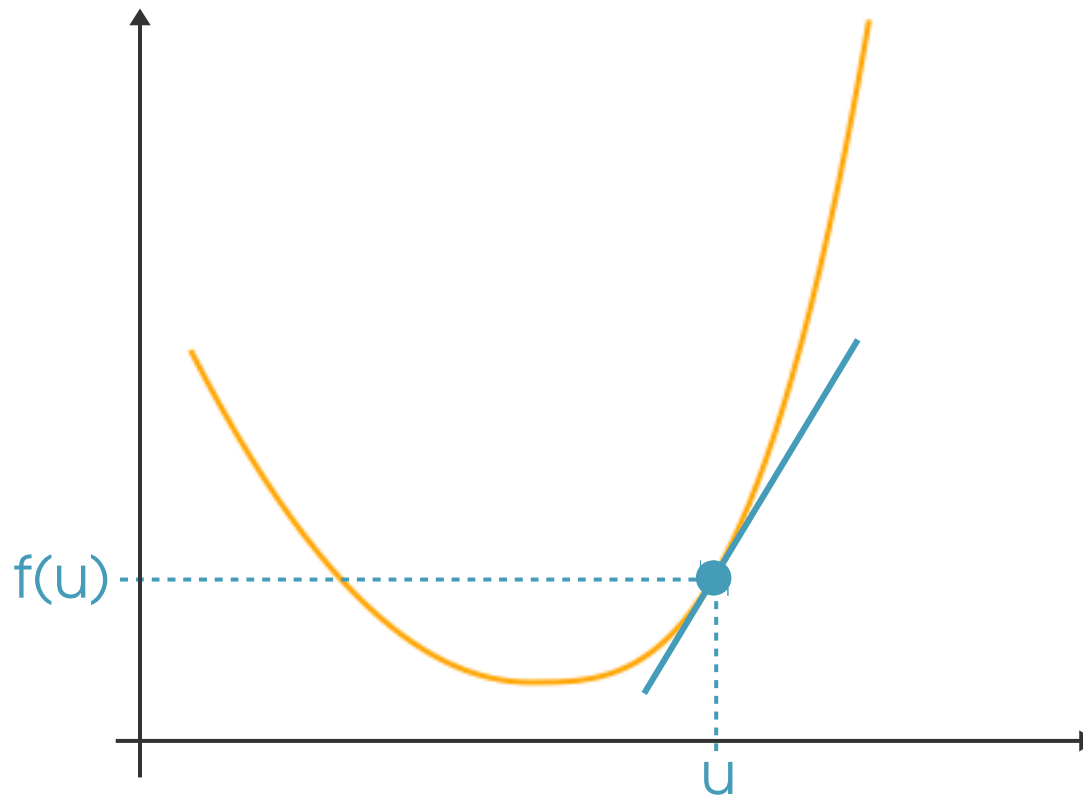
- Suppose  $f$  differentiable
- **Given the first-order characterization of convex functions, how can we solve  $\min_{u \in \text{dom}(f)} f(u)$  ?**

**Set the gradient of  $f$  to 0**

- But what if  $\nabla f(u) = 0$  cannot be solved?

# Gradient descent

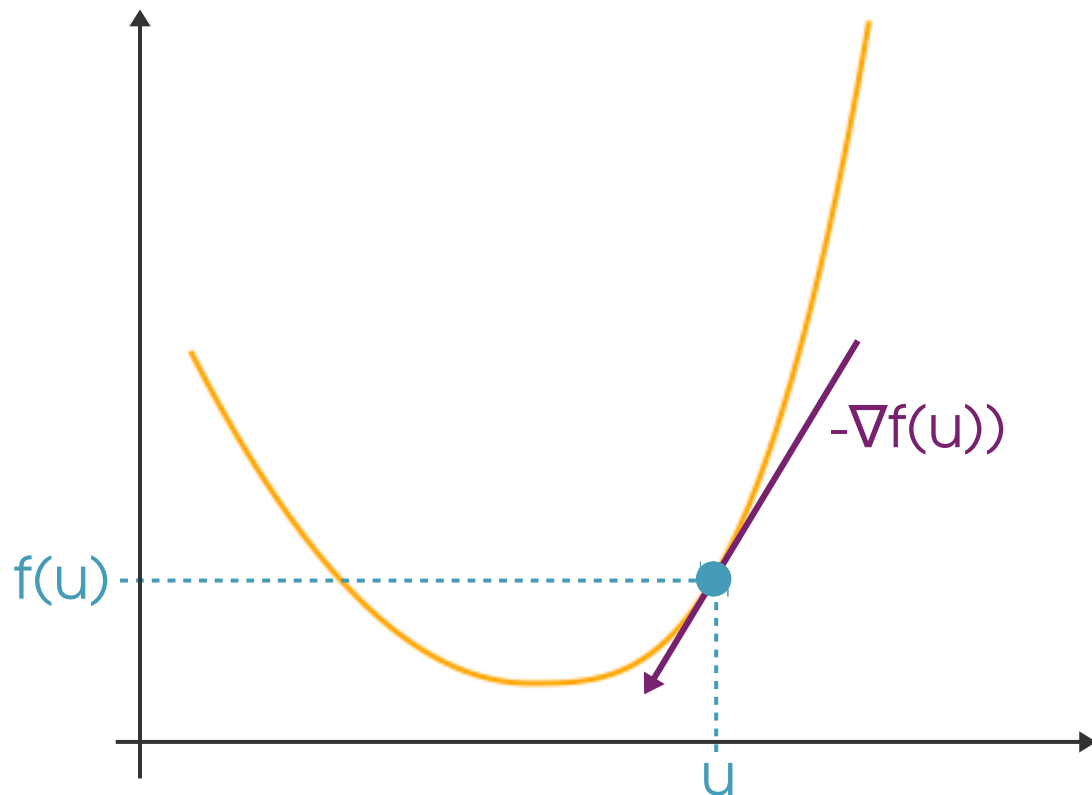
- Start from a random point  $u$ .
- **How do I get closer to the solution?** 



# Gradient descent

- Start from a random point  $u$ .
- **How do I get closer to the solution?**
- Follow the **opposite of the gradient**.

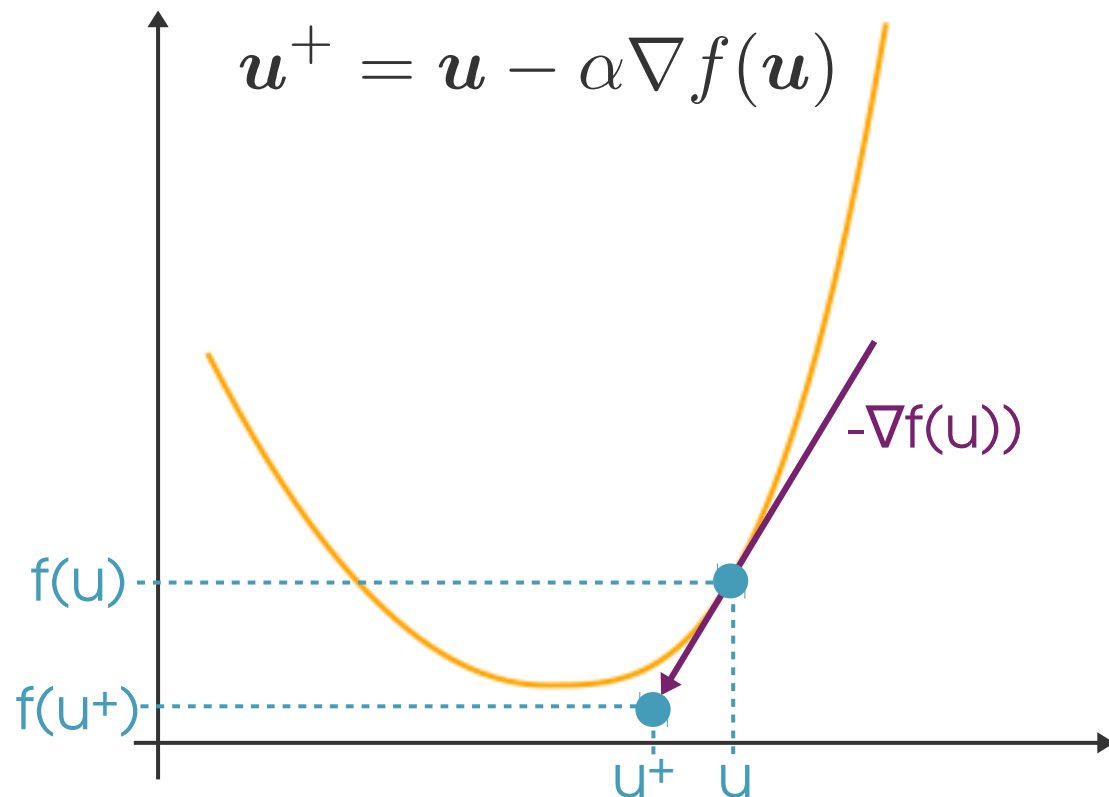
The gradient indicates the direction of steepest increase.



# Gradient descent

- Start from a random point  $\mathbf{u}$ .
- **How do I get closer to the solution?**
- Follow the **opposite of the gradient**.

The gradient indicates the direction of steepest increase.





# Gradient descent algorithm

- Choose an initial point  $\mathbf{u}^{(0)} \in \mathbb{R}^n$
- Repeat for  $k=1, 2, 3, \dots$ 
$$\mathbf{u}^{(k)} = \mathbf{u}^{(k-1)} - \alpha_k \nabla f(\mathbf{u}^{(k-1)})$$
- Stop at some point

# Gradient descent algorithm

- Choose an initial point  $\mathbf{u}^{(0)} \in \mathbb{R}^n$
- Repeat for  $k=1, 2, 3, \dots$

$$\mathbf{u}^{(k)} = \mathbf{u}^{(k-1)} - \boxed{\alpha_k} \nabla f(\mathbf{u}^{(k-1)})$$

step size

- Stop at some point  
stopping criterion

# Gradient descent algorithm

- Choose an initial point  $\mathbf{u}^{(0)} \in \mathbb{R}^n$
- Repeat for  $k=1, 2, 3, \dots$

$$\mathbf{u}^{(k)} = \mathbf{u}^{(k-1)} - \boxed{\alpha_k} \nabla f(\mathbf{u}^{(k-1)})$$

step size

- Stop at some point

stopping criterion

Usually: stop when  $\|\nabla f(\mathbf{u}^{(k)})\|_2 < \epsilon$      $\epsilon = 10^{-\text{sthg}}$



$$\nabla f(\mathbf{u}) = 0 \Leftrightarrow \mathbf{u} \text{ minimizes } f$$

# Gradient descent algorithm

- Choose an initial point  $\mathbf{u}^{(0)} \in \mathbb{R}^n$
- Repeat for  $k=1, 2, 3, \dots$

$$\mathbf{u}^{(k)} = \mathbf{u}^{(k-1)} - \boxed{\alpha_k} \nabla f(\mathbf{u}^{(k-1)})$$

step size

- If the step size is too big



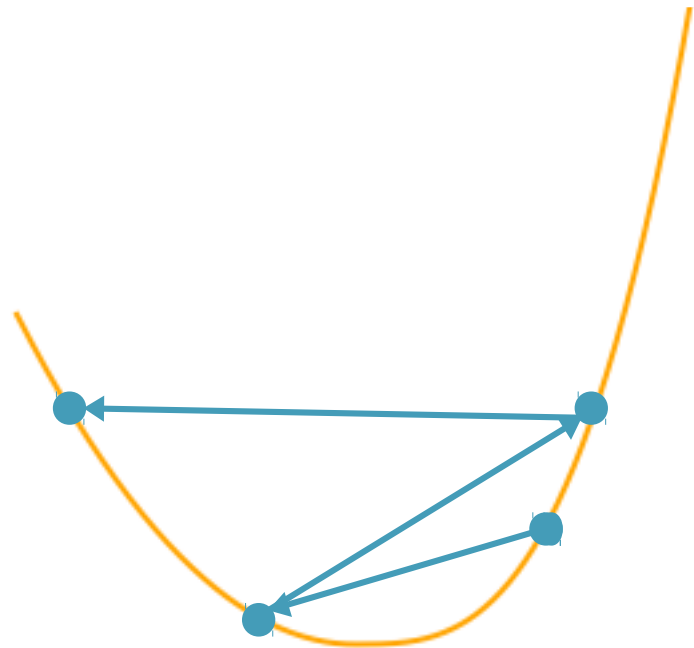
# Gradient descent algorithm

- Choose an initial point  $\mathbf{u}^{(0)} \in \mathbb{R}^n$
- Repeat for  $k=1, 2, 3, \dots$

$$\mathbf{u}^{(k)} = \mathbf{u}^{(k-1)} - \boxed{\alpha_k} \nabla f(\mathbf{u}^{(k-1)})$$

step size

- If the step size is too big,



# Gradient descent algorithm

- Choose an initial point  $\mathbf{u}^{(0)} \in \mathbb{R}^n$
- Repeat for  $k=1, 2, 3, \dots$

$$\mathbf{u}^{(k)} = \mathbf{u}^{(k-1)} - \boxed{\alpha_k} \nabla f(\mathbf{u}^{(k-1)})$$

step size

- If the step size is too big, the search might diverge

# Gradient descent algorithm

- Choose an initial point  $\mathbf{u}^{(0)} \in \mathbb{R}^n$
- Repeat for  $k=1, 2, 3, \dots$

$$\mathbf{u}^{(k)} = \mathbf{u}^{(k-1)} - \boxed{\alpha_k} \nabla f(\mathbf{u}^{(k-1)})$$

step size

- If the step size is too big, the search might diverge
- If the step size is too small, 

# Gradient descent algorithm

- Choose an initial point  $\mathbf{u}^{(0)} \in \mathbb{R}^n$
- Repeat for  $k=1, 2, 3, \dots$

$$\mathbf{u}^{(k)} = \mathbf{u}^{(k-1)} - \boxed{\alpha_k} \nabla f(\mathbf{u}^{(k-1)})$$

step size

- **If the step size is too big**, the search might diverge
- **If the step size is too small**, the search might take a very long time



# Gradient descent algorithm

- Choose an initial point  $\mathbf{u}^{(0)} \in \mathbb{R}^n$
- Repeat for  $k=1, 2, 3, \dots$

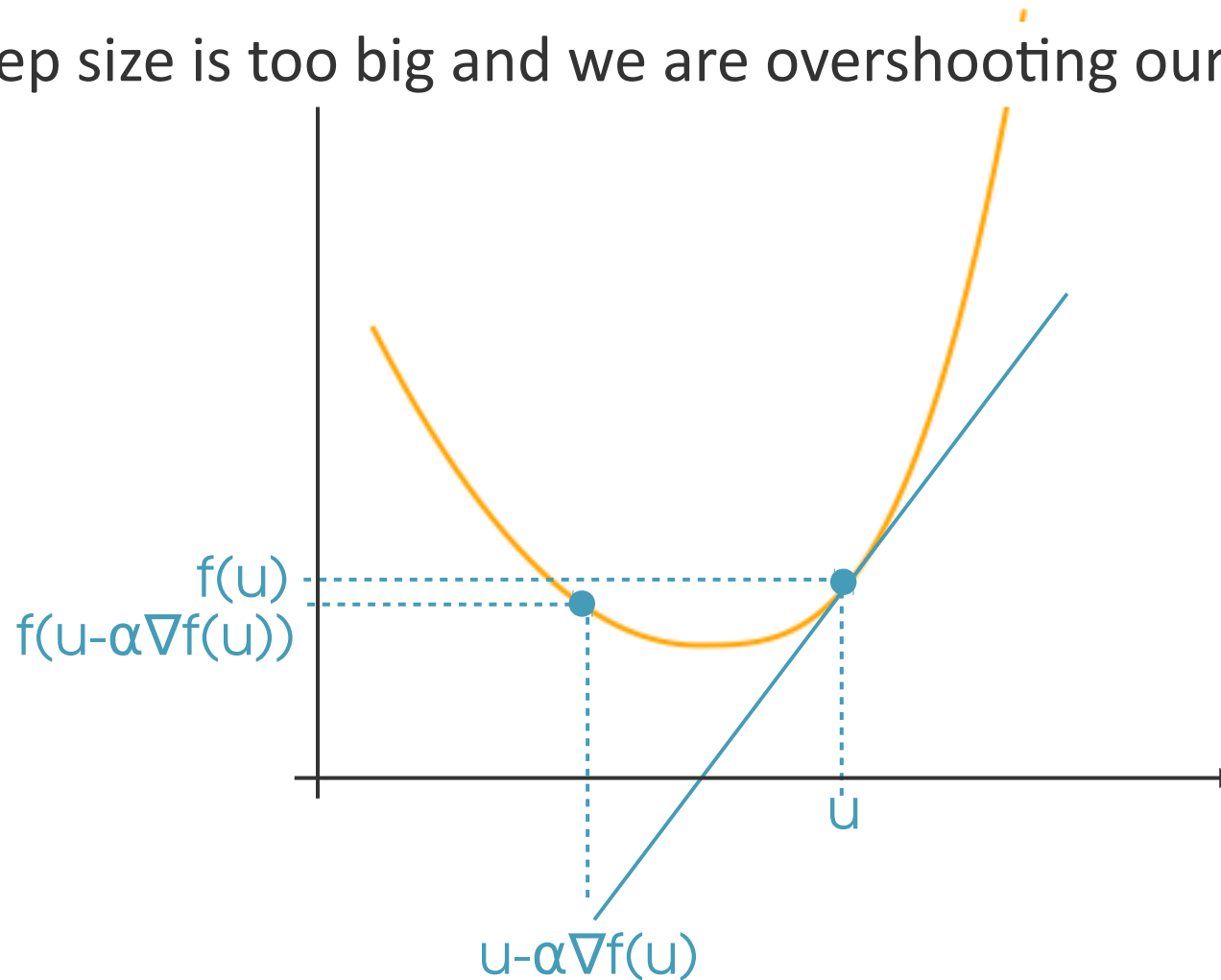
$$\mathbf{u}^{(k)} = \mathbf{u}^{(k-1)} - \boxed{\alpha_k} \nabla f(\mathbf{u}^{(k-1)})$$

step size

- **If the step size is too big**, the search might diverge
- **If the step size is too small**, the search might take a very long time
- **Backtracking line search** makes it possible to choose the step size **adaptively**.

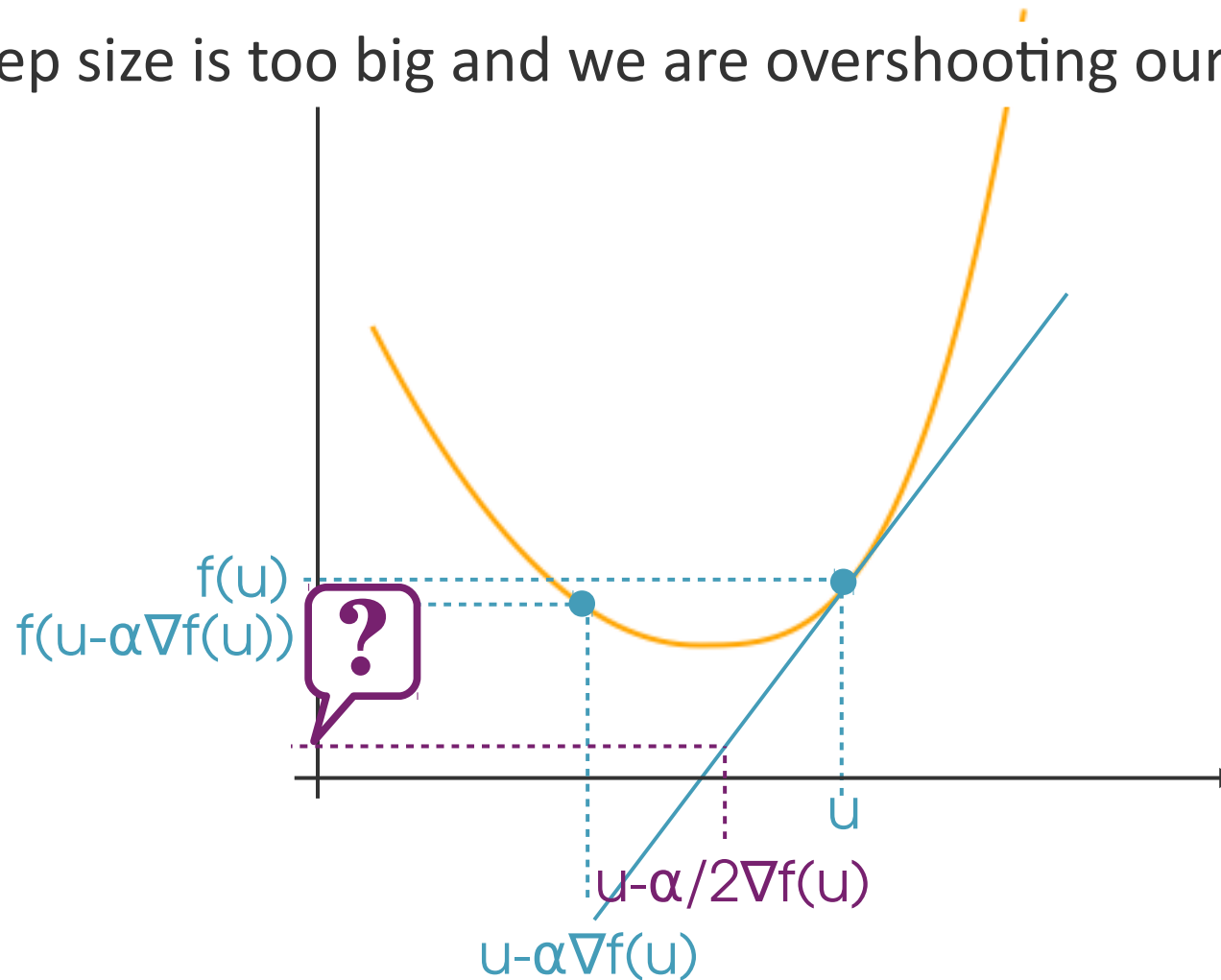
# BLS: shrinking needed

The step size is too big and we are overshooting our goal.



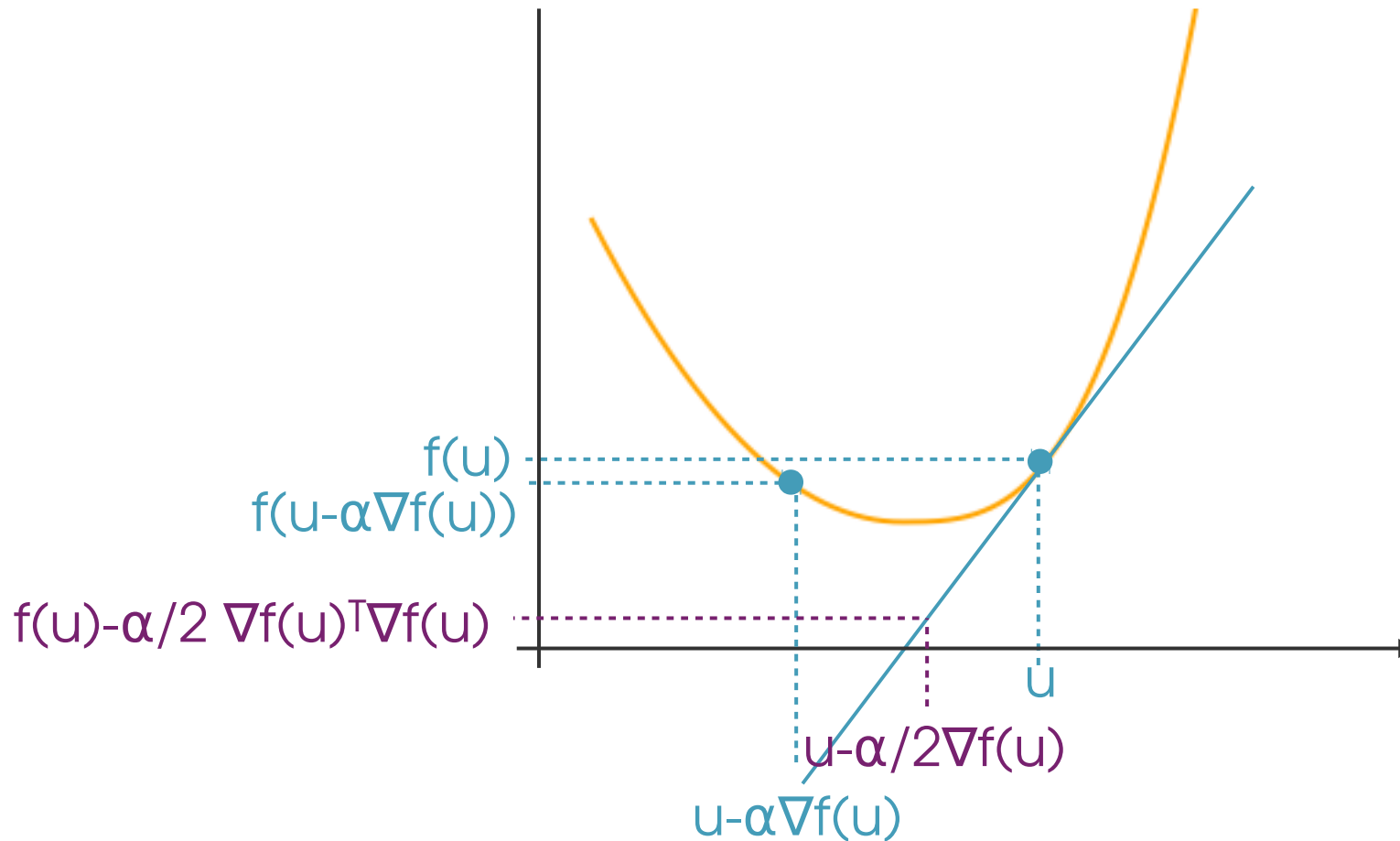
# BLS: shrinking needed

The step size is too big and we are overshooting our goal.



# BLS: shrinking needed

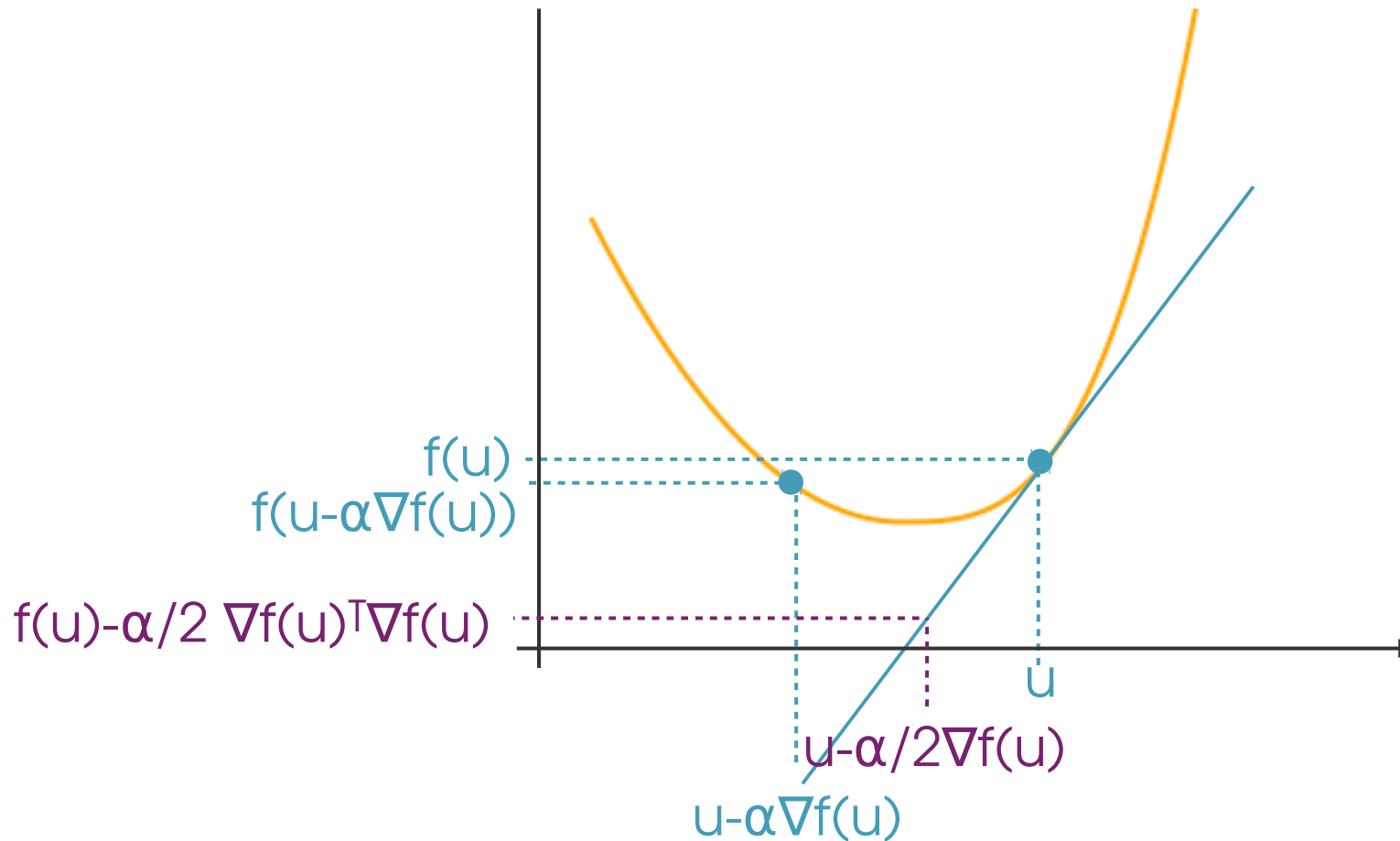
The step size is too big and we are overshooting our goal.



# BLS: shrinking needed

$$f(u - \alpha \nabla f(u)) > f(u) - \frac{\alpha}{2} \nabla f(u)^T \nabla f(u)$$

The step size is too big and we are overshooting our goal.

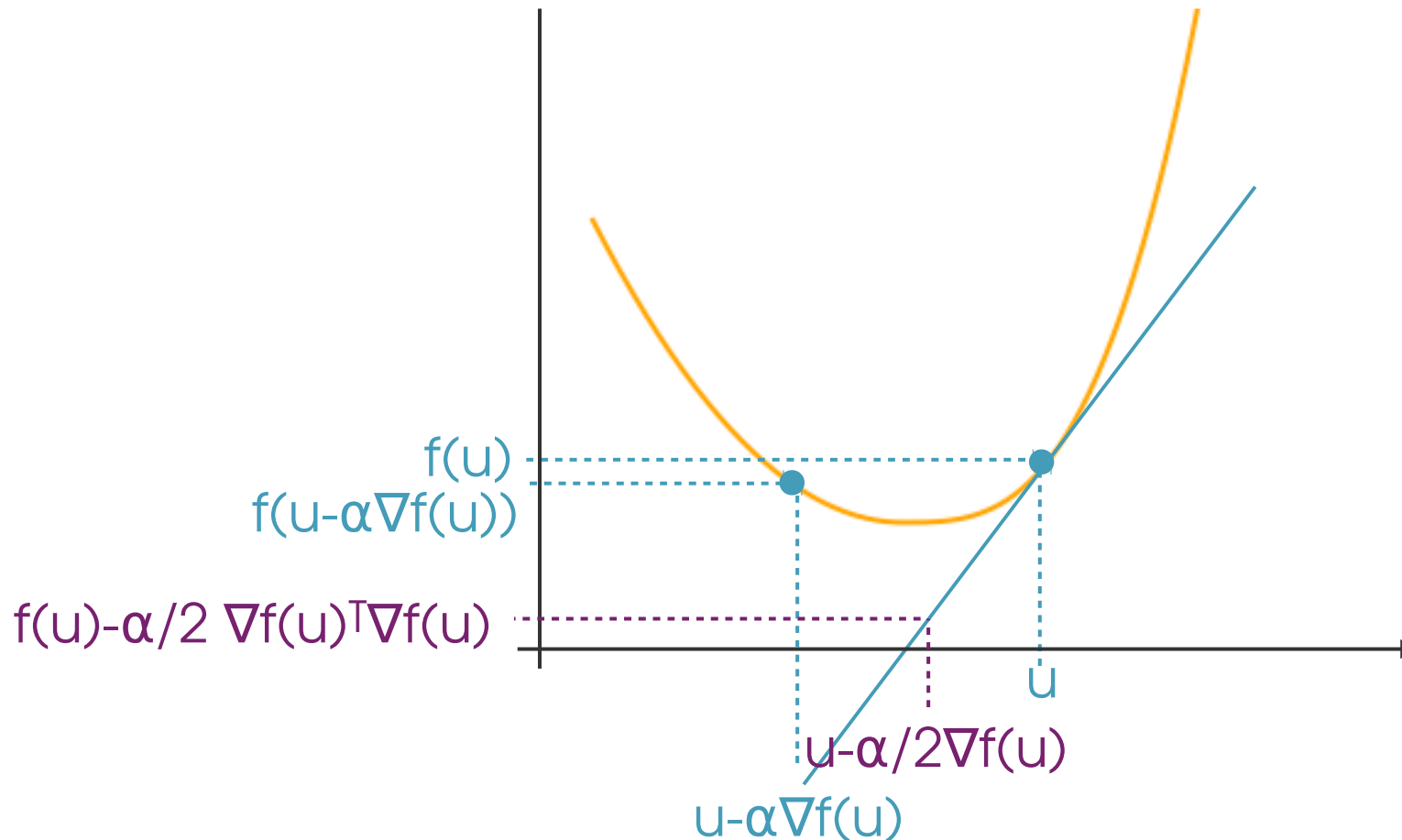


# BLS: shrinking needed

$$f(\mathbf{u}^{(k-1)} - \alpha_{k-1} \nabla f(\mathbf{u}^{(k-1)})) > f(\mathbf{u}^{(k-1)}) - \frac{1}{2} \alpha_{k-1} \|\nabla f(\mathbf{u}^{(k-1)})\|_2^2$$

$$f(\mathbf{u} - \alpha \nabla f(\mathbf{u})) > f(\mathbf{u}) - \alpha/2 \nabla f(\mathbf{u})^\top \nabla f(\mathbf{u})$$

The step size is too big and we are overshooting our goal.

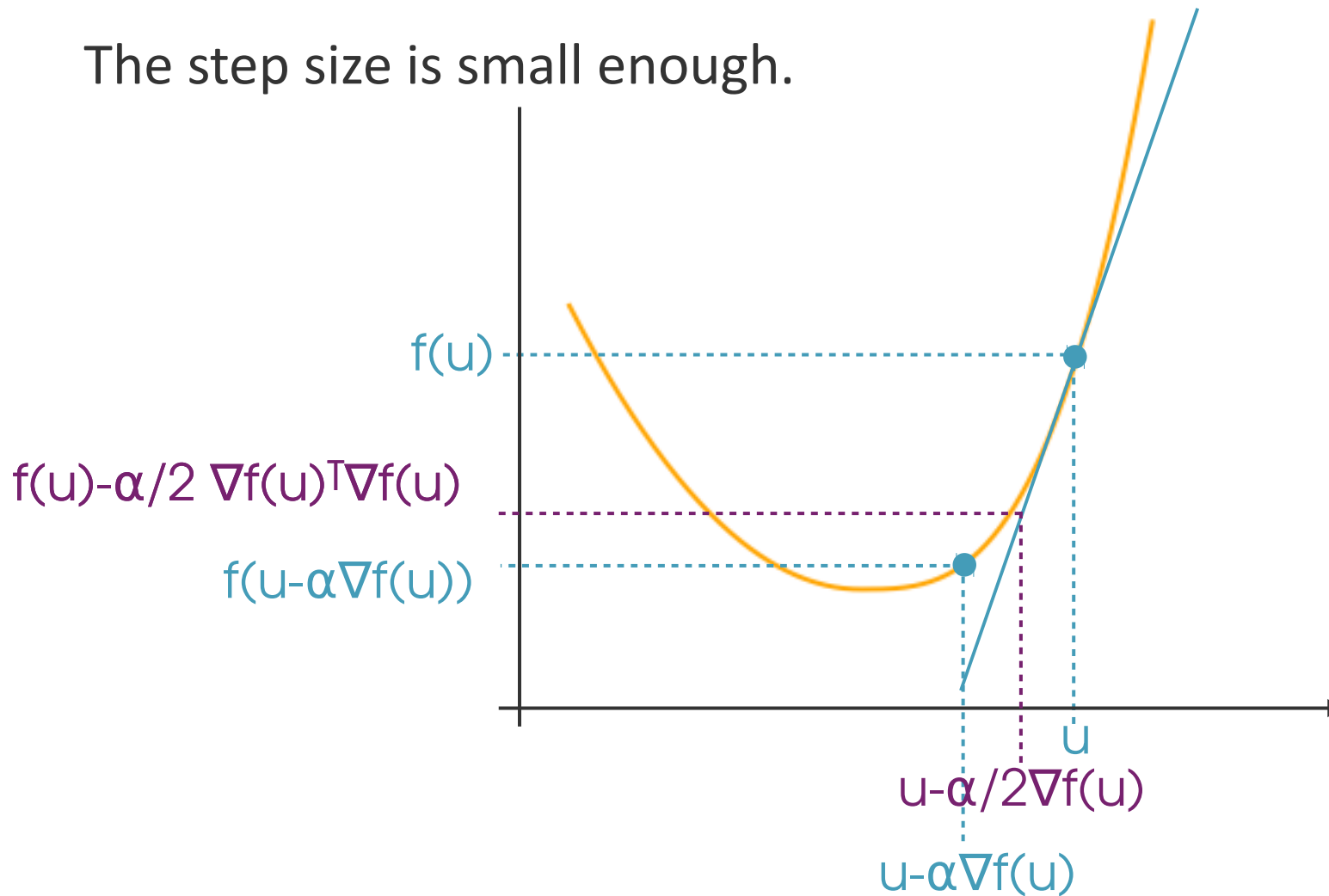


# BLS: no shrinking needed

$$f(\mathbf{u}^{(k-1)} - \alpha_{k-1} \nabla f(\mathbf{u}^{(k-1)})) \leq f(\mathbf{u}^{(k-1)}) - \frac{1}{2} \alpha_{k-1} \|\nabla f(\mathbf{u}^{(k-1)})\|_2^2$$

$$f(\mathbf{u} - \alpha \nabla f(\mathbf{u})) \leq f(\mathbf{u}) - \frac{\alpha}{2} \nabla f(\mathbf{u})^\top \nabla f(\mathbf{u})$$

The step size is small enough.



# Backtracking line search

- **Shrinking parameter**  $0 < \beta < 1$ , **initial step size**  $\alpha_0$
- Choose an initial point  $\mathbf{u}^{(0)} \in \mathbb{R}^n$
- Repeat for  $k=1, 2, 3, \dots$ 
  - If  $f(\mathbf{u}^{(k-1)} - \alpha_{k-1} \nabla f(\mathbf{u}^{(k-1)})) > f(\mathbf{u}^{(k-1)}) - \frac{1}{2} \alpha_{k-1} \|\nabla f(\mathbf{u}^{(k-1)})\|_2^2$   
**shrink the step size:**  $\alpha_k = \beta \alpha_{k-1}$
  - Else:  $\alpha_k = \alpha_{k-1}$
  - Update:  $\mathbf{u}^{(k)} = \mathbf{u}^{(k-1)} - \alpha_k \nabla f(\mathbf{u}^{(k-1)})$
- Stop when  $\|\nabla f(\mathbf{u}^{(k)})\|_2 < \epsilon$      $\epsilon = 10^{-\text{sthg}}$



# Newton's method

$$\mathbf{u}^{(k)} = \mathbf{u}^{(k-1)} - \alpha_k \nabla f(\mathbf{u}^{(k-1)})$$

- Suppose  $f$  is twice derivable
- Second-order Taylor's expansion:

$$f(\mathbf{v}) \approx f(\mathbf{u}) + \nabla f(\mathbf{u})^\top (\mathbf{v} - \mathbf{u}) + \frac{1}{2}(\mathbf{v} - \mathbf{u})^\top \nabla^2 f(\mathbf{u})(\mathbf{v} - \mathbf{u})$$

- Minimize in  $\mathbf{v}$  instead of in  $\mathbf{u}$

# Newton's method

$$\mathbf{u}^{(k)} = \mathbf{u}^{(k-1)} - \alpha_k \nabla f(\mathbf{u}^{(k-1)})$$

- Suppose  $f$  is twice derivable
- Second-order Taylor's expansion:

$$f(\mathbf{v}) \approx \boxed{f(\mathbf{u}) + \nabla f(\mathbf{u})^\top (\mathbf{v} - \mathbf{u}) + \frac{1}{2}(\mathbf{v} - \mathbf{u})^\top \nabla^2 f(\mathbf{u})(\mathbf{v} - \mathbf{u})}$$

$g(\mathbf{v})$

- Minimize in  $\mathbf{v}$  instead of in  $\mathbf{u}$  

# Newton's method

$$\mathbf{u}^{(k)} = \mathbf{u}^{(k-1)} - \alpha_k \nabla f(\mathbf{u}^{(k-1)})$$

- Suppose  $f$  is twice derivable
- Second-order Taylor's expansion:

$$f(\mathbf{v}) \approx \boxed{f(\mathbf{u}) + \nabla f(\mathbf{u})^\top (\mathbf{v} - \mathbf{u}) + \frac{1}{2}(\mathbf{v} - \mathbf{u})^\top \nabla^2 f(\mathbf{u})(\mathbf{v} - \mathbf{u})}$$

$g(\mathbf{v})$


- Minimize in  $\mathbf{v}$ :  $\nabla_{\mathbf{v}} g(\mathbf{v}) = \nabla f(\mathbf{u}) + \nabla^2 f(\mathbf{u})(\mathbf{v} - \mathbf{u})$

$$\nabla_{\mathbf{v}} g(\mathbf{v}) = 0 \Rightarrow \mathbf{v} - \mathbf{u} = -(\nabla^2 f(\mathbf{u}))^{-1} \nabla f(\mathbf{u})$$

$$\boxed{\alpha_k = \left( \nabla^2 f(\mathbf{u}^{(k-1)}) \right)^{-1}}$$

# Newton CG (conjugate gradient)

$$\mathbf{u}^{(k)} = \mathbf{u}^{(k-1)} - \left( \nabla^2 f(\mathbf{u}^{(k-1)}) \right)^{-1} \nabla f(\mathbf{u}^{(k-1)})$$

- Computing the inverse of the Hessian is **computationally intensive**.
- Instead, compute  $\nabla^2 f(\mathbf{u}^{(k-1)})$  and  $\nabla f(\mathbf{u}^{(k-1)})$  and solve  $\nabla^2 f(\mathbf{u}^{(k-1)})\delta_k = \nabla f(\mathbf{u}^{(k-1)})$  for  $\delta_k$
- What is the new update rule? 

# Newton CG (conjugate gradient)

$$\mathbf{u}^{(k)} = \mathbf{u}^{(k-1)} - \left( \nabla^2 f(\mathbf{u}^{(k-1)}) \right)^{-1} \nabla f(\mathbf{u}^{(k-1)})$$

- Computing the inverse of the Hessian is **computationally intensive**.
- Instead, compute  $\nabla^2 f(\mathbf{u}^{(k-1)})$  and  $\nabla f(\mathbf{u}^{(k-1)})$  and solve  $\nabla^2 f(\mathbf{u}^{(k-1)})\delta_k = \nabla f(\mathbf{u}^{(k-1)})$  for  $\delta_k$
- **New update rule:**

$$\mathbf{u}^{(k)} = \mathbf{u}^{(k-1)} - \delta_k$$

# Newton CG (conjugate gradient)

$$\mathbf{u}^{(k)} = \mathbf{u}^{(k-1)} - \left( \nabla^2 f(\mathbf{u}^{(k-1)}) \right)^{-1} \nabla f(\mathbf{u}^{(k-1)})$$

- Computing the inverse of the Hessian is **computationally intensive**.
- Instead, compute  $\nabla^2 f(\mathbf{u}^{(k-1)})$  and  $\nabla f(\mathbf{u}^{(k-1)})$  and solve  $\nabla^2 f(\mathbf{u}^{(k-1)})\delta_k = \nabla f(\mathbf{u}^{(k-1)})$  for  $\delta_k$

This is a problem of the form  $A\mathbf{x} - \mathbf{b} = 0 \quad A \succeq 0$

# Newton CG (conjugate gradient)

$$\mathbf{u}^{(k)} = \mathbf{u}^{(k-1)} - \left( \nabla^2 f(\mathbf{u}^{(k-1)}) \right)^{-1} \nabla f(\mathbf{u}^{(k-1)})$$

- Computing the inverse of the Hessian is **computationally intensive**.
- Instead, compute  $\nabla^2 f(\mathbf{u}^{(k-1)})$  and  $\nabla f(\mathbf{u}^{(k-1)})$  and solve  $\nabla^2 f(\mathbf{u}^{(k-1)}) \delta_k = \nabla f(\mathbf{u}^{(k-1)})$  for  $\delta_k$

This is a problem of the form  $A\mathbf{x} - \mathbf{b} = 0$

$$A \succeq 0$$



# Newton CG (conjugate gradient)

$$\mathbf{u}^{(k)} = \mathbf{u}^{(k-1)} - \left( \nabla^2 f(\mathbf{u}^{(k-1)}) \right)^{-1} \nabla f(\mathbf{u}^{(k-1)})$$

- Computing the inverse of the Hessian is **computationally intensive**.
- Instead, compute  $\nabla^2 f(\mathbf{u}^{(k-1)})$  and  $\nabla f(\mathbf{u}^{(k-1)})$  and solve  $\nabla^2 f(\mathbf{u}^{(k-1)})\delta_k = \nabla f(\mathbf{u}^{(k-1)})$  for  $\delta_k$

This is a problem of the form  $A\mathbf{x} - \mathbf{b} = 0$   $A \succeq 0$

**Second-order characterization of convex functions**



# Newton CG (conjugate gradient)

$$\mathbf{u}^{(k)} = \mathbf{u}^{(k-1)} - \left( \nabla^2 f(\mathbf{u}^{(k-1)}) \right)^{-1} \nabla f(\mathbf{u}^{(k-1)})$$

- Computing the inverse of the Hessian is **computationally intensive**.
- Instead, compute  $\nabla^2 f(\mathbf{u}^{(k-1)})$  and  $\nabla f(\mathbf{u}^{(k-1)})$  and solve  $\nabla^2 f(\mathbf{u}^{(k-1)})\delta_k = \nabla f(\mathbf{u}^{(k-1)})$  for  $\delta_k$

This is a problem of the form  $A\mathbf{x} - \mathbf{b} = 0 \quad A \succeq 0$

Solve using the **conjugate gradient method**.

# Conjugate gradient method

$$\text{Solve } Ax - b = 0 \quad A \succeq 0$$

- Idea: build a set of **A-conjugate vectors** (basis of  $\mathbb{R}^n$ )

$$\{v_1, v_2, \dots, v_n\} : v_i^\top A v_j = 0 \quad \forall i \neq j$$

- **Initialisation:**  $v_0 = r_0 = b - Ax^{(0)}$

- At step t:

- **Update rule:**  $x^{(t)} = x^{(t-1)} + \alpha_t v_{t-1}$
- **residual**  $r_t = b - Ax^{(t)}$
- $v_t = r_t + \beta_t v_{t-1}$

$$\alpha_t = \frac{r_{t-1}^\top r_{t-1}}{v_{t-1}^\top A v_{t-1}}$$

$$\beta_t = \frac{r_t^\top r_t}{r_{t-1}^\top r_{t-1}}$$

ensures

$$\begin{aligned} v_t^\top A v_{t-1} &= 0 \\ v_t^\top A v_i &= 0 \quad \forall i < t \end{aligned}$$

- **Convergence:**

$$r_i^\top r_j = 0 \quad \forall i \neq j \text{ hence } r_n = 0$$

# Conjugate gradient method

Prove  $v_t^\top A v_{t-1} = 0$



Given

– **Initialisation:**  $v_0 = r_0 = b - Ax^{(0)}$

– At step t:

• **Update rule:**  $x^{(t)} = x^{(t-1)} + \alpha_t v_{t-1}$

• **residual**  $r_t = b - Ax^{(t)}$

•  $v_t = r_t + \beta_t v_{t-1}$

and assuming  $r_t^\top r_{t-1} = 0$

$$\alpha_t = \frac{r_{t-1}^\top r_{t-1}}{v_{t-1}^\top A v_{t-1}}$$

$$\beta_t = \frac{r_t^\top r_t}{r_{t-1}^\top r_{t-1}}$$

**Prove**  $\mathbf{v}_t^\top A \mathbf{v}_{t-1} = 0$

Given

$$\alpha_t = \frac{\mathbf{r}_{t-1}^\top \mathbf{r}_{t-1}}{\mathbf{v}_{t-1}^\top A \mathbf{v}_{t-1}}$$

– **Initialisation:**  $\mathbf{v}_0 = \mathbf{r}_0 = \mathbf{b} - A\mathbf{x}^{(0)}$

– **Update rule:**  $\mathbf{x}^{(t)} = \mathbf{x}^{(t-1)} + \alpha_t \mathbf{v}_{t-1}$

– **residual**  $\mathbf{r}_t = \mathbf{b} - A\mathbf{x}^{(t)}$

$$\beta_t = \frac{\mathbf{r}_t^\top \mathbf{r}_t}{\mathbf{r}_{t-1}^\top \mathbf{r}_{t-1}}$$

–  $\mathbf{v}_t = \mathbf{r}_t + \beta_t \mathbf{v}_{t-1}$  and assuming  $\mathbf{r}_t^\top \mathbf{r}_{t-1} = 0$

By definition,  $\mathbf{r}_t = \mathbf{b} - A\mathbf{x}^{(t)}$  and  $\mathbf{x}^{(t)} = \mathbf{x}^{(t-1)} + \alpha_t \mathbf{v}_{t-1}$ . Hence  $\mathbf{r}_t = \mathbf{b} - A(\mathbf{x}^{(t-1)} + \alpha_t \mathbf{v}_{t-1})$  and

$$\mathbf{r}_t = \mathbf{r}_{t-1} - \alpha_t A \mathbf{v}_{t-1}. \quad (1)$$

By definition,  $\mathbf{v}_t = \mathbf{r}_t + \beta_t \mathbf{v}_{t-1}$  and  $\beta_t = \frac{\mathbf{r}_t^\top \mathbf{r}_t}{\mathbf{r}_{t-1}^\top \mathbf{r}_{t-1}}$ . Hence  $\mathbf{v}_t^\top A \mathbf{v}_{t-1} = \mathbf{r}_t^\top A \mathbf{v}_{t-1} + \beta_t \mathbf{v}_{t-1}^\top A \mathbf{v}_{t-1}$  and therefore

$$\mathbf{v}_t^\top A \mathbf{v}_{t-1} = \mathbf{r}_t^\top A \mathbf{v}_{t-1} + \frac{\mathbf{r}_t^\top \mathbf{r}_t}{\mathbf{r}_{t-1}^\top \mathbf{r}_{t-1}} \mathbf{v}_{t-1}^\top A \mathbf{v}_{t-1}.$$

Because  $\alpha_t = \frac{\mathbf{r}_{t-1}^\top \mathbf{r}_{t-1}}{\mathbf{v}_{t-1}^\top A \mathbf{v}_{t-1}}$ ,

$$\mathbf{v}_t^\top A \mathbf{v}_{t-1} = \mathbf{r}_t^\top A \mathbf{v}_{t-1} + \frac{1}{\alpha_t} \mathbf{r}_t^\top \mathbf{r}_t.$$

From (1),  $A \mathbf{v}_{t-1} = \frac{1}{\alpha_t} (\mathbf{r}_{t-1} - \mathbf{r}_t)$ , and therefore

$$\mathbf{v}_t^\top A \mathbf{v}_{t-1} = \frac{1}{\alpha_t} \mathbf{r}_t^\top (\mathbf{r}_{t-1} - \mathbf{r}_t) + \frac{1}{\alpha_t} \mathbf{r}_t^\top \mathbf{r}_t = \frac{1}{\alpha_t} \mathbf{r}_t^\top \mathbf{r}_{t-1} = 0.$$

This is true if and only if we have shown that  $\mathbf{r}_t^\top \mathbf{r}_{t-1} = 0$ .

# Conjugate gradient method

Prove  $r_t^\top r_{t-1} = 0$  and conclude the proof 

Given

– **Initialisation:**  $v_0 = r_0 = b - Ax^{(0)}$

– At step t:

• **Update rule:**  $x^{(t)} = x^{(t-1)} + \alpha_t v_{t-1}$

• **residual**  $r_t = b - Ax^{(t)}$

•  $v_t = r_t + \beta_t v_{t-1}$

$$\alpha_t = \frac{r_{t-1}^\top r_{t-1}}{v_{t-1}^\top A v_{t-1}}$$

$$\beta_t = \frac{r_t^\top r_t}{r_{t-1}^\top r_{t-1}}$$

**Prove**  $\mathbf{r}_t^\top \mathbf{r}_{t-1} = 0$

Given

$$\alpha_t = \frac{\mathbf{r}_{t-1}^\top \mathbf{r}_{t-1}}{\mathbf{v}_{t-1}^\top A \mathbf{v}_{t-1}}$$

$$\beta_t = \frac{\mathbf{r}_t^\top \mathbf{r}_t}{\mathbf{r}_{t-1}^\top \mathbf{r}_{t-1}}$$

- **Initialisation:**  $\mathbf{v}_0 = \mathbf{r}_0 = \mathbf{b} - A\mathbf{x}^{(0)}$
- **Update rule:**  $\mathbf{x}^{(t)} = \mathbf{x}^{(t-1)} + \alpha_t \mathbf{v}_{t-1}$
- **residual**  $\mathbf{r}_t = \mathbf{b} - A\mathbf{x}^{(t)}$
- $\mathbf{v}_t = \mathbf{r}_t + \beta_t \mathbf{v}_{t-1}$

From (1),

$$\begin{aligned} \mathbf{r}_t^\top \mathbf{r}_{t-1} &= \mathbf{r}_{t-1}^\top \mathbf{r}_{t-1} - \alpha_t (A\mathbf{v}_{t-1})^\top \mathbf{r}_{t-1} \\ &= \mathbf{r}_{t-1}^\top \mathbf{r}_{t-1} - \frac{\mathbf{r}_{t-1}^\top \mathbf{r}_{t-1}}{\mathbf{v}_{t-1}^\top A \mathbf{v}_{t-1}} \mathbf{v}_{t-1}^\top A^\top \mathbf{r}_{t-1} = (\mathbf{r}_{t-1}^\top \mathbf{r}_{t-1}) \left( 1 - \frac{\mathbf{v}_{t-1}^\top A^\top \mathbf{r}_{t-1}}{\mathbf{v}_{t-1}^\top A \mathbf{v}_{t-1}} \right). \end{aligned}$$

Because  $A \succeq 0$ ,  $A = A^\top$  and hence

$$\mathbf{r}_t^\top \mathbf{r}_{t-1} = (\mathbf{r}_{t-1}^\top \mathbf{r}_{t-1}) \left( 1 - \frac{\mathbf{v}_{t-1}^\top A \mathbf{r}_{t-1}}{\mathbf{v}_{t-1}^\top A \mathbf{v}_{t-1}} \right). \quad (2)$$

- If  $t = 1$ , because  $\mathbf{r}_0 = \mathbf{v}_0$ ,  $\mathbf{r}_1^\top \mathbf{r}_0 = 0$ . From the previous proof, we now have that  $\mathbf{v}_1^\top A \mathbf{v}_0 = 0$ .
- For  $t = 2$ , by definition  $\mathbf{v}_{t-1} = \mathbf{r}_{t-1} + \beta_{t-1} \mathbf{v}_{t-2}$  and we can replace  $\mathbf{r}_{t-1}$  accordingly to find that

$$\mathbf{v}_{t-1}^\top A \mathbf{r}_{t-1} = \mathbf{v}_{t-1}^\top A \mathbf{v}_{t-1} - \beta_{t-1} \mathbf{v}_{t-1}^\top A \mathbf{v}_{t-2}.$$

Because we have shown that  $\mathbf{v}_1^\top A \mathbf{v}_0 = 0$ , we conclude from (2) that  $\mathbf{r}_t^\top \mathbf{r}_{t-1} = 0$ .

- We can iterate this procedure to show alternatively that  $\mathbf{v}_t^\top A \mathbf{v}_{t-1} = 0$  and  $\mathbf{r}_t^\top \mathbf{r}_{t-1} = 0$  for all values of  $t$ .

# Quasi-Newton methods

$$\mathbf{u}^{(k)} = \mathbf{u}^{(k-1)} - \delta_k \quad \boxed{\nabla^2 f(\mathbf{u}^{(k-1)})} \delta_k = \nabla f(\mathbf{u}^{(k-1)})$$

- What if the Hessian is unavailable / expensive to compute at each iteration?

- **Approximate the inverse Hessian:**  $\delta_k = \boxed{W_{k-1}} \nabla f(\mathbf{u}^{(k-1)})$   
update  $W_k \approx (\nabla^2 f(\mathbf{u}^{(k)}))^{-1}$  iteratively

- **Conditions:**

- $W_k \succeq 0$

- **Secant equation:**  $\boxed{\nabla f(\mathbf{u}) - \nabla f(\mathbf{v}) \approx \nabla^2 f(\mathbf{u})(\mathbf{u} - \mathbf{v})}$

$$\Rightarrow W_k \left( \nabla f(\mathbf{u}^{(k)}) - \nabla f(\mathbf{u}^{(k-1)}) \right) = \mathbf{u}^{(k)} - \mathbf{u}^{(k-1)}$$

- **Initialization:** Identity  $W_0 = I_n$

# Quasi-Newton methods

$$\mathbf{u}^{(k)} = \mathbf{u}^{(k-1)} - \delta_k \quad \boxed{\nabla^2 f(\mathbf{u}^{(k-1)})} \delta_k = \nabla f(\mathbf{u}^{(k-1)})$$

- What if the Hessian is unavailable / expensive to compute at each iteration?

- **Approximate the inverse Hessian:**  $\delta_k = \boxed{W_{k-1}} \nabla f(\mathbf{u}^{(k-1)})$   
update  $W_k \approx (\nabla^2 f(\mathbf{u}^{(k)}))^{-1}$  iteratively

- **BFGS:** Broyden-Fletcher-Goldfarb-Shanno

$$W_k = W_{k-1} - \frac{\mathbf{s}_k \mathbf{y}_k^\top W_{k-1} + W_{k-1} \mathbf{y}_k \mathbf{s}_k^\top}{\langle \mathbf{y}_k, \mathbf{s}_k \rangle} + \left( 1 + \frac{\langle \mathbf{y}_k, W_{k-1} \mathbf{y}_k \rangle}{\langle \mathbf{y}_k, \mathbf{s}_k \rangle} \right) \frac{\mathbf{s}_k \mathbf{s}_k^\top}{\langle \mathbf{y}_k, \mathbf{s}_k \rangle}$$

$$\mathbf{s}_k = \mathbf{u}^{(k)} - \mathbf{u}^{(k-1)} \quad \mathbf{y}_k = \nabla f(\mathbf{u}^{(k)}) - \nabla f(\mathbf{u}^{(k-1)})$$

- **L-BFGS:** Limited memory variant

Do not store the full matrix  $W_k$ .



# Stochastic gradient descent

- For  $f : \mathbf{u} \mapsto \sum_{i=1}^m h_i(\mathbf{u})$

- **Gradient descent:**

$$\mathbf{u}^{(k)} = \mathbf{u}^{(k-1)} - \alpha_k \sum_{i=1}^m \nabla h_i(\mathbf{u}^{(k-1)})$$

- **Stochastic gradient descent:**

$$\mathbf{u}^{(k)} = \mathbf{u}^{(k-1)} - \alpha_k \nabla h_{i_k}(\mathbf{u}^{(k-1)})$$

- **Cyclic:** cycle over 1, 2, ..., m, 1, 2, ..., m, ...
- **Randomized:** chose  $i_k$  uniformly at random in  $\{1, 2, \dots, m\}$ .

# Coordinate Descent

- For  $f : \mathbb{R}^n \rightarrow \mathbb{R}$   
$$u \mapsto g(u) + \sum_{i=1}^n h_i(u_i)$$
  - $g$ : convex and differentiable
  - $h_i$ : convex  $\Rightarrow$  the **non-smooth** part of  $f$  is **separable**.
- Minimize **coordinate by coordinate**:
  - Initialisation:  $u^{(0)} \in \mathbb{R}^n$
  - For  $k=1, 2, \dots$ :
    - $u_1^{(k)} \in \arg \min_{u_1} f(\boxed{u_1}, u_2^{(k-1)}, \dots, u_n^{(k-1)})$
    - $u_2^{(k)} \in \arg \min_{u_2} f(u_1^{(k)}, \boxed{u_2}, \dots, u_n^{(k-1)})$
    - $\dots$
    - $u_n^{(k)} \in \arg \min_{u_n} f(u_1^{(k)}, u_2^{(k)}, \dots, \boxed{u_n})$

# Coordinate Descent

- For  $f : \mathbb{R}^n \rightarrow \mathbb{R}$   

$$u \mapsto g(u) + \sum_{i=1}^n h_i(u_i)$$
  - $g$ : convex and differentiable
  - $h_i$ : convex  $\Rightarrow$  the **non-smooth** part of  $f$  is **separable**.

- Minimize **coordinate by coordinate**:

- Initialisation:  $u^{(0)} \in \mathbb{R}^n$

- For  $k=1, 2, \dots$ :
 

$$u_1^{(k)} \in \arg \min_{u_1} f(\boxed{u_1}, u_2^{(k-1)}, \dots, u_n^{(k-1)})$$

$$u_2^{(k)} \in \arg \min_{u_2} f(\boxed{u_1^{(k)}}, \boxed{u_2}, \dots, u_n^{(k-1)})$$

...

$$u_n^{(k)} \in \arg \min_{u_n} f(\boxed{u_1^{(k)}}, \boxed{u_2^{(k)}}, \dots, \boxed{u_n})$$

# Coordinate Descent

- For  $f : \mathbb{R}^n \rightarrow \mathbb{R}$   
$$u \mapsto g(u) + \sum_{i=1}^n h_i(u_i)$$
  - $g$ : convex and differentiable
  - $h_i$ : convex  $\Rightarrow$  the **non-smooth** part of  $f$  is **separable**.
- Minimize **coordinate by coordinate**:
  - Initialisation:  $u^{(0)} \in \mathbb{R}^n$
  - For  $k=1, 2, \dots$ :
    - $u_1^{(k)} \in \arg \min_{u_1} f(u_1, u_2^{(k-1)}, \dots, u_n^{(k-1)})$
    - $u_2^{(k)} \in \arg \min_{u_2} f(u_1^{(k)}, u_2, \dots, u_n^{(k-1)})$

## Variants:

- re-order the coordinates randomly
- Proceed by blocks of coordinates (2 or more at a time)

# Summary: Unconstrained convex optimization

If  $f$  is **differentiable**

- Set its gradient to zero
- If hard to solve: **gradient descent**  $\mathbf{u}^{(k)} = \mathbf{u}^{(k-1)} - \alpha_k \nabla f(\mathbf{u}^{(k-1)})$

Setting the learning rate:

- **Backtracking Line Search** (adapt heuristically to avoid “overshooting”)
- **Newton’s method**: Suppose  $f$  **twice differentiable**
  - $\alpha_k = \left( \nabla^2 f(\mathbf{u}^{(k-1)}) \right)^{-1}$
  - If the Hessian is hard to invert, compute  $\delta_k = \alpha_k \nabla f(\mathbf{u}^{(k-1)})$  by solving  $\nabla^2 f(\mathbf{u}^{(k-1)}) \delta_k = -\nabla f(\mathbf{u}^{(k-1)})$  by the **conjugate gradient method**
  - If the Hessian is hard to compute, approximate the inverse Hessian with a **quasi-Newton method** such as **BFGS** (**L-BFGS**: less memory)
- If  $f$  is separable: **stochastic gradient descent**
- If the non-smooth part of  $f$  is separable: **coordinate descent**.

# Constrained convex optimization

# Constrained convex optimization

- **Convex optimization program/problem:**

$$\min_{\mathbf{u} \in D} f(\mathbf{u})$$

subject to  $g_i(\mathbf{u}) \leq 0, i = 1, \dots, m$

$$h_j(\mathbf{u}) = 0, j = 1, \dots, r$$

- $f$  is **convex**
- $g_i, i = 1, \dots, m$  are **convex**
- $h_j, j = 1, \dots, r$  are **affine**  $h_j : \mathbf{u} \mapsto \mathbf{a}_j^\top \mathbf{u} + b_j$
- The **feasible set** is convex

$$\mathcal{C} = \{\mathbf{v} : \mathbf{v} \in D; g_i(\mathbf{v}) \leq 0, i = 1, \dots, m; h_j(\mathbf{v}) = 0, j = 1, \dots, r\}$$

# Lagrangian

$$\min_{\mathbf{u} \in D} f(\mathbf{u})$$

subject to  $g_i(\mathbf{u}) \leq 0, i = 1, \dots, m$

$$h_j(\mathbf{u}) = 0, j = 1, \dots, r$$

- **Lagrangian:**  $L : \mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R}^r \rightarrow \mathbb{R}$

$$\mathbf{u}, \boxed{\boldsymbol{\alpha}, \boldsymbol{\beta}} \mapsto f(\mathbf{u}) + \sum_{i=1}^m \boxed{\alpha_i} g_i(\mathbf{u}) + \sum_{j=1}^r \boxed{\beta_j} h_j(\mathbf{u})$$

= Lagrange multipliers

= dual variables



# Lagrange dual function

- Lagrangian:  $L : \mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R}^r \rightarrow \mathbb{R}$

$$u, \alpha, \beta \mapsto f(u) + \sum_{i=1}^m \alpha_i g_i(u) + \sum_{j=1}^r \beta_j h_j(u)$$

- Lagrange dual function:

$$Q : \mathbb{R}^m \times \mathbb{R}^r \rightarrow \mathbb{R}$$
$$\alpha, \beta \mapsto \inf_{u \in D} L(u, \alpha, \beta)$$

**Infimum** = the greatest value  $x$   
such that  $x \leq L(u, \alpha, \beta)$

- $Q$  is **concave** (independently of the convexity of  $f$ )

# Lagrange dual function

- $Q : \mathbb{R}^m \times \mathbb{R}^r \rightarrow \mathbb{R}$

$$\alpha, \beta \mapsto \inf_{\mathbf{u} \in D} f(\mathbf{u}) + \sum_{i=1}^m \alpha_i g_i(\mathbf{u}) + \sum_{j=1}^r \beta_j h_j(\mathbf{u})$$

- Q is **concave** (independently of the convexity of f)

Consider  $(\alpha_1, \beta_1), (\alpha_2, \beta_2)$  and  $0 \leq \lambda \leq 1$ . Set  $\alpha = \lambda\alpha_1 + (1 - \lambda)\alpha_2$  and  $\beta = \lambda\beta_1 + (1 - \lambda)\beta_2$ .  
 $Q(\alpha, \beta) = \inf_{\mathbf{u} \in D} f(\mathbf{u}) + \alpha^\top \mathbf{g}(\mathbf{u}) + \beta^\top \mathbf{h}(\mathbf{u})$ ,  
where  $\mathbf{g}(\mathbf{u}) = (g_1(\mathbf{u}), g_2(\mathbf{u}), \dots, g_m(\mathbf{u}))$  and  $\mathbf{h}(\mathbf{u}) = (h_1(\mathbf{u}), h_2(\mathbf{u}), \dots, h_r(\mathbf{u}))$ . Hence


$$\begin{aligned} Q(\alpha, \beta) &= \inf_{\mathbf{u} \in D} f(\mathbf{u}) + \lambda\alpha_1^\top \mathbf{g}(\mathbf{u}) + (1 - \lambda)\alpha_2^\top \mathbf{g}(\mathbf{u}) + \lambda\beta_1^\top \mathbf{h}(\mathbf{u}) + (1 - \lambda)\beta_2^\top \mathbf{h}(\mathbf{u}) \\ &= \inf_{\mathbf{u} \in D} \lambda (f(\mathbf{u}) + \alpha_1^\top \mathbf{g}(\mathbf{u}) + \beta_1^\top \mathbf{h}(\mathbf{u})) + (1 - \lambda) (f(\mathbf{u}) + \alpha_2^\top \mathbf{g}(\mathbf{u}) + \beta_2^\top \mathbf{h}(\mathbf{u})). \end{aligned}$$

This last equality holds because  $f(\mathbf{u}) = \lambda f(\mathbf{u}) + (1 - \lambda)f(\mathbf{u})$ . Hence

$$\begin{aligned} Q(\alpha, \beta) &\geq \lambda \inf_{\mathbf{u} \in D} (f(\mathbf{u}) + \alpha_1^\top \mathbf{g}(\mathbf{u}) + \beta_1^\top \mathbf{h}(\mathbf{u})) + (1 - \lambda) \inf_{\mathbf{u} \in D} (f(\mathbf{u}) + \alpha_2^\top \mathbf{g}(\mathbf{u}) + \beta_2^\top \mathbf{h}(\mathbf{u})) \\ &\geq \lambda Q(\alpha_1, \beta_1) + (1 - \lambda) Q(\alpha_2, \beta_2). \end{aligned}$$

# Lagrange dual function

- The dual function gives a **lower bound** on our solution


Let  $p^* = \min_{\mathbf{u} \in \mathcal{C}} f(\mathbf{u})$   **feasible set**

Then for any  $\boldsymbol{\alpha} \in \mathbb{R}_+^m$   $\alpha_1 \geq 0, \alpha_2 \geq 0, \dots, \alpha_m \geq 0$   
 $\boldsymbol{\beta} \in \mathbb{R}^r$

$$Q(\boldsymbol{\alpha}, \boldsymbol{\beta}) \leq p^*$$

Consider any feasible point  $\mathbf{u}$ . Then  $g_1(\mathbf{u}) \leq 0, g_2(\mathbf{u}) \leq 0, \dots, g_m(\mathbf{u}) \leq 0$  and  $h_1(\mathbf{u}) = h_2(\mathbf{u}) = \dots = h_r(\mathbf{u}) = 0$ . By definition,  $L(\mathbf{u}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = f(\mathbf{u}) + \sum_{i=1}^m \alpha_i g_i(\mathbf{u}) + \sum_{j=1}^r \beta_j h_j(\mathbf{u})$ . Because the second term is negative ( $\alpha_i \geq 0, g_i(\mathbf{u}) \leq 0$ ) and the third one is zero ( $h_j(\mathbf{u}) = 0$ ),  $L(\mathbf{u}, \boldsymbol{\alpha}, \boldsymbol{\beta}) \leq f(\mathbf{u})$  for every feasible point  $\mathbf{u}$ . Hence  $\inf_{\mathbf{u} \in \mathcal{D}} L(\mathbf{u}, \boldsymbol{\alpha}, \boldsymbol{\beta}) \leq \min_{\mathbf{u} \in \mathcal{C}} f(\mathbf{u})$ .

# Weak duality

- $Q(\alpha, \beta) \leq p^*$  for any  $\alpha \in \mathbb{R}_+^m, \beta \in \mathbb{R}^r$
- What is the **best lower bound** on  $p^*$  we can get? 

# Weak duality

- $Q(\alpha, \beta) \leq p^*$  for any  $\alpha \in \mathbb{R}_+^m, \beta \in \mathbb{R}^r$
- What is the **best lower bound** on  $p^*$  we can get?

$$\max_{\alpha, \beta} Q(\alpha, \beta)$$

subject to  $\alpha_1 \geq 0, \alpha_2 \geq 0, \dots, \alpha_m \geq 0$

# Weak duality

- $Q(\alpha, \beta) \leq p^*$  for any  $\alpha \in \mathbb{R}_+^m, \beta \in \mathbb{R}^r$
- What is the **best lower bound** on  $p^*$  we can get?

$$\max_{\alpha, \beta} Q(\alpha, \beta)$$

subject to  $\alpha_1 \geq 0, \alpha_2 \geq 0, \dots, \alpha_m \geq 0$

**Lagrange dual problem**

# Weak duality

- $Q(\alpha, \beta) \leq p^*$  for any  $\alpha \in \mathbb{R}_+^m, \beta \in \mathbb{R}^r$
- What is the **best lower bound** on  $p^*$  we can get?

$$\max_{\alpha, \beta} Q(\alpha, \beta)$$

subject to  $\alpha_1 \geq 0, \alpha_2 \geq 0, \dots, \alpha_m \geq 0$

**Lagrange dual problem**

- Optimal values  $\alpha^*, \beta^*$  of  $\alpha, \beta$  are called **dual optimal** or **optimal Lagrange multipliers**.
- Original optimization problem = **primal**
- The dual is a **convex optimization problem** (even if the primal is not!)

# Weak duality

- Let  $d^*$  be the solution to the dual problem

$$d^* = \max_{\alpha, \beta} Q(\alpha, \beta)$$

subject to  $\alpha_1 \geq 0, \alpha_2 \geq 0, \dots, \alpha_m \geq 0$

- Because for every dual admissible  $\alpha, \beta$   $Q(\alpha, \beta) \leq p^*$

$$d^* \leq p^*$$

**Weak duality**  
(always holds)



# Strong duality & Slater's conditions

- **Strong duality:**  $d^* = p^*$ 
  - Does not hold in general
  - But often holds for convex optimization problems

# Strong duality & Slater's conditions

- **Strong duality:**  $d^* = p^*$ 
  - Does not hold in general
  - But often holds for convex optimization problems
- **Constraint qualifications:** conditions under which strong duality holds (in addition to convexity)

# Strong duality & Slater's conditions

- **Strong duality:**  $d^* = p^*$ 
  - Does not hold in general
  - But often holds for convex optimization problems
- **Constraint qualifications:** conditions under which strong duality holds (in addition to convexity)
- In particular: **Slater's conditions:**
  - If the primal is **convex** and there exists at least one **strictly feasible point** (i.e. the inequalities hold strictly), then strong duality holds
  - Strict inequalities only need to hold for **non-affine** constraints.

# Karush-Kuhn-Tucker conditions

- Suppose  $f, g_i, h_j$  differentiable + strong duality

$$f(u^*) = Q(\alpha^*, \beta^*) \quad \text{[strong duality]}$$

# Karush-Kuhn-Tucker conditions

- Suppose  $f, g_i, h_j$  differentiable + strong duality

$$f(\mathbf{u}^*) = Q(\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*) \quad \text{[strong duality]}$$

$$= \inf_{\mathbf{u} \in D} f(\mathbf{u}) + \sum_{i=1}^m \alpha_i^* g_i(\mathbf{u}) + \sum_{j=1}^r \beta_j^* h_j(\mathbf{u}) \quad \text{[definition of Q]}$$

# Karush-Kuhn-Tucker conditions

- Suppose  $f, g_i, h_j$  differentiable + strong duality

$$f(\mathbf{u}^*) = Q(\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*) \quad \text{[strong duality]}$$

$$= \inf_{\mathbf{u} \in D} f(\mathbf{u}) + \sum_{i=1}^m \alpha_i^* g_i(\mathbf{u}) + \sum_{j=1}^r \beta_j^* h_j(\mathbf{u}) \quad \text{[definition of } Q\text{]}$$

$$\leq f(\mathbf{u}^*) + \sum_{i=1}^m \alpha_i^* g_i(\mathbf{u}^*) + \sum_{j=1}^r \beta_j^* h_j(\mathbf{u}^*) \quad \text{[definition of inf]}$$

# Karush-Kuhn-Tucker conditions

- Suppose  $f, g_i, h_j$  differentiable + strong duality

$$f(\mathbf{u}^*) = Q(\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*) \quad \text{[strong duality]}$$

$$= \inf_{\mathbf{u} \in D} f(\mathbf{u}) + \sum_{i=1}^m \alpha_i^* g_i(\mathbf{u}) + \sum_{j=1}^r \beta_j^* h_j(\mathbf{u}) \quad \text{[definition of } Q]$$

$$\leq f(\mathbf{u}^*) + \sum_{i=1}^m \alpha_i^* g_i(\mathbf{u}^*) + \sum_{j=1}^r \beta_j^* h_j(\mathbf{u}^*) \quad \text{[definition of inf]}$$

What is the sign of this expression?



# Karush-Kuhn-Tucker conditions

- Suppose  $f, g_i, h_j$  differentiable + strong duality

$$f(\mathbf{u}^*) = Q(\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*) \quad [\text{strong duality}]$$

$$= \inf_{\mathbf{u} \in D} f(\mathbf{u}) + \sum_{i=1}^m \alpha_i^* g_i(\mathbf{u}) + \sum_{j=1}^r \beta_j^* h_j(\mathbf{u}) \quad [\text{definition of } Q]$$

$$\leq f(\mathbf{u}^*) + \boxed{\sum_{i=1}^m \alpha_i^* g_i(\mathbf{u}^*)} + \boxed{\sum_{j=1}^r \beta_j^* h_j(\mathbf{u}^*)} \quad [\text{definition of inf}]$$

$\mathbf{u}^* \text{ feasible} \Rightarrow h_j = 0$

$$\mathbf{u}^* \text{ feasible} \Rightarrow g_i \leq 0$$

$$\boldsymbol{\alpha}^* \text{ feasible} \Rightarrow \alpha_i \geq 0$$



# Karush-Kuhn-Tucker conditions

- Suppose  $f, g_i, h_j$  differentiable + strong duality

$$f(\mathbf{u}^*) = Q(\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*) \quad \text{[strong duality]}$$

$$= \inf_{\mathbf{u} \in D} f(\mathbf{u}) + \sum_{i=1}^m \alpha_i^* g_i(\mathbf{u}) + \sum_{j=1}^r \beta_j^* h_j(\mathbf{u}) \quad \text{[definition of } Q]$$

$$\leq f(\mathbf{u}^*) + \boxed{\sum_{i=1}^m \alpha_i^* g_i(\mathbf{u}^*)} + \boxed{\sum_{j=1}^r \beta_j^* h_j(\mathbf{u}^*)} \quad \text{[definition of inf]}$$

$\mathbf{u}^*$  feasible  $\Rightarrow h_j = 0$

$$\mathbf{u}^* \text{ feasible} \Rightarrow g_i \leq 0$$

$$\leq f(\mathbf{u}^*). \quad \boldsymbol{\alpha}^* \text{ feasible} \Rightarrow \alpha_i \geq 0$$

# Karush-Kuhn-Tucker conditions

- Suppose  $f, g_i, h_j$  differentiable + strong duality

$$f(\mathbf{u}^*) = Q(\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*) \quad [\text{strong duality}]$$

$$= \inf_{\mathbf{u} \in D} f(\mathbf{u}) + \sum_{i=1}^m \alpha_i^* g_i(\mathbf{u}) + \sum_{j=1}^r \beta_j^* h_j(\mathbf{u}) \quad [\text{definition of } Q]$$


$$\leq f(\mathbf{u}^*) + \boxed{\sum_{i=1}^m \alpha_i^* g_i(\mathbf{u}^*)} + \boxed{\sum_{j=1}^r \beta_j^* h_j(\mathbf{u}^*)} \quad [\text{definition of inf}]$$

$\mathbf{u}^*$  feasible  $\Rightarrow h_j = 0$

$$\mathbf{u}^* \text{ feasible} \Rightarrow g_i \leq 0$$

$$\boldsymbol{\alpha}^* \text{ feasible} \Rightarrow \alpha_i \geq 0$$

$$\leq f(\mathbf{u}^*).$$

- Hence 

# Karush-Kuhn-Tucker conditions

- Suppose  $f, g_i, h_j$  differentiable + strong duality

$$f(\mathbf{u}^*) = Q(\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*) \quad \text{[strong duality]}$$

$$= \inf_{\mathbf{u} \in D} f(\mathbf{u}) + \sum_{i=1}^m \alpha_i^* g_i(\mathbf{u}) + \sum_{j=1}^r \beta_j^* h_j(\mathbf{u}) \quad \text{[definition of } Q]$$

$$\leq f(\mathbf{u}^*) + \boxed{\sum_{i=1}^m \alpha_i^* g_i(\mathbf{u}^*)} + \boxed{\sum_{j=1}^r \beta_j^* h_j(\mathbf{u}^*)} \quad \text{[definition of inf]}$$

$\mathbf{u}^*$  feasible  $\Rightarrow h_j = 0$

$$\mathbf{u}^* \text{ feasible} \Rightarrow g_i \leq 0$$

$$\leq f(\mathbf{u}^*). \quad \boldsymbol{\alpha}^* \text{ feasible} \Rightarrow \alpha_i \geq 0$$

- Hence all above inequalities are equalities

# Karush-Kuhn-Tucker conditions

- Suppose  $f, g_i, h_j$  differentiable + strong duality

$$f(\mathbf{u}^*) = Q(\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*) \quad [\text{strong duality}]$$

$$= \inf_{\mathbf{u} \in D} f(\mathbf{u}) + \sum_{i=1}^m \alpha_i^* g_i(\mathbf{u}) + \sum_{j=1}^r \beta_j^* h_j(\mathbf{u}) \quad [\text{definition of } Q]$$

$$= f(\mathbf{u}^*) + \boxed{\sum_{i=1}^m \alpha_i^* g_i(\mathbf{u}^*)} + \boxed{\sum_{j=1}^r \beta_j^* h_j(\mathbf{u}^*)} \quad [\text{definition of inf}]$$

$\mathbf{u}^* \text{ feasible} \Rightarrow h_j = 0$

$$\mathbf{u}^* \text{ feasible} \Rightarrow g_i \leq 0$$

$$= f(\mathbf{u}^*). \quad \boldsymbol{\alpha}^* \text{ feasible} \Rightarrow \alpha_i \geq 0$$

# Karush-Kuhn-Tucker conditions

- Suppose  $f, g_i, h_j$  differentiable + strong duality

$$\begin{aligned}
 f(u^*) &= Q(\alpha^*, \beta^*) && \text{[strong duality]} \\
 &= \inf_{u \in D} \overbrace{f(u) + \sum_{i=1}^m \alpha_i^* g_i(u) + \sum_{j=1}^r \beta_j^* h_j(u)}^{L(u, \alpha^*, \beta^*)} && \text{[definition of } Q\text{]} \\
 &= f(u^*) + \underbrace{\sum_{i=1}^m \alpha_i^* g_i(u^*)}_{\substack{\mathbf{u}^* \text{ feasible} \Rightarrow g_i \leq 0}} + \underbrace{\sum_{j=1}^r \beta_j^* h_j(u^*)}_{\substack{\mathbf{u}^* \text{ feasible} \Rightarrow h_j = 0}} && \text{[definition of inf]} \\
 &= f(u^*). && \substack{\mathbf{u}^* \text{ feasible} \Rightarrow g_i \leq 0 \\ \alpha^* \text{ feasible} \Rightarrow \alpha_i \geq 0}
 \end{aligned}$$

- $f(u^*) = L(u^*, \alpha^*, \beta^*) = \inf_{u \in D} L(u, \alpha^*, \beta^*)$

# Karush-Kuhn-Tucker conditions

- Suppose  $f, g_i, h_j$  differentiable + strong duality

$$f(u^*) = Q(\alpha^*, \beta^*) \quad \text{[strong duality]}$$

$$= \inf_{u \in D} \left[ f(u) + \sum_{i=1}^m \alpha_i^* g_i(u) + \sum_{j=1}^r \beta_j^* h_j(u) \right] \quad \text{[definition of } Q \text{]}$$

$$= f(u^*) + \sum_{i=1}^m \alpha_i^* g_i(u^*) + \sum_{j=1}^r \beta_j^* h_j(u^*) \quad \text{[definition of inf]}$$

$u^*$  feasible  $\Rightarrow h_j = 0$

$u^*$  feasible  $\Rightarrow g_i \leq 0$

$$= f(u^*). \quad \alpha^* \text{ feasible} \Rightarrow \alpha_i \geq 0$$

- $f(u^*) = L(u^*, \alpha^*, \beta^*) = \inf_{u \in D} L(u, \alpha^*, \beta^*)$

$$\nabla_u L(u^*, \alpha^*, \beta^*) = 0$$

stationarity

# Karush-Kuhn-Tucker conditions

- Suppose  $f, g_i, h_j$  differentiable + strong duality

$$\begin{aligned}
 f(\mathbf{u}^*) &= Q(\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*) && \text{[strong duality]} \\
 &= \inf_{\mathbf{u} \in D} \left[ f(\mathbf{u}) + \sum_{i=1}^m \alpha_i^* g_i(\mathbf{u}) + \sum_{j=1}^r \beta_j^* h_j(\mathbf{u}) \right] && \text{[definition of } Q\text{]} \\
 &= f(\mathbf{u}^*) + \sum_{i=1}^m \alpha_i^* g_i(\mathbf{u}^*) + \sum_{j=1}^r \beta_j^* h_j(\mathbf{u}^*) && \text{[definition of inf]} \\
 &= f(\mathbf{u}^*). && \begin{array}{l} \mathbf{u}^* \text{ feasible} \Rightarrow h_j = 0 \\ \mathbf{u}^* \text{ feasible} \Rightarrow g_i \leq 0 \\ \boldsymbol{\alpha}^* \text{ feasible} \Rightarrow \alpha_i \geq 0 \end{array}
 \end{aligned}$$

## complementary slackness

- $\alpha_i^* g_i(\mathbf{u}^*) = 0$ 
  - $\alpha_i > 0 \Rightarrow g_i(\mathbf{u}^*) = 0$
  - $g_i(\mathbf{u}^*) < 0 \Rightarrow \alpha_i = 0$

# Karush-Kuhn-Tucker conditions

- Let's sum up all of our conditions:

- **Primal feasibility:**  $g_i(\mathbf{u}^*) \leq 0 \quad i = 1, \dots, m$

$$h_j(\mathbf{u}^*) = 0 \quad j = 1, \dots, r$$

- **Dual feasibility:**  $\alpha_i^* \geq 0 \quad i = 1, \dots, m$

- **Complementary slackness:**  $\alpha_i^* g_i(\mathbf{u}^*) = 0 \quad i = 1, \dots, m$

- **Stationarity:** 
$$\nabla_u f(\mathbf{u}^*) + \sum_{i=1}^m \alpha_i^* \nabla_u g_i(\mathbf{u}^*) + \sum_{j=1}^r \beta_j^* \nabla_u h_j(\mathbf{u}^*) = 0.$$



# Karush-Kuhn-Tucker conditions

- Let's sum up all of our conditions:

- Primal feasibility:**  $g_i(\mathbf{u}^*) \leq 0 \quad i = 1, \dots, m$   
 $h_j(\mathbf{u}^*) = 0 \quad j = 1, \dots, r$
- Dual feasibility:**  $\alpha_i^* \geq 0 \quad i = 1, \dots, m$
- Complementary slackness:**  $\alpha_i^* g_i(\mathbf{u}^*) = 0 \quad i = 1, \dots, m$
- Stationarity:** 
$$\nabla_u f(\mathbf{u}^*) + \sum_{i=1}^m \alpha_i^* \nabla_u g_i(\mathbf{u}^*) + \sum_{j=1}^r \beta_j^* \nabla_u h_j(\mathbf{u}^*) = 0.$$

**Karush-Kuhn-Tucker (KKT) conditions**

# Karush-Kuhn-Tucker conditions

- Let's sum up all of our conditions:

- Primal feasibility:**  $g_i(\mathbf{u}^*) \leq 0 \quad i = 1, \dots, m$   
 $h_j(\mathbf{u}^*) = 0 \quad j = 1, \dots, r$
- Dual feasibility:**  $\alpha_i^* \geq 0 \quad i = 1, \dots, m$
- Complementary slackness:**  $\alpha_i^* g_i(\mathbf{u}^*) = 0 \quad i = 1, \dots, m$
- Stationarity:** 
$$\nabla_u f(\mathbf{u}^*) + \sum_{i=1}^m \alpha_i^* \nabla_u g_i(\mathbf{u}^*) + \sum_{j=1}^r \beta_j^* \nabla_u h_j(\mathbf{u}^*) = 0.$$

**Karush-Kuhn-Tucker (KKT) conditions**

- For **convex optimization problems**, any  $(\mathbf{u}, \alpha, \beta)$  that verify the KKT conditions are **optimal**.

# Karush-Kuhn-Tucker conditions

- For **convex optimization problems**, any  $(\mathbf{u}, \boldsymbol{\alpha}, \boldsymbol{\beta})$  that verify the KKT conditions are **optimal**.

Consider  $\bar{\mathbf{u}}, \bar{\boldsymbol{\alpha}}, \bar{\boldsymbol{\beta}}$  that verify the KKT conditions.

*Primal feasibility* implies that  $\bar{\mathbf{u}}$  is feasible.

$L(\mathbf{u}, \bar{\boldsymbol{\alpha}}, \bar{\boldsymbol{\beta}})$  is convex in  $\mathbf{u}$ . Indeed  $L(\mathbf{u}, \bar{\boldsymbol{\alpha}}, \bar{\boldsymbol{\beta}}) = f(\mathbf{u}) + \sum_{i=1}^m \bar{\alpha}_i g_i(\mathbf{u}) + \sum_{j=1}^r \bar{\beta}_j h_j(\mathbf{u})$ , and  $f$  and  $g_i$  are convex,  $h_j$  is affine, and *dual feasibility* implies  $\bar{\alpha}_i \geq 0$ .

Because  $L(\mathbf{u}, \bar{\boldsymbol{\alpha}}, \bar{\boldsymbol{\beta}})$  is convex in  $\mathbf{u}$ , *stationarity* implies that  $\bar{\mathbf{u}}$  minimizes  $L(\mathbf{u}, \bar{\boldsymbol{\alpha}}, \bar{\boldsymbol{\beta}})$ , hence  $Q(\bar{\boldsymbol{\alpha}}, \bar{\boldsymbol{\beta}}) = L(\bar{\mathbf{u}}, \bar{\boldsymbol{\alpha}}, \bar{\boldsymbol{\beta}})$ . Because of *complementary slackness*,  $\sum_{i=1}^m \bar{\alpha}_i g_i(\bar{\mathbf{u}}) = 0$ , and because of *primal feasibility*,  $\sum_{j=1}^r \bar{\beta}_j h_j(\bar{\mathbf{u}}) = 0$ . Therefore,  $Q(\bar{\boldsymbol{\alpha}}, \bar{\boldsymbol{\beta}}) = f(\bar{\mathbf{u}})$ .

Let  $p^*$  be the optimal value of the primal, and  $d^*$  the optimal value of the dual. By definition of the optimal,  $Q(\bar{\boldsymbol{\alpha}}, \bar{\boldsymbol{\beta}}) \leq d^*$ . By weak duality,  $d^* \leq p^*$ . Therefore  $f(\bar{\mathbf{u}}) = Q(\bar{\boldsymbol{\alpha}}, \bar{\boldsymbol{\beta}}) \leq d^* \leq p^*$ . Because  $p^*$  is the optimal value of the primal, this implies  $f(\bar{\mathbf{u}}) = p^*$  and hence all the above inequalities are equalities:  $f(\bar{\mathbf{u}}) = p^* = d^* = Q(\bar{\boldsymbol{\alpha}}, \bar{\boldsymbol{\beta}})$ . Hence  $\bar{\mathbf{u}}, \bar{\boldsymbol{\alpha}}, \bar{\boldsymbol{\beta}}$  are optimal.

# Geometric interpretation

$$\min_{\boldsymbol{u} \in D} f(\boldsymbol{u}) \quad \text{subject to } g(\boldsymbol{u}) \leq 0$$

# Geometric interpretation

$$\min_{\mathbf{u} \in D} f(\mathbf{u}) \quad \text{subject to } g(\mathbf{u}) \leq 0$$

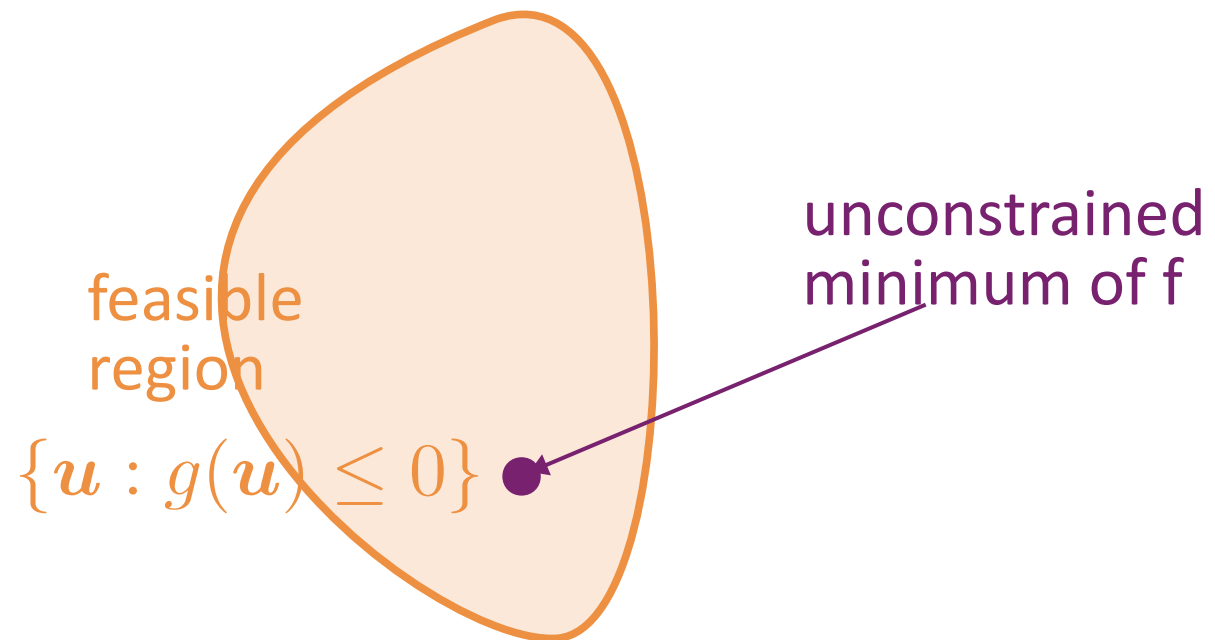
unconstrained  
minimum of  $f$



# Geometric interpretation

$$\min_{\mathbf{u} \in D} f(\mathbf{u}) \quad \text{subject to } g(\mathbf{u}) \leq 0$$

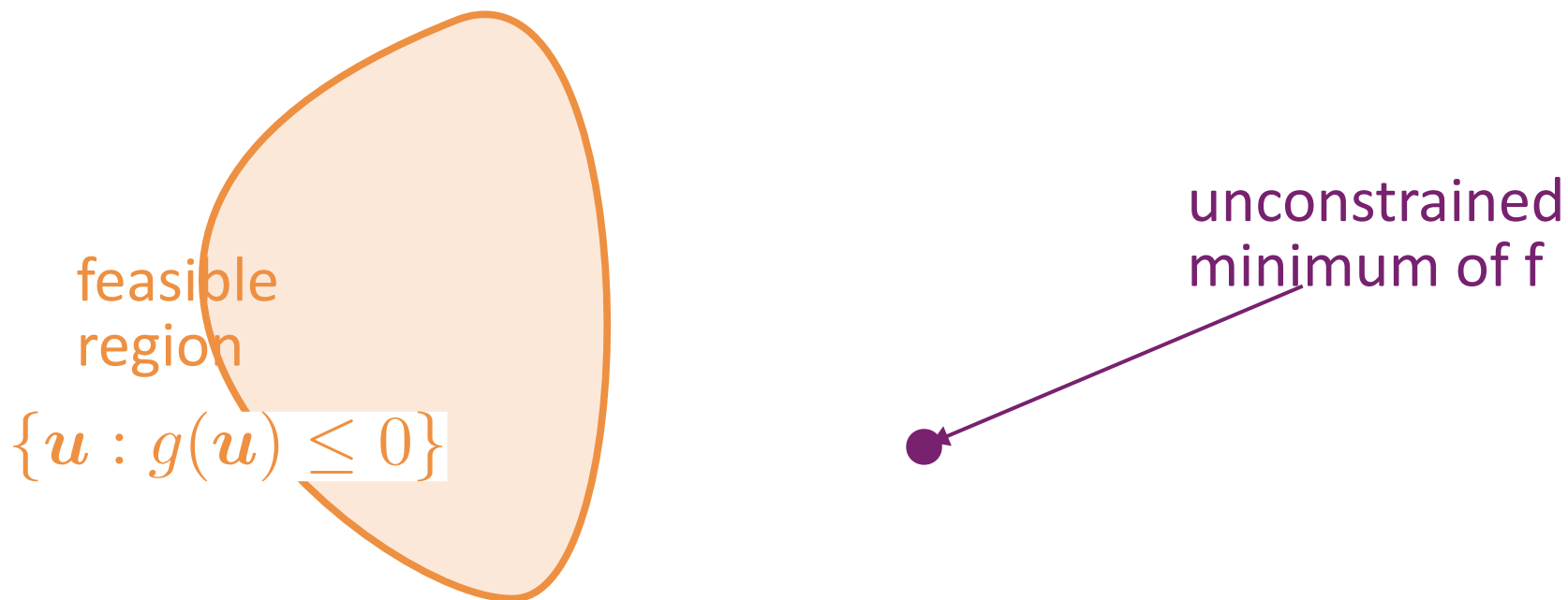
- **Case 1:** the unconstrained minimum lies in the feasible region



# Geometric interpretation

$$\min_{\mathbf{u} \in D} f(\mathbf{u}) \quad \text{subject to } g(\mathbf{u}) \leq 0$$

- **Case 1:** the unconstrained minimum lies in the feasible region
- **Case 2:** it does not.

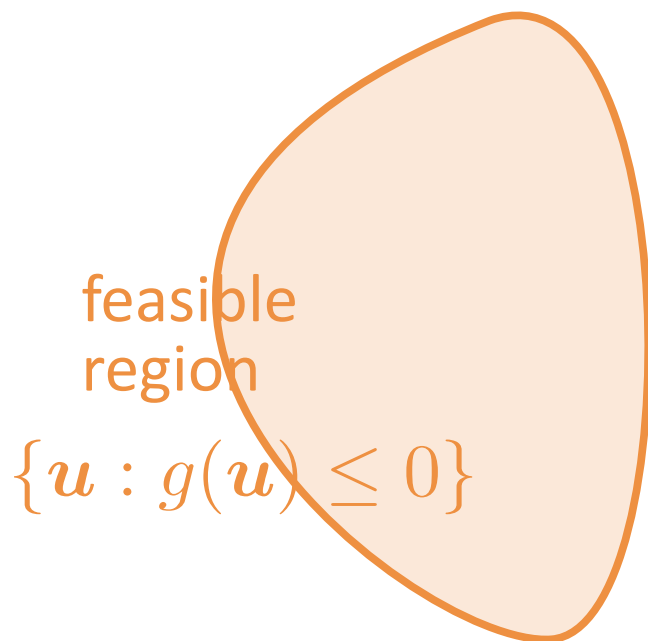


# Geometric interpretation

$$\min_{\mathbf{u} \in D} f(\mathbf{u}) \quad \text{subject to } g(\mathbf{u}) \leq 0$$

- **Case 1:** the unconstrained minimum lies in the feasible region
- **Case 2:** it does not.

SOLUTION 



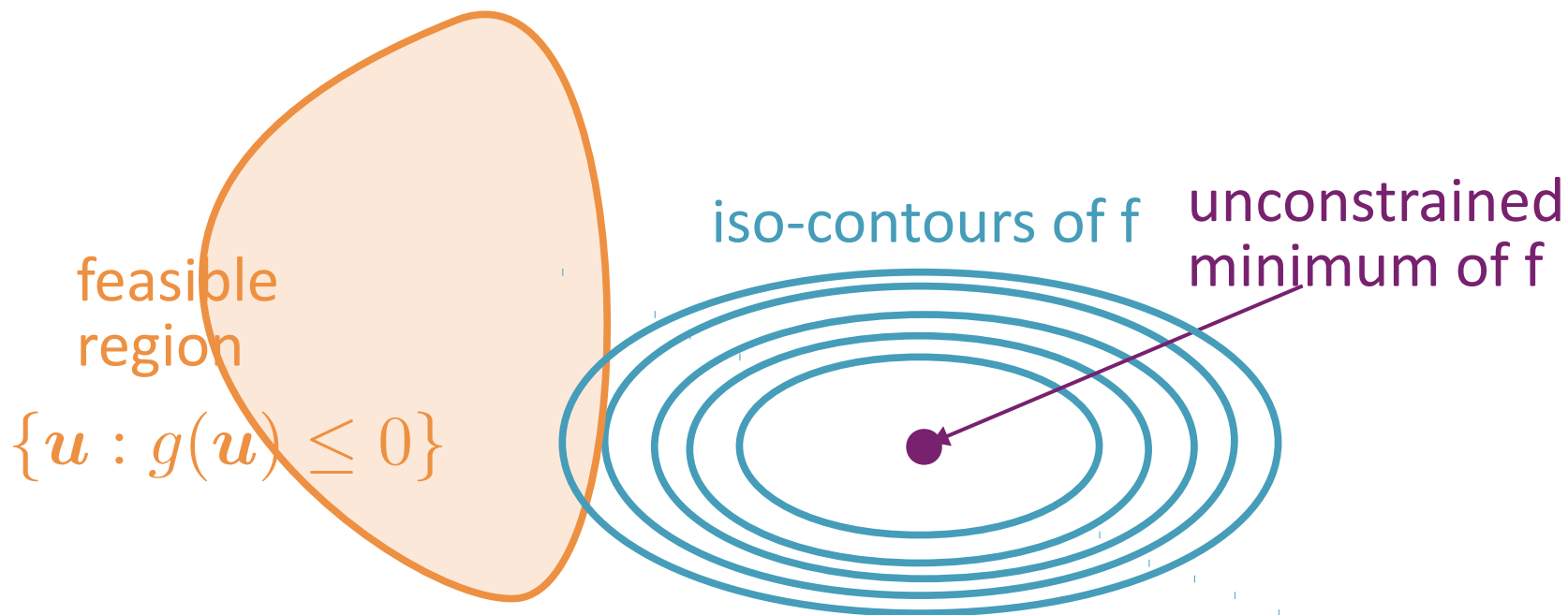


# Geometric interpretation

$$\min_{\mathbf{u} \in D} f(\mathbf{u}) \quad \text{subject to } g(\mathbf{u}) \leq 0$$

- **Case 1:** the unconstrained minimum lies in the feasible region
- **Case 2:** it does not.

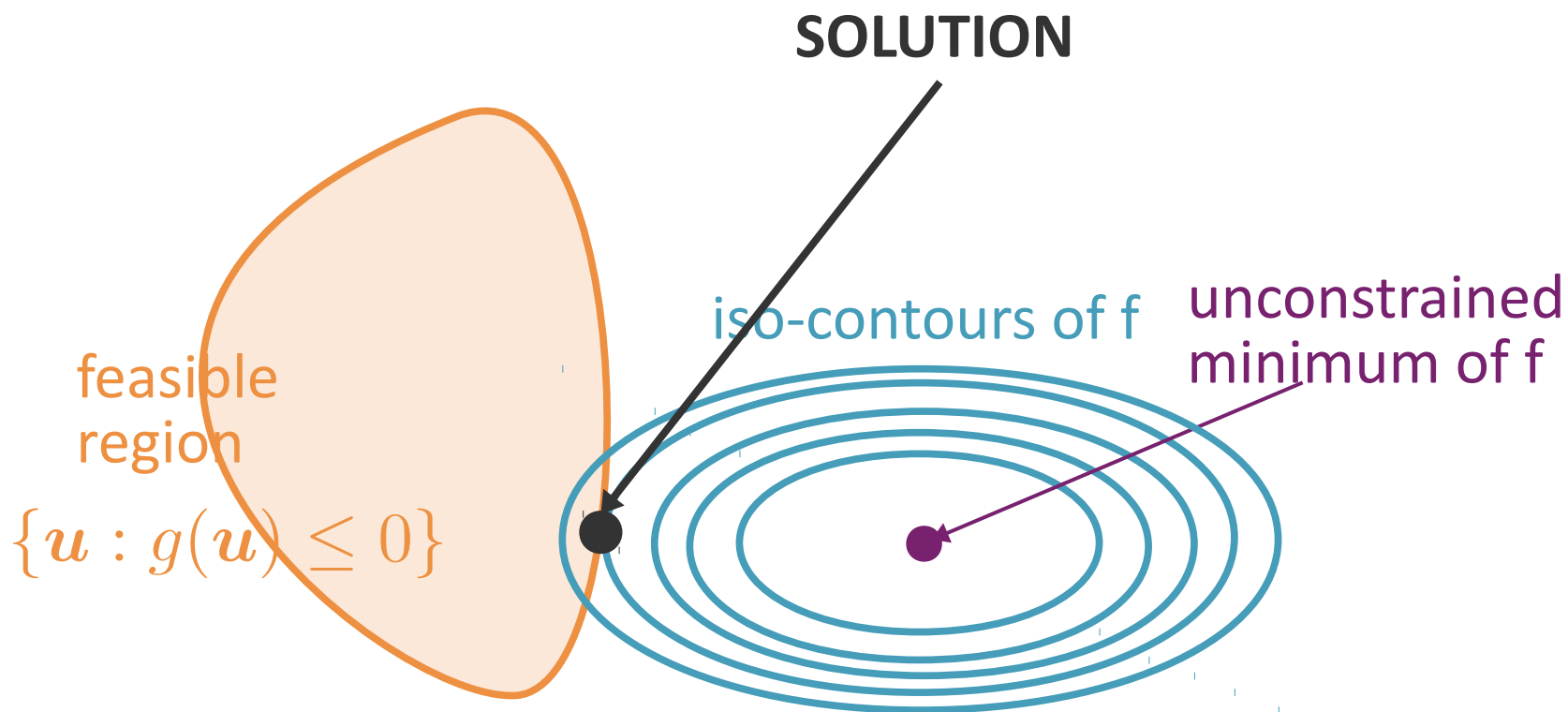
SOLUTION 



# Geometric interpretation

$$\min_{\mathbf{u} \in D} f(\mathbf{u}) \quad \text{subject to } g(\mathbf{u}) \leq 0$$

- **Case 1:** the unconstrained minimum lies in the feasible region
- **Case 2:** it does not.

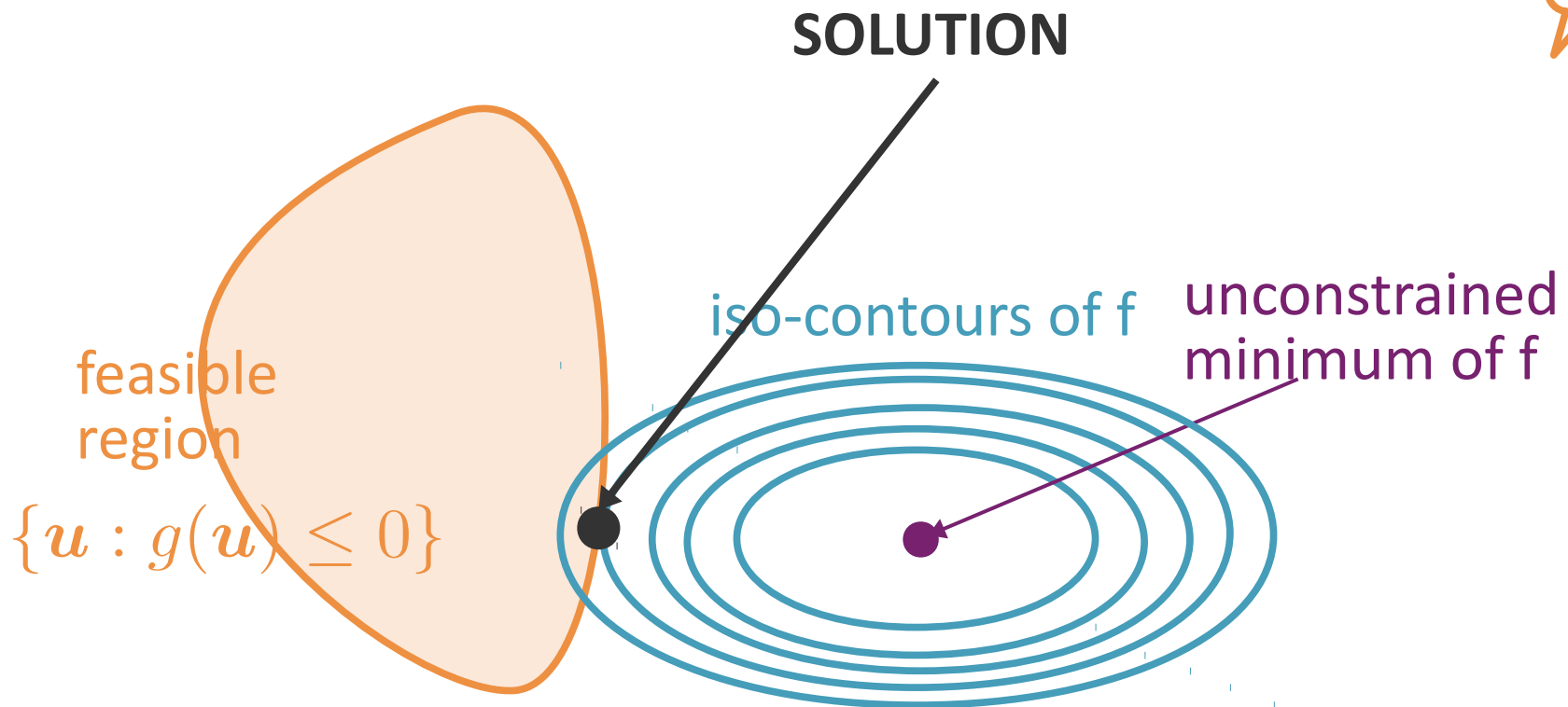


# Geometric interpretation

$$\min_{\mathbf{u} \in D} f(\mathbf{u}) \quad \text{subject to } g(\mathbf{u}) \leq 0$$

- **Case 1:** the unconstrained minimum lies in the feasible region
- **Case 2:** it does not.

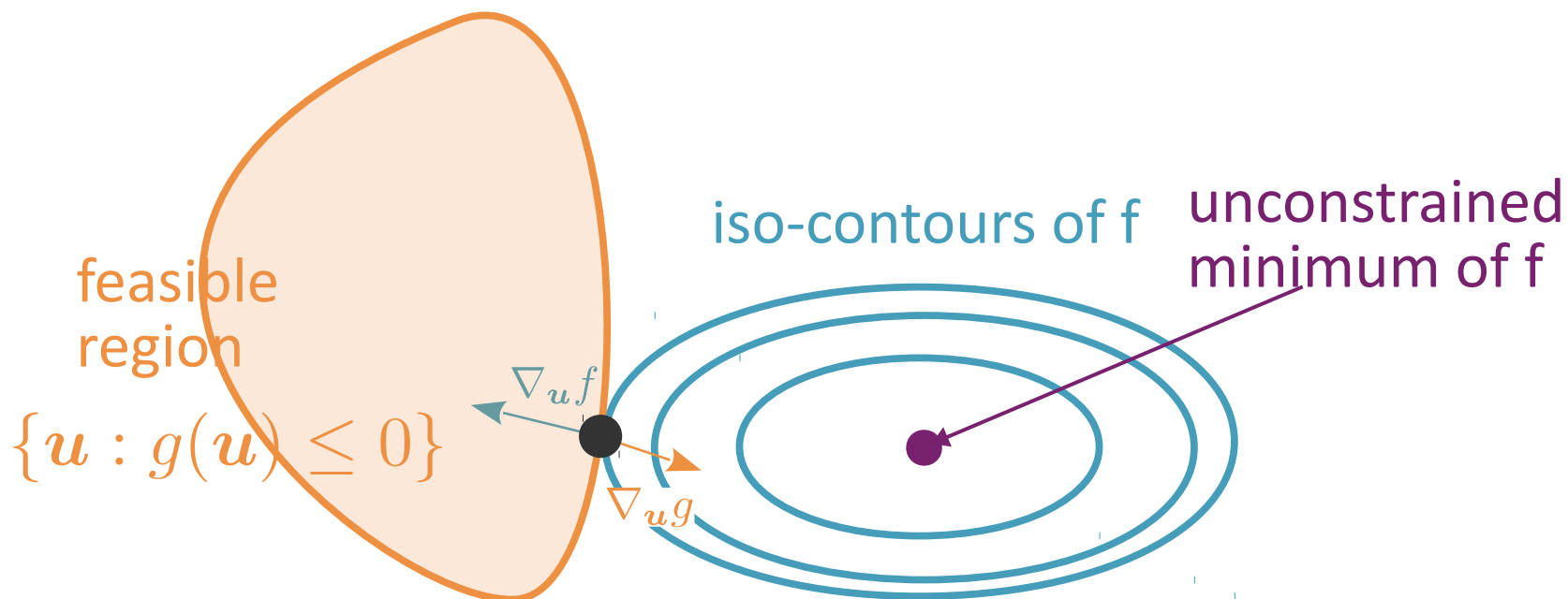
What can you say about the gradients of  $f$  and  $g$ ?



# Geometric interpretation

$$\min_{\mathbf{u} \in D} f(\mathbf{u}) \quad \text{subject to } g(\mathbf{u}) \leq 0$$

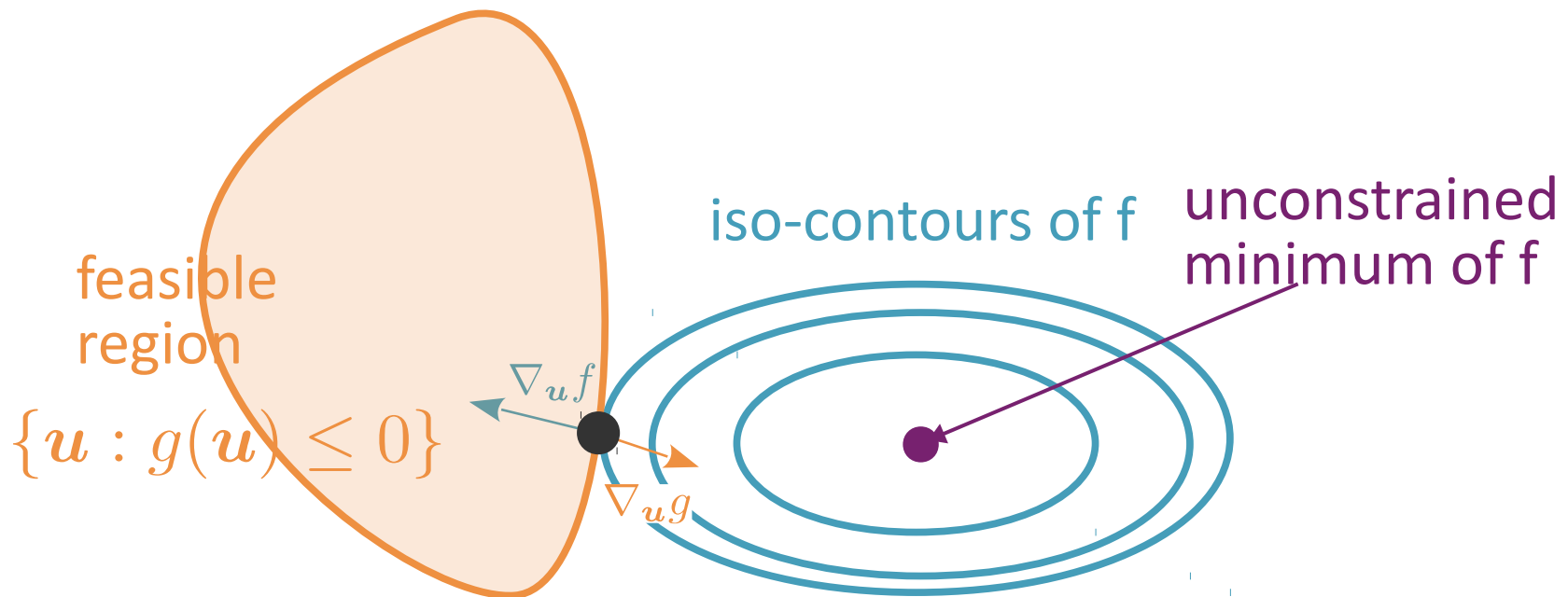
- **Case 1:** the unconstrained minimum lies in the feasible region
- **Case 2:** it does not. The solution lies where the iso-contours of  $f$  meet the feasible region.



# Geometric interpretation

$$\min_{\mathbf{u} \in D} f(\mathbf{u}) \quad \text{subject to } g(\mathbf{u}) \leq 0$$

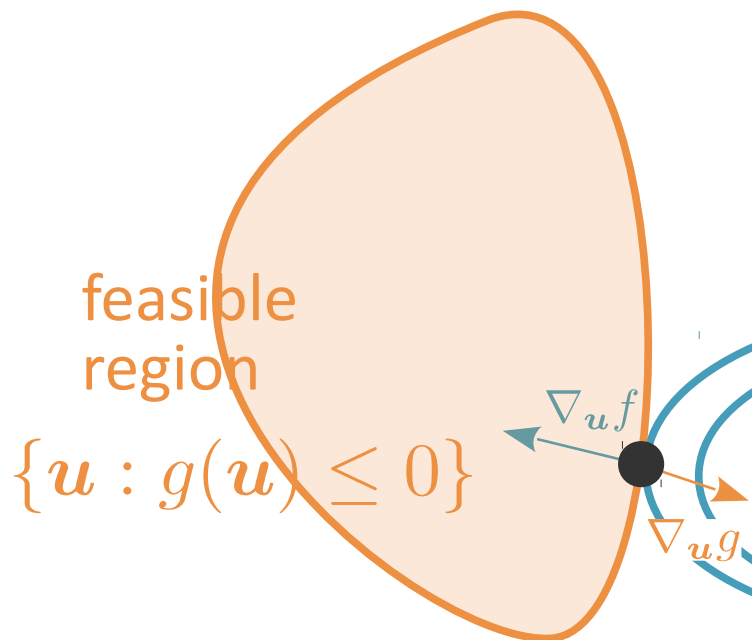
- **Case 1:** the unconstrained minimum lies in the feasible region  $\nabla f(\mathbf{u}^*) = 0$  and  $g(\mathbf{u}^*) \leq 0$
- **Case 2:** it does not. The solution lies where the iso-contours of  $f$  meet the feasible region.



# Geometric interpretation

$$\min_{\mathbf{u} \in D} f(\mathbf{u}) \quad \text{subject to } g(\mathbf{u}) \leq 0$$

- **Case 1:** the unconstrained minimum lies in the feasible region  $\nabla f(\mathbf{u}^*) = 0$  and  $g(\mathbf{u}^*) \leq 0$
- **Case 2:** it does not. The solution lies where the iso-contours of  $f$  meet the feasible region.



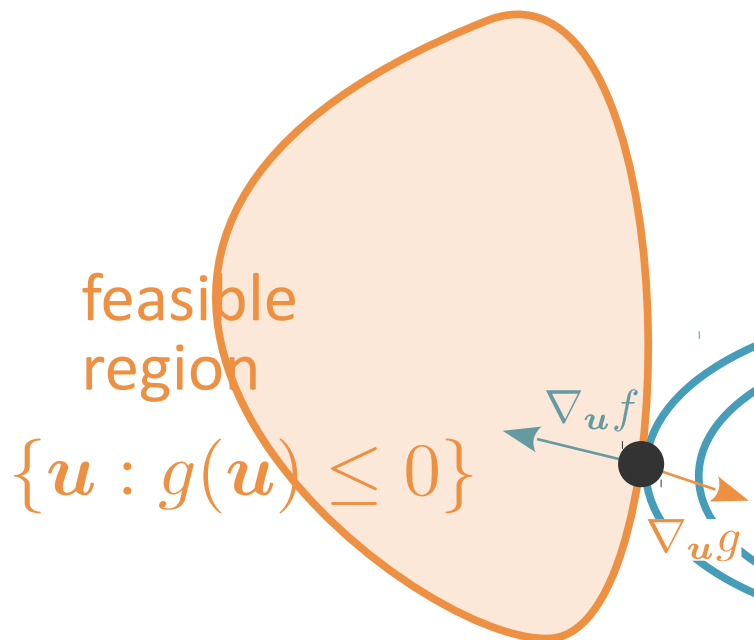
- The gradients of  $f$  and  $g$  are parallel, of opposite directions

$$\nabla f(\mathbf{u}^*) = -\alpha \nabla g(\mathbf{u}^*) \quad \alpha \geq 0$$

# Geometric interpretation

$$\min_{\mathbf{u} \in D} f(\mathbf{u}) \quad \text{subject to } g(\mathbf{u}) \leq 0$$

- **Case 1:** the unconstrained minimum lies in the feasible region  $\nabla f(\mathbf{u}^*) = 0$  and  $g(\mathbf{u}^*) \leq 0$
- **Case 2:** it does not. The solution lies where the iso-contours of  $f$  meet the feasible region.



- The gradients of  $f$  and  $g$  are parallel, of opposite directions  
$$\nabla f(\mathbf{u}^*) = -\alpha \nabla g(\mathbf{u}^*) \quad \alpha \geq 0$$
- The solution lies at the border of the feasible space  
$$g(\mathbf{u}^*) = 0$$

# Geometric interpretation

$$\min_{\mathbf{u} \in D} f(\mathbf{u}) \quad \text{subject to } g(\mathbf{u}) \leq 0$$

- **Case 1:**  $\nabla f(\mathbf{u}^*) = 0$  and  $g(\mathbf{u}^*) \leq 0$
- **Case 2:**  $\nabla f(\mathbf{u}^*) = -\alpha \nabla g(\mathbf{u}^*) \quad \alpha \geq 0$   
and  $g(\mathbf{u}^*) = 0$
- Can be summarized as:  $\nabla f(\mathbf{u}^*) + \alpha \nabla g(\mathbf{u}^*) = 0$   
and  $\alpha g(\mathbf{u}^*) = 0$ 
  - Either  $\alpha = 0$  and  $g(\mathbf{u}^*) \leq 0$  (case 1)
  - or  $g(\mathbf{u}^*) = 0$  (case 2).



# Geometric interpretation

$$\min_{\mathbf{u} \in D} f(\mathbf{u}) \quad \text{subject to } g(\mathbf{u}) \leq 0$$

- **Case 1:**  $\nabla f(\mathbf{u}^*) = 0$  and  $g(\mathbf{u}^*) \leq 0$
- **Case 2:**  $\nabla f(\mathbf{u}^*) = -\alpha \nabla g(\mathbf{u}^*) \quad \alpha \geq 0$   
and  $g(\mathbf{u}^*) = 0$

- Can be summarized as:  $\nabla f(\mathbf{u}^*) + \alpha \nabla g(\mathbf{u}^*) = 0$



and  $\alpha g(\mathbf{u}^*) = 0$

- Either  $\alpha = 0$  and  $g(\mathbf{u}^*) \leq 0$  (case 1)
- or  $g(\mathbf{u}^*) = 0$  (case 2).


# Geometric interpretation

$$\min_{\mathbf{u} \in D} f(\mathbf{u}) \quad \text{subject to } g(\mathbf{u}) \leq 0$$

- **Case 1:**  $\nabla f(\mathbf{u}^*) = 0$  and  $g(\mathbf{u}^*) \leq 0$
- **Case 2:**  $\nabla f(\mathbf{u}^*) = -\alpha \nabla g(\mathbf{u}^*) \quad \alpha \geq 0$   
and  $g(\mathbf{u}^*) = 0$  stationarity
- Can be summarized as:  $\nabla f(\mathbf{u}^*) + \alpha \nabla g(\mathbf{u}^*) = 0$   
and  $\alpha g(\mathbf{u}^*) = 0$ 
  - Either  $\alpha = 0$  and  $g(\mathbf{u}^*) \leq 0$  (case 1)
  - or  $g(\mathbf{u}^*) = 0$  (case 2).

# Geometric interpretation

$$\min_{\mathbf{u} \in D} f(\mathbf{u}) \quad \text{subject to } g(\mathbf{u}) \leq 0$$

- **Case 1:**  $\nabla f(\mathbf{u}^*) = 0$  and  $g(\mathbf{u}^*) \leq 0$
- **Case 2:**  $\nabla f(\mathbf{u}^*) = -\alpha \nabla g(\mathbf{u}^*) \quad \alpha \geq 0$   
and  $g(\mathbf{u}^*) = 0$  stationarity
- Can be summarized as:  $\nabla f(\mathbf{u}^*) + \alpha \nabla g(\mathbf{u}^*) = 0$   
and  $\alpha g(\mathbf{u}^*) = 0$  
  - Either  $\alpha = 0$  and  $g(\mathbf{u}^*) \leq 0$  (case 1)
  - or  $g(\mathbf{u}^*) = 0$  (case 2).

# Geometric interpretation

$$\min_{\mathbf{u} \in D} f(\mathbf{u}) \quad \text{subject to } g(\mathbf{u}) \leq 0$$


- **Case 1:**  $\nabla f(\mathbf{u}^*) = 0$  and  $g(\mathbf{u}^*) \leq 0$
- **Case 2:**  $\nabla f(\mathbf{u}^*) = -\alpha \nabla g(\mathbf{u}^*) \quad \alpha \geq 0$   
and  $g(\mathbf{u}^*) = 0$  stationarity
- Can be summarized as:  $\nabla f(\mathbf{u}^*) + \alpha \nabla g(\mathbf{u}^*) = 0$   
and  $\alpha g(\mathbf{u}^*) = 0$  complementary slackness
  - Either  $\alpha = 0$  and  $g(\mathbf{u}^*) \leq 0$  (case 1)
  - or  $g(\mathbf{u}^*) = 0$  (case 2).

# Quadratic Programs

- Special case of convex optimization problems where
  - $f$  is **quadratic**

$$f : u \mapsto \frac{1}{2} u^\top Q u + b^\top u + c \quad Q \succeq 0 \quad b \in \mathbb{R}^n \quad c \in \mathbb{R}$$

- $g_i$  and  $h_j$  are **affine**  $u \mapsto c^\top u + d \quad c \in \mathbb{R}^n \quad d \in \mathbb{R}$

- The feasible set is 

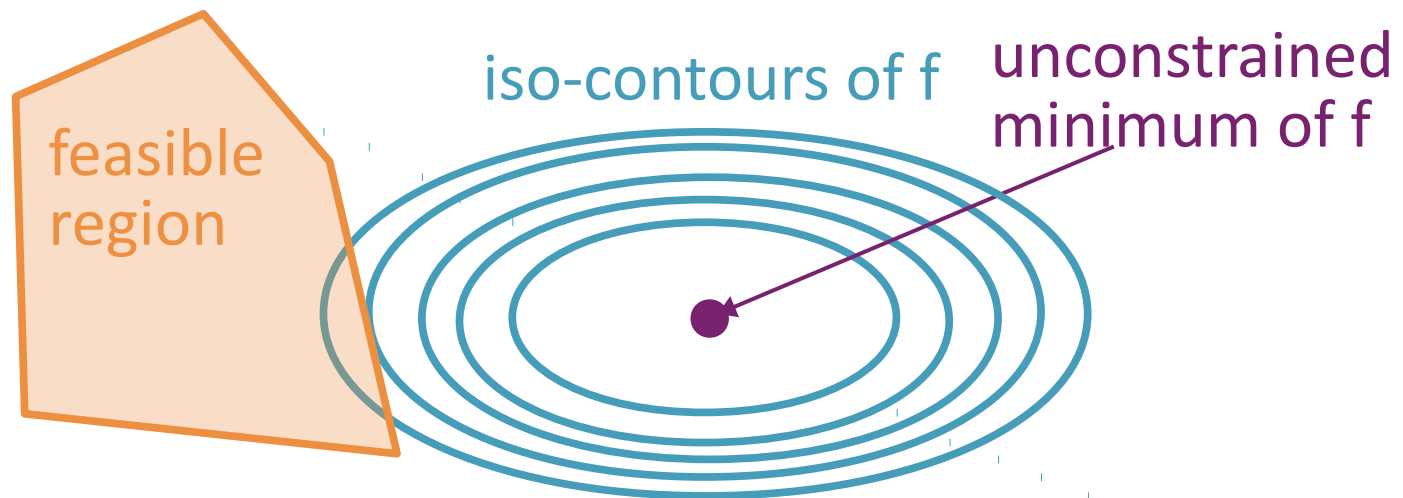
# Quadratic Programs

- Special case of convex optimization problems where
  - $f$  is **quadratic**

$$f : u \mapsto \frac{1}{2} u^\top Q u + b^\top u + c \quad Q \succeq 0 \quad b \in \mathbb{R}^n \quad c \in \mathbb{R}$$

- $g_i$  and  $h_j$  are **affine**  $u \mapsto c^\top u + d \quad c \in \mathbb{R}^n \quad d \in \mathbb{R}$

- The feasible set is a **polyhedron**.



# Quadratic Programs

- Many methods can be used to solve QPs, for example
  - Interior point methods
  - Active set methods
- Many solvers implement them
  - CPLEX
  - CVXOPT
  - CGAL and more.

# Slack variables

- Replace the inequality constraints

$$g_i(\mathbf{u}) \leq 0 \Leftrightarrow \exists s_i \geq 0 : g_i(\mathbf{u}) + s_i = 0$$

- $s_i =$  **slack variable**.

$$\min_{\mathbf{u} \in D} f(\mathbf{u})$$

subject to  $g_i(\mathbf{u}) + s_i = 0, i = 1, \dots, m$

$$s_i \geq 0, i = 1, \dots, m$$

$$h_j(\mathbf{u}) = 0, j = 1, \dots, r$$



# Summary

- We often try to formulate **machine learning** problems as **convex optimization** problems

$$\begin{array}{ll} \min_{\mathbf{u} \in D} f(\mathbf{u}) \\ \text{subject to } g_i(\mathbf{u}) \leq 0, i = 1, \dots, m \\ h_j(\mathbf{u}) = 0, j = 1, \dots, r \end{array}$$

- If **f differentiable**:  $\nabla f(\mathbf{u}) = 0 \Leftrightarrow \mathbf{u}$  minimizes  $f$
- **Unconstrained** convex optimization problems can be solved by **gradient descent**.  
  
**Flavors:** Backtracking line search, Newton's methods, BFGS, stochastic gradient descent.
- **Constrained** convex optimization problems can be solved in **dual space** via the **Lagrangian**.

# References

- *Convex optimization*. S. Boyd and L. Vandenberghe.  
<https://web.stanford.edu/~boyd/cvxbook/>
  - **Convex sets**: Chapter 2.1
  - **Convex functions**: Chapter 3.1.1 – 3.1.5, 3.2
  - **Convex optimization problems**: 4.1.1 – 4.1.2 + 4.2.2
  - **Unconstrained minimization**: 9.1.1 + 9.2 – 9.3 (gradient descent) + 9.5 (Newton)
  - **QP**: 4.4.1 + 5.1 (Lagrange) + 5.2 (Duality) + 5.3.2 (Slater) + 5.5.3 (KKT)
  - **Slack variables**: 4.1.3
  - Also see the Bibliography section at the end of each chapter.
- To go further
  - *Numerical Optimization*. J. Bonnans, J. Gilbert, C. Lemaréchal, C. Sagastizábal. **Quasi-Newton methods**: 4.3 – 4.4.
  - *Stochastic gradient descent tricks*. L. Bottou (2012).  
<http://leon.bottou.org/publications/pdf/tricks-2012.pdf>
  - *Coordinate Descent Algorithms*. S. Wright (2015).  
<https://arxiv.org/abs/1502.04759>

# Homework

- **By Monday (Oct 2<sup>nd</sup>)**

Visit <http://tinyurl.com/ma2823-2017>

Download and read **the complete syllabus**.

**Set up your computer for the labs.**

- **By Friday (Oct 6<sup>th</sup>)**

Download, solve and turn in **HW 1**.

**See you on Monday, 8:30am in Amphi sc.046!**