# Preliminaries on Optimization

Adaptive and Cooperative Algorithms (ECE 457A)

ECE, MME, and MSCI Departments,
University of Waterloo, ON, Canada
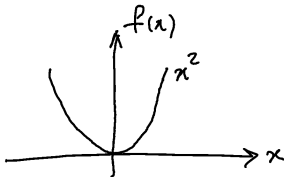
Course Instructor: Benyamin Ghojogh
Fall 2023

**What is Optimization?**

# Optimization problem

- Consider a function representing some cost. We call it **cost function** or **objective function**.
- We want to **minimize** or **maximize** this objective function.
- Examples:
  - Example for **minimization**: the cost function can be the error of some airplane structure from the perfect aerodynamic structure.
  - Example for **maximization**: the objective function can be the profit of the company.
  - **All life** is optimization!
  - **All machine learning** in artificial intelligence is optimization!
- The variables of the objective function are called the **objective variables** or **decision variables** or **optimization variables**.
- Example:

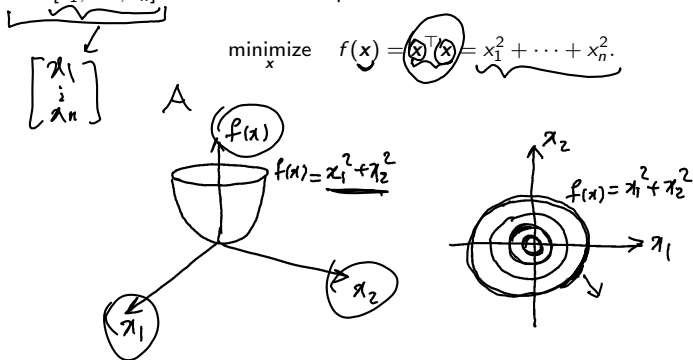$$\underset{x}{\text{minimize}} \quad f(x) = x^2.$$

# Univariate and multivariate optimization problems

- The optimization problem can be **univariate**, meaning that the optimization problem has only one scalar variable. Example:

$$\underset{x}{\text{minimize}} \quad f(x) = x^2.$$

- The optimization problem can be **multivariate**, meaning that the optimization problem has several scalar variables $\{x_1, \ldots, x_n\}$. These variables can be combined into a vector $\boldsymbol{x} = [x_1, \ldots, x_n]^\top$ or matrix. Example:

$$\underset{x}{\text{minimize}} \quad f(\boldsymbol{x}) = \boldsymbol{x}^\top \boldsymbol{x} = x_1^2 + \cdots + x_n^2.$$
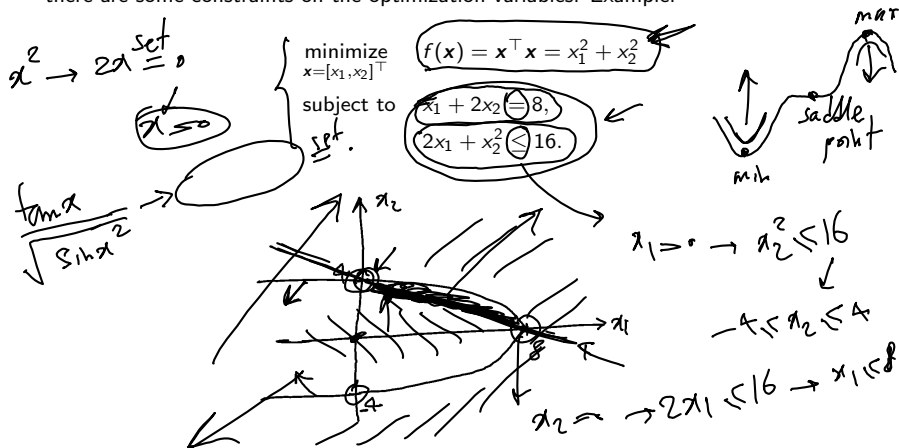
# Unconstrained and constrained problems

- The optimization problem can be **unconstrained**, meaning that we simply optimize a function only. Example:
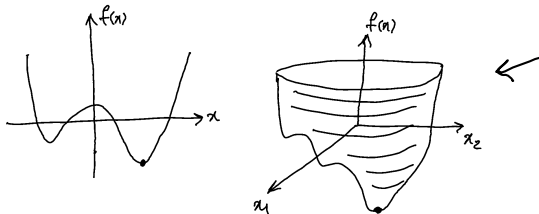
$$\underset{x}{\text{minimize}} \quad f(x) = x^\top x.$$

- The optimization problem can be **constrained**, meaning that we optimize a function while there are some constraints on the optimization variables. Example:

$$\underset{x=[x_1,x_2]^\top}{\text{minimize}} \quad f(x) = x^\top x = x_1^2 + x_2^2$$

$$\text{subject to} \quad x_1 + 2x_2 = 8,$$

$$2x_1 + x_2^2 \leq 16.$$

$$x^2 \rightarrow 2x \overset{set}{=} 0$$

$$x = 0$$

$$\frac{\tan x}{\sqrt{\sin x^2}}$$

$$x_1 > 0 \rightarrow x_2^2 \leq 16$$

$$-4 \leq x_2 \leq 4$$

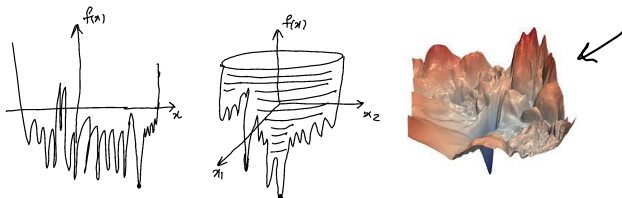$$x_2 = 0 \rightarrow 2x_1 \leq 16 \rightarrow x_1 \leq 8$$

# Optimization versus search

- If the objective problem is **simple enough**, we can solve it using **classic optimization** methods. We will learn important classic methods.
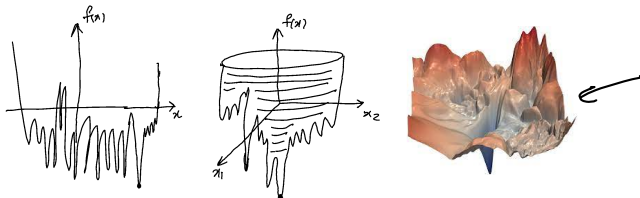


- If the objective function is **complicated** or if we have **too many constraints**, we can use **search** for finding a good solution.
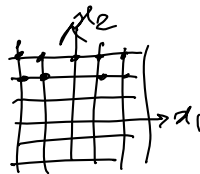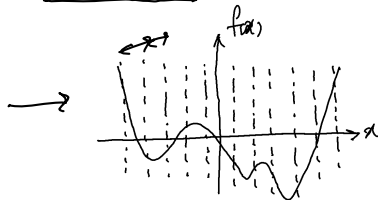
# When to use search-based (metaheuristic) optimization

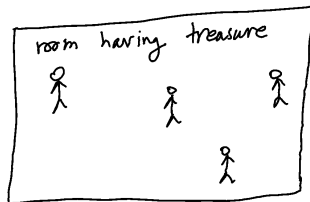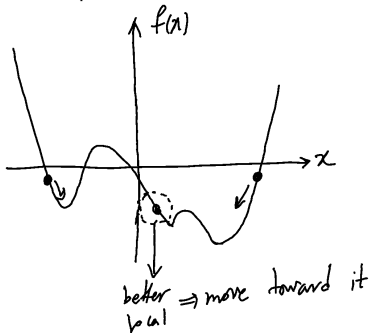- When we have a **complicated** (highly con-convex) optimization landscape:



- or when the **gradient** of function is **hard to compute**,
- or when the function is not known but it works as a **black-box**, i.e., it outputs a value for each input fed to it.
- In these cases, we need:
  - either **non-convex optimization**,
  - or **search-based** optimization (**metaheuristic** optimization).

# Search for optimization

- We can do **grid search** or **brute-force search**.



- Or we can **search wisely** by **metaheuristic optimization**. We will learn several important metaheuristic optimization methods.

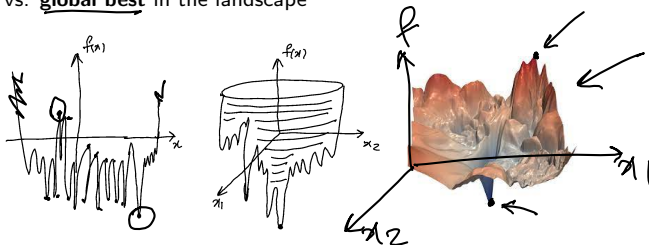**Preliminaries on Metaheuristic Optimization**

# Heuristic and Metaheuristic Methods

- When the problem is complicated, a **heuristic** method approximates its solution.
- It gives a **good enough guess** of the solution to the problem, but that you may not really know **how good** it is.
- **Heuristics** are often **problem-dependent**, i.e., you define a heuristic for a given problem.
- **Metaheuristics** are **problem-independent** techniques that can be applied to a broad range of problems.
- **Metaheuristic optimization** methods can solve various complicated optimization problems using wise search.
- **Metaheuristic** methods are considered as a family of methods in **soft computing**.

# Exploration vs. Exploitation

Some things need to be defined:

- **Fitness** vs. **cost**: We usually minimize the cost function but maximize the fitness function. Fitness and cost can be converted to each other by changing maximization to minimization or vice versa.

- **Optimization landscape**: the optimization cost/fitness function

- **Local best** vs. **global best** in the landscape



- **Exploitation**: local search around the solution because the global optimum might be close to the current solution.

- **Exploration**: search far away from the solution (explore the landscape) because the global optimum might be far away from the current solution. It helps not to get **stuck in local optimum**.

$$\sum_{i=1}^{n} \left( \hat{y}_i - y_i \right)^2$$

$W_1$
$W_2$
$W_3$

min

cost    loss

min
$\{ w_1, \dots w_{1000} \}$
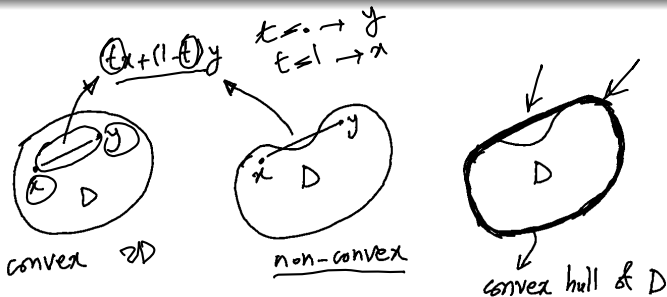
**Preliminaries on Sets and Norms**

# Convex set

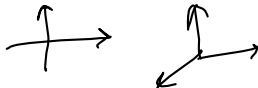## Definition (Convex set and convex hull)

A set $\mathcal{D}$ is a convex set if it completely contains the line segment between any two points in the set $\mathcal{D}$:

$$\forall x, y \in \mathcal{D}, 0 \leq t \leq 1 \implies tx + (1-t)y \in \mathcal{D}.$$

The convex hull of a (not necessarily convex) set $\mathcal{D}$ is the smallest convex set containing the set $\mathcal{D}$. If a set is convex, it is equal to its convex hull.

# Inner product

## Definition (Inner product of vectors)

Consider two vectors $\boldsymbol{x} = [x_1, \ldots, x_d]^\top \in \mathbb{R}^d$ and $\boldsymbol{y} = [y_1, \ldots, y_d]^\top \in \mathbb{R}^d$. Their **inner product**, also called **dot product**, is:

$$\langle \boldsymbol{x}, \boldsymbol{y} \rangle = \boldsymbol{x}^\top \boldsymbol{y} = \sum_{i=1}^{d} x_i\, y_i.$$

## Definition (Inner product of matrices)

We also have inner product between matrices $\boldsymbol{X}, \boldsymbol{Y} \in \mathbb{R}^{d_1 \times d_2}$. Let $\boldsymbol{X}_{ij}$ denote the $(i, j)$-th element of matrix $\boldsymbol{X}$. The inner product of $\boldsymbol{X}$ and $\boldsymbol{Y}$ is:

$$\langle \boldsymbol{X}, \boldsymbol{Y} \rangle = \mathbf{tr}(\boldsymbol{X}^\top \boldsymbol{Y}) = \sum_{i=1}^{d_1} \sum_{j=1}^{d_2} \boldsymbol{X}_{i,j}\, \boldsymbol{Y}_{i,j},$$
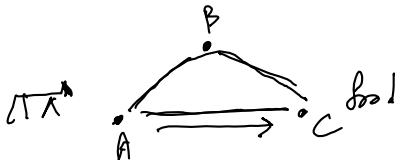
where $\mathbf{tr}(.)$ denotes the trace of matrix.

# Norm

## Definition (Norm)

A function $\|\cdot\| : \mathbb{R}^d \to \mathbb{R}$, $\|\cdot\| : x \mapsto \|x\|$ is a **norm** if it satisfies:

1. $\|x\| \geq 0, \forall x$
2. $\|ax\| = |a|\|x\|, \forall x$ and all scalars $a$
3. $\|x\| = 0$ if and only if $x = 0$
4. Triangle inequality: $\|x + y\| \leq \|x\| + \|y\|$.

## Important norms for vectors

Some important norms for a vector $\boldsymbol{x} = [x_1, \ldots, x_d]^\top$ are as follows.

- The $\ell_p$ **norm** is:
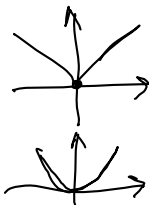
$$\|\boldsymbol{x}\|_p := \left(|x_1|^p + \cdots + |x_d|^p\right)^{1/p},$$

where $p \geq 1$ and $|.|$ denotes the absolute value.

- Two well-known $\ell_p$ norms are $\ell_1$ **norm** and $\ell_2$ **norm** (also called the **Euclidean norm**) with $p = 1$ and $p = 2$, respectively:

$$\|\boldsymbol{x}\|_1 := |x_1| + \cdots + |x_d| = \sum_{i=1}^{d} |x_i|,$$

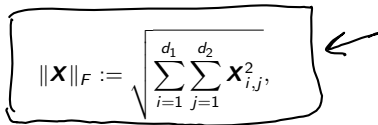$$\|\boldsymbol{x}\|_2 := \sqrt{x_1^2 + \cdots + x_d^2} = \sqrt{\sum_{i=1}^{d} x_i^2}$$

- The $\ell_\infty$ **norm**, also called the **infinity norm**, the **maximum norm**, or the **Chebyshev norm**, is:

$$\|\boldsymbol{x}\|_\infty := \max\{|x_1|, \ldots, |x_d|\}.$$

# Important norms for matrices

Some important norms for a matrix $\boldsymbol{X} \in \mathbb{R}^{d_1 \times d_2}$ are as follows.

- The formulation of the **<u>Frobenius norm</u>** for a matrix is similar to the formulation of $\ell_2$ norm for a vector:

$$\|\boldsymbol{X}\|_F := \sqrt{\sum_{i=1}^{d_1} \sum_{j=1}^{d_2} \boldsymbol{X}_{i,j}^2},$$

where $\boldsymbol{X}_{ij}$ denotes the $(i,j)$-th element of $\boldsymbol{X}$.
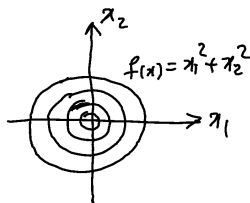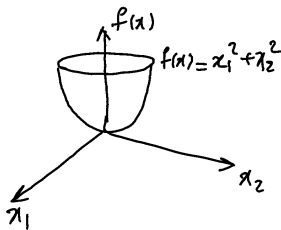
# Quadratic forms using norms

For $\boldsymbol{x} \in \mathbb{R}^d$ and $\boldsymbol{X} \in \mathbb{R}^{d_1 \times d_2}$, we have:

$$\star \quad \underbrace{\|\boldsymbol{x}\|_2^2 = \boldsymbol{x}^\top \boldsymbol{x}}_{} = \langle \boldsymbol{x}, \boldsymbol{x} \rangle = \sum_{i=1}^{d} x_i^2,$$

$$\star \quad \|\boldsymbol{X}\|_F^2 = \underbrace{\text{tr}(\boldsymbol{X}^\top \boldsymbol{X})}_{} = \langle \boldsymbol{X}, \boldsymbol{X} \rangle = \sum_{i=1}^{d_1} \sum_{j=1}^{d_2} \boldsymbol{X}_{i,j}^2,$$

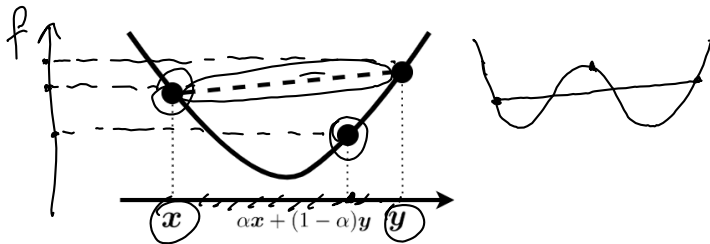which are convex and in quadratic forms.

**Preliminaries on Functions**

# Convex function

## Definition (Convex function)

A function $f(.)$ with domain $\mathcal{D}$ is convex if:

$$f(\alpha\boldsymbol{x} + (1-\alpha)\boldsymbol{y}) \leq \alpha f(\boldsymbol{x}) + (1-\alpha)f(\boldsymbol{y}), \quad \forall \boldsymbol{x}, \boldsymbol{y} \in \mathcal{D}, \tag{1}$$

where $\alpha \in [0, 1]$.



If $\leq$ is changed to $\geq$ in Eq. (1), the function is *concave*.

# Convex function

## Definition (Convex function)

If the function $f(.)$ is differentiable, it is convex if:

$$f(\boldsymbol{x}) \geq f(\boldsymbol{y}) + \nabla f(\boldsymbol{y})^\top (\boldsymbol{x} - \boldsymbol{y}), \quad \forall \boldsymbol{x}, \boldsymbol{y} \in \mathcal{D}. \tag{2}$$
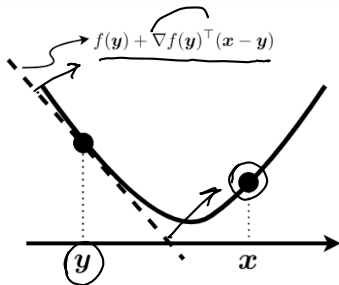


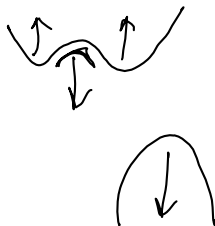If $\geq$ is changed to $\leq$ in Eq. (2), the function is *concave*.

# Convex function

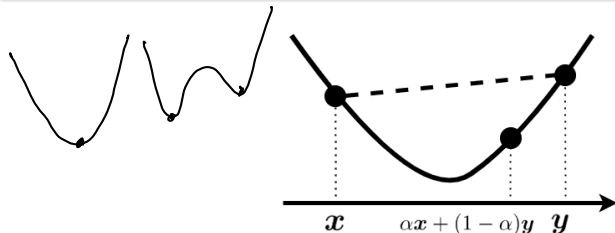## Definition (Convex function)

If the function $f(.)$ is twice differentiable, it is convex if its second-order derivative is positive semi-definite:

$$\nabla^2 f(\boldsymbol{x}) \succeq \boldsymbol{0}, \quad \forall \boldsymbol{x} \in \mathcal{D}. \tag{3}$$
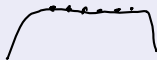


If $\succeq$ is changed to $\preceq$ in Eq. (3), the function is *concave*.

**Preliminaries on Optimization**

# Min, max, sup, inf

## Definition (Minimum, maximum, infimum, and supremum)

A **minimum** and **maximum** of a function $f : \mathbb{R}^d \to \mathbb{R}$, $f : \boldsymbol{x} \mapsto f(\boldsymbol{x})$, with domain $\mathcal{D}$, are defined as:

$$\min_{\boldsymbol{x}} f(\boldsymbol{x}) \leq f(\boldsymbol{y}), \ \forall \boldsymbol{y} \in \mathcal{D},$$
$$\max_{\boldsymbol{x}} f(\boldsymbol{x}) \geq f(\boldsymbol{y}), \ \forall \boldsymbol{y} \in \mathcal{D},$$

respectively.

The minimum and maximum of a function belong to the range of function.

# Min, max, sup, inf

## Definition (Infimum and supremum)

**Infimum** and **supremum** are the lower-bound and upper-bound of function, respectively:

$$\inf_x f(\boldsymbol{x}) := \max\{z \in \mathbb{R} \mid z \leq f(\boldsymbol{x}), \forall \boldsymbol{x} \in \mathcal{D}\},$$
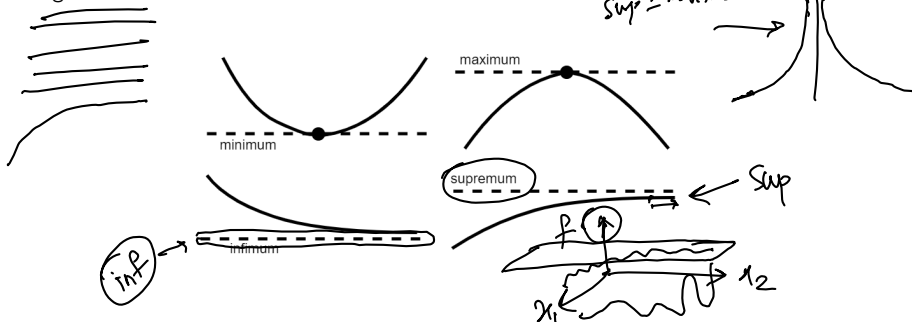$$\sup_x f(\boldsymbol{x}) := \min\{z \in \mathbb{R} \mid z \geq f(\boldsymbol{x}), \forall \boldsymbol{x} \in \mathcal{D}\}.$$

*(handwritten annotations: "lower bound", "upper bound")*

Depending on the function, the infimum and supremum of a function may or may not belong to the range of function.

*(handwritten annotations: "Sup = max = ∞", "Sup", "f", "$x_1$", "$x_2$", "inf")*



maximum

minimum

supremum

infimum

# Local and global minimizers

∃ : there exists

∀ : for all

## Definition (Local minimizer)

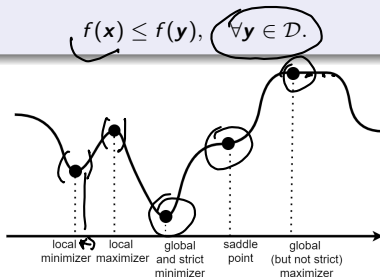A point $x \in \mathcal{D}$ is a **local minimizer** of function $f(.)$ if and only if:

$$\exists \epsilon > 0 : \forall y \in \mathcal{D}, \|y - x\|_2 \leq \epsilon \Longrightarrow f(x) \leq f(y), \tag{4}$$

meaning that in an $\epsilon$-neighborhood of $x$, the value of function is minimum at $x$.

## Definition (Global minimizer)

A point $x \in \mathcal{D}$ is a **global minimizer** of function $f(.)$ if and only if:

$$f(x) \leq f(y), \quad \forall y \in \mathcal{D}. \tag{5}$$



local minimizer    local maximizer    global and strict minimizer    saddle point    global (but not strict) maximizer

# Minimizer in convex function

### Lemma (Minimizer in convex function)

*In a **convex function**, any local minimizer is a global minimizer. In other words, in a convex function, there exists only **one** local minimum value which is the global minimum value.*

### Proof.

Proof can be found in the appendix of the tutorial [1]. □

As an imagination, a convex function is like a multi-dimensional bowl with only one minimum value (it may have several local minimizers but with the same minimum values).

strongly convex        convex (but not strongly convex)
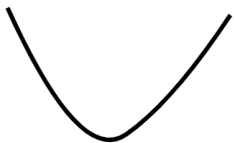
# Minimizer in convex function

## Lemma (Gradient of a convex function at the minimizer point)

*When the function $f(.)$ is convex and differentiable, a point $x^*$ is a minimizer if and only if:*

$$\nabla f(x^*) = 0.$$ $\longrightarrow$ first order condition

## Proof.

Proof can be found in the appendix of the tutorial [1]. □
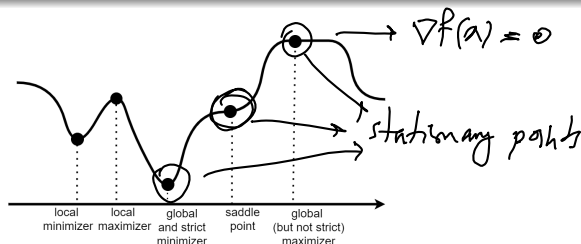


strongly convex          convex (but not strongly convex)

# Stationary, extremum, and saddle points

## Definition (Stationary, extremum, and saddle points)

- In a general (not-necessarily-convex) function $f(.)$, a point $\boldsymbol{x}^*$ is a **stationary** if and only if $\nabla f(\boldsymbol{x}^*) = \boldsymbol{0}$.
- By passing through a **saddle point**, the sign of the second derivative flips to the opposite sign.
- **Minimizer** and **maximizer** points (locally or globally) minimize and maximize the function, respectively.
- A **saddle point** is neither minimizer nor maximizer, although the gradient at a saddle point is zero.
- Both minimizer and maximizer are also called the **extremum points**.
- A stationary point can be either a minimizer, a maximizer, or a saddle point of function.



local
minimizer

local
maximizer

global
and strict
minimizer

saddle
point

global
(but not strict)
maximizer

$\nabla f(a) = 0$

stationary points

# First-order optimality condition

## Lemma (First-order optimality condition [2, Theorem 1.2.1])

If $\boldsymbol{x}^*$ is a local minimizer for a differentiable function $f(.)$, then:

$$\nabla f(\boldsymbol{x}^*) = \boldsymbol{0}. \tag{6}$$

Note that if $f(.)$ is convex, this equation is a necessary and sufficient condition for a minimizer.

## Proof.

Proof can be found in the appendix of the tutorial [1]. $\qquad\square$

## Note

If setting the derivative to zero, $\nabla f(\boldsymbol{x}^*) = \boldsymbol{0}$, gives a closed-form solution for $\boldsymbol{x}^*$, the optimization is done. Otherwise, we should solve it iteratively by either classic optimization or metaheuristic optimization.

# Arguments of optimization
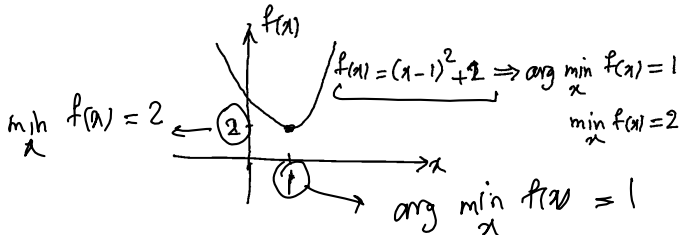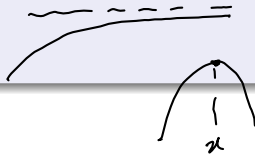
## Definition (Arguments of minimization and maximization)

In the domain of function, the point which minimizes (resp. maximizes) the function $f(.)$ is the argument for the minimization (resp. maximization) of function.

The minimizer and maximizer of function are denoted by

$$\arg\min_{x} f(x), \text{ and}$$

$$\arg\max_{x} f(x),$$

respectively.

# Converting optimization problems
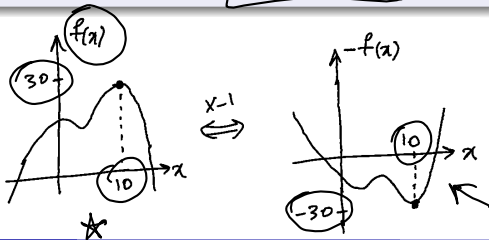
$$P(Y|X) = \frac{P(X|Y) P(Y)}{P(X)}$$

## Converting max to min and vice versa

We can convert convert maximization to minimization and vice versa:

$$\text{maximize } f(x) = -\text{minimize } (-f(x)),$$
$$\text{minimize } f(x) = -\text{maximize } (-f(x)).$$

We can have similar conversions for the arguments of maximization and minimization but as the sign of optimal value of function is not important in argument, we do not have the negative sign before maximization and minimization:

$$\arg\max_x f(x) = \arg\min_x (-f(x)),$$
$$\arg\min_x f(x) = \arg\max_x (-f(x)).$$

# Converting optimization problems

## Converting max to min and vice versa

We can convert convert maximization to minimization and vice versa using the reciprocal of cost function:
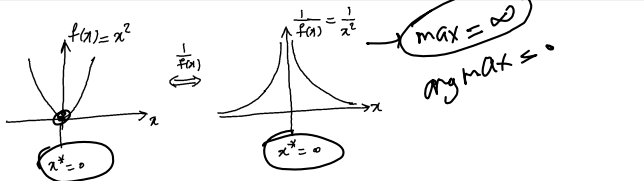
$$\underset{x}{\text{maximize}}\ f(x) = \frac{1}{\underset{x}{\text{minimize}}\ \frac{1}{f(x)}},$$

$$\underset{x}{\text{minimize}}\ f(x) = \frac{1}{\underset{x}{\text{maximize}}\ \frac{1}{f(x)}}.$$

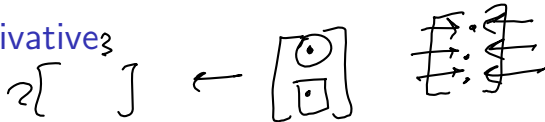We can have similar conversions for the arguments of maximization and minimization:

$$\arg\max_x f(x) = \arg\min_x \frac{1}{f(x)},$$

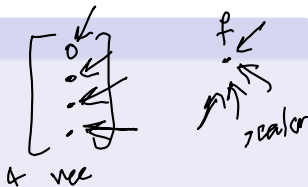$$\arg\min_x f(x) = \arg\max_x \frac{1}{f(x)}.$$

**Preliminaries on Derivatives**

# Dimensionality of derivatives

- Consider a function $f : \mathbb{R}^{d_1} \to \mathbb{R}^{d_2}$, $f : \boldsymbol{x} \mapsto f(\boldsymbol{x})$.
- Derivative of function $f(\boldsymbol{x}) \in \mathbb{R}^{d_2}$ with respect to (w.r.t.) $\boldsymbol{x} \in \mathbb{R}^{d_1}$ has dimensionality $(d_1 \times d_2)$.
- This is because tweaking every element of $\boldsymbol{x} \in \mathbb{R}^{d_1}$ can change every element of $f(\boldsymbol{x}) \in \mathbb{R}^{d_2}$. The $(i,j)$-th element of the $(d_1 \times d_2)$-dimensional derivative states the amount of change in the $j$-th element of $f(\boldsymbol{x})$ resulted by changing the $i$-th element of $\boldsymbol{x}$.
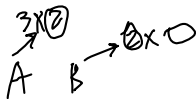
## Examples

- The derivative of a scalar w.r.t. a scalar is a scalar.
- The derivative of a scalar w.r.t. a vector is a vector.
- The derivative of a scalar w.r.t. a matrix is a matrix.
- The derivative of a vector w.r.t. a vector is a matrix.
- The derivative of a vector w.r.t. a matrix is a rank-3 tensor.
- The derivative of a matrix w.r.t. a matrix is a rank-4 tensor.

## Dimensionality of derivative



In more details:

- If the function is $f : \mathbb{R} \to \mathbb{R}, f : x \mapsto f(x)$, the derivative $(\partial f(x)/\partial x) \in \mathbb{R}$ is a scalar because changing the scalar $x$ can change the scalar $f(x)$.
- If the function is $f : \mathbb{R}^d \to \mathbb{R}, f : \boldsymbol{x} \mapsto f(\boldsymbol{x})$, the derivative $(\partial f(\boldsymbol{x})/\partial \boldsymbol{x}) \in \mathbb{R}^d$ is a vector because changing every element of the vector $\boldsymbol{x}$ can change the scalar $f(\boldsymbol{x})$.
- If the function is $f : \mathbb{R}^{d_1 \times d_2} \to \mathbb{R}, f : \boldsymbol{X} \mapsto f(\boldsymbol{X})$, the derivative $(\partial f(\boldsymbol{X})/\partial \boldsymbol{X}) \in \mathbb{R}^{d_1 \times d_2}$ is a matrix because changing every element of the matrix $\boldsymbol{X}$ can change the scalar $f(\boldsymbol{X})$.
- If the function is $f : \mathbb{R}^{d_1} \to \mathbb{R}^{d_2}, f : \boldsymbol{x} \mapsto f(\boldsymbol{x})$, the derivative $(\partial f(\boldsymbol{x})/\partial \boldsymbol{x}) \in \mathbb{R}^{d_1 \times d_2}$ is a matrix because changing every element of the vector $\boldsymbol{x}$ can change every element of the vector $f(\boldsymbol{x})$.
- If the function is $f : \mathbb{R}^{d_1 \times d_2} \to \mathbb{R}^{d_3}, f : \boldsymbol{X} \mapsto f(\boldsymbol{X})$, the derivative $(\partial f(\boldsymbol{X})/\partial \boldsymbol{X})$ is a $(d_1 \times d_2 \times d_3)$-dimensional tensor because changing every element of the matrix $\boldsymbol{X}$ can change every element of the vector $f(\boldsymbol{X})$.
- If the function is $f : \mathbb{R}^{d_1 \times d_2} \to \mathbb{R}^{d_3 \times d_4}, f : \boldsymbol{X} \mapsto f(\boldsymbol{X})$, the derivative $(\partial f(\boldsymbol{X})/\partial \boldsymbol{X})$ is a $(d_1 \times d_2 \times d_3 \times d_4)$-dimensional tensor because changing every element of the matrix $\boldsymbol{X}$ can change every element of the matrix $f(\boldsymbol{X})$.

# Gradient, Jacobian, and Hessian

$\triangleq$

### Definition (Gradient)

Consider a function $f : \mathbb{R}^d \to \mathbb{R}$, $f : x \mapsto f(x)$. In optimizing the function $f$, the derivative of function w.r.t. its variable $x$ is called the **gradient**, denoted by:

$$\nabla f(x) := \frac{\partial f(x)}{\partial x} \in \mathbb{R}^d.$$

### Definition (Hessian)

Consider a function $f : \mathbb{R}^d \to \mathbb{R}$, $f : x \mapsto f(x)$. The second derivative of function w.r.t. to its derivative is called the **Hessian** matrix, denoted by:

$$B_{ij} = B_{ji}$$
$$B = B^\top$$

$$B = \nabla^2 f(x) := \frac{\partial^2 f(x)}{\partial x^2} \in \mathbb{R}^{d \times d}.$$

The Hessian matrix is symmetric. If the function is convex, its Hessian matrix is positive semi-definite.

$$\frac{\partial f}{\partial x} \in \mathbb{R}^d = g \rightarrow \frac{\partial g \rightarrow \mathbb{R}^d}{\partial x \rightarrow \mathbb{R}^d} \in \mathbb{R}^{d \times d}$$

# Gradient, Jacobian, and Hessian

$$\begin{bmatrix} \lambda_1 \\ \vdots \\ \lambda_{d_1} \end{bmatrix} \qquad \begin{bmatrix} f_1 \\ f_2 \\ \vdots \\ f_{d_2} \end{bmatrix}$$
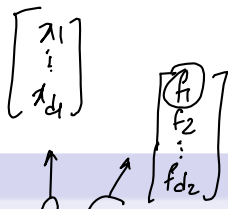
## Definition (Jacobian)

If the function is multi-dimensional, i.e., $f : \mathbb{R}^{d_1} \to \mathbb{R}^{d_2}$, $f : \boldsymbol{x} \mapsto f(\boldsymbol{x})$, the gradient becomes a matrix:

$$\boldsymbol{J} := \Big[\frac{\partial f}{\partial x_1}, \ldots, \frac{\partial f}{\partial x_{d_1}}\Big]^\top = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_{d_2}}{\partial x_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_1}{\partial x_{d_1}} & \cdots & \frac{\partial f_{d_2}}{\partial x_{d_1}} \end{bmatrix} \in \mathbb{R}^{d_1 \times d_2},$$

where $\boldsymbol{x} = [x_1, \ldots, x_{d_1}]^\top$ and $f(\boldsymbol{x}) = [f_1, \ldots, f_{d_2}]^\top$.

This matrix derivative is called the **Jacobian** matrix.

**Optimization Problems**

# General optimization problem

Consider the function $f : \mathbb{R}^d \to \mathbb{R}$, $f : \boldsymbol{x} \mapsto f(\boldsymbol{x})$. Let the domain of function be $\mathcal{D}$ where $\boldsymbol{x} \in \mathcal{D}, \boldsymbol{x} \in \mathbb{R}^d$.

## Definition (Unconstrained optimization)

**Unconstrained** minimization of a cost function $f(.)$:

$$\underset{\boldsymbol{x}}{\text{minimize}} \quad f(\boldsymbol{x}),$$

where $\boldsymbol{x}$ is called the **optimization variable** and the function $f(.)$ is called the **objective function** or the **cost function**.

# General optimization problem

## Definition (Constrained optimization)

**Constrained** optimization problem where we want to minimize the function $f(x)$ while satisfying $m_1$ inequality constraints and $m_2$ equality constraint:

$$\begin{cases} \underset{x}{\text{minimize}} & f(x) \\ \text{subject to} & y_i(x) \leq 0, \ i \in \{1, \ldots, m_1\}, \\ & h_i(x) = 0, \ i \in \{1, \ldots, m_2\}. \end{cases}$$

$f(x)$ is the **objective function**, every $y_i(x) \leq 0$ is an **inequality constraint**, and every $h_i(x) = 0$ is an **equality constraint**.

$$y_1(x) \leq 0$$
$$y_2(x) \leq \cdot$$
$$\vdots$$
$$y_m(x) \leq 0$$

$$h_1(x) = \circ$$
$$\vdots$$
$$h_{m_2}(x) = \cdot$$

## Note

If some of the inequality constraints are not in the form $y_i(x) \leq 0$, we can restate them as:

$$y_i(x) \geq 0 \implies -y_i(x) \leq 0,$$
$$y_i(x) \leq c \implies y_i(x) - c \leq 0.$$

$$y(x)$$

Therefore, all inequality constraints can be written in the form $y_i(x) \leq 0$.

# General optimization problem

Example:

$$
\begin{aligned}
\underset{x}{\text{minimize}} \quad & x_1 + 3x_2^2 \\
\text{subject to} \quad & 2x_1 - 10x_2 \leq 5, \\
& -2x_1 + 5x_2 \geq 3, \\
& 4x_1 + 10x_2 = 6,
\end{aligned}
$$

*(handwritten annotations: minimize $x$; minimize $\{x_1, x_2\}$; $[x_1, x_2]^T$; minimize $x$)*

can be converted to:

$$
\begin{aligned}
\underset{x}{\text{minimize}} \quad & x_1 + 3x_2^2 \\
\text{subject to} \quad & 2x_1 - 10x_2 - 5 \leq 0, \\
& 2x_1 - 5x_2 + 3 \leq 0, \\
& 4x_1 + 10x_2 - 6 = 0.
\end{aligned}
$$

*(handwritten annotations: $\rightarrow m_1 = 2$; $\rightarrow m_2 = 1$)*
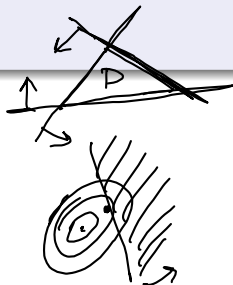
# Feasible point

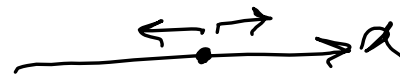## Definition (Feasible point)

The point **x** for the optimization problem:

$$\begin{cases} \underset{\mathbf{x}}{\text{minimize}} & f(\mathbf{x}) \\ \text{subject to} & y_i(\mathbf{x}) \leq 0, \ i \in \{1, \ldots, m_1\}, \\ & h_i(\mathbf{x}) = 0, \ i \in \{1, \ldots, m_2\}, \end{cases}$$

is feasible if:

$$\mathbf{x} \in \mathcal{D}, \text{ and}$$
$$y_i(\mathbf{x}) \leq 0, \quad \forall i \in \{1, \ldots, m_1\}, \text{ and}$$
$$h_i(\mathbf{x}) = 0, \quad \forall i \in \{1, \ldots, m_2\}.$$

$x \in \mathbb{R}$

line $\rightarrow$ ~~point~~ dot

$x \in \mathbb{R}^2$

line $\rightarrow$ line

$x \in \mathbb{R}^3$
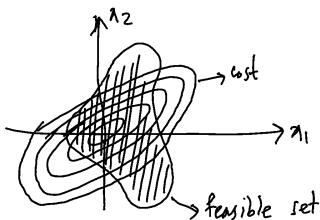
line $\rightarrow$ plane

hyperplane

# Constrained optimization with the feasible set

## Definition (Constrained optimization)

The **constrained** optimization problem can also be stated as:

$$\min_{\boldsymbol{x}} \quad f(\boldsymbol{x})$$
$$\text{subject to} \quad \boldsymbol{x} \in \mathcal{S},$$

where $\mathcal{S}$ is the feasible set of constraints.

**Converting Constrained Optimization to Unconstrained Optimization**

# Converting Constrained Optimization to Unconstrained Optimization

- Consider the following constrained optimization problem:

$$
\begin{cases}
\underset{x}{\text{minimize}} & f(x) \\
\text{subject to} & g(x) \leq 0.
\end{cases}
\tag{7}
$$

  This is a **constrained hard penalty** because it kills the problem if it is not satisfied (it can never happen as it would be **infeasible**).

- We can convert it to an unconstrained regularized problem:

$$
\underset{x}{\text{minimize}} \quad f(x) + \lambda g(x),
\tag{8}
$$

  where $\lambda > 0$ is the regularization parameter. This is a **soft penalty** which tolerates some violence of the constraint. The solution of this problem is approximately similar to the solution of the constrained hard penalty.

$$
\left\{ \begin{array}{l} \min\limits_{\boxed{x}} \ f(x) \\ \\ \text{s.t.} \quad \cancel{\otimes} \ y_1(x) \cancel{\phi} \leq 0 \\ \qquad\quad y_2(x) \leq 0 \\ \qquad\quad h(x) = 0 \end{array} \right\}
$$

$$
\min_{x, \lambda_1, \lambda_2, \nu} \ \mathscr{L} = \boxed{f(x)} + \boxed{\lambda_1 \, y_1(x)} + \boxed{\lambda_2 \, y_2(x)} + \boxed{\nu \, h(a)}
$$

Lagrangian

Lagrange

multiplier

$$
\frac{\partial \mathscr{L}}{\partial x} \overset{\text{set}}{=} 0
$$

$$
\frac{\partial \mathscr{L}}{\partial \lambda_1} \overset{\text{set}}{=} 0
$$

$$
\frac{\partial \mathscr{L}}{\partial \lambda_2} \overset{\text{set}}{=} 0
$$

dual variables

regularisation params

# Converting Constrained Optimization to Unconstrained Optimization

$$\lim_{\lambda \to \infty} f(\lambda) + \lambda g(\lambda) \qquad g(\lambda) \leq 0$$

- We can also convert it to an **unconstrained** problem using indicator function. In optimization, **indicator function** $\mathbb{I}(.)$ is zero if its condition is satisfied and is infinite otherwise.

$$\mathbb{I}(\boldsymbol{x} \in \mathcal{S}) = \begin{cases} 0 & \text{if } \boldsymbol{x} \in \mathcal{S} \\ \infty & \text{if } \boldsymbol{x} \notin \mathcal{S}. \end{cases} \tag{9}$$

The problem can become regularized using the indicator function:

$$\underset{\boldsymbol{x}}{\text{minimize}} \quad f(\boldsymbol{x}) + \lambda \mathbb{I}(g(\boldsymbol{x}) \leq 0), \tag{10}$$
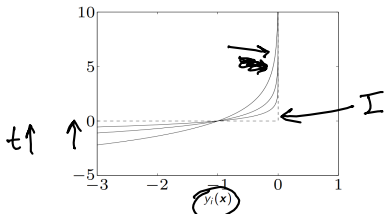
where $\lambda > 0$ is the regularization parameter. This is a **unconstrained hard penalty** because it kills the problem if it is not satisfied (it blows the regularized cost function to infinity). The solution of this problem is exactly the same as the solution of the constrained hard penalty.

# Converting Constrained Optimization to Unconstrained Optimization

- In practice, it is not possible to implement the indicator function in computer as it has infinity and also it is not smooth (differentiable). Therefore, we can use approximations of the indicator function by **barrier functions**. One of the barrier functions is logarithm, named the **logarithmic barrier** or **log barrier** in short. It approximates the indicator function by:

Newton's method

log-barrier method

$$\mathbb{I}(y_i(\mathbf{x}) \le 0) \approx -\frac{1}{t} \log(-y_i(\mathbf{x})), \tag{11}$$

where $t > 0$ (usually a large number such as $t = 10^6$) and the approximation becomes more accurate by $t \to \infty$.



The solution of the problem, which uses barrier functions as approximations of the indicator function, is **approximately** similar to the solution of the hard penalty. This is the technique that the **interior point method** uses.

# References

[1] B. Ghojogh, A. Ghodsi, F. Karray, and M. Crowley, "KKT conditions, first-order and second-order optimization, and distributed optimization: Tutorial and survey," *arXiv preprint arXiv:2110.01858*, 2021.

[2] Y. Nesterov, *Lectures on convex optimization*, vol. 137.
Springer, 2018.