# Lecture 1.    Introduction

Mathematical optimization problem

$$\text{minimize} \quad f(x)$$

$$\text{subject to} \quad x \in \Omega$$

where   $f: \mathbb{R}^n \to \mathbb{R}$   objective function   目标函数

$x = (x_1, \ldots, x_n)^T \in \mathbb{R}^n$.   optimization variables.

$\Omega \subseteq \mathbb{R}^n$   feasible set / constraint set

- $x \in \Omega$   feasible  可行   infeasible o.w.

- $\Omega$   specified by constraint functions $g_1, \ldots g_m$.

$$\min_{x} \quad f(x)$$

$$\text{s.t.} \quad g_i(x) \leq 0, \quad i = 1, \ldots, m.$$

$x^*$: optimal solution.   $f(x)$ achieves optimal.

$$x^* = \arg\min f(x).$$

Remark:  maximizing $f(x)$ is equivalent to minimizing $-f(x)$

why linear and convex?

In general optimization problems are very difficult to solve.

# Knapsack problem.    背包问题

$n$ types of knapsacks.

$i^{th}$ type    carry $a_i$ pencils   and $b_i$ books.   costs. $c_i$

A pencils and B books in total.

Goal : spend least money .

$$\min \quad \sum_i c_i x_i$$

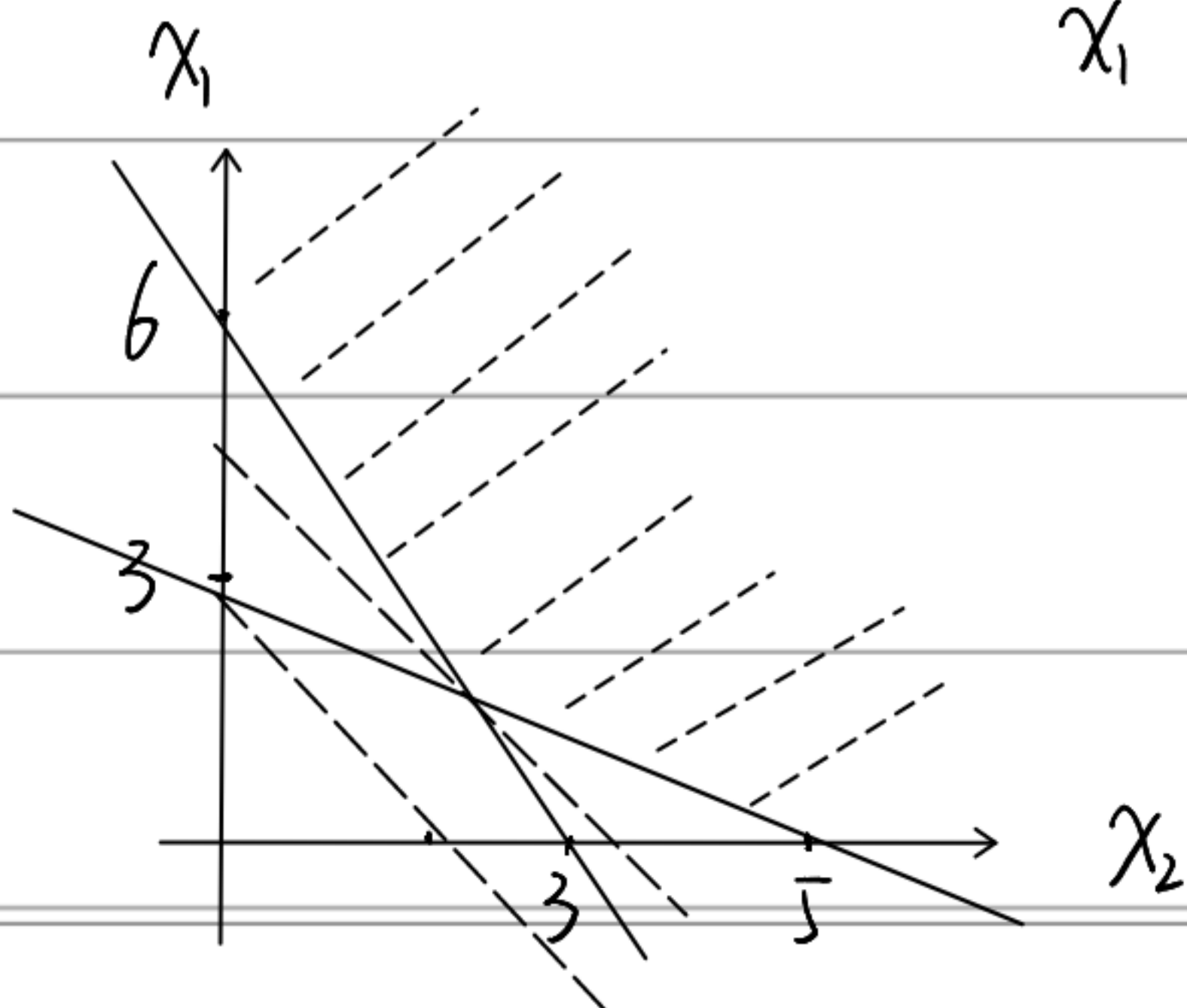$$\text{s.t.} \quad \sum_i a_i x_i \geqslant A$$

$$\sum_i b_i x_i \geqslant B.$$

How to solve it ?

if only 2 types .

$$\min \quad 10 x_1 + 15 x_2.$$

$$\text{s.t.} \quad 5 x_1 + 3 x_2 \geqslant 15$$

$$x_1 + 2 x_2 \geqslant 6$$



more types ?

simplex algorithm.

primal dual

# Data fitting.

| T | V |
|---|---|
| 30 | 1.011 |
| 40 | 1.019 |
| 50 | 1.032 |
| 60 | 1.041 |
| ... | ... |

$$\underset{\downarrow}{V} \quad \underset{\downarrow}{T}$$

$$y = kx + b.$$

what are the two coefficient

k and b ?

## Least squares method.

given n measurements.

$(x_1, y_1), \cdots (x_n, y_n).$

assume $i^{th}$ error denoted by $\varepsilon_i$.

least squares criterion

minimize $\sum \varepsilon_i = \sum (y_i - kx_i - b)^2$

## Geometric explanation : projection

$$\begin{cases} 30k + b = 1.011 \\ 40k + b = 1.019 \\ 50k + b = 1.032 \end{cases}$$

$$k \begin{bmatrix} 30 \\ 40 \\ 50 \end{bmatrix} + b \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = ?$$

projection of $\begin{bmatrix} 1.011 \\ 1.019 \\ 1.032 \end{bmatrix}$ onto the subspace spanned by

minimize $\left\| \begin{bmatrix} 30 & 1 \\ 40 & 1 \\ 50 & 1 \end{bmatrix} \begin{bmatrix} k \\ b \end{bmatrix} - \begin{bmatrix} 1.011 \\ 1.019 \\ 1.032 \end{bmatrix} \right\|_2^2$ $\begin{bmatrix} 30 \\ 40 \\ 50 \end{bmatrix}$ and $\begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$

Norm. inner product. distance to hyperplane $\qquad$ not necessary $\mathbb{R}^n$

Fourier

IP: an inner product $\langle \cdot, \cdot \rangle$ is a function $S \times S \to \mathbb{R}$. s.t.

1. nonnegative $\qquad \langle x, x \rangle \geqslant 0$. $\qquad = 0$ iff $x = 0$

2. symmetric $\qquad \langle x, y \rangle = \langle y, x \rangle$

3. linearity $\qquad \langle sx, y \rangle = s\langle x, y \rangle$ $\qquad$ (homogeneity)

$$\langle x+y, z \rangle = \langle x, z \rangle + \langle y, z \rangle \qquad \text{(additivity)}$$

if $\langle x, y \rangle = 0$ then $x$ and $y$ are called orthogonal. $\frac{1}{\perp}$ 交.

Euclidean inner product space: $\langle x, y \rangle = x^T y = \sum x_i y_i$

$\| \cdot \|$.

Norm: a norm is a function $\mathbb{R}^n \to \mathbb{R}$. s.t.

1. nonnegative $\qquad \| x \| \geqslant 0$ $\qquad = 0$ iff $x = 0$.

2. positive homogeneity $\qquad \| a x \| = |a| \, \| x \|$

3. triangle inequality $\qquad \| x + y \| \leq \| x \| + \| y \|$.

$L^P$-norm. or $p$-norm for real $p \geqslant 1$.

$$\| x \|_p = \left( |x_1|^P + |x_2|^P + \cdots + |x_n|^P \right)^{1/P}.$$

in particular. 1-norm : $\| x \|_1 = \sum_i |x_i|$

(Euclidean norm) 2 : $\| x \| \overset{\Delta}{=} \| x \|_2 = \sqrt{\langle x, x \rangle} = \sqrt{\sum_i x_i^2}$

default

$\infty$ - norm :     $\|x\|_\infty = \max \{ |x_1|, |x_2|, \ldots, |x_n| \}$.

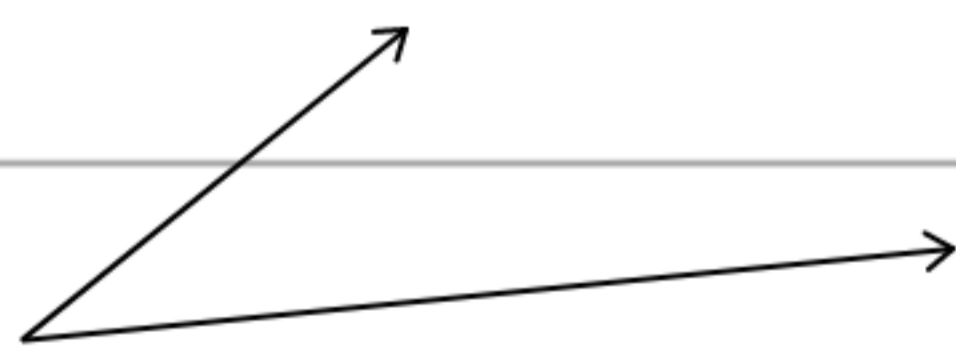1 norm

Manhattan
distance.

Cauchy — Schwarz inequality

$$|\langle u, v \rangle|^2 \leq \langle u, u \rangle \cdot \langle v, v \rangle \qquad \text{or.}$$

$$|\langle u, v \rangle| \leq \|u\| \, \|v\|.$$

$n$-dimensional Euclidean space. ($\mathbb{R}^n$)

$$\left( \sum_i u_i v_i \right)^2 \leq \left( \sum_i u_i^2 \right) \left( \sum_i v_i^2 \right).$$
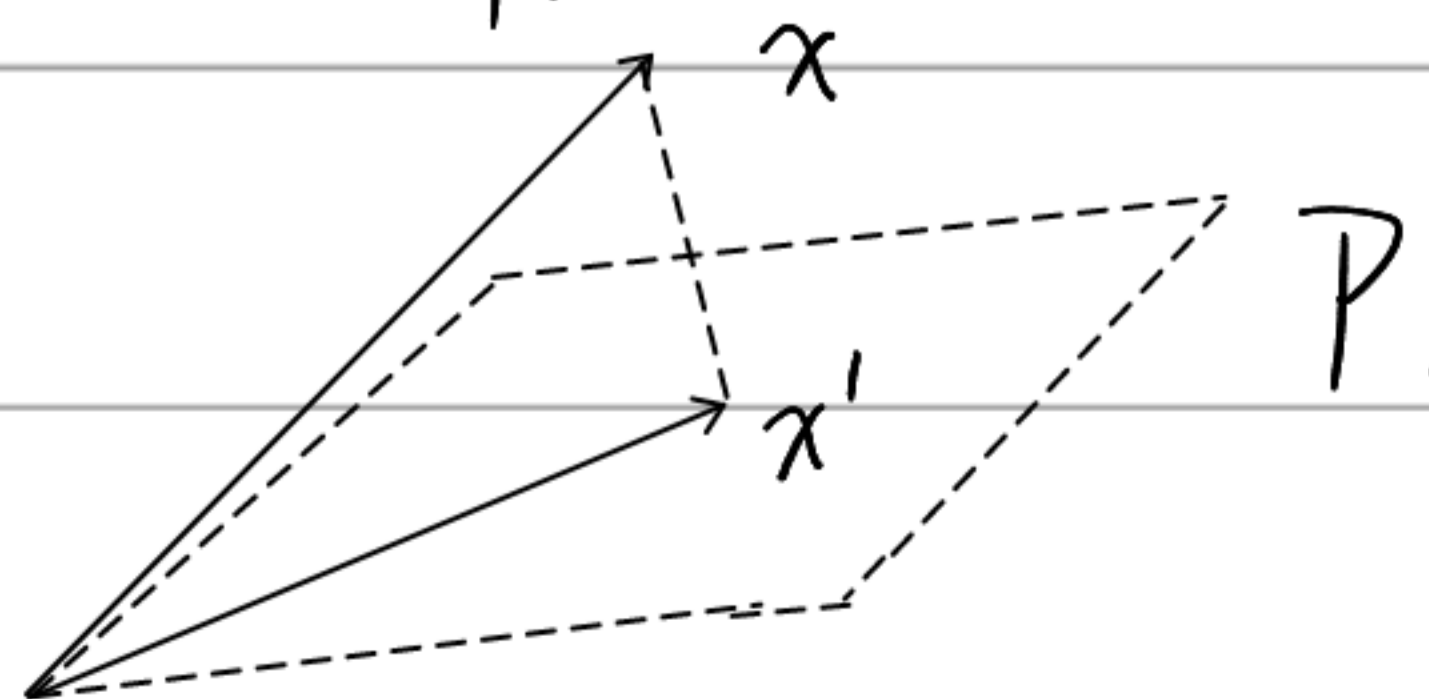
Geometric explanation ·    projection.

Distance to hyperplane.

hyperplane $P$:  $w^T x + b = 0$        $w \perp P$.

orthogonal projection

$(x - x') \perp P$    $w^T x' + b = 0.$

$$x - x' = r \cdot w \quad \text{for some } r \in \mathbb{R}.$$

$$w^T (x - r \cdot w) + b = 0 \implies r = \frac{w^T x + b}{w^T w} \longleftarrow \|w\|^2$$

distance from $x$ to $P$ is.

$$\min_{y \in P} \| x - y \| = \| x - x' \| = \| r w \| = \frac{|w^T x + b|}{\|w\|}$$

## Linear least squares regression.

given $m$ measurements $(x_1, y_1) \cdots (x_m, y_m)$.

assume that $\quad y = w^T x + b$.

The least squares regression. is to compute the following OPT.

$$\min_{w \in \mathbb{R}^n, \, b \in \mathbb{R}} \sum_i (Xw + b \cdot \mathbf{1} - y)^2$$

$$= \min_{w \in \mathbb{R}^n, \, b \in \mathbb{R}} \| Xw + b \cdot \mathbf{1} - y \|$$

where $X = (x_1, \ldots x_n)^T \in \mathbb{R}^{m \times n}$. $y = (y_1, \ldots y_m)^T \in \mathbb{R}^m$

$e \perp$ hyperspace spanned by column vector of $X$.
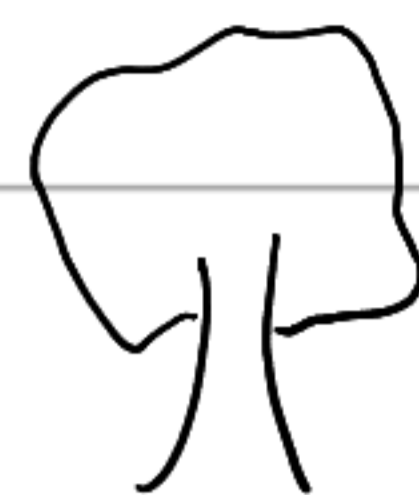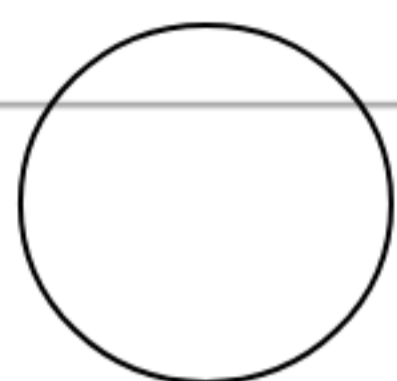
$$y - e \; \cancel{\in} \qquad X^T e = 0$$

suppose $y - e = X \hat{w}$ thus we have.

$$X^T (y - X \hat{w}) = 0.$$

$$X^T X \hat{w} = X^T y$$

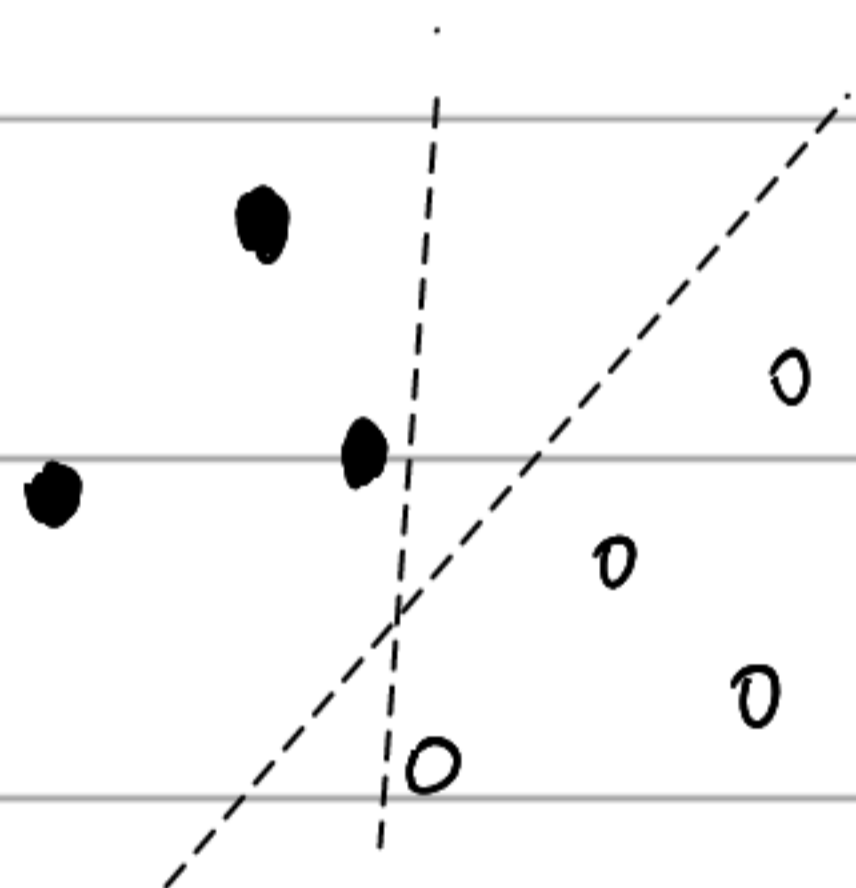# Classification and support vector machine.

classify

given $m$ data points. $(x_1, y_1) \cdots (x_m, y_m)$.

classfier is a function s.t.

$$\begin{cases} f(x_i) > 0 & \text{iff } y_i = +1 \\ f(x_i) < 0 & \text{iff } y_i = -1 \end{cases} \iff y_i \, f(x_i) > 0.$$

linear classifier : $f(x) = w^T x + b$.

which one is better?

against noise

maximize the minimum
distance to the hyper plane.

Support vector machine : linear classifer with max margin

$$\max_{w, b} \min_{1 \leq i \leq m} \frac{|w^T x_i + b|}{\|w\|}.$$

$$\text{s.t.} \quad y_i (w^T x_i + b) > 0.$$

since $y_i = \text{sgn}(w^T x_i + b)$. $|w^T x_i + b| = y_i (w^T x_i + b)$.

$\forall \alpha > 0$. $\tilde{w} = \alpha w$. $\tilde{b} = \alpha b$ also feasible.

choosing $\alpha$ properly. s.t. $\min\limits_{1 \leq i \leq m} y_i (\tilde{w}^T x_i + \tilde{b}) = 1.$

$$\max \quad \frac{1}{\|\tilde{w}\|}$$

$$\text{s.t.} \quad y_i (\tilde{w}^T x_i + \tilde{b}) \geq 1.$$

which is equivalent to.

$$\min \quad \frac{1}{2} \|\tilde{w}\|^2$$

$$\text{s.t.} \quad y_i (\tilde{w}^T x_i + \tilde{b}) \geq 1.$$


Global optima and local optima.

$$\min\limits_{x \in X} f(x). \qquad \text{let } x^* \triangleq \arg\min\limits_{x \in X} f(x).$$

$x^*$ is a global minimum if $f(x^*) \leq f(x)$

global maximum.

global optima may not exist.

$-$ $f(x) = x.$ $X = \mathbb{R}.$ $\inf f(x) = -\infty.$

$\sim$ $f(x) = \frac{1}{x}$ $X = \mathbb{R}_{>0}$ $\inf f(x) = 0.$

when will global optima exist ?

Continuous functions on compact sets have global optima.