FoDA . Clustering

L22 .
         Voronoi Diagrams

# What is Clustering?

**Input**    Set of objects  $X = \{x_1, x_2, \ldots x_m\}$

   Distance  $D: X \times X \to \mathbb{R}^+$

(this class: $X \subset \mathbb{R}^d$ , $D(x_1, x_2) = \|x_1 - x_2\|$    *Euclidean*
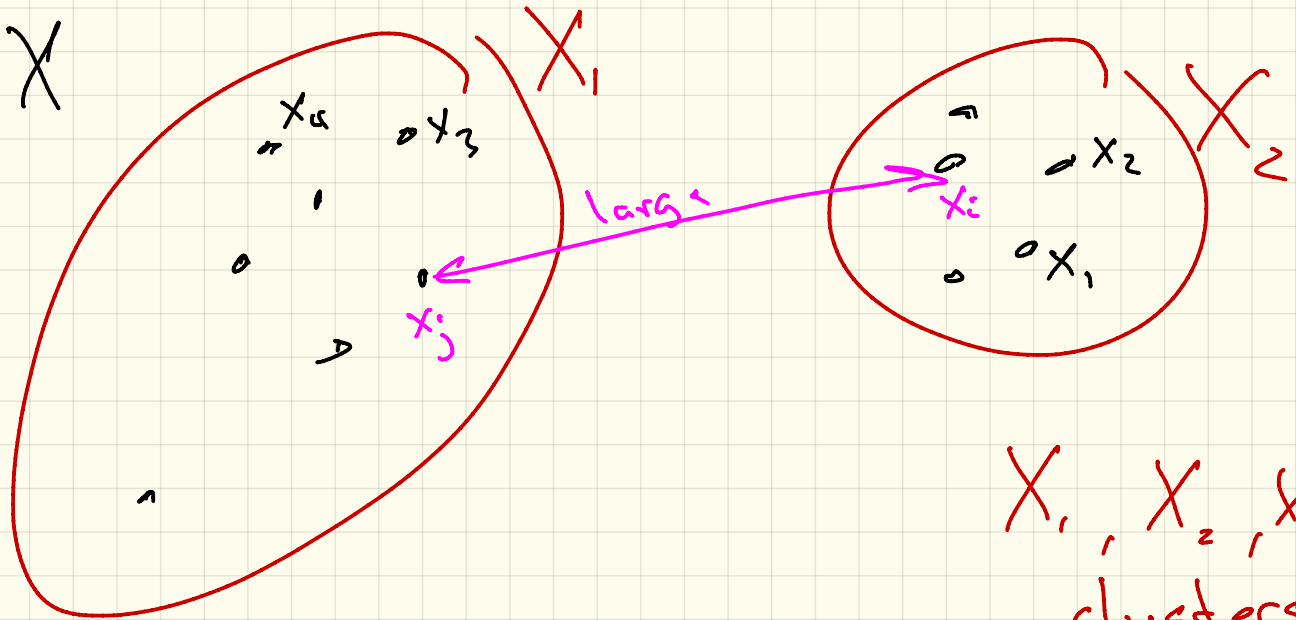
↙ usually (part of input)

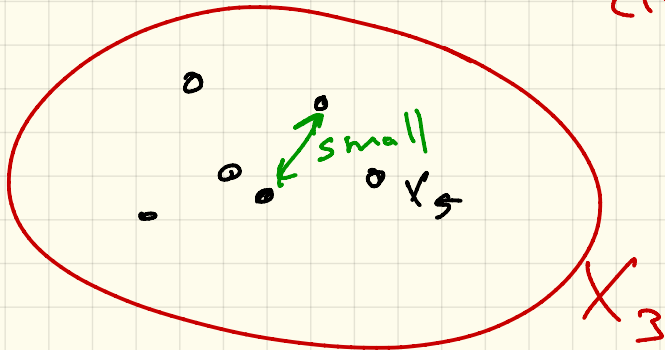if ill-defined ⟼ trouble

**Goal**

$\mathbb{R}$  subsets  $\{X_1, X_2, \ldots X_R\}$

$X_i \subset X$
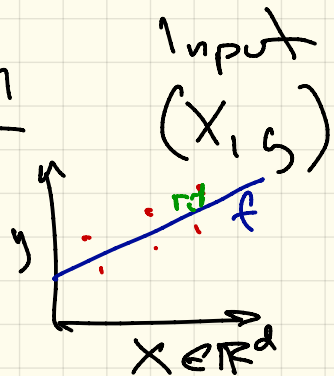
$x_i, x_i' \in X \Rightarrow D(x_i, x_i')$ small

$x_i \in X \; x_j \in X_j \; i \neq j \Rightarrow D(x_i, x_j)$ large

$X$

$X_1$

$x_4$   $x_3$

large

$x_i$

$X_2$

$x_2$

$x_1$

$x_j$

$X_1, X_2, X_3$

clusters

$k = 3$

small

$x_5$

$X_3$

# Regression

Input
$$(X_i, y)$$
$$\rightarrow f(x_i \in X) = \hat{y}_i$$



measure $\quad r_i = y_i - \hat{y}_i$

$$\min \; \frac{1}{i} \, r_i^2$$

# Dimensionality Reduction
## PCA

$$A = X \subset \mathbb{R}^d \rightarrow B \subset \mathbb{R}^d$$
$$\text{rank} = k$$



$$\min \; \sum_i r_i^2$$

$$\| A - A_k \|_F^2 = \| A - \pi_B (A) \|_F^2$$
$$= \sum_{i=1}^n \| a_i - \pi_B (a_i) \|^2$$
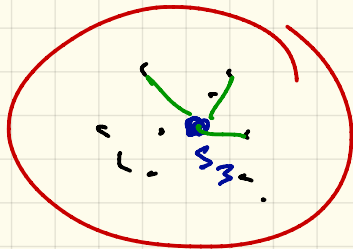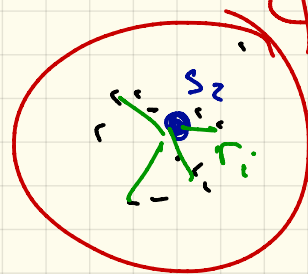
# Assignment-based Clustering

$X \subset \mathbb{R}^d$, $\quad D : \|\cdot - \cdot\|$ $\qquad$ k clusters

$$D(x_i, x_j) = \|x_i - x_j\|$$

Model $\quad S = \{s_1, s_2, \ldots, s_k\} \subset \mathbb{R}^d$

$$\phi_S(x) = \overset{arg}{\underset{s_i \in S}{min}} \|x - s_i\|$$

maps to closest site

$$r_i = \|x_i - s_j\|$$

$$= \|x_i - \phi_S(x_i)\|$$

# Post Office Problem

$$\phi_S(x) = \underset{s_i \in S}{\arg\min} \|x - s_i\|$$

$$\phi_S : \mathbb{R}^d \to S$$

$|S| = k$

$d = 1$    Solvable in $O(\log k)$ time

# Voronoi Diagram

$d = 2$

$\phi_S$

$S_2$

$S_1$

$\{x \in \mathbb{R}^d \text{ s.t. } \|S_2 - x\| = \|S_5 - x\| = \|S_6 - x\|$

$\leq \|S_j - x\| \; j \neq 2, 5, 6\}$

$\{x \in \mathbb{R}^2 \mid \phi_S(x) = S_5\}$

$= V_{2,5,6}$

$e_{1,2} = \{x \in \mathbb{R}^2 \text{ s.t.}$

$\|x - S_1\| = \|x - S_2\|$

$\leq \|x - S_j\|$

$j \neq 1, 2\}$

$S_4$

$S_6$

Voronoi cell of $S_5$

$S_3$

$S_5$

# Voronoi Diagram in $\mathbb{R}^2$

→ "Complexity" is $O(k)$
      # vertices & # edges

→ Compute in $O(k \log k)$ time

→ Solve $\phi_S(x)$ in $O(\log k)$ time

**Bad News** not true for $d \geq 3$
complexity in $\mathbb{R}^3, \mathbb{R}^4$   $O(k^2)$
complexity in $\mathbb{R}^d$   $O\left(k^{\lceil d/2 \rceil}\right)$ ← curse of dimensionality

$$\phi_S(x) = \underset{s_j \in S}{\arg\min} \, \| s_j - x \|$$

$$S = \{ s_1, \ldots s_k \}$$

in high-d

⇓

very high

complexity

---

0. $s = s_1$
   $m = \| x - s_1 \|$
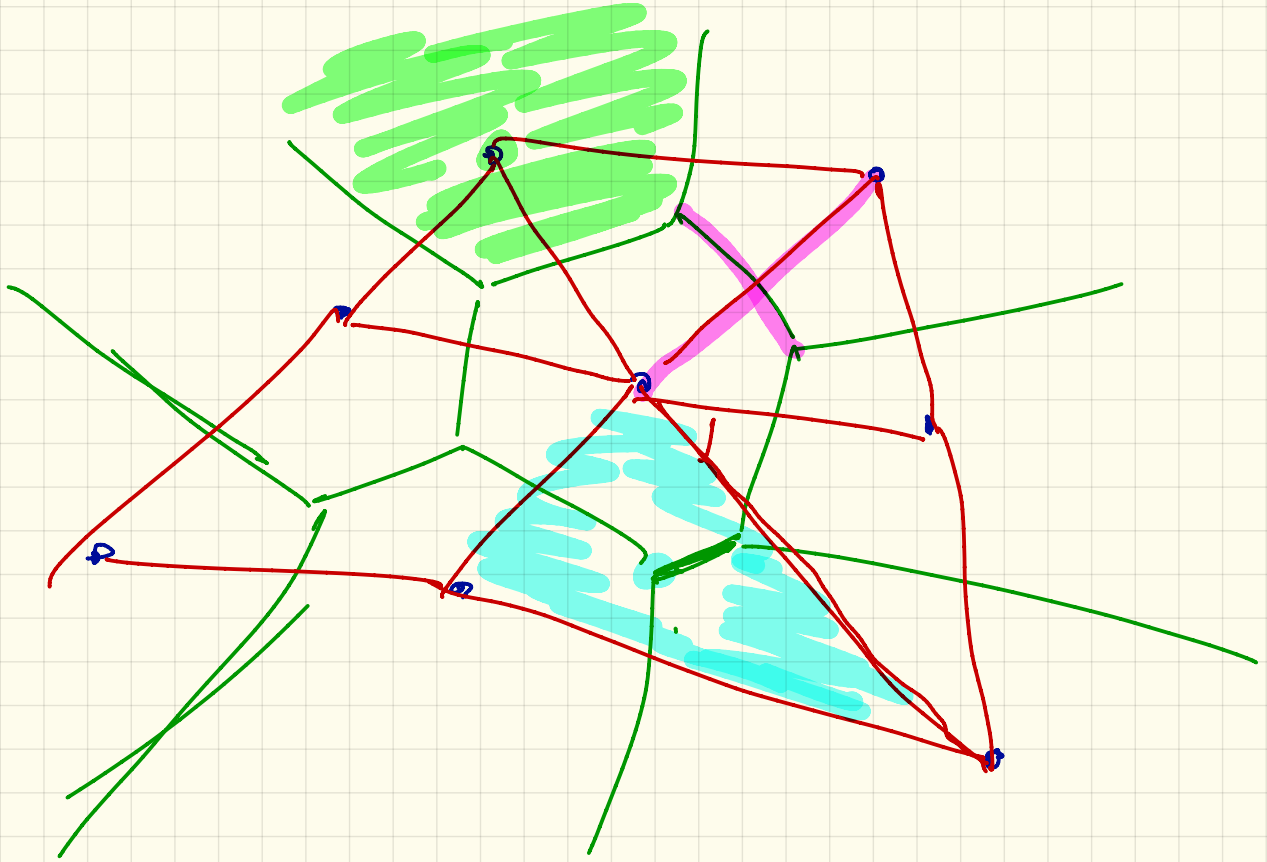1. for $i = 2$ to $k$
      if $m > \| x - s_i \|$
         $s = s_i$
         $m = \| x - s_i \|$
2. return $s$

---

$O(k)$ time

# Delaunay Triangulation

Input $X \subset \mathbb{R}^d$, $k$

---

Assignment - Based Clustering

Goal ~~Find~~ $S = \{s_1, \ldots s_k\} \subset \mathbb{R}^d$

$$\underset{S}{\text{minimize}} \quad f\left\{ \|x_i - \phi_S(x_i)\| \right\}$$

find $S = \{s_1 \ldots s_k\}$   $\quad f = \sum_i r_i \quad \leftarrow$ k-median
to try to                        $\quad\quad\quad\quad\quad\quad\quad \leftarrow$ k-mediod $(S \subset X)$
minimize $\sum_i \|x_i - \phi_S(x_i)\|^2$   $\quad f = \sum_i r_i^2 \quad \leftarrow$ k-means
$S$

$\quad\quad\quad\quad\quad\quad\quad\quad f = \max r_i \quad \leftarrow$ k-center

# Dim Reduction (PCA)

## K-means clustering

| Dim Reduction (PCA) | K-means clustering |
|---|---|
| $k$ basis functions $V_B = \{v_1, v_2, \dots v_k\}$ $\|v_i\|$ $\langle v_i, v_j \rangle = 0$ | $k$ sites $S = \{s_1, s_2, \dots s_k\}$ |

$\longleftarrow$ Find $\longrightarrow$

### Assignment

$$\pi_B(x) = \underset{b \in B}{\arg\min} \|b - x\|$$
projection

### Assignment

$$\phi_S(x) = \underset{s_j \in S}{\arg\min} \|x - s_j\|$$
NN

### Goal
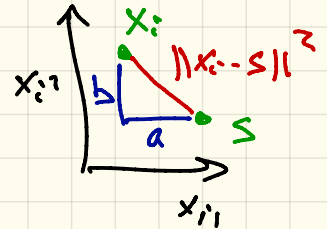
$$\sum_i r_i^2 = \sum_i \|x_i - \pi_B(x_i)\|^2$$

### Goal

$$\sum_i r_i^2 = \sum_{i=1}^{n} \|x_i - \phi_S(x_i)\|^2$$

Set $\quad X \in \{x_1, \dots x_n\} \qquad x_i \in \mathbb{R}^d$

Find $\quad s \quad$ minimize

$$\sum_{i=1}^{n} \| s - x_i \|^2$$

$$\| x_i - s \|^2 = \sum_{j=1}^{d} (x_{ij} - s_j)^2$$



$x_i$

$\| x_i - s \|^2$

$x_{i2}$ $\quad b$

$a$ $\quad s$

$x_{i1}$

$\| x_i - s \|^2 = a^2 + b^2$

Soln: $\quad s = \frac{1}{n} \sum_{i=1}^{n} x_i$

$$s_j = \frac{1}{n} \sum_{i=1}^{n} x_{ij}$$