

Unsupervised Learning

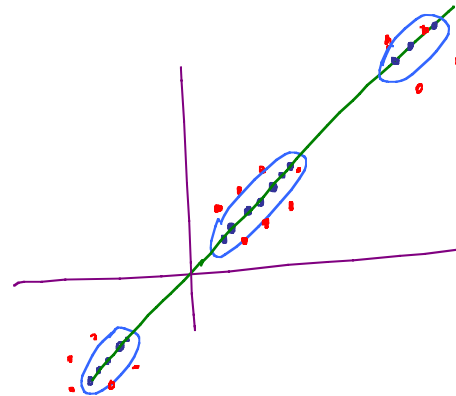
- ↳ Representation learning

- ↳ PCA

- ↳ kernel PCA

- ↳

- ↳ Clustering (today)

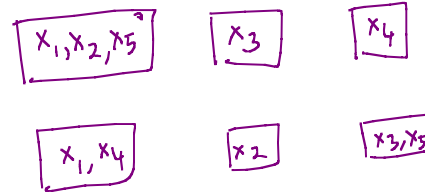


"Clustered together"

Goal: $\{x_1, \dots, x_n\} \quad x_i \in \mathbb{R}^d$

Partition the data k different clusters

Example: $\{x_1, x_2, x_3, x_4, x_5\} \quad \underline{k=3}$



3^5 possibilities

$x_1, x_2, \dots, x_n \leftarrow$ DATA POINTS

$z_1, z_2, \dots, z_n \leftarrow$ CLUSTER INDICATOR

$z_i \in \{1, \dots, k\}$

Question: Given a cluster assignment, how good is it?

$$F(z_1, \dots, z_n) = \sum_{i=1}^n \|x_i - \mu_{z_i}\|_2^2$$

↪ Mean/average of z_i^{th} cluster

$$\mu_k = \frac{\sum_{i=1}^n x_i \mathbb{1}(z_i = k)}{\sum_{i=1}^n \mathbb{1}(z_i = k)}$$

$$\mathbb{1}(u) = \begin{cases} 1 & \text{if } u \text{ is true} \\ 0 & \text{o/w} \end{cases}$$

Example:

$$\begin{array}{ccccc} \textcircled{x_1} & \textcircled{x_2} & \textcircled{x_3} & \textcircled{x_4} & \textcircled{x_5} \\ z_1 = -1 & z_2 = 2 & z_3 = -1 & z_4 = -1 & z_5 = 2 \end{array}$$

$k=2$

$$M_1 = \frac{x_1 + x_3 + x_4}{3} \quad ; \quad M_2 = \frac{x_2 + x_5}{2}$$

Goal

$$\min_{\{z_1, \dots, z_n\}} \sum_{i=1}^n \|x_i - \mu_{z_i}\|^2$$

Too many possibilities! (k^n)

NP-HARD

LLYOD'S ALGORITHM / K-MEANS ALGORITHM

INITIALIZATION

$$\overset{\text{ITERATION}}{\downarrow} z_1^0, z_2^0, \dots, z_n^0 \in \{1, \dots, k\}$$

→ UNTIL CONVERGENCE

- COMPUTE MEANS

$$\forall k \quad \mu_k^t = \frac{\sum_{i=1}^n x_i \mathbb{1}(z_i^t = k)}{\sum_{i=1}^n \mathbb{1}(z_i^t = k)}$$

- RE-ASSIGNMENT STEP

$$\forall i \quad z_i^{t+1} = \arg \min_k \|x_i - \mu_k^t\|_2^2$$

mean of
if the current
assignment is
smallest, then
don't reassign

FACT: LLOYD'S ALGORITHM CONVERGES. ← Good news.

- Converged solution may not be "optimal"
 - But produces "reasonable" clusters in practice.
-

QUESTIONS

- CONVERGENCE?
- NATURE OF CLUSTERS?
- INITIALIZATION?
- CHOICE OF K ?