# Introduction

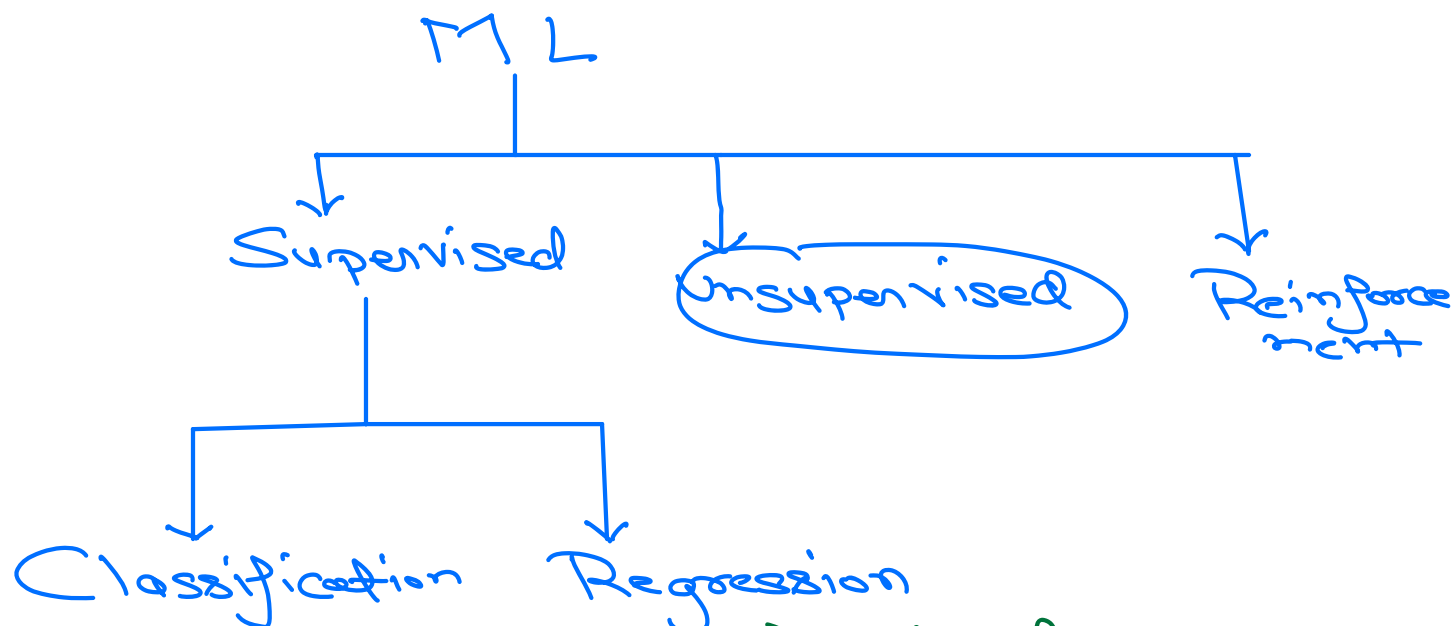## Unsupervised ML

- Kmeans Clustering
- Hierarchical Clustering
- Gaussian Mixture Models
- Outlier/Novelty detection Algos
- PCA / T-SNE / U-map

---

# Topics

- Intro to Unsupervised ML
- Case Study: Customer Segmentation
- Clustering
- Dunn Index
- K-means Intro
- Mathematical formulation of K-means
- Lloyd's Algo
- Implementation of Lloyd's Algo
- Determining K
- Home-work

## Intro to Unsupervised ML

ML

├── Supervised
│   ├── Classification
│   └── Regression
├── Unsupervised
└── Reinforcement

Classification $\Rightarrow \{x_i, y_i \; ; \; x_i \in \mathbb{R}^d$

(Features) (Labels)

$$y_i \in \{0,1,2 \dots n\}$$

Binary $\Rightarrow y_i \in \{0,1\}$

Regression $\Rightarrow \{x_i, y_i \; ; \; x_i \in \mathbb{R}^d$

(Features) (Labels)

$$y_i \in \mathbb{R}$$

* Unsupervised $\Rightarrow$

$$D \rightarrow \{x_i \; ; \; x_i \in \mathbb{R}^d\}$$

→ Create Meaningful Groups of Customers who have Similar Behaviour → Create Clusters

Clustering
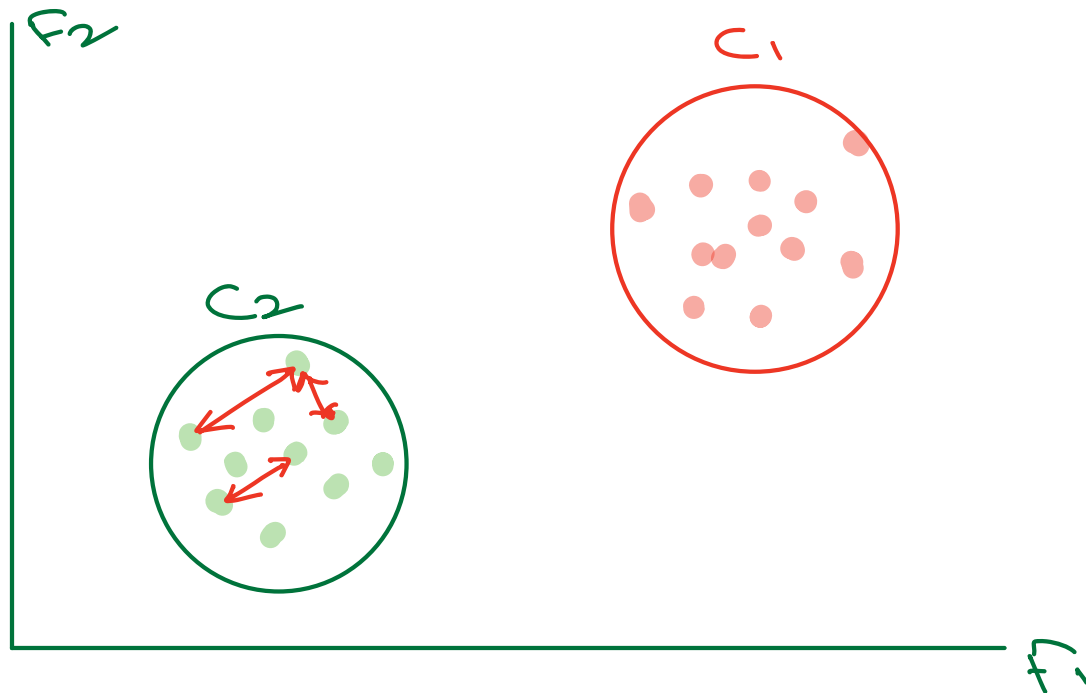


* n-Cluster → Hyperparameter
* Every point will belong to one of the Clusters
* Data points in a Cluster are Close Each Other

# Evaluating Clusters

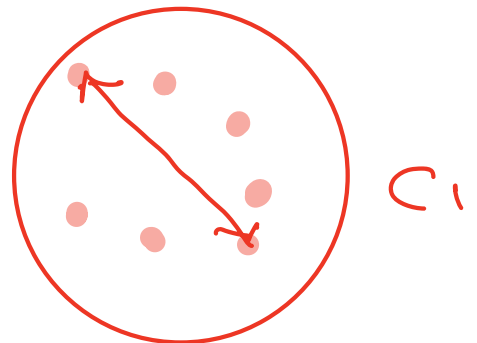① No Ground Truth, Hence we Can't Classification



* Intra Cluster Distance

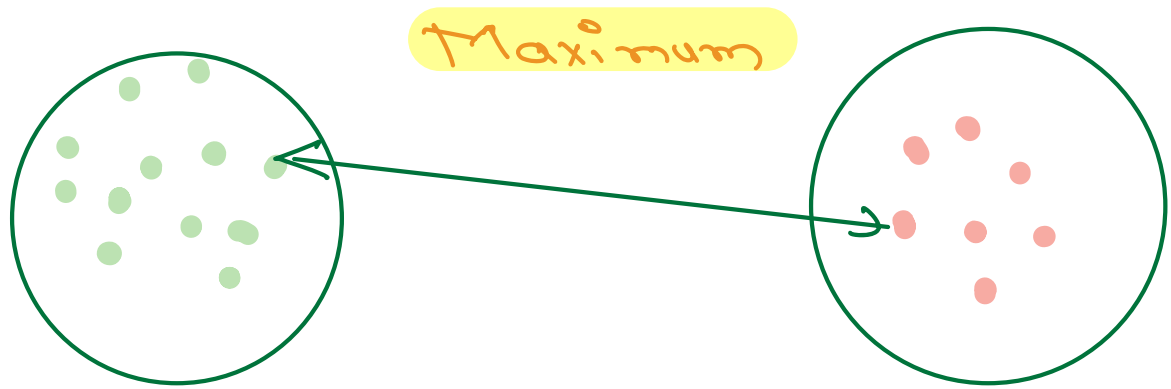① Avg distance among All Pairs

or

① Distance of Farthest Pair

or

① Closest Points Distance



Q : Intra-Cluster : Minimum

* Inter Cluster Distance



Maximum

* Avg of All pairs
* Distance among farthest pair
* Distance among Closest pair

Summary

Intra Cluster Distance Minimum
Inter Cluster Distance Maximum

* Methods to Calculate Distance
  1) Euclidean Distance → Low D
  2) Manhattan Distance → Moderate D
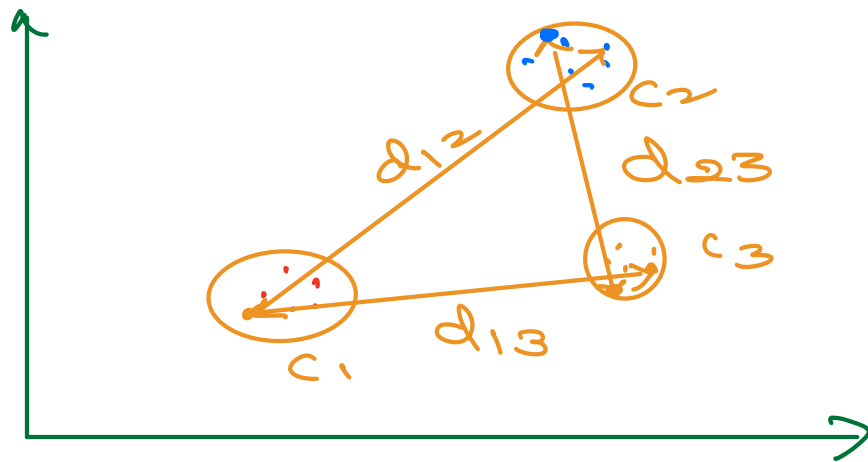  3) Cosine Similarity → High D

$$D_i \ni \frac{\min_{i,j} \text{ distance }(i,j)}{\max_k \text{ distance'}(k)}$$

* distance $(i,j)$ : Inter-Cluster Distance

Farthest point in cluster $i$ and Cluster $j$

* distance'k : Intra Cluster Distance farthest pair



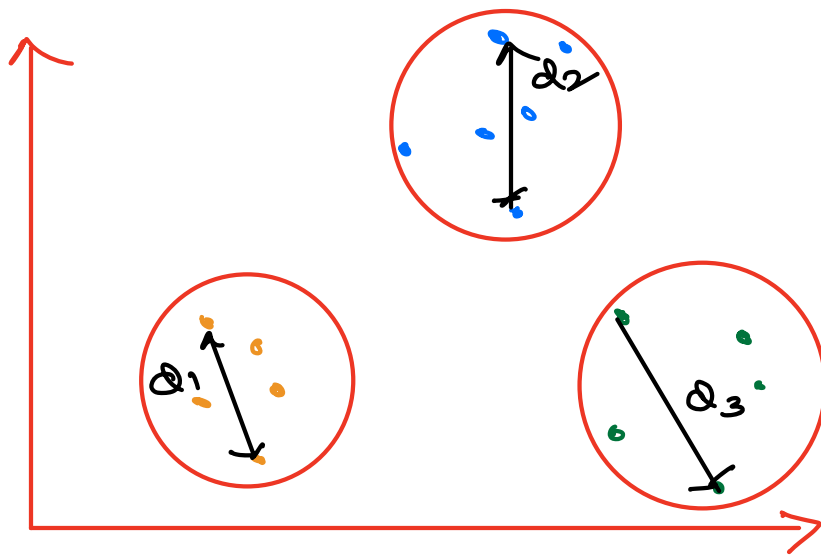$d(i,j) \ni \qquad d_{13}, d_{12}, d_{23}$

$d_{23} \leftarrow$ numerator

$$(d_1, d_2, d_3)$$

$\downarrow$

maximum

$\downarrow$

$d_3$

Dunn-index $\Rightarrow$ $\dfrac{d_{23} \leftarrow \text{High}}{d_3 \leftarrow \text{Low}}$

$\downarrow$

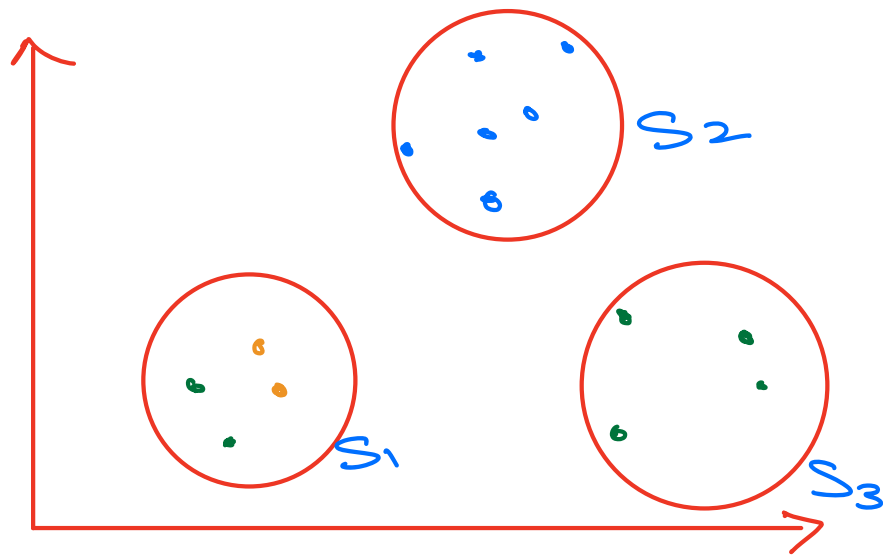- Higher Value of Dunn-index is Better

- Range $\Rightarrow [0, \infty)$

$$\boxed{K\text{-means}}$$

$\downarrow$

$n\text{-Clusters (No of Clusters)}$

$\downarrow$

Similar to K-nn

$$x_i \in \mathbb{R}^d$$

$n$ inputs
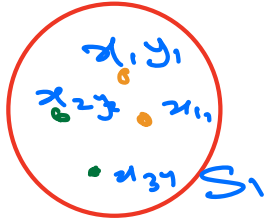
$d$ features



$S_n \ni D$

$$S_1 \cup S_2 \cup S_3 \ni D \ni S_3$$

$$S_1 \cap S_2 \cap S_3 \ni \emptyset$$

$$S_1 \cap S_2 \ni \emptyset \ni S_2 \cap S_3$$

* K-means is a Centroid Based algorithm where

$$C_1 \ni \frac{x_1 + x_2 + x_3 \cdots x_n \in S_1}{len(S_1)}$$

$$C_1 \ni \frac{x_1 + x_2 + x_3 + x}{len(S_1)} \quad x_i^2$$

$$\frac{y_1 + y_2 + y_3 + y_4}{len(S_1)} \quad y_i^2$$



$$C_1 \Rightarrow x_i^2, y_i^2$$

$$\boxed{C_i \ni \frac{1}{|S_i|} \sum_{x_j \in S_i} x_j}$$

$j^{th}$ vector belonging to $i^{th}$ Set

Size / Cardinality of $i^{th}$ Set

# Objective Function

C₁, C₂, C₃ --- Cₙ

→ max ( Inter Cluster Distance )

and

→ min ( Intra Cluster Distance )

Difficult to Solve with Traditional Optimization

* Approximate Algo:

**Lloyd's Algorithm**

→ Simple and Easy to Calculate

→ Doesn't Guarantee the Best Solution