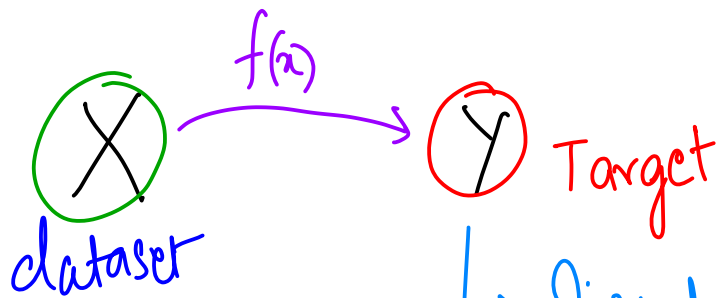


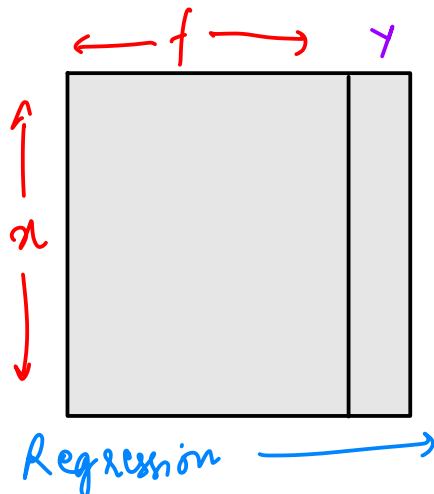
K MEANS

- ① Unsupervised learning
- ② Higher dimensional analyse/visualisation
- ③ GMM
- ④ Anomaly detection



↳ Discrete value [Classification]

↳ Real Continuous value [Regression]

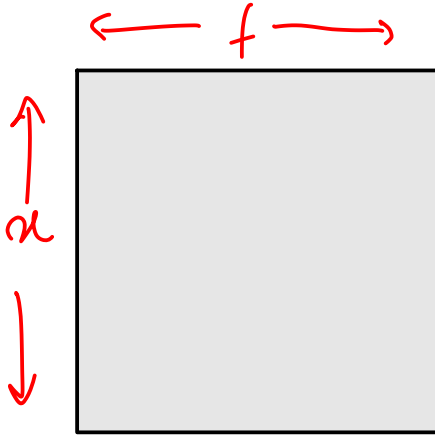


Supervised

↳ Classification

$$\mathcal{D} = [(\hat{x}_i, \hat{y}_i), x_i \in \mathbb{R}^d, y_i \in \{0, 1\}]$$

$$\mathcal{D} = [(\hat{x}_i, \hat{y}_i), x_i \in \mathbb{R}^d, y_i \in \mathbb{R}^r]$$

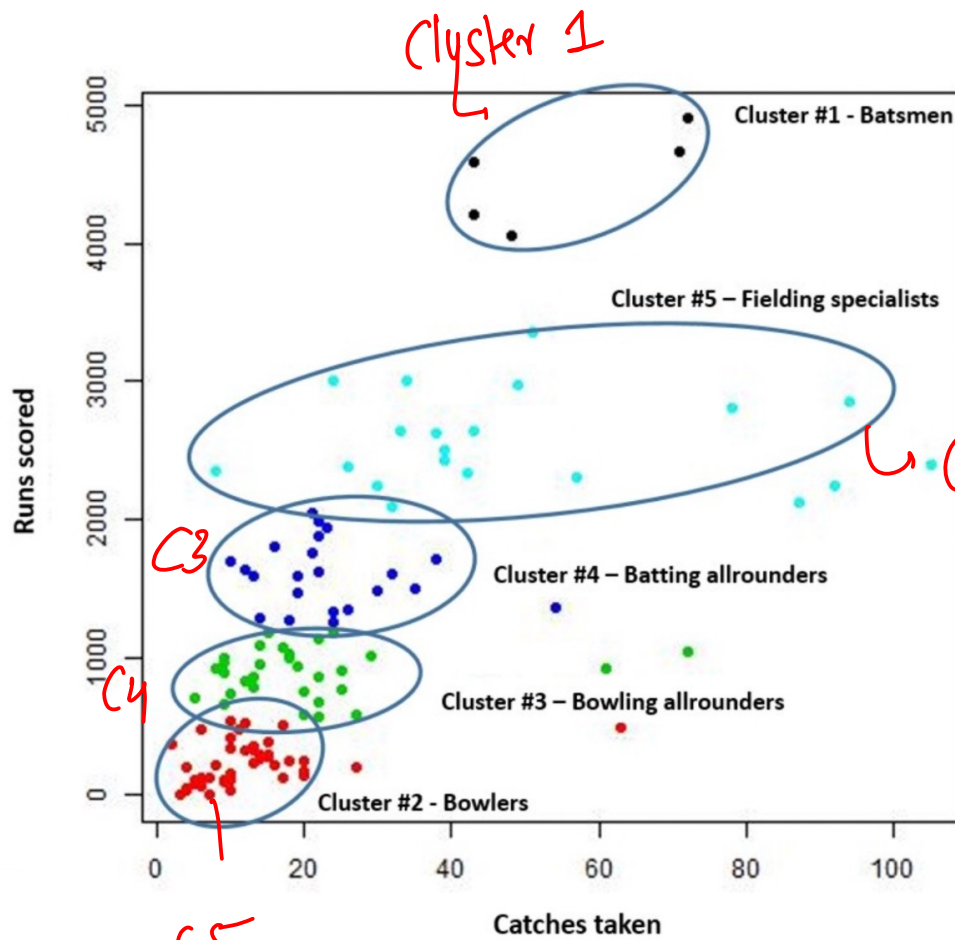


Unsupervised

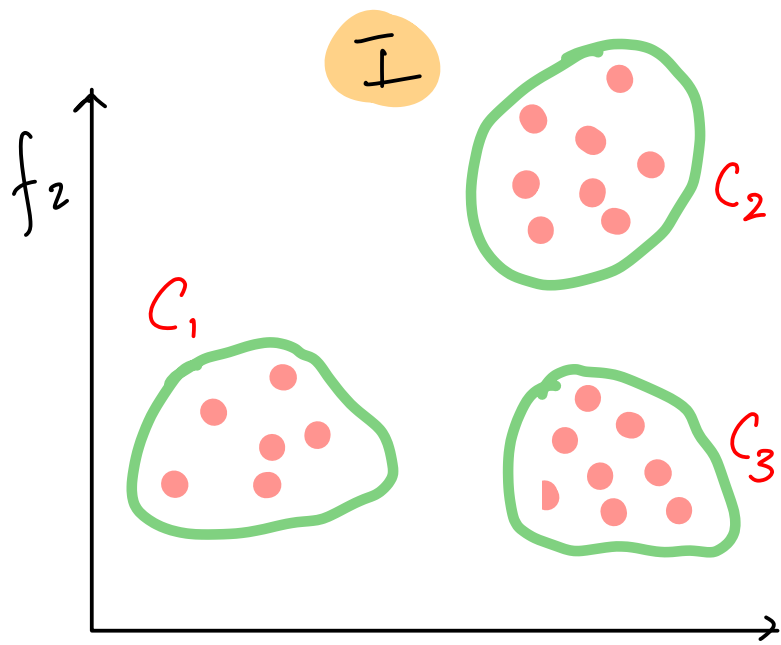
Amazon wants to segment their customers so that they can provide relevant and similar items to their customers Which algorithm do you think can be used here?

8 users have participated

A	Regression	0%
B	Classification	13%
✓ C	Clustering	87%
D	None of the above	0%

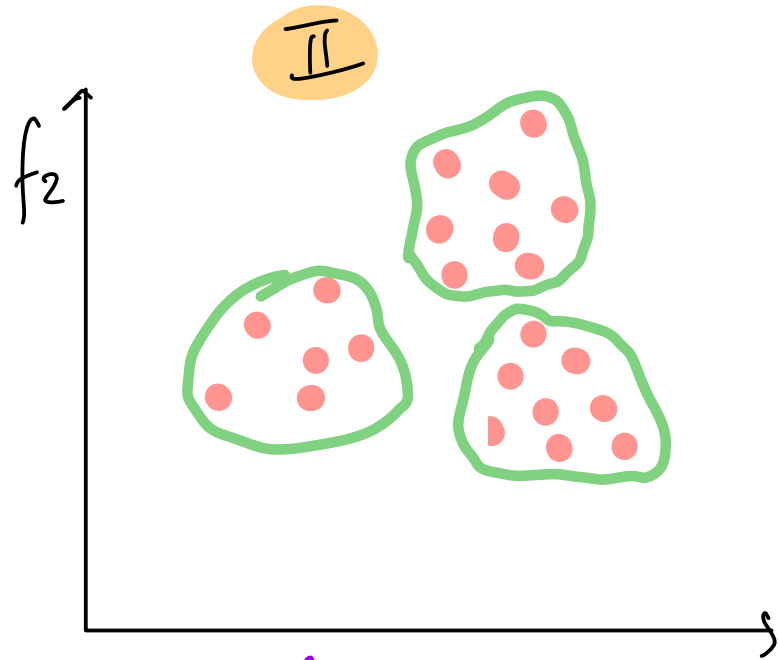


Clustering



① Avg distance b/w points within grps should be less.

② Avg distance b/w points in one group to another group should be more.



→ diff. algorithms

→ Good or Bad.

"Cluster output should make biz sense."

$C_1, C_2, f_1, f_2$

① Inter cluster distance

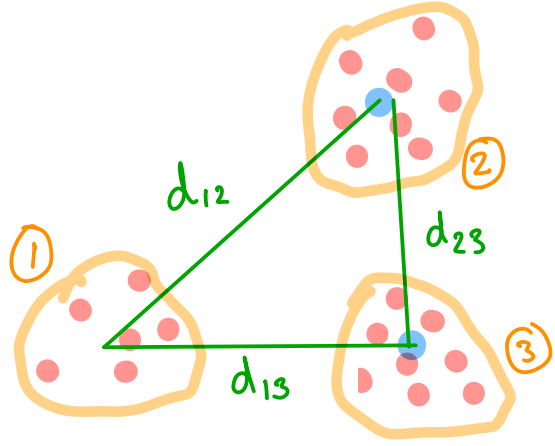
b/w different clusters

② Intra cluster distance

within cluster

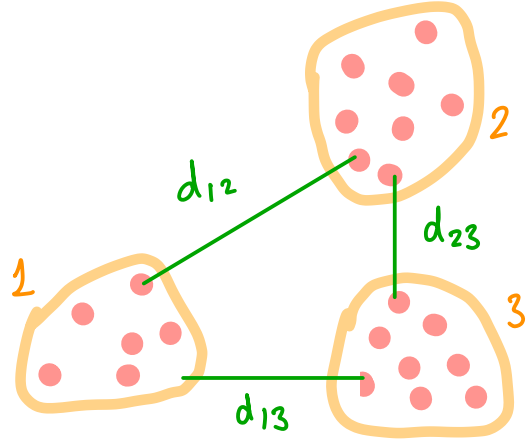
# Inter cluster distance

I



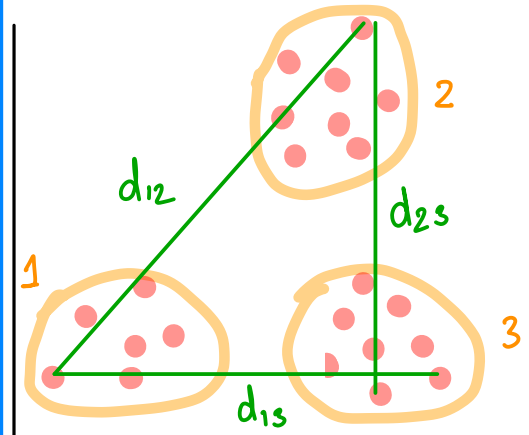
Distances b/w  
average values of  
cluster point.

II



distance b/w closest  
point from the  
cluster.

III

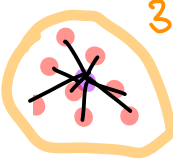
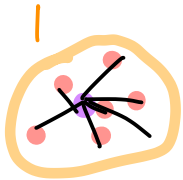
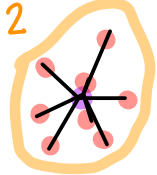


Distance b/w  
farthest point from  
cluster.



# Intra Cluster Distances

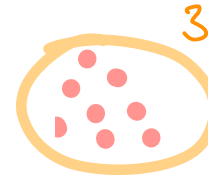
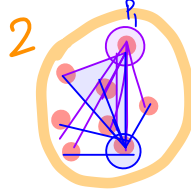
$$\frac{1}{n} \sum d_i$$



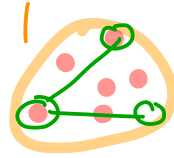
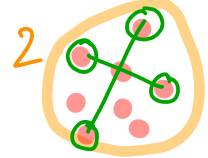
Avg distance w.r t  
some central point

Avg

$P_1 \rightarrow$  Every other  
point in  
cluster



Avg distance b/w  
all pairs of points



distance b/w farthest  
points within cluster.

$$\mathcal{D} = \left\{ (x_i)_{i=1}^n \cdot x_i \in \mathbb{R}^d \right\}$$

inter cluster distance ( $\uparrow$  high)

intra cluster distance ( $\downarrow$  low)

Ways of finding distances ??

\* Euclidean  $\Rightarrow$  low - dims

\* Manhattan  $\Rightarrow$  low - medium dims

\* Cosine <sub>text</sub>  $\Rightarrow$  high dims.

If only one Inter or Intra distance is given, Can we judge how good or bad the clusters are?

It will not work.

We need both.

---

What metric can be used?

\* Biz Case.

---

Raymond. → S, M, L, XL, XXL

would you create 100  
Clusters?

Metric  $\Rightarrow$  **DUNN-index**  $\Rightarrow$  metric used to evaluate clustering algorithm.

$\hookrightarrow$  OBJECTIVE: identify cluster with small variance  
b/w cluster members and are well separated.

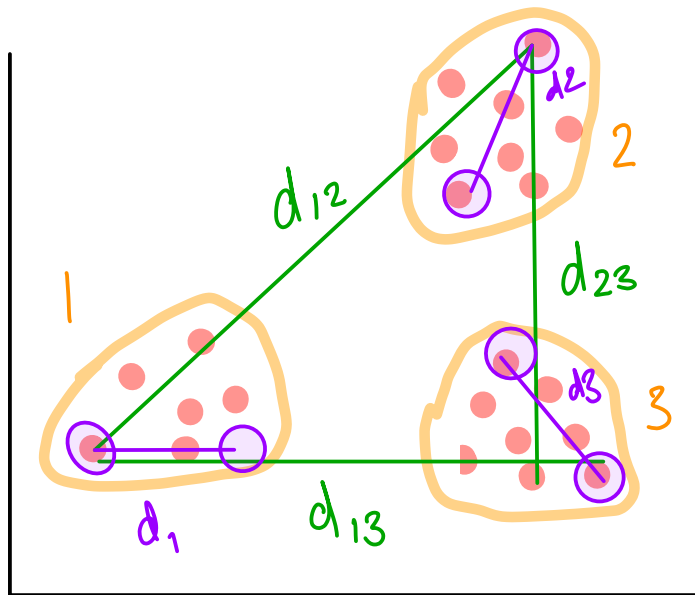
$$D = \frac{\min_{i,j} \text{distance}(i,j)}{\max_k \text{distance}(k)}$$

$D(i, j) \Rightarrow$  Inter cluster distance

$\hookrightarrow$  Distance b/w farthest point of the cluster  $C_1$  &  $C_2$ .

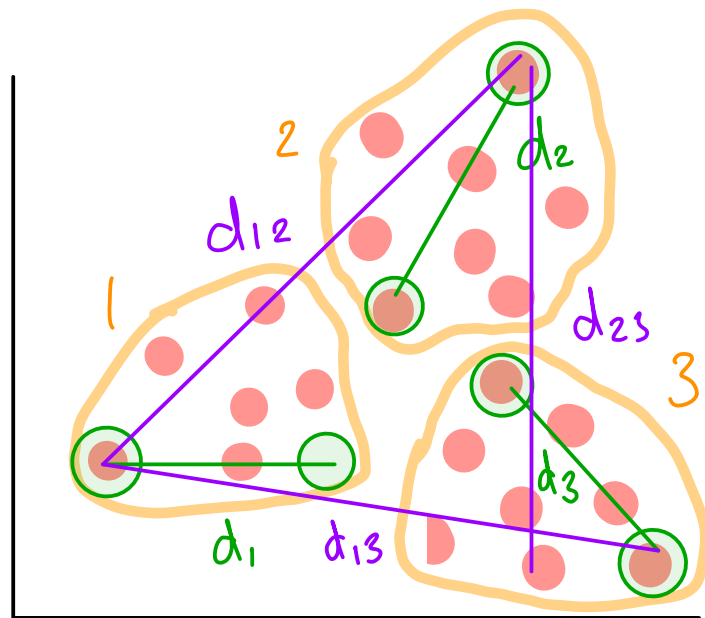
$D(k) \Rightarrow$  Intra cluster distance

$\hookrightarrow$  Distance b/w farthest points within cluster  $k^{th}$



$$d(i,j) \mid \min(d_{12}, d_{13}, \underline{d_{23}})$$

$$d(k) \mid \max(d_1, d_2, \underline{d_3})$$



$$d'(k) \mid \max(d_1, d_2, \underline{d_3})$$

$$d(i,j) \mid \min(d_{12}, d_{13}, \underline{d_{23}})$$

k clusters

$$\underline{QVNN \text{ Index}} = Q = \frac{\uparrow \min_{i,j} d(i,j)}{\downarrow \max_k d(k)} \Rightarrow \begin{array}{l} \text{Inter cluster} \\ \text{Intra cluster} \end{array}$$

$$Q_1 > Q_2$$

↑  
good

$$\frac{\text{Inter } 10}{\text{Intra } 5}$$

①

$$\frac{\text{Inter } 10}{\text{Intra } 10}$$

②

**In the Dunn Index formula, what does "distance(i, j)" represent?**

8 users have participated

- |   |   |     |
|---|---|-----|
| A | The average distance between all points in cluster i and all points in cluster j. | 13% |
| B | The distance between the centroids of cluster i and cluster j.                    | 0%  |
| C | The distance between the closest points in cluster i and cluster j.               | 12% |
| D | The distance between the farthest points in cluster i and cluster j.              | 75% |



**We have two clustering algorithms, and if we have to choose one of them, the algorithms chosen should have:**

8 users have participated

- |   |   |                           |      |
|---|---|---------------------------|------|
| ✓ | A | higher dunn-index         | 100% |
|   | B | lower dunn-index          | 0%   |
|   | C | Dunn-index doesn't matter | 0%   |