

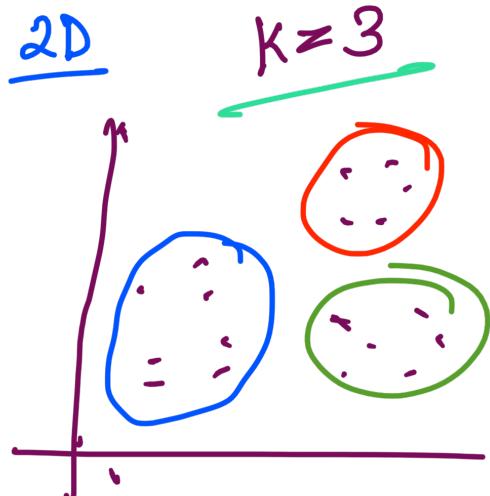
Day 8

1. K-means \leftarrow UnSupervised ✓
 2. KNN \leftarrow Supervised. ✓
 3. Overfitting
 4. Underfitting
- } concepts

1. K-means

$K \leftarrow \# \text{ of clusters}$

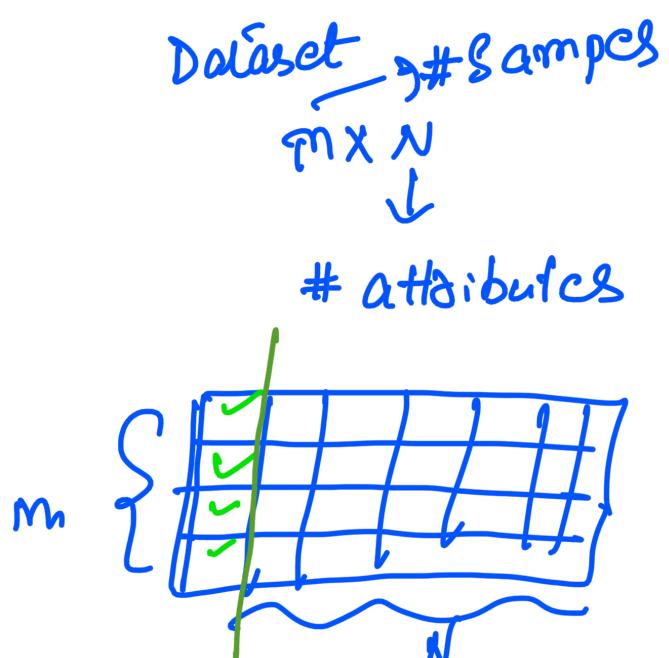
$K = 2$



Example

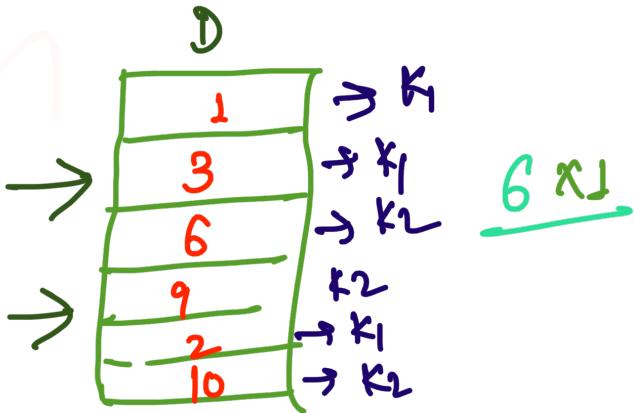
ID

$m \times 1$



$m = \# \text{ samples}$

$s \leftarrow \text{feature}$



1. $K = 2$ ←
2. $\underline{K_1} = 3 \quad \underline{K_2} = 9$
- 3.

		K_1	K_2
D_1	1	$(1-3)^2$ $= 4$	$(1-9)^2$ $= 64$
D_2	3	x	x
D_3	6	$(6-3)^2$ $= 9$	$(6-9)^2$ $= 9$
D_4	9	x	x
D_5	2	$(2-3)^2$ $= 1$	$(2-9)^2$
D_6	10	$(10-3)^2$	$(10-9)^2$

$$\begin{aligned} \underline{\underline{K_1}} &= \{ \underline{1}, \underline{3}, \underline{2} \} \\ \underline{\underline{K_2}} &= \{ 6, 9, 10 \} \end{aligned}$$

8th iteration

1. Centroid
2. Mean

$$\begin{aligned} K_1 \text{ Centroid} &= 2 \\ \text{Mean} &= 2 \end{aligned}$$

K_2

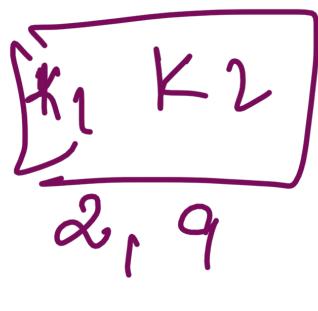
$$K_1 = \{ 1, 3, \underline{2} \}$$

$$\frac{1+5}{2} \quad K_1 = 3 \quad K_2 = 9$$

$$K_2 = \{ 6, 9, 10 \}$$

$$\frac{2+9}{2} \quad K_1 = 2 \quad K_2 = 9$$

$$\begin{aligned} K_1 &= \{ 1, 2, 3 \} \\ K_2 &= \{ 6, 9, 10 \} \end{aligned}$$



$$D_1 = \underline{(1-2)^2} \quad (2-9)^2$$

$$D_2 \quad \underline{(3-2)^2} \quad (3-9)^2$$

Iteration: S₁₀₀

Centroid, K , is not changing
Cluster assignment is not changed

Steps(k-means):

1. Select the number of clusters (K)
2. Select the no. of data points = K
3. Assign each datapoint to one of the ' K '

3.1. Distance

✓ - L₂ norm $(a-b)^2$ } Research area

- L₁ norm $a-b$

- $|a-b|$

- 3.2 Find the min^m and assign it to that ' K '

4. Find the cluster centroid.

5. Repeat Step 3 until convergence.

Convergence

- k' cluster centroid is not changed
- assignment of data points are not changed.

Dataset

$m \times 2$

2 jealous

Cluster assignment

$$k_1 = [x_{N,1}, x_{N,2}]$$

$$k_2 = [x'_{N',1}, x'_{N',2}]$$

	$x_{1,1}$	$x_{1,2}$
m	$x_{2,1}$	$x_{2,2}$
:		
	$x_{m,1}$	$x_{m,2}$

Euclidean distance - (ED)

1st point: $\sqrt{(x_{1,1} - x_{N,1})^2 + (x_{1,2} - x_{N,2})^2}$

1st point, $(1, 2)$

k_1 $(3, 4)$

$N < m$
$N' < m$
$N \neq N'$
$N, N' \subseteq m$

$$(ED) = \sqrt{(1-3)^2 + (2-4)^2} = \sqrt{8}$$

K-Means

- Popular
- Social media Analysis.
- Healthcare $s_1 s_2 s_3 \rightarrow c_1 \rightarrow d_1$
- Segmentation
- market analysis

Applications

S = Symptom

D = Disease

Major Research Problem $(k) ?$

Unsupervised → It will not have class label.
How to decide the value $\underline{\underline{k}}$?

Supervised

Classification
class labels
discrete

Regression
class labels
continuous.

Unsupervised

→ Clustering
→ K means.
(pattern)

Association.
→ grouping.
→ Hierarchical

- Advantage of using an association:
- The value of K need not to be given.

Hierarchical Clustering -

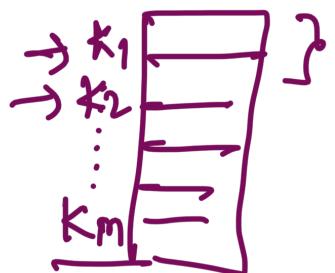
- ✓ → Agglomerative \times (A)
→ Divisive. (top down approach)

(A) - Consider the entire dataset
 $\# \text{datapoint} = \# K$

bottom-up
end up with 1 cluster

initial = K centroids

Step 1 : Assign each data point $m \times 1$ as cluster centroid

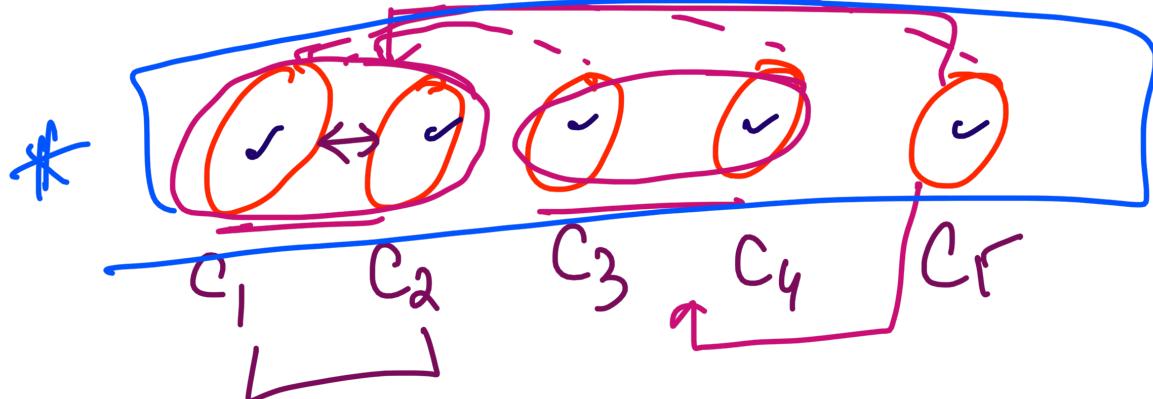
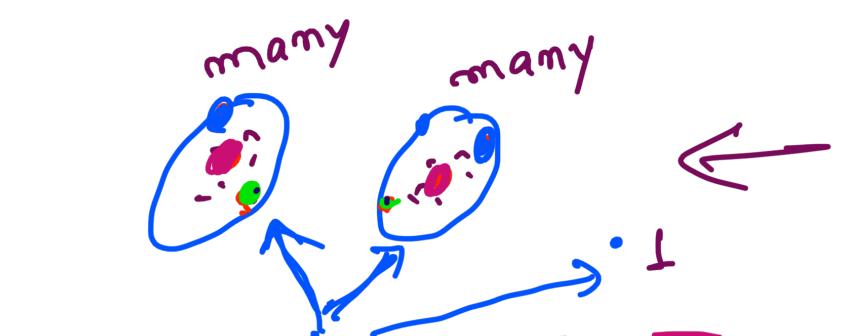


Step 2: $(k, k_5) \rightarrow (k-1)$

↳ centroid

Step 3: $(k_2, k_6) \rightarrow (k-2)$

↳ remaining.



Step 4: Find the distance between clusters

→ Single linkage ✓

→ shortest distance between two clusters

→ Complete linkage: ✓

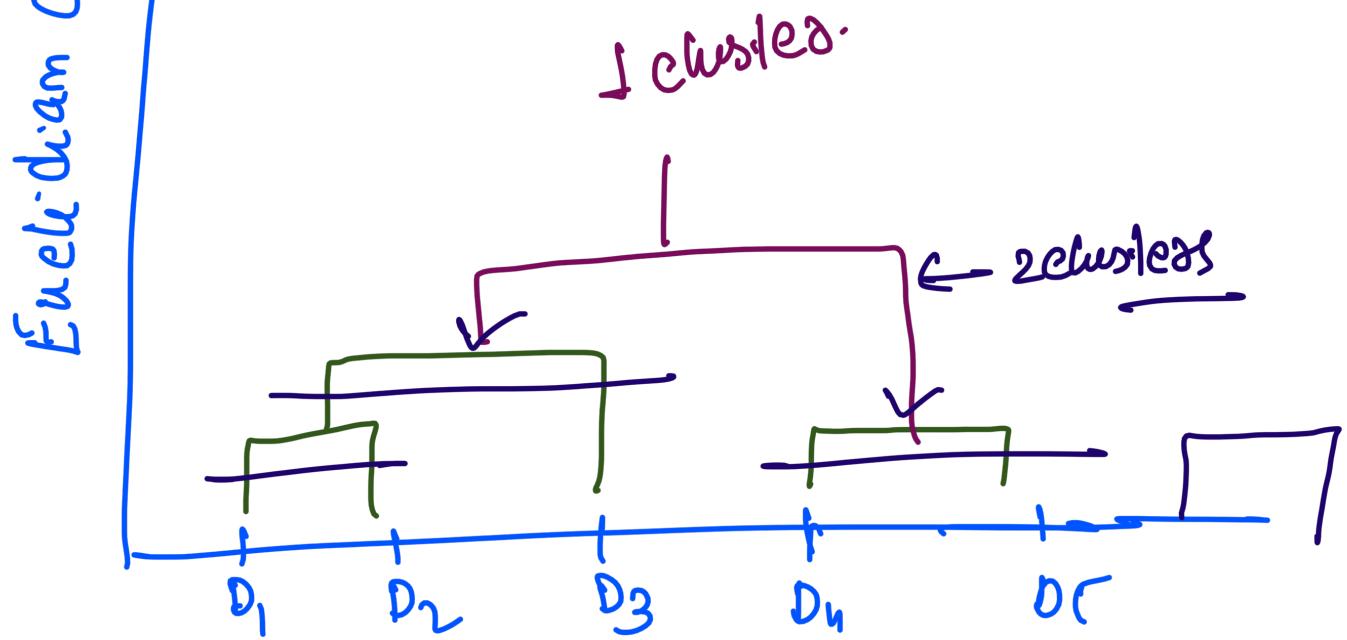
→ Farthest distance

→ Average/Centroid linkage ✓

* Euclidean distance.

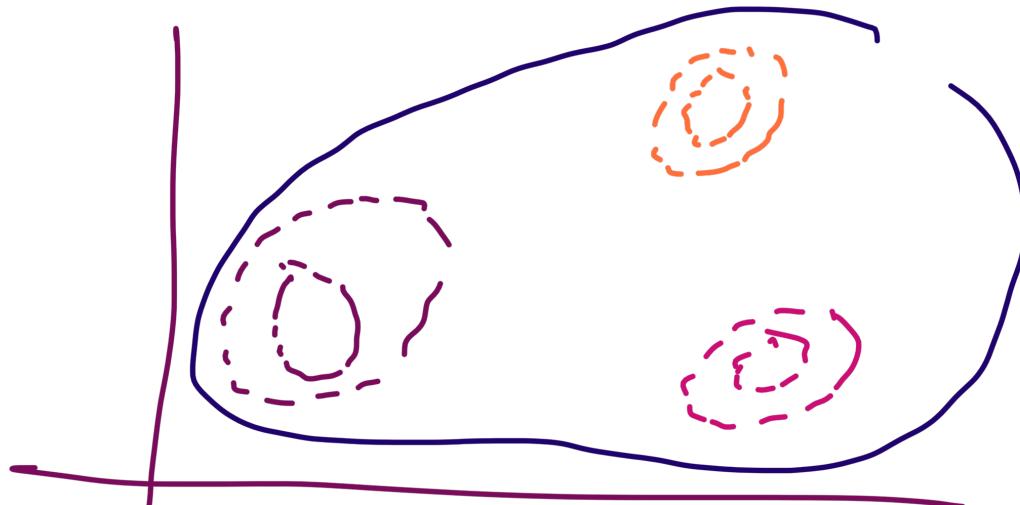
Representation

Dendrogram



Data points

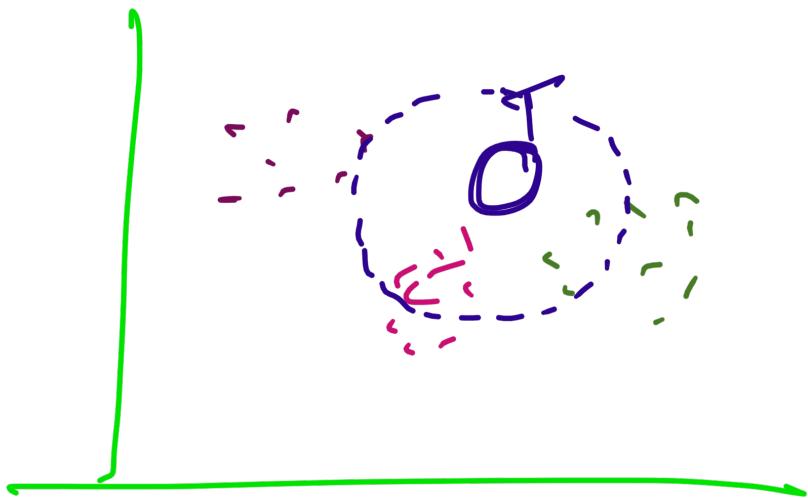
→ Cluster Similarity.



K NN

Nearest
Neighbour

Supervised



$K \leftarrow$ datapoints

$K = 5 \leftarrow$
T if will take L_2 distance
with 5 datapoint which falls

nearer to c'

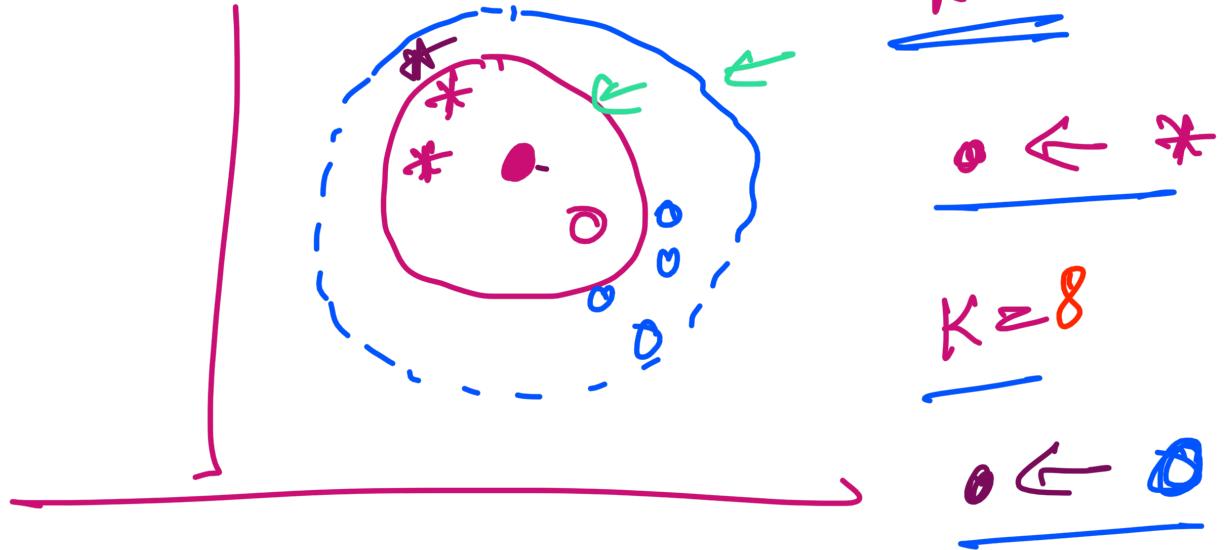
$$(T - D_1)^2 = a$$

T
 $D_1 \rightarrow @$
 $D_2 \rightarrow B$
 $\} \rightarrow C_1$

$D_3 \rightarrow C$

$D_4 \rightarrow D$

$D_5 \rightarrow E$



Lazy algorithm

1. Overfitting
2. Underfitting.

Training (good)

Dataset \rightarrow Algorithm \rightarrow Model.

Learning

Testing

Inference if not getting good results

Overshooting

2. Underfitting

Training (bad)
Testing (bad)

model is
not
leaving
any point

