

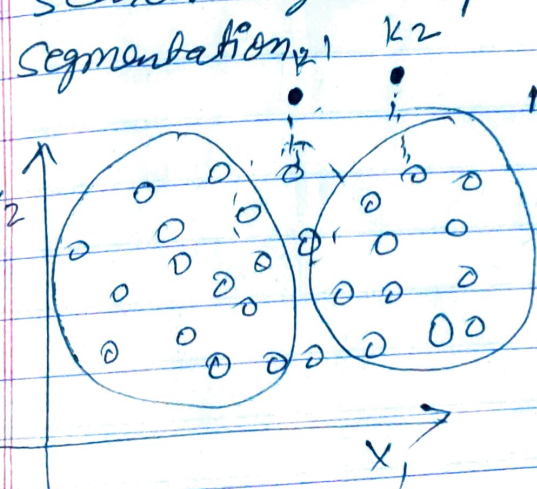
Clustering

In clustering we divide datasets in certain groups & train that dataset.
(based on similarity)

① K-Mean ② Hierarchical clustering ③ DBSCAN

K-mean

Appⁿ - Cost segmentation, Data Analysis, Anomaly detection
Search engine optimization, speech/image segmentation

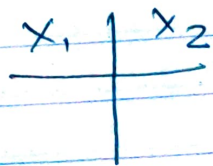


WCSS

Within cluster
Sumation of Square.

We find K using
WCSS.

K=2 drop take 2 random points.



Associate each object with
dataset centroid shifts such
that clusters are formed.

after this if point D_i comes it will calculate
distance to centroid of each cluster
whichever is less than Cluster D_i is from.

How to take K value?

$\left\{ \begin{array}{l} \text{K mean} \\ \text{K mean ++} \\ \text{mini batch K mean} \end{array} \right.$

$WCSS = \text{Within Cluster Summation of Squares}$
 n clusters.

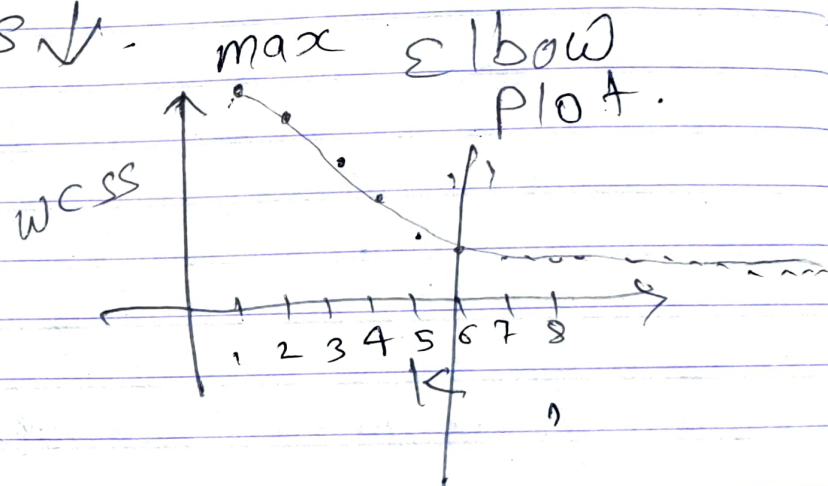
$$WCSS = \sum \sum (C_i - P_c)^2$$

\rightarrow (inertia) \Rightarrow wcss

$$K = n \Rightarrow WCSS = 0$$

$$K = 1 \Rightarrow WCSS = \text{max}$$

$K \uparrow WCSS \downarrow$



Drawback: Allocation
 K mean of centroid.

(To overcome)
 K mean ++

dispersion

1. Always take ^{distance b/w two or} centroid (minimum).

2. Triangle rule.

3. Provide dataset into batches (Mini batch K mean)

* By default `kmeans++` is used in code.

Hierarchical

Agglomerative. - n data points
↳ n clusters.
Divisive
↳ ex - `kmean`.

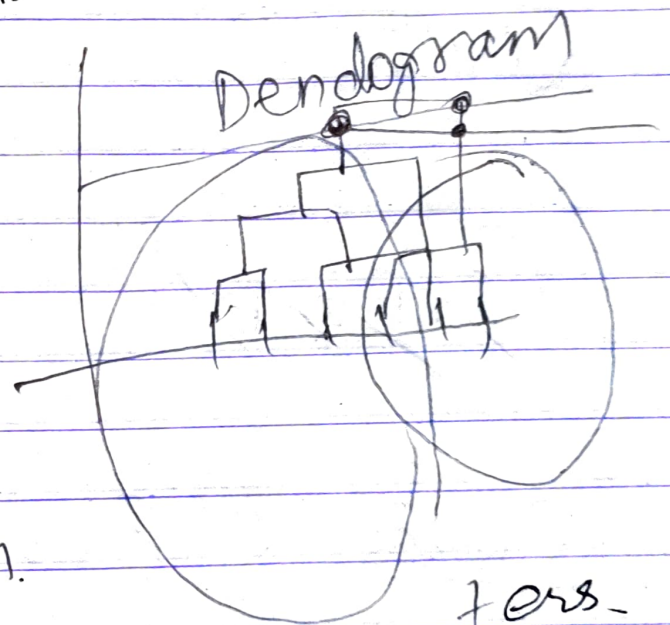
↳ Hierarchical clustering

Bottoms up approach.

more distance
vertical

⇒ Take that
line
for cluster
formation.

(mean of it)



2 clusters.
on dendrogram
approach.



DBSCAN (Best Clustering).

Density Based Spatial Clustering
of Appⁿ with Noise.

Epsilon - radius of circle (for a data point)

Min Points - min no. of datapoints inside circle
to make cluster.

Core points - point taken to form cluster.

Border points - point which is part of some
core point cluster but not in

Noise - not inside any cluster or
(outliers). can't form any cluster.

~~EPS is bad~~

Drawback: ① EPS very high, EIPS very low

② Min point very high.

Validation of clustering.

Clusters are supposed to be created in such a way such that it must be having very high inter classing similarity & very low intraclass similarity.

Rand index

$$= \frac{\text{Total Agree}}{\text{Total Agree} + \text{Total Disagree}}$$

Order status | @lenovo
com.

Jaccard coefficient = $\left(\frac{SS}{SS + SD + DS} \right)$ $SS \rightarrow$ ground truth.

Entropy

Purity

$SS \rightarrow 1$
 \Rightarrow for a point agreeing to ground truth.

Silhouette coefficient.

cohesion

separation.

$DD \rightarrow 0$
 \Rightarrow Point don't belong to ground truth as well as algo.

(SD →

Cluster Prediction on

			S	D
			1	0
S	GT	1	SS	SD
D	0	0	DS	DD

$$\text{Rand Index} = \frac{SS + DD}{SS + DS + SD + DD}$$

$$\text{Entropy} = -P_i \log P_i$$

$$= -\sum P_i \log(P_i)$$

Purity = Total % age of dataset correctly placed in cluster.

= Purity of Cluster i

$$P_i = \max(P_{ij})$$

→ i th class in j th cluster

$$P(\text{cluster})^{\text{whole}} = \sum \frac{m_i}{n} p_i$$

Not in script learn.

Silhouette score $S(x) = \frac{b(x) - a(x)}{\max\{a(x), b(x)\}}$

$a(x) \rightarrow$ average distance of dataset from all other pts in same cluster

should be high.

$b(x) \rightarrow$ avg. dist. of x from all points in another cluster

$$Sc = \frac{1}{N} \sum S(x)$$

0

$a(x) > b(x)$
Wrong