

Data Science & AI



Machine Learning



Supervised Learning

Lecture No.- 06



By- Krish Naik Sir

Recap of Previous Lecture



Topic

Topic

Topic

Topic

Topic

Feature Transformation
naive baye's

Topics to be Covered



Topic

Random forest classifier

Topic

KNN classifier & regressor

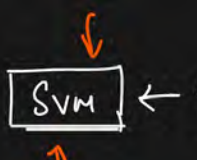
Topic

K-means

Topic

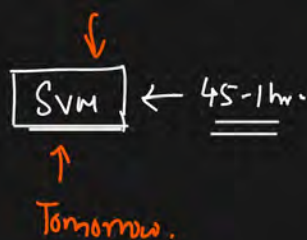
Topic

Machine Learning

- ① Random Forest Classifier And Regressor [Bagging]
 - ② KNN Classifier And Regressor [KNN].
 - ③ K Means ✓
 - ④ Hierarchical Mean ✓
 - ⑤ DBSCAN ✓
 - ⑥ Silhouette Scoring. ✓
- 

SVM

↑
Tomorrow.



размеры



Sequentiell y

Assignment

$[ROC - AUC]$

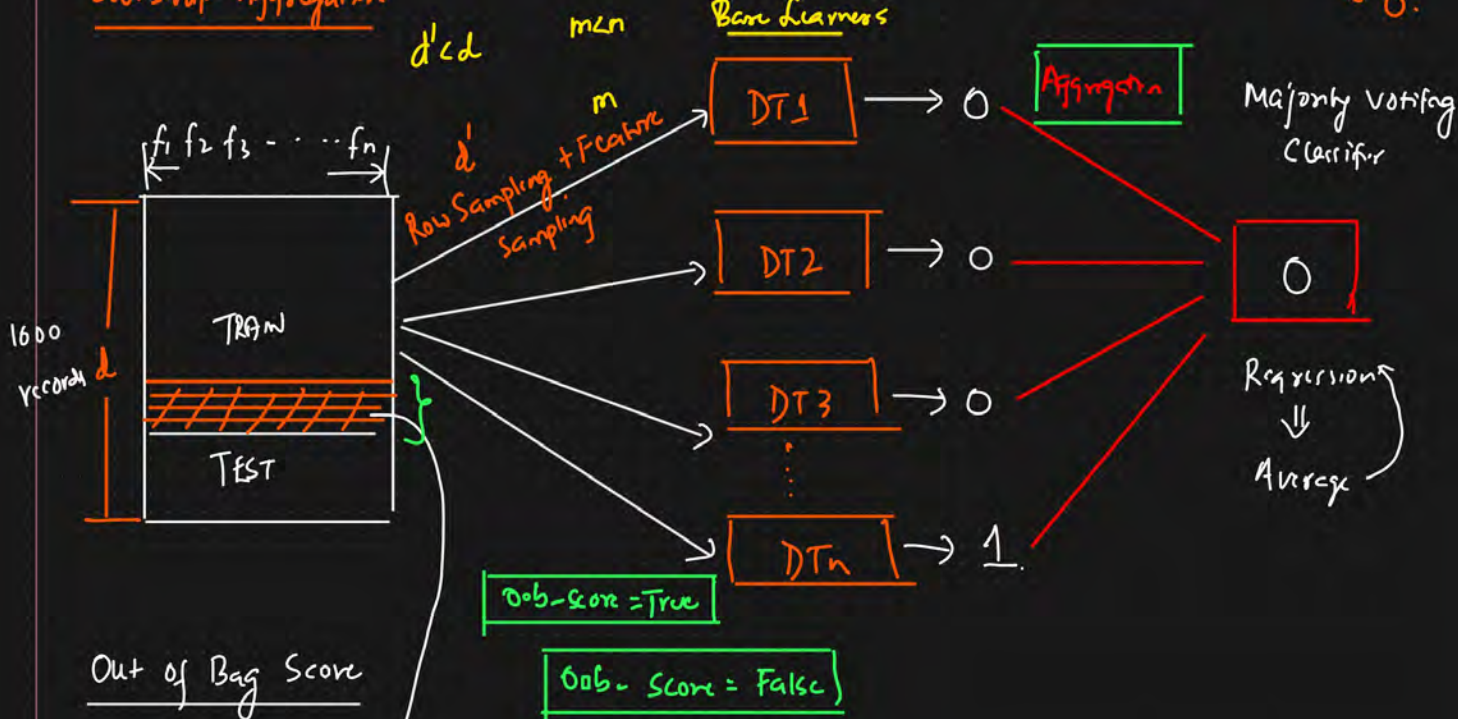
Curve

- ## ① Random Forest Classifier And Regressor [Bagging]

Bagging

Decision Tree \rightarrow Overfitting \rightarrow Low Bias } Low Bias
High Variance } Low Variance

Bootstrap Aggregation



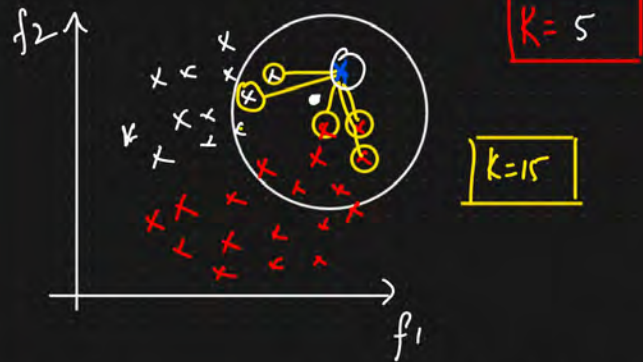
Out of Bag Data \Rightarrow Validation DATA \rightarrow Cross Validation
 \downarrow
 validation accuracy \rightarrow {OOB-score}

② KNN Classifier And Regressor [K Nearest Neighbor]

① Classifier

② Regressor

① Classification



f_1 f_2 O/p (Yes/No)

① We have to initialize the k value
 $K=1, 2, 3, 4, \dots \Rightarrow$ hyperparameter

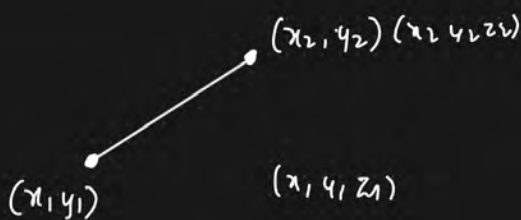
② Find the K Nearest Neighbor from the new Test data

③ Maximum no. of points is in n th category

Test data \rightarrow Red category point.

In Regression we consider the Average of all data points.

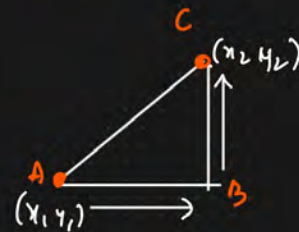
① Eucledian Distance



$$\text{Distance} = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

$$\text{Distance} = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2 + (z_2 - z_1)^2}$$

② Manhattan Distance



$$\Rightarrow AB + BC$$

$$|x_2 - x_1| + |y_2 - y_1|$$

Application :-

Air Traffic Control

Eucledian Distance

Manhattan



Manhattan Distance

A.

B

C

Categorical \rightarrow Numerical

$f_1 \rightarrow f_1'$

Yes

No

Yes.

One hot Encoding

Yes

1

0

1

No

0

1

0

Location

New York

Florida

Texas.

New Y

1

0

0

F

0

1

0

T

0

0

1

Variants of KNN

Disadvantage:

Time Complexity is more to search Nearest Neighbour.

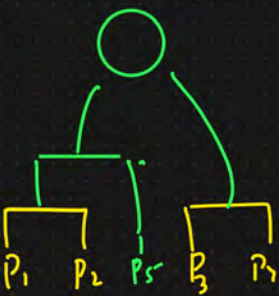
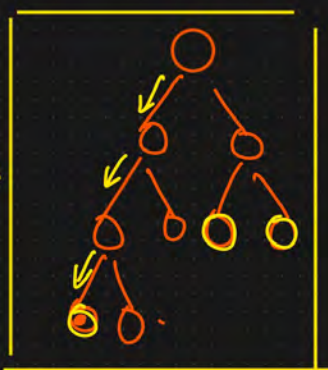
Brut search

$O(n)$

$k=3$

$O(\log n)$

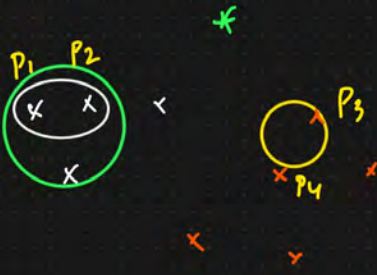
\Rightarrow



$O(\log n)$

Variants of KNN

- ① KD Tree \Rightarrow K Dimension Tree
- ② Ball tree \leftarrow





Feature Scaling ?? \Rightarrow Distance Formula

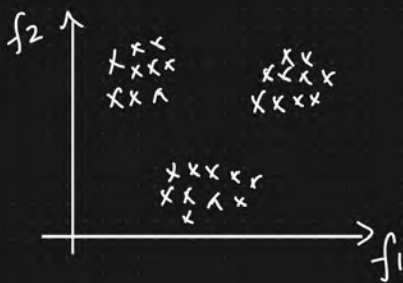


- ① Impact of Outliers?
- ② Feature Scaling?

③ Unsupervised Machine Learning

- ① K Means
- ② Hierarchical Clustering
- ③ DBSCAN Clustering.

① K Means Clustering



\Rightarrow



K=3 K=4

Steps

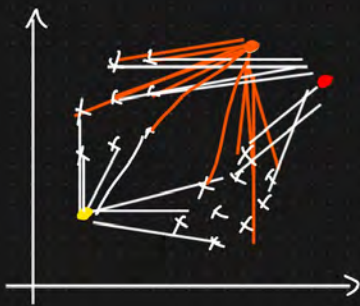
K=?

- ① Initialize some K Centroids.
- ② Points nearest to the centroid will be grouped.
- ③ Move the centroid \rightarrow Mean

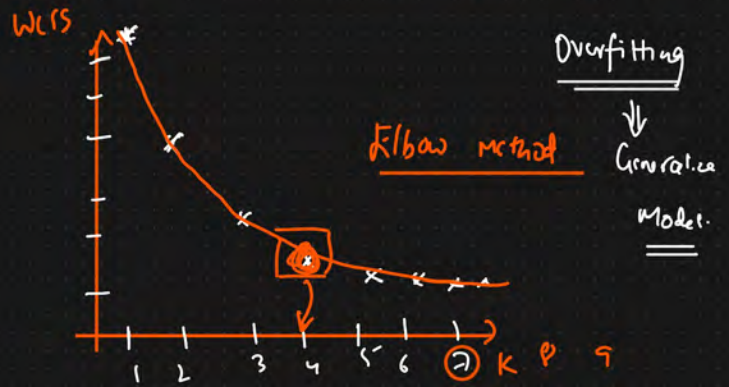
How do we select the K value

Initialize K=1 to 20

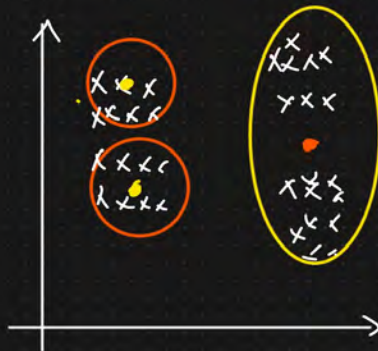
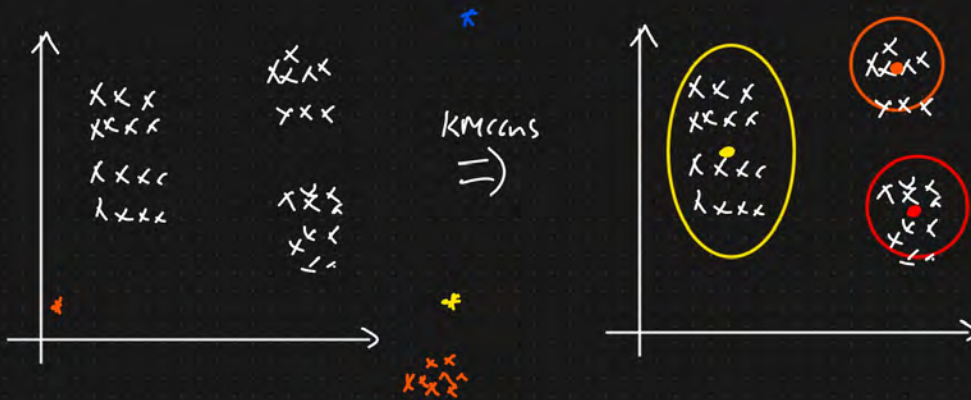
Within Cluster Sum of square



$$WCSS = \sum_{i=1}^K \left(\text{Distance between points to the nearest centroid} \right)^2$$



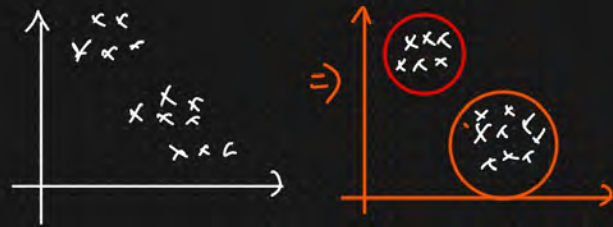
Random Initialization Trap (K Means++)



④ Hierarchical Clustering

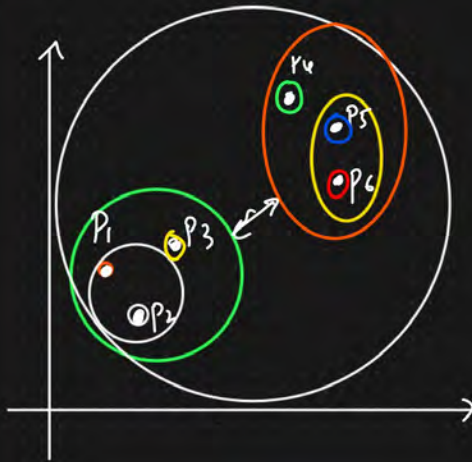
Two Types

- ① Agglomerative
 - ② Divisive
- } \Rightarrow Geometric Intuition

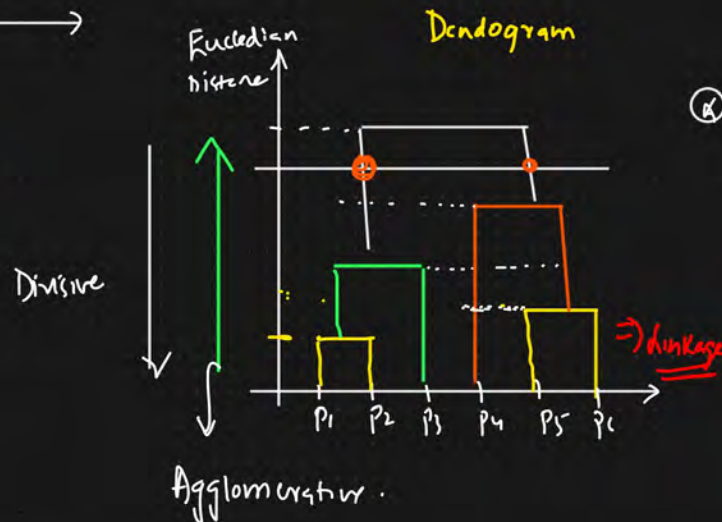


Steps

- ① For each point initially we consider it as a separate cluster
- ② Find the nearest point and create a new cluster



2 centroids



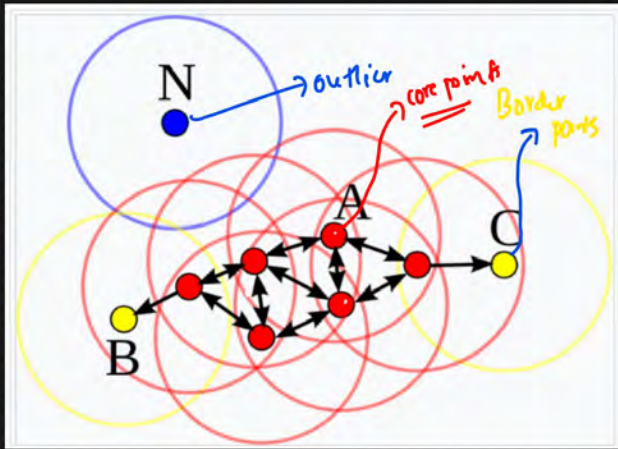
- ④ Select the longest vertical such that no horizontal line passes through it

K Means Vs Hierarchical Clustering

- ① Data size
 - \rightarrow Huge \Rightarrow K Means
 - \rightarrow Small \Rightarrow Hierarchical clustering

- ② Types of Data
 - \rightarrow Numerical \Rightarrow K Means
 - \rightarrow Variety of data \Rightarrow Hierarchical Means

① Density Based Spatial Clustering . [DBSCAN clustering]



- Core point
 - border points
 - Noise / outliers .
- } Non linear Clustering

Hyperparameter

- ① $minpts = 4$
- ② $\epsilon = \text{radius}$.

① Core point

No. of points within the ϵ should be $\geq \underline{\underline{minpts}}$

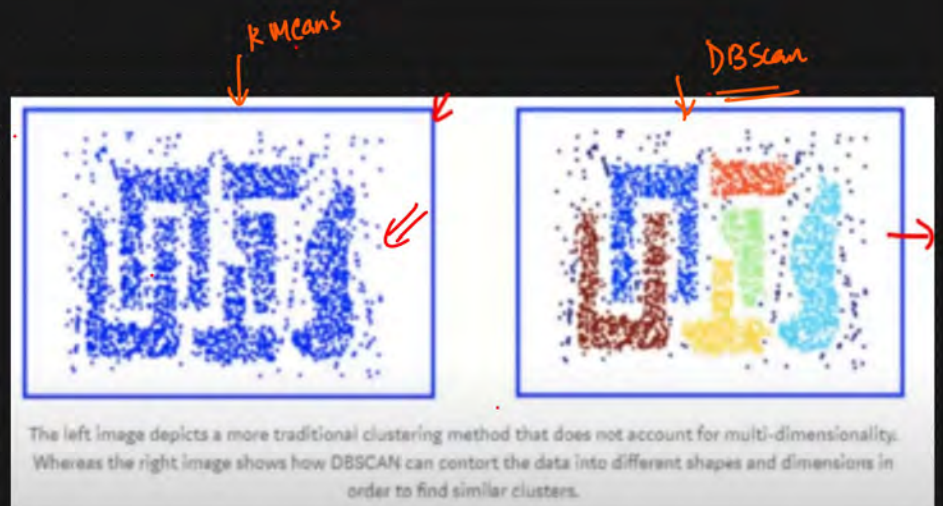
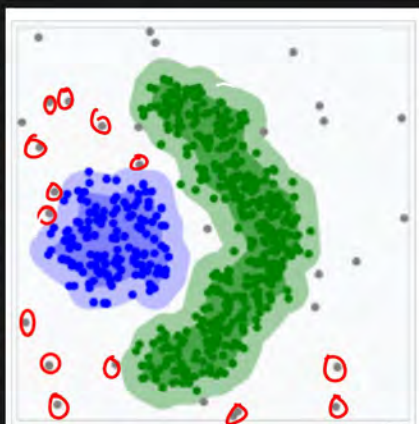


② Border point

No. of data points within the radius ϵ but is less than $\underline{\underline{minpts = 4}}$.



③ Outliers [DBSCAN is Robust to Outliers]



* Silhouette Clustering → Performance metrics of clustering

①

Assume the data have been clustered via any technique, such as **k-medoids** or **k-means**, into k clusters.

For data point $i \in C_I$ (data point i in the cluster C_I), let

$$a(i) = \frac{1}{|C_I| - 1} \sum_{j \in C_I, i \neq j} d(i, j)$$

be the mean distance between i and all other data points in the same cluster, where $|C_I|$ is the number of points belonging to cluster C_I , and $d(i, j)$ is the distance between data points i and j in the cluster C_I (we divide by $|C_I| - 1$ because we do not include the distance $d(i, i)$ in the sum). We can interpret $a(i)$ as a measure of how well i is assigned to its cluster (the smaller the value, the better the assignment).

$a(i) =$



②

We then define the mean dissimilarity of point i to some cluster C_J as the mean of the distance from i to all points in C_J (where $C_J \neq C_I$).

For each data point $i \in C_I$, we now define

$$b(i) = \min_{J \neq I} \frac{1}{|C_J|} \sum_{j \in C_J} d(i, j)$$

$b(i) \Rightarrow$



③

We now define a *silhouette* (value) of one data point i

$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}, \text{ if } |C_I| > 1$$

and

$$s(i) = 0, \text{ if } |C_I| = 1$$

Which can be also written as:

$$s(i) = \begin{cases} 1 - \frac{a(i)}{b(i)}, & \text{if } a(i) < b(i) \\ 0, & \text{if } a(i) = b(i) \\ \frac{b(i)}{a(i)} - 1, & \text{if } a(i) > b(i) \end{cases}$$

From the above definition it is clear that

$$-1 \leq s(i) \leq 1$$

≤ 1

Clustering is
not done
well

→ -1 to 0

[-ve]

↓

THANK - YOU