

Q1 - Overview of Clustering

Clustering is the use of unsupervised techniques for grouping similar objects. In ML, unsupervised refers to the pattern of finding hidden structure within unlabeled data.

- clustering is a method often used for exploratory analysis of the data.
 - In clustering there are no predictions made ~~paths~~ further.
 - Clustering methods find the similarities b/w objects according to the object attributes and group similar objects into clusters.
 - cluster techniques are utilized in marketing, economics and various branches of science.
- A popular clustering method is K-means.

Q2) ^{Imp} K means \rightarrow Given a collection of objects each with n measurable attributes, K-means is an analytical techniques that for a chosen values of K , identifies K clusters of objects based on the objects' proximity of the center of the K groups.

Algorithm:-

- 1) Selects the no of clusters $K=3$.
- 2) Randomly select 3 distinct data points
- 3) Mean distance b/w its 1st pt & selected 3 clusters

"K means is a clustering algorithm whose main goal is to group similar elements or data points into clusters"

Note:- "K" in K-means represent the number of clusters"

Algorithm:-

Steps:-

1. Select the number K to decide the number of clusters.
2. Select random K points as centroids.
3. Assign each data point to their closest centroid, which will form the predefined K clusters.
4. Calculate the variance & place a new centroid for each cluster.
5. Repeat the step 3.
6. If any reassignment occurs then go to step 4, else to finish.
7. Stop.

Example for K -means.

10 marks.

2022 Q) Suppose we have 4 different medicines with their weight index & pH as mentioned below, making use of K means algorithm group the below objects into 2 groups. use $c_1 = (1, 1)$ & $c_2 = (2, 1)$ as initial clusters.

Objects	weight	pH	$c_1(x_1, y_1)$	$c_2(x_2, y_2)$	clusters
A1	2	1	0	1	C_1
B2	2	1	1	0	C_2
C4	4	3	3.60	2.82	C_2
D5	5	4	5	4.24	C_2

$$d(P_1, P_2) = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

$$C_1 = (A1)$$

$$C_2 = (B2, C4, D5)$$

What is Regression Analysis?

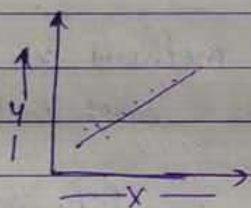
⇒ they are used to determine the relationship b/w a dataset's dependent & independent variables. (x, y)

- It is widely used when the dependent & independent variables are linked in a linear or non-linear fashion. & target variable has a set of continuous values.

Example - the best way to examine the relationship b/w sales & advertising expenditures.

1. Linear Regression:-

Linear Regression is an analytical technique used to model the relationship b/w a dependent variable (y) and an independent variable (x).



It employs a regression line, also known as best-fit line.

The line equation is defined as

$$y = c + m \cdot x + e$$

c = intercept, m = slope of line,

e = error term.

2) Logistic Regression. List the reasons to choose b/w the regression algo & explain the same.

⇒ when choosing b/w regression algorithms, several factors should be considered based on the specific problem at hand & the character of the dataset.

Some reasons to choose linear regression:-

1. Interpretability:- The it provides a clear pos' interpretation of the relship b/w the variables. Coefficients of the linear regression equation indicate the magnitude & direction.
2. Linearity assumption:- Linear Regr assumes a linear rel b/w the independent variables & the target variable.
3. Simplicity and Speed:- it is computationally efficient and relatively simple to implement compared to other regression algo.
4. Limited training data:- it performs well even when the dataset is small or the number of training instances is limited.
5. Baseline model:- it is a good baseline model to compare against more complex algo.

Q) Explain how least squared methods are used in minimizing errors in linear regression?

- If the data shows a linear relationship between two variables, it results in a least-squares regression line.
- This minimizes the vertical distance from the data points to the regression line.
 - The term least squares is used b'coz it is the smallest sum of squares of error, which is also called the variance.

8) In K means algorithm, have to choose the appropriate value of k and discuss the mathematical basis for the same.

→ Elbow Method.

Elbow method is one of the most popular ways to find the optimal number of clusters.

- This method uses the concept of WCSS value.

- WCSS: within cluster sum of squares, which defines the total variations within a cluster.

Formula is :-

for $k=2$

$$SSE_1 = \sum \text{distance}(x - c_1)^2$$

$$SSE_2 = \sum \text{distance}(x - c_2)^2$$

$$SSE = SSE_1 + SSE_2$$

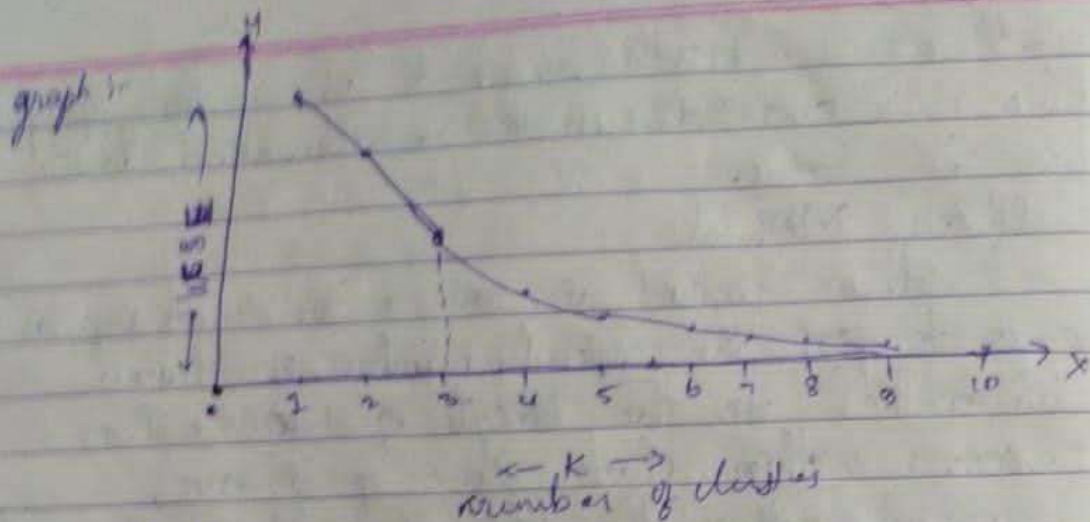
$$- WCSS = \sum_{P_i \text{ in cluster}_1} \text{distance}(P_i, c_1)^2 + \sum_{P_i \text{ in } c_2} \text{distance}(P_i, c_2)^2 + \dots$$

$\sum_{P_i \text{ in cluster}_1} \text{distance}(P_i, c_1)^2$:- it is sum of the square of the distance b/w each data point and its centroid within a cluster & the same for other terms.

- To measure the distance b/w data points & centroid, we can use any method such as Euclidean distance or Manhattan distance.

- To find optimal value, elbow method follows the steps:-

1. It executes the K-means clustering on a given dataset for different k values (ranges from 1-10)
2. For each value of k, calculate the WCSS value.
3. Plots a curve b/w calculated WCSS values & the number of clusters k
4. The shape of bend or a point of the plot look like an arm, that point is considered as the best value of k



(K) $K=3$ As K value increase, the SSE value decrease.

Q) Illustrate the use of Box-Jenkins Methodology for time series analysis.

→ A time series consists of an ordered sequence of equally spaced values over time.

Example: Monthly unemployment rates, daily website visits, or stock prices every second.

- A time series can consist of the following components:

- Trend
- Seasonality
- Cyclic
- Random.

1. Trend: refers to the long-term movement in a time series. It indicates whether the observation values are increasing or decreasing over time.

eg: increase in sales over months.

2. Seasonality: it describes the fixed, periodic fluctuation in the observation over time.

- it's often related to the calendar.

eg: monthly retail sales can fluctuate over the year due to the weather and holidays.

3) Cyclic: it refers to a periodic fluctuation, but ^{one} that is not as fixed as in case of a seasonally component.

Eg:- retail sales are influenced by the general state of the economy. Thus, retail sales time series can often follow the lengthy boom-bust cycles of the economy.

4) Random :- component is what remains, after accounting for the other 3 components.

Although noise is certainly part of this random component, there is also some underlying structure to this random component that need to be forecast future values.

- Developed by Georg Box - The Box-Jenkins methodology for time series analysis involves the 3 main steps

1. Condition data and select a model.

- Identify & account for any trends or seasonality in the time series

- Examine the remaining time series & determine a suitable model.

2. Estimate the model parameters.

3. Assess the model and return to steps, if necessary

Q.1 Explain the k nearest Neighbour algorithm.
List out the modeling assumptions to be made while using KNN algorithm.

- \Rightarrow k -nearest Neighbour is one of the simplest ML algo based on supervised learning technique.
- \Rightarrow KNN assumes the similarity b/w the new cases and available cases & put the new cases into the category that is most similar to the available categories.
- \Rightarrow KNN is used for the classification problems.
- \Rightarrow KNN is a non-parametric algo, means it does not make any assumption on underlying data.
- \Rightarrow it is called lazy learner algo, bcz it does not learn from the training set immediately, it performs an action on the dataset.

Eg:- Suppose we have 2 image of animal that looks similar to cat & dog, but we want to know its a cat or dog, so here we can use KNN algo as it work on similarity measure.

- It will find the similar features of the new data set to the cats & dog images & based on the most similar features it will put it in either cat or dog category.

