

Course: PMIT-6307; Class Test-03; Marks-10; Time-30 Minutes

- Ques. 1. Illustrate the process of the 5-nearest neighbor's algorithm with an example. [4]
 2. Consider a situation where you have 1000 fruits which are either 'banana' or 'apple' or 'other'. The possible attributes of any fruit are: [Long, Sweet, Yellow]. The training dataset & summary of the training dataset will look like Table-1 & Table-2.
 Now, consider a case where you're given that a fruit is long, sweet and yellow; and you need to predict what type of fruit it is [using Naive Bayes]. [6]

Table 1: Training Dataset

Long	Sweet	Yellow	Fruit
0	0	1	Apple
1	0	1	Banana
0	1	0	Apple
1	1	1	Other
..

Table 2: Summary

Type	Long	Not Long	Sweet	Not sweet	Yellow	Not Yellow	Total
Banana	400	100	350	150	450	50	500
Apple	0	300	150	150	300	0	300
Other	100	100	150	50	50	150	200
Total	500	500	650	350	800	200	1000

- Ques. 3. Why do most researcher prefer SVM classification? How does SVM work?
 4. What is the class imbalance problem? Why is it a problem? How mitigate this problem?
 5. Define Scalability
 6. Define hierarchical clustering with example. What is the most useful technique?
 7. Write the steps of k-means clustering, what is the difference of clustering problem?
 8. Are some limitation of k-means clustering? How can you overcome this problem?
 9. What do you mean by clustering? What are the different types of cluster?
 10. How a k-means classification preferred? How do you decide number of k?
 11. In density based clustering, how do you decide minPoint and Eps?

$$\begin{aligned}
 (12) \sqrt{(61-165)^2 + (61-65)^2} &= 5.66 \\
 (13) \sqrt{(61-168)^2 + (61-63)^2} &= 7.07 \\
 (14) \sqrt{(61-168)^2 + (61-63)^2} &= 7.28 \\
 (15) \sqrt{(61-168)^2 + (61-66)^2} &= 8.60 \\
 (16) \sqrt{(61-170)^2 + (61-63)^2} &= 9.22 \\
 (17) \sqrt{(61-170)^2 + (61-64)^2} &= 9.49 \\
 (18) \sqrt{(61-170)^2 + (61-68)^2} &= 11.40
 \end{aligned}$$

Step 1: Calculate similarity based on distance function using Euclidean Distance

Step 2: Find K-Nearest Neighbors; Let K=5, Searches for the 5 customer closest to house A. Of 4 of them had medium size and one is large size. Then best guess for house A is medium T shirt size.

Nearest Neighbor Classifier

Let
 $K=5$

Given that,

Customer name = Monica.

Height = 161 cm.

Weight = 61 kg

T-shirt size = ?

Calculate Euclidean Distance
in all attributes:

- (1) $\sqrt{(161-158)^2 + (61-58)^2} = 4.24$
- (2) $\sqrt{(161-158)^2 + (61-59)^2} = 3.61$
- (3) $\sqrt{(161-158)^2 + (61-63)^2} = 3.61$
- (4) $\sqrt{(161-160)^2 + (61-59)^2} = 1.41$
- (5) $\sqrt{(161-160)^2 + (61-67)^2} = 1.41$
- (6) $\sqrt{(161-163)^2 + (61-69)^2} = 2.24$
- (7) $\sqrt{(161-163)^2 + (61-61)^2} = 2.00$
- (8) $\sqrt{(161-160)^2 + (61-68)^2} = 3.16$
- (9) $\sqrt{(161-163)^2 + (61-64)^2} = 3.61$
- (10) $\sqrt{(161-165)^2 + (61-61)^2} = 4.00$
- (11) $\sqrt{(161-165)^2 + (61-62)^2} = 4.12$
- (12) $\sqrt{(161-165)^2 + (61-65)^2} = 5.66$
- (13) $\sqrt{(161-168)^2 + (61-62)^2} = 7.07$
- (14) $\sqrt{(161-168)^2 + (61-63)^2} = 7.28$
- (15) $\sqrt{(161-168)^2 + (61-66)^2} = 8.60$
- (16) $\sqrt{(161-170)^2 + (61-63)^2} = 9.22$
- (17) $\sqrt{(161-170)^2 + (61-64)^2} = 9.49$
- (18) $\sqrt{(161-170)^2 + (61-68)^2} = 11.40$

Step 1: Calculate similarity based on distance function using
Euclidean Distance

Step 2: Find K-Nearst Neighbors : Let $K=5$, Search for
the 5 customer closest to Monica. Of 4 of them had
medium(M) size as large(L) size. Then best guess for
Monica as medium T-shirt size.

#	Height	Weight	Size	Distance	Posit.
01	158	58	M	4.24	
02	158	59	M	3.61	
03	158	63	M	3.61	
04	160	59	M	1.41	01
05	160	60	M	1.42	02
06	163	60	M	2.24	04
07	163	61	M	2.00	03
08	160	64	L	3.16	05
09	163	64	L	3.61	
10	165	61	L	7.00	
11	165	62	L	4.12	
12	165	65	L	5.66	
13	168	62	L	7.07	
14	168	63	L	7.28	
15	168	66	L	8.60	
16	170	63	L	9.22	
17	170	67	L	9.49	
18	170	68	L	11.40	

Naïve Bayes Classifier (Bayes Theorem).

Frequency Table:-

Weather	No	Yes
Overscast	5	0
Rainy	2	2
Sunny	3	2
Total.	10	4

Likelihood Table:-

Weather	No	Yes
Overscast	0	5
Rainy	2	2
Sunny	2	3
All	$\frac{4}{14}$	$\frac{10}{14} = \frac{5}{7}$

Here, $P(\text{Sunny}) = 0.35$, $P(\text{Yes}) = 0.71$, $P(\text{Sunny}/\text{Yes}) = \frac{5}{10} = 0.5$

$$P(\text{Yes}|\text{Sunny}) = \frac{0.3 \times 0.71}{0.35} = 0.60$$

Here, $P(\text{Sunny}) = 0.35$, $P(\text{No}) = 0.29$, $P(\text{Sunny}/\text{No}) = \frac{2}{4} = 0.5$

$$P(\text{No}|\text{Sunny}) = \frac{0.5 \times 0.29}{0.35} = 0.41$$

Here, $P(\text{Yes}|\text{Sunny}) > P(\text{No}|\text{Sunny})$.

Hence, a sunny day, Player can play the game.

→ Vertical Orient.

CRM

→ Horizontal Orient.

prob

$$P(\text{Sunny}/\text{Yes}) = \frac{\text{Sunny Yes}}{\text{Yes Total}}$$

$$= \frac{3}{10}$$

$$= 0.3$$

$$P(\text{Sunny}/\text{No}) = \frac{\text{Sunny No}}{\text{No Total}}$$

$$= \frac{2}{4}$$

$$= 0.5$$

$$P(\text{Yes}|\text{Sunny}) = \frac{P(\text{Sunny}/\text{Yes}) \times P(\text{Yes})}{P(\text{Sunny})} = \frac{0.3 \times 0.71}{0.35} = 0.60$$

$$P(\text{No}|\text{Sunny}) = \frac{P(\text{Sunny}/\text{No}) \times P(\text{No})}{P(\text{Sunny})} = \frac{0.5 \times 0.29}{0.35} = 0.41$$

$$\dots \rightarrow P_{\text{PT-9}}.$$

Naïve Bayes Classifier (Bay's Theorem).

Frequency Table:-

Weather	Yes	No
Overcast	5	0
Rainy	2	2
Sunny	3	2
Total	10	4

Likelihood Table:-

Weather	No	Yes
Overcast	0	5
Rainy	2	2
Sunny	2	3
All	4/14	10/14 = 0.7
	= 0.29	

here, $P_{\text{Sunny}} = 0.35$, $P_{\text{Yes}} = 0.71$, $P_{\text{Rainy}}/(\text{No}) = \frac{2}{10} = 0.2$

$$P(\text{Yes}|\text{Sunny}) = \frac{0.3 \times 0.71}{0.35} = 0.60$$

here, $P_{\text{Sunny}} = 0.35$, $P_{\text{No}} = 0.29$, $P_{\text{Sunny}/(\text{No})} = \frac{3}{4} = 0.75$

$$P(\text{No}|\text{Sunny}) = \frac{0.5 \times 0.29}{0.35} = 0.41$$

Here, $P(\text{Yes}|\text{Sunny}) > P(\text{No}|\text{Sunny})$.

hence, a sunny day, P can play the game.

→ Vertical column

Chances

Horizontal margin
prob.

$$P_{\text{Sunny}/(\text{Yes})} = \frac{\text{Sunny Yes}}{\text{Total}}$$

$$= \frac{3}{10}$$

$$= 0.3$$

$$= \frac{3}{10}$$

$$P(\text{Yes}|\text{Sunny}) = \frac{P_{\text{Sunny}/(\text{Yes})} \times P_{\text{Yes}}}{P_{\text{Sunny}}} = \frac{0.3 \times 0.71}{0.35} \rightarrow 0.60$$

$$P(\text{No}|\text{Sunny}) = \frac{P_{\text{Sunny}/(\text{No})} \times P_{\text{No}}}{P_{\text{Sunny}}} = \frac{0.5 \times 0.29}{0.35} = 0.41$$

$$P_{\text{Sunny}/(\text{No})} = \frac{\text{Sunny No}}{\text{Total}}$$

$$= \frac{2}{4}$$

$$= 0.5$$

Example

outlook

	Ys	No	P(Ys)	P(No)
Sunny	2	3	2/9	3/5
Cloudy	4	0	4/9	0/5
Rainy	3	2	3/9	2/5
Total	9	5	100%	100%

Temp.

	Ys	No	P(Ys)	P(No)
Fair	2	2	2/9	2/5
Normal	4	2	4/9	2/5
Cool	3	1	3/9	1/5
Total	9	5	100%	100%

Wend.

	Ys	No	P(Ys)	P(No)
Fair	2	2	2/9	2/5
Temp	3	3	3/9	3/5
Total	9	5	100%	100%

Play (Today).

Play	P(Yes)/P(No)
Yes	9/14
No	5/14
Total	14

Day	outlook	Temp.	Humidity	Wind
0	Rain	Fair	High	False
-1	Pain	Fair	High	True
2	Cloudy	Hot	High	F
3	Sunny	Wild	High	F
4	Sunny	Cloud	Normal	F
5	Sunny	Cool	Normal	T
6	Cloudy	Cool	Normal	T
7	Pain	Wild	High	F
8	Rain	Cool	Normal	F
9	Sunny	Wild	Normal	F
10	Rain	Wild	Normal	T
11	Cloudy	Wild	High	T
12	Cloudy	Fair	Normal	F
13	Sunny	Wild	High	T

Probability of playing golf is given by:-

$$P(\text{Yes Today}) = \frac{P(\text{Sunny} | \text{outlook}) P(\text{Fair} | \text{Temp}) P(\text{Normal} | \text{Humidity}) P(\text{No} | \text{Wind}))}{P(\text{today})}$$

$$= \frac{2}{9} \cdot \frac{2}{9} \cdot \frac{6}{9} \cdot \frac{6}{9} \cdot \frac{9}{14} \approx 0.01411$$

$$P(\text{Not today}) = \frac{P(\text{Sunny} | \text{outlook} | \text{No}) P(\text{Fair} | \text{Temp} | \text{No}) P(\text{Normal} | \text{Humidity} | \text{No}) P(\text{No} | \text{Wind} | \text{No})}{P(\text{today})}$$

$$\approx \frac{3}{5} \cdot \frac{2}{5} \cdot \frac{1}{5} \cdot \frac{2}{5} \cdot \frac{5}{14} = 0.0068$$

So, prediction that golf would be played is "Yes"

CLUSTERING

1. What are the different types of clustering methods used in business Intelligence?

Clustering is an undirected technique used in data mining for identifying several hidden patterns in the data without coming up with any specific hypothesis. The reason behind using clustering is to identify similarities between certain objects and make a group of similar ones.

There are two different types of clustering, which are hierarchical and non-hierarchical methods. Non-hierarchical Clustering In this method, the dataset containing N objects is divided into M clusters. In business intelligence, the most widely used non-hierarchical clustering technique is K-means. Hierarchical Clustering In this method, a set of nested clusters are produced. In these nested clusters, every pair of objects is further nested to form a large cluster until only one cluster remains in the end.

The different types of clustering methods below.

1. Density-Based Clustering
2. DBSCAN (Density-Based Spatial Clustering of Applications with Noise)
3. OPTICS (Ordering Points to Identify Clustering Structure)
4. HDBSCAN (Hierarchical Density-Based Spatial Clustering of Applications with Noise)
5. Hierarchical Clustering
6. Fuzzy Clustering
7. Partitioning Clustering
8. PAM (Partitioning Around Medoids)
9. Grid-Based Clustering

2. When is Clustering used?

The primary function of clustering is to perform segmentation, whether it is store, product, or customer. Customers and products can be clustered into hierarchical groups based on different attributes. Another usage of the clustering technique is seen for detecting anomalies like fraud transactions. Here, a cluster with all the good transactions is detected and kept as a sample. This is said to be a normal cluster. Whenever something is out of the line from this cluster, it comes under the suspect section. This method is found to be really useful in detecting the presence of abnormal cells in the body. Other than that, clustering is widely used to break down large datasets to create smaller data groups. This enhances the efficiency of assessing the data.

3. What are the advantages of Clustering?

Clustering is said to be more effective than a random sampling of the given data due to several reasons. The two major advantages of clustering are:
Requires fewer resources A cluster creates a group of fewer resources from

the entire sample. Due to this, there is a lesser requirement of resources as compared to random sampling. Random sampling will require travel and administrative expenses, but this is not the case over here. Feasible option Here, every cluster determines an entire set of the population as homogeneous groups are created from the entire population. With this, it becomes easy to include more subjects in a single study.

4. Why do we use clustering in ML?

Cluster analysis is usually used to classify data into structures that are more easily understood and manipulated. It is an unsupervised machine learning task.

5. What is the difference between clustering and classification in ML?

Classifying the input labels basis on the class labels is classification. On the other hand, the process of grouping basis the similarity without taking help from class labels is known as clustering.

6. Why clustering is better than classification?

The machine learns from the existing data in clustering because the need for multiple pieces of training is not required. Grouping is done on similarities as it is unsupervised learning. Classification on the contrary is complex because it is a supervised type of learning and requires training on the data sets.

7. What is hierarchical clustering?

It is a form of clustering algorithm that produces 1 to n clusters, where n represents the number of observations in a data set. There are two types of hierarchical clustering, divisive (top-down) and agglomerative (bottom-up).

8. Clustering Data Mining techniques: Hierarchical Clustering

When you're on a quest to find data pieces and map them according to cluster probability, the Hierarchical Clustering method works like a **charm**. Now, the mapped data pieces may belong to a Cluster with distinct qualities in terms of multidimensional scaling, cross-tabulation, or quantitative relationships among data variables in several aspects.

Considering how to find a single cluster after merging the various clusters while retaining the hierarchy of the attributes on which they are classed in mind? The stages of the Hierarchical Clustering method mentioned below can be used to accomplish this:

- Begin by picking the data points and clustering them according to the hierarchy.
- Are you considering how the clusters will be interpreted? With a Top-down or Bottom-up method, a Dendrogram may be utilized to comprehend the hierarchy of clusters properly.
- Clusters are merged until only one remains, and we may use a variety of metrics to determine how close the clusters are when merging them,

such as Euclidean distance, Manhattan distance, or Mahalanobis distance.

- For the time being, the process has ended because the intersection point has been discovered and well mapped on the dendrogram.

What is Clustering in Data Mining?

In clustering, a group of different data objects is classified as similar objects. One group means a cluster of data. Data sets are divided into different groups in the cluster analysis, which is based on the similarity of the data.

Data clustering, also called data segmentation, aims to partition a collection of data into a predefined number of subsets (or clusters) that are optimal in terms of some predefined criterion function. Data clustering is a fundamental and enabling tool that has a broad range of applications in many areas.

What is clustering and its types?

Clustering itself can be categorized into two types viz. Hard Clustering and Soft Clustering. In hard clustering, one data point can belong to one cluster only. But in soft clustering, the output provided is a probability likelihood of a data point belonging to each of the pre-defined numbers of clusters.

Density-Based Clustering

In this method, the clusters are created based upon the density of the data points which are represented in the data space. The regions that become dense due to the huge number of data points residing in that region are considered as clusters.

DBSCAN (Density-Based Spatial Clustering of Applications with Noise)

DBSCAN groups data points together based on the distance metric. It follows the criterion for a minimum number of data points. It can discover clusters of different shapes and sizes from a large amount of data, which is containing noise and outliers. It takes two parameters – *eps* and *minimum points*. Eps indicates how close the data points should be to be considered as neighbors. The criterion for minimum points should be completed to consider that region as a dense region.

5) Clustering Data Mining techniques: Density-Based Spatial Clustering

When it comes to discovering clusters in bigger geographical databases, the Density-based Spatial Clustering Algorithm with Noise (or DBSCAN) is a superior alternative to K Means when it comes to cross-examining the density of its data points. It's also more appealing and efficient than CLARANS, which stands for Clustering LARge ApplicatioNS via Medoid-based partitioning approach.

The DBSCAN Clustering algorithm approach is beneficial and comparable to the mean-shift density-based Clustering algorithm.

DBSCAN's method starts with an unvisited data point and uses distance (Epsilon) to extract the neighborhood before designating the point as visited. If two points are within a certain distance of each other, they have termed neighbors. Following are the steps of DBSCAN:

- When enough points (based on minPoints) are found, the clustering process begins, using the current data point as the initial point in the new cluster. If there aren't enough points, the algorithm flags it as visited and classifies it as noise defect clustering.
- The initial point in the new cluster utilizes the same distance to define its neighborhood, resulting in a clustered point neighborhood, and the process continues for every additional cluster point added to the group. This process is repeated until all data points have been labeled and visited.
- After all of the data points in the neighborhood have been visited, a fresh unvisited data point is chosen for clustering. As a result, all data points are labeled as noise or clustered under the visited label.

In density based clustering , How do you decide min Point an EPS

The eps should be chosen based on the distance of the dataset (we can use a k-distance graph to find it), but in general small eps values are preferable. **minPoints:** As a general rule, a minimum minPoints can be derived from a number of dimensions (D) in the data set, as $\text{minPoints} \geq D + 1$.

Parameters: The DBSCAN algorithm basically requires 2 parameters:

eps: specifies how close points should be to each other to be considered a part of a cluster. It means that if the distance between two points is lower or equal to this value (eps), these points are considered neighbors.

minPoints: the minimum number of points to form a dense region. For example, if we set the minPoints parameter as 5, then we need at least 5 points to form a dense region.

Parameter estimation:

The parameter estimation is a problem for every data mining task. To choose good parameters we need to understand how they are used and have at least a basic previous knowledge about the data set that will be used.

eps: if the eps value chosen is too small, a large part of the data will not be clustered. It will be considered outliers because don't satisfy the number of points to create a dense region. On the other hand, if the value that was chosen is too high, clusters will merge and the majority of objects will be in the same cluster. The eps should be chosen based on the distance of the dataset (we can use a k-distance graph to find it), but in general small eps values are preferable.

minPoints: As a general rule, a minimum minPoints can be derived from a number of dimensions (D) in the data set, as $\text{minPoints} \geq D + 1$. Larger values are usually better for data sets with noise and will form more significant clusters.

The minimum value for the minPoints must be 3, but the larger the data set, the larger the minPoints value that should be chosen.

Why should we use DBSCAN?

The DBSCAN algorithm should be used to find associations and structures in data that are hard to find manually but that can be relevant and useful to find patterns and predict trends.

Clustering methods are usually used in biology, medicine, social sciences, archaeology, marketing, characters recognition, management systems and so on. Let's think in a practical use of DBSCAN. Suppose we have an e-commerce and we want to improve our sales by recommending relevant products to our customers. We don't know exactly what our customers are looking for but based on a data set we can predict and recommend a relevant product to a specific customer. We can apply the DBSCAN to our data set (based on the e-commerce database) and find clusters based on the products that the users have bought. Using this clusters we can find similarities between customers, for example, the customer A have bought 1 pen, 1 book and 1 scissors and the customer B have bought 1 book and 1 scissors, then we can recommend 1 pen to the customer B. This is just a little example of use of DBSCAN, but it can be used in a lot of applications in several areas.

How can we easily implement DBSCAN?

As I already wrote (tip: don't believe in everything I write) the DBSCAN is a well-known algorithm, therefore, you don't need to worry about implement it yourself. You can use one of the libraries/packages that can be found on the internet. Here is a list of links that you can find the DBSCAN implementation: I also have developed an application (in Portuguese) to explain how DBSCAN works in a didactically way. The application was written in C++ and you can find it on [Github](#).

What is a real life example of DBSCAN clustering?

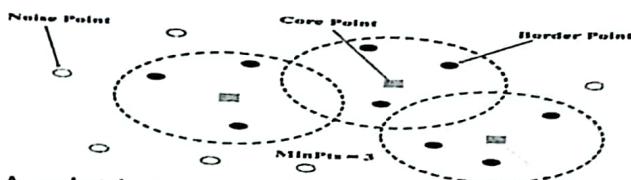
Using this clusters we can find similarities between customers, for example, if customer A has bought a pen, a book and one pair scissors, while customer B purchased a book and one pair of scissors, then you could recommend a pen to customer B.

The basic principle of dbscan clustering

The principle of DBSCAN is to find the neighborhoods of data points exceeds certain density threshold. The density threshold is defined by two parameters: the radius of the neighborhood (eps) and the minimum number of neighbors/data points (minPts) within the radius of the neighborhood. With these two thresholds in mind, DBSCAN starts from a random point to find its first density neighborhood. Then from its first density neighborhood, it starts to find the second density neighborhood. If the second density neighborhood exists, DBSCAN will merge the first and second density

neighborhoods to become a bigger density neighborhood. The same process keep going on until its neighborhood cannot meet the threshold anymore.

Pictorial example of Core Point, Noise Point and Border Point-Explain.



A point is a core point if it has more than a specified number of points (MinPts) within Eps—These are points that are at the interior of a cluster.
A border point has fewer than MinPts within Eps, but is in the neighborhood of a core point.

A noise point is any point that is not a core point nor a border point.

What is DBSCAN useful for?

DBSCAN is a clustering method that is used in machine learning to separate clusters of high density from clusters of low density.

Advantages. DBSCAN does not require one to specify the number of clusters in the data a priori, as opposed to k-means. DBSCAN can find arbitrarily-shaped clusters. It can even find a cluster completely surrounded by (but not connected to) a different cluster.

Density-Based Clustering refers to one of the most popular unsupervised learning methodologies used in model building and machine learning algorithms. The data points in the region separated by two clusters of low point density are considered as noise.

HDBSCAN (Hierarchical Density-Based Spatial Clustering of Applications with Noise)

HDBSCAN is a density-based clustering method that extends the DBSCAN methodology by converting it to a hierarchical clustering algorithm.

Hierarchical Clustering

Hierarchical Clustering groups (Agglomerative or also called as Bottom-Up Approach) or divides (Divisive or also called as Top-Down Approach) the clusters based on the distance metrics.

In agglomerative clustering, initially, each data point acts as a cluster, and then it groups the clusters one by one. This comes under in one of the most sought-after clustering methods.

Divisive is the opposite of Agglomerative, it starts off with all the points into one cluster and divides them to create more clusters. These algorithms create a distance matrix of all the existing clusters and perform the linkage between

the clusters depending on the criteria of the linkage. The clustering of the data points is represented by using a dendrogram. There are different types of linkages: –

- o **Single Linkage:** – In single linkage the distance between the two clusters is the shortest distance between points in those two clusters.
- o **Complete Linkage:** – In complete linkage, the distance between the two clusters is the farthest distance between points in those two clusters.
- **Average Linkage:** – In average linkage the distance between the two clusters is the average distance of every point in the cluster with every point in another cluster.

K-Means Clustering: – K-Means clustering is one of the most widely used algorithms. It partitions the data points into k clusters based upon the distance metric used for the clustering. The value of 'k' is to be defined by the user. The distance is calculated between the data points and the centroids of the clusters.

K-means clustering is a type of unsupervised learning used when you have unlabeled data (i.e., data without defined categories or groups). This algorithm aims to find groups in the data, with the number of groups represented by the variable K. In this clustering method, the number of clusters found from the data is denoted by the letter 'K.'

The data point which is closest to the centroid of the cluster gets assigned to that cluster. After an iteration, it computes the centroids of those clusters again and the process continues until a pre-defined number of iterations are completed or when the centroids of the clusters do not change after an iteration. It is a very computationally expensive algorithm as it computes the distance of every data point with the centroids of all the clusters at each iteration. This makes it difficult for implementing the same for huge data sets.

2) Clustering Data Mining Techniques: K-Means Clustering

After determining the centroid value between two data points, the K-Means Clustering Algorithm repeatedly discovers the k number of clusters. It is rather useful to compute Cluster Centroids with their vector quantization observations, by virtue of which data points with changeable characteristics may be brought to Clustering.

As the clustering process accelerates, a large amount of unlabeled real-world data will become more efficient as it is split into clusters of different forms and densities. Have you ever considered how the centroid distance is calculated? Take a look at the K-means steps stated below:

- First, decide on the number of clusters that will differ in shape and density. Let's call that number k , and its value can be anything from 3,4 to anything else.
- You may now assign data points to the number of the cluster. The centroid distance is then calculated using the least squared Euclidean distance once the data point and cluster have been chosen.
- The data point resembles the cluster if it is substantially closer to the centroid distance; otherwise, it does not.
- Iteratively compute the centroid distances with the selected data point until you find the largest number of clusters made up of related data points. When assured convergence (a point where data points are well clustered) is reached, the algorithm stops clustering.

What are the Applications of Data Mining Clustering Techniques?

- Clustering can assist marketers identify unique groups in their consumer bases and describe them based on purchase behaviors in the business world.
- It may be used in biology to create plant and animal taxonomies to classify genes with similar functions.
- Clustering can also assist in identifying regions of land usage in an earth observation database, as well as groupings of motor insurance customers with a high average claim cost.
- Cluster analysis may be used as a standalone data mining function to obtain insight into data distribution, notice the features of each cluster, and focus on a specific group of clusters for further study.

What are the Requirements of Clustering Data Mining Techniques?

- **Scalability:** Many clustering techniques work well on small data sets with less than 200 data objects, however, a huge database might include millions of objects. Clustering on a subset of a big dataset might result in skewed findings. Clustering methods that are highly scalable are required.
- **Usability and interpretability:** Users anticipate interpretable, thorough, and usable clustering findings. As a result, clustering may require unique semantic interpretations and applications. It's crucial to investigate how the application aim influences Clustering Data Mining technique selection.
- **High dimensionality:** A database or a data warehouse can have several dimensions or properties. Many clustering algorithms excel

at dealing with low-dimensional data (two or three dimensions). Human eyes are capable of assessing clustering quality in up to three dimensions. Clustering data items in a high-dimensional space may be difficult, especially when the data is sparse and heavily skewed (misleading data).

- **Constraint-based clustering:** Clustering may be required in real-world applications due to a variety of restrictions. Assume you're in charge of selecting locations for a certain number of new automatic cash dispensing machines (ATMs) in a city. You may decide this by clustering households while taking into account limits such as the city's waterways, highway networks, and client needs per area. Finding groupings of data with appropriate clustering behavior that fulfill stated requirements is a difficult issue.

Limitations of K-means clustering

- difficult to choose the number of clusters, k
- cannot be used with arbitrary distances
- sensitive to scaling – requires careful preprocessing
- does not produce the same result every time
- clever initialization helps finding better results
- sensitive to outliers (squared errors emphasize outliers)
- cluster sizes can be quite unbalanced (e.g., one-element outlier clusters)

Process of overcome the Limitations of K-means clustering

A common way to avoid this problem is to use the K-means++ algorithm, which selects the initial centers based on a probability distribution that favors points that are far from each other. This can reduce the number of iterations and improve the quality of the final clusters.

How a KNN Classification

KNN classifier is a machine learning algorithm used for classification and regression problems. It works by finding the K nearest points in the training dataset and uses their class to predict the class or value of a new data point.

What is an example of KNN classification?

With the help of KNN algorithms, we can classify a potential voter into various classes like "Will Vote", "Will not Vote", "Will Vote to Party 'Congress'", "Will Vote to Party 'BJP'". Other areas in which KNN algorithm can be used are Speech Recognition, Handwriting Detection, Image Recognition and Video Recognition.

Why use KNN classification?

The KNN algorithm can compete with the most accurate models because it makes highly accurate predictions. Therefore, you can use the KNN algorithm for applications that require high accuracy but that do not require a human-readable model. The quality of the predictions depends on the distance measure.

Is KNN best for classification?

The main advantage of KNN over other algorithms is that KNN can be used for multiclass classification. Therefore if the data consists of more than two labels or in simple words if you are required to classify the data in more than two categories then KNN can be a suitable algorithm.

Step of KNN Classification

The K-NN working can be explained on the basis of the below algorithm:

Step-1: Select the number K of the neighbors.

Step-2: Calculate the Euclidean distance of K number of neighbors.

Step-3: Take the K nearest neighbors as per the calculated Euclidean distance.

Why do we use SVM? Why do most researcher perform SVM Classification.

SVMs are used in applications like handwriting recognition, intrusion detection, face detection, email classification, gene classification, and in web pages. This is one of the reasons we use SVMs in machine learning. It can handle both classification and regression on linear and non-linear data.

What do you mean by OLAP?

Online analytical processing (OLAP) is software technology you can use to analyze business data from different points of view. Organizations collect and store data from multiple data sources, such as websites, applications, smart meters, and internal systems.

What is an example of an OLAP?



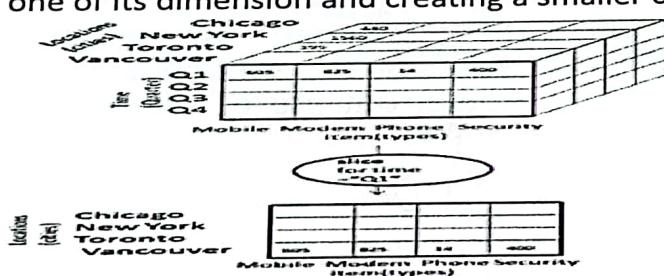
For example, a user can request that data be analyzed to display a spreadsheet showing all of a company's beach ball products sold in Florida in the month of July, compare revenue figures with those for the same products in September and then see a comparison of other product sales in Florida in the same time period.

What is Slicing and Dicing?

Slicing and Dicing refers to a way of segmenting, viewing and comprehending data in a database. Large blocks of data is cut into smaller segments and the process is repeated until the correct level of detail is achieved for proper analysis. Therefore slicing and dicing presents the data in new and diverse perspectives and provides a closer view of it for analysis. For example a report is showing annual performance of a particular product. If we want to view the quarterly performance, we can use slicing and dicing strategy to drill down to the quarterly level.

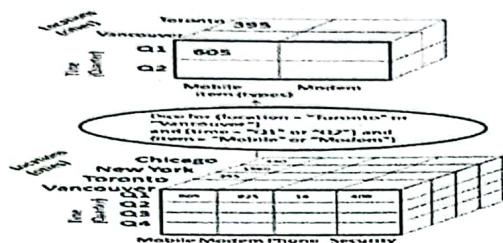
OLAP (Online Analytical Processing) is a computer process that enables user to select and extract data from different viewpoints. Slicing and Dicing term is generally used in OLAP databases that presents data to the end user in multidimensional cube format like a 3D spreadsheet (called an OLAP cube). The OLAP cube comprises numeric facts called measures which are categorized by dimensions.

Slicing refers to selecting a subset of the cube by choosing a single value for one of its dimension and creating a smaller cube with one less dimension.



In the above diagram, the slice is performed for dimension "time" using the variable "Q1" (criteria time). It generates a new cube with 2 dimensions. Thus slice operation produced a sliced OLAP cube by enabling the user to choose specific value for one of its dimension.

Dicing on the other hand generates a subcube by picking two or more values from multiple dimensions of the cube. The cube is rotated independent of its dimension.



In the above diagram, the dice operation is based on the following selection criteria involving all the three variables:

- Location = "Vancouver" or "Toronto"
- Time = "Q1" or "Q2"
- Item = " Modem" or "Mobile"

Thus dicing and slicing allows detailed analysis of data from different angles.

Hence, this concludes the definition of Slicing and Dicing along with its overview. This article has been researched & authored by the [Business Concepts Team](#). It has been reviewed & published by the MBA Skool Team. The content on MBA Skool has been created for educational & academic purpose only. Browse the definition and meaning of more similar terms. The Management Dictionary covers over 1800 business concepts from 5 categories.

why EDA is necessary for data mining

The main purpose of EDA is to help look at data before making any assumptions. It can help identify obvious errors, as well as better understand patterns within the data, detect outliers or anomalous events, find interesting relations among the variables.

Why is EDA important in data mining?

In data mining, Exploratory Data Analysis (EDA) is an approach to analyzing datasets to summarize their main characteristics, often with visual methods. EDA is used for seeing what the data can tell us before the modeling task.

Why should we perform EDA?

Exploratory data analysis can help detect obvious errors, identify outliers in datasets, understand relationships, unearth important factors, find patterns within data, and provide new insights.

Why do we use cross-validation in classification?

The purpose of cross-validation is to test the ability of a machine learning model to predict new data. It is also used to flag problems like overfitting or selection bias and gives insights on how the model will generalize to an independent dataset.

Cross-validation is a technique that helps you assess the accuracy and generalizability of your predictive models. It involves splitting your data into multiple subsets and training and testing your model on different combinations of them.

Advantages of the learning curve model

Strategic planning to improve the output of employees or even whole departments. Motivating company staff by creating a culture of ongoing learning and progress-tracking. Identifying trends that can be used for more accurate forecasting and better business decisions.

The ROC curve is used to assess the overall diagnostic performance of a test and to compare the performance of two or more diagnostic tests. It is also used to select an optimal cut-off value for determining the presence or absence of a disease.

Advantages of the ROC curve

You say "Typically, the optimal cut-off value for a screening test is the one closest to the top left hand corner [of the ROC curve]." One could expand on that and say that the optimal point varies according to the pre-test probability (or prevalence) of the condition and the relative disadvantages (financial and/or health outcomes) of false positive and false negative test results. To me, the advantage of keeping the whole ROC curve is that the correct optimal point can be chosen for each case rather than losing information by making a somewhat arbitrary choice of optimal cut-off point.

Professional Masters in IT (PMIT-6307); Class Test-01; Marks-15; Time-30 Minutes

1. Discuss whether or not each of the following activities is a data mining task. If the answer is yes, then also specify which one of the task categories it will belong to: 4
 - a. An image analyst obtains some new images and wants to automatically detect the number of distinct objects in the image. He does not have any prior information about these objects.
 - b. By looking at a CT scan, a doctor wants to identify if a patient has cancer or not. There are a lot of labelled CT scans that the doctor will use for making the decision.
 - c. Predicting the future stock price of a company using historical records. → Yes,
 - d. Monitoring the heart rate of a patient for abnormalities.
2. Write short note on following application in terms of goals and approaches: 5
 - a. Loan Default Prediction
 - b. Credit Card Fraud Detection
3. Briefly describe three data quality problems.
4. Is temperature an interval or ratio measurement statistically? Justify your answer. 3

5. What is data mining & its application
6. Difference b/w descriptive & predictive data mining
7. Example of Descriptive data mining.
8. Difference b/w Classification & clustering.
9. What are the sequential step involved in classification?
10. What is the association rule discovery with example.
11. What is data mining in a machine learning?
12. Why do we need study more on data mining?
13. Challenges of implementation data mining.
14. Explain reasons for data mining.
15. Scalability define - Diff. b/w Data mining & database.
16. Categorization of Attribute
17. Same above - 4.
18. List the different types of attributes.

Q) Discuss whether or not each of the following activities is a data mining task. If your answer is yes, then also specify which one of the task categories it will belong to.

a) An image analysis obtain some new images and want to automatically detect the number of distinct object in the image. He does not have any prior information about these objects.

Ans: Yes.

b) By looking at a CT scan, a doctor wants to identify if a patient has cancer or not. There are a lot of labeled CT scans that the doctor will use for making a decision.

Goal: predict brain tumor from available

Approach: MRI report collecting and analyzing objects. The MRI report data is grouped together to detect brain tumor.

c) Predicting the future stock price of company based on historical record.

Yes, we would attempt to create a model that can predict the current value of the stock price. This is an example of the area of data mining known as predictive modeling. We could use regression for this modelling, although researchers generally feel there is a wide variety of techniques for predicting time series.

d) Monitoring the heart rate of a patient for abnormality.

Yes, we would build a model of the normal behavior of heart rate and raise an alarm when unusual heart behavior occurred. This would involve the area of data mining known as anomaly detection. This could also be considered as a classification problem if we have examples of both normal and abnormal heart behavior.

(B) Briefly describe three Data quality problems.

Ans:- Data has become the heart of all Business across the world. Organizations heavily rely on data assets for the decision making process, but unfortunately most of the clean, accurate and quality data doesn't exist. Data is impacted by numerous factors that deteriorate its quality. Data quality refers to the measurement of the current state of data against traits like completeness, accuracy, reliability, relevance and timeliness. While the data quality issue indicate the presence of defect that harm the above mentioned traits. Data is beneficial only, if it's of high quality, some of the consequences of poor quality data are as follows:-

- ① Human Error ② Data Duplication ③ Inconsistent Data
- ④ Inaccurate and Missing Data ⑤ Using the wrong formula
- ⑥ Data Overload ⑦ Data DownTime ⑧ Hidden Data ⑨ Outdated Data
- ⑩ Data illiteracy ⑪ Noise and outliers ⑫ False data ⑬ Duplicate data

Human Error:- Even with all the automation data is still typed on various web interfaces, hence there is a high possibility of typographical mistakes, leading to inaccurate data. This data entry can be done both by customer and the employee. Customer may write the correct data into the wrong data field. Similarly, employee may take a mistake while handling or migrating the data.

Data Duplication:- Nowadays, data comes from multiple channels giving rise to duplicate data when merged. It results in multiple variations of the same record providing skewed analytical results and incorrect insights. The budget is also wasted on this duplicate records. We can make use of Data Duplicate Tools to find similar type of records and fix them.

Inaccurate & Missing Data:- Inaccurate data can seriously impact the decision making holding the business to achieve their goals. It is tough to detect because the former will be caused by the spelling mistake.

Classification

(2) Write short notes on following application in terms of goals and approaches:-

- ① Loan default prediction
- ② Credit card fraud detection.

Loan default prediction:-

Goal:- Predict loan defaulter-

Approach:-

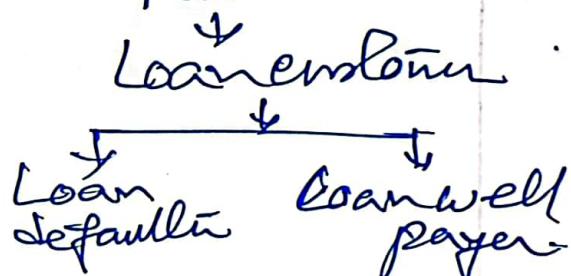
- i) Analysis bank transaction statement
- ii) Regularly payment / irregular payment history
- iii) Irregular payment history continuously or sudden period.
- iv) If sudden period it is not detected.
- v) If irregular payment history in continuous to detect loan defaulter.

Credit Card Fraud Detection:-

Goal:- Predict fraudulent entry in credit card transaction.

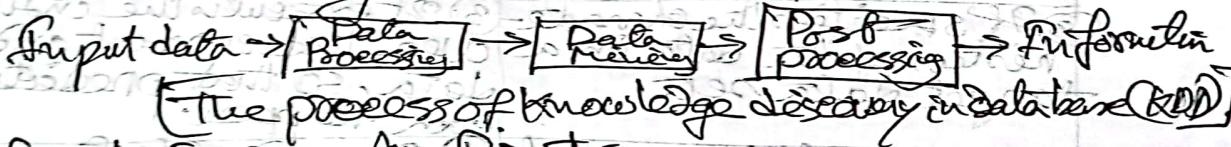
Approach:-

- i) Use credit card transaction and the information on its accounts held as all item
- ii) Previously Transaction History as fraud or fair.
This from the class attributes.
- iii) Learn a model for the class of Transaction pattern.
- iv) Use this data



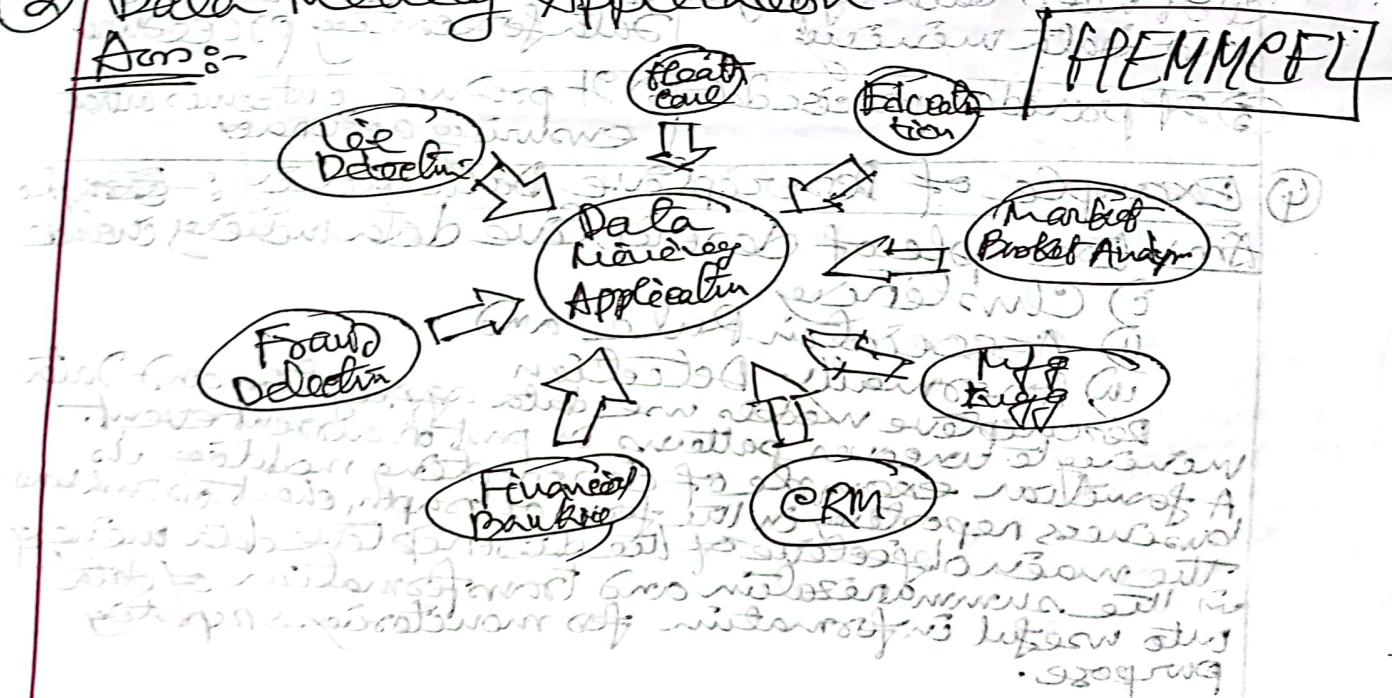
Q) What is Data mining?

Ans:- The process of extracting information to identify patterns, trends and useful data that would allow the business the business to take the data-driven decision from huge sets of data is called Data mining.



② **Data Mining Application**

Ans:-



③ Difference Between Descriptive & Predictive Data Mining

Descriptive Data Mining

- ① Descriptive mining is usually used to provide correlation, cross-tabulation frequency etc.

Predictive Data Mining

- ① The term Predictive means to predict something. So predictive data mining is the analysis done to predict the future event or other data or trend.

- ② It is based on the reactive approach.

- ② It is based on the proactive approach.

- ③ It specifies the characteristics of the data in a target data set.

- ③ It executes the induction over the current and past data, so that prediction can happen.

- ④ It need data aggregation and data mining

- ④ It need statistics and data forecasting procedure.

- ⑤ It provides precise data

- ⑤ It produce outcomes after ensuring accuracy

④ Example of Descriptive Data Mining:

Ans: Examples of descriptive data mining include

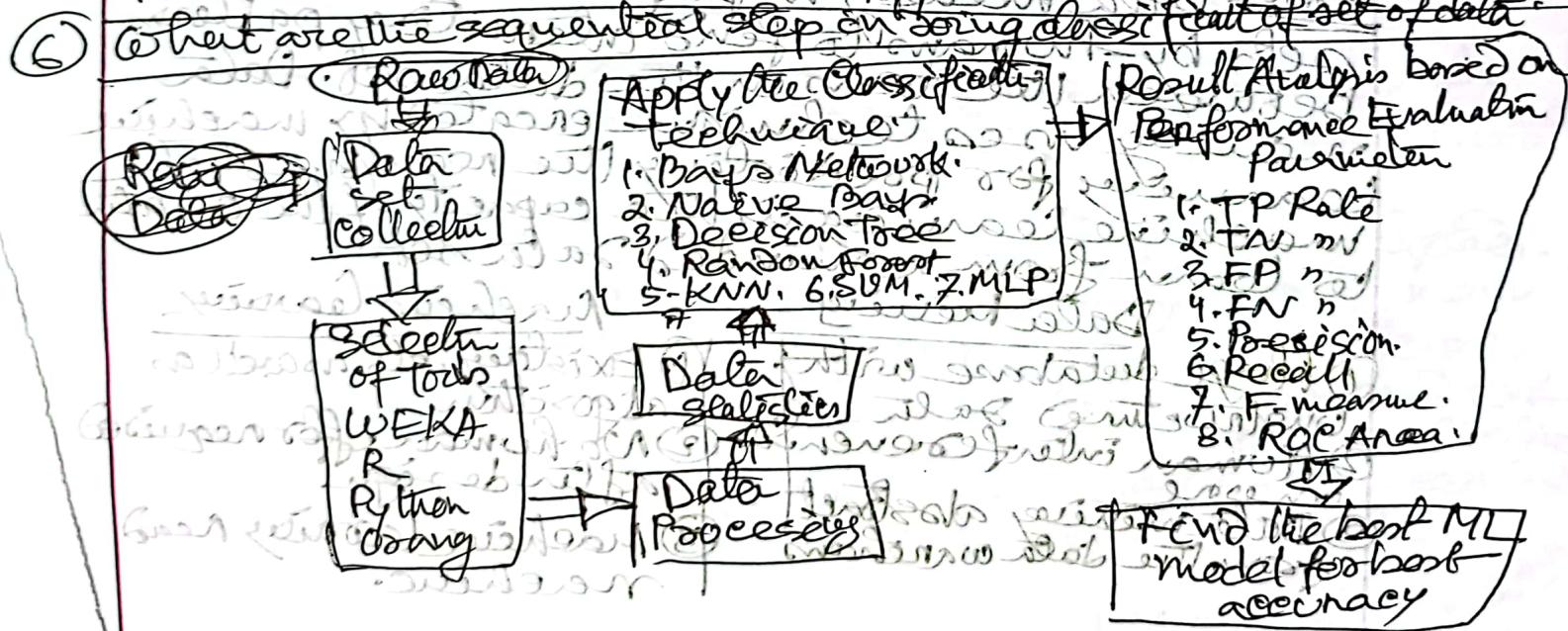
i) Clustering

ii) Association Rule and

iii) Anomaly Detection

Descriptive models use data aggregation and data mining to uncover patterns in past or current event. A familiar example of descriptive modeling is business reporting in the form of graph, chart and dashboard. The main objective of the descriptive data mining is the summarization and transformation of data into useful information for monitoring & reporting purpose.

	Classification	Clustering
①	Classification is used to label data.	① Clustering is used to group similar data instance together.
②		
③	The no. of classes is known.	③ The no. of classes is unknown.
④	Training data is required.	④ No training data is required.
⑤	Process on the training data, the classification model is used to classify future instance into already defined classes.	⑤ Clustering is used to make sense of existing data.
⑥	Popular algorithm for classification include Naive Bayes classifier, Decision Tree and Random Forests.	⑥ Popular algorithm used for clustering include k-means, Means Shift clustering and Density Based Spatial Clustering of Applications with Noise.



Q) What is the Associate Rule Discovery with example?

Ans- Association rule learning is a rule-based machine learning method for discovering interesting relations between large database. It is intended to identify strong rules discovered in database using some measures of interestingness.

Associate rules are useful data mining for analyzing and forecasting customer behaviour. For example, peanut butter and jelly are frequently purchased together because a lot of people like to make PB&J sandwiches.

Q) In data mining is a machine learning?

Ans- Data mining is performed on certain data sets by humans to find interesting patterns between the items in the data set. Data mining uses techniques created by machine learning for predicting the results. So here's machine learning is the capacity of the computer to learn from a intended data set.

Data mining - Machine Learning

- | | |
|---|--|
| ① Huge database with
structured data | ① Existing data as well as
algorithm |
| ② Human intervention
more | ② No human effort required
after design |
| ③ Data mining abstract
from the data仓库 | ③ Machine learning need
machines |

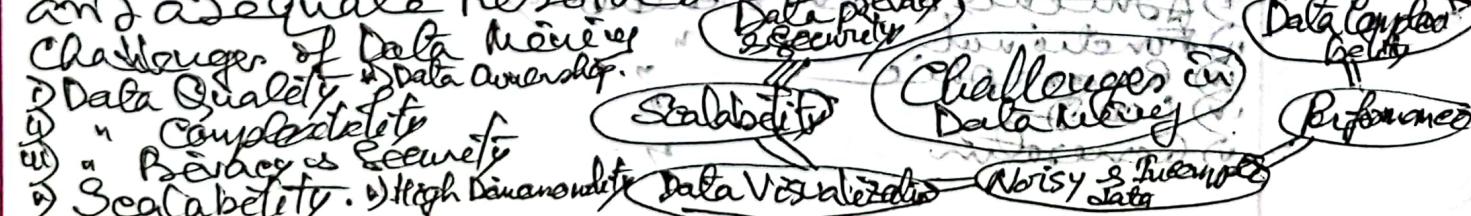
why do we need study more on data mining.

- Ans:- We need more study on data mining in the following reasons:-
- i) The data mining techniques enables organization to obtain knowledge-based data.
 - ii) Data mining enables organizations to make iterative modifications at operation and production.
 - iii) Coupled with other statistical data applications, data mining is a cost-efficient.
 - iv) Data mining helps the decision-making process of an organization.
 - v) It facilitates the automated discovery of hidden patterns as well as the prediction of trends and behaviors.
 - vi) It can be induced in the new system as well as the existing platforms.
 - vii) It is a quick process that makes it easy for new users to analyze enormous amounts of data in a short time.

Ques 16)

Challenges of implementation data mining, Explain.

Ans:- Data mining is very powerful, it faces many challenges during its execution. Various challenges could be related to performance, data, methods and tools etc. The process of data mining becomes effective when the challenges or problems are correctly recognized and adequate resolved.



Q1

Explain reasons for Data affected.

Ans: There are various types of reasons for affected or loss of data as follows:-

- i) Hardware failure.
- ii) Human Errors.
- iii) Malware.
- iv) Misplacing the Device or the data.
- v) Software causes.
- vi) Natural Disasters.
- vii) Hackers.
- viii) Miscellaneous Issues.

Q2

What is Scalability? Scalability is the property of a system handle a growing amount of work. The definition of software system specifies that, this may be done by adding resources to the system.

In an economic context, a scalable business model implies that company can increase sales given increase resources. In example: a package delivery system is scalable because more packages can be delivered by adding more delivery vehicles. However, if all packages had to first pass through a single warehouse (for renting), the system would not be scalable, because one warehouse can handle only a limited number of packages.

Scalability can be measured over multiple dimensions:-

- i) Administrative Scalability
- ii) Functional Scalability
- iii) Geographic Scalability
- iv) Load Scalability
- v) Generational Scalability

Categorisation of Attributes :-

Quantitative Data :- Quantitative means we can count it and its numerical.

- i) Discrete Data - Involved whole number. Example - People in class, Finger in your hand, No. of children etc.
- ii) Continuous Data - It can be divided up as much as we want, and measure to many decimal places. Example - Weight of the car, Temperature, Speed of an Aeroplane etc.

Qualitative Data :- Qualitative means we can't count it and its not numerical.

- i) Nominal Data - Of used to label variables without any quantitative values. Example (=, +).
 - a) Male/Female (Sex)
 - b) Black/White
 - c) Nationalities
 - d) Name of people
- ii) Ordinal Data - Of used for measure of satisfaction, happiness, etc. Example: How likely are you recommend our services (likely, very likely, unlikely).
- iii) Interval Data - Of allows to measure standard deviation and central tendencies. Example - Temperature in degree celsius. 20°C warmer than 10°C and the difference between 10°C and 0°C is also 10°C . The difference between 20°C and 10°C is 10°C . The difference between 10°C and 0°C is also 10°C .
- iv) Ratio Data - Ratio data is very similar interval data, except zero means none. For ratio data, it is not possible to have negative values. Example - Height is not possible to have negative height. If an object height is zero, then there is no object. This is different than something like temperature. Both 0°C and -5°C are completely valid and meaningful temperature ($*, /$). Current age, mass, length, Temperature Kelvin

(Q)

Is temperature an interval or ratio measurement statistically? Justify your answer.

Ans: Temperature can be either an interval or ratio of attributes measurement statistically, depending on its measurement scale. When measured on the Kelvin scale, a temperature of 2 is, in a physically meaningful way, twice that of a temperature of 1. This is not true when temperature is measured on either the Celsius or Fahrenheit scale, as an interval scale. Because zero is not the lowest possible temperature. In the Kelvin scale, a ratio scale, zero represents a total lack of thermal energy.

Fahrenheit temperature is not physically meaningful.

(S)

list the different types of attributes with properties.

Ans: Different types of attributes with their properties follow:

- i) Nominal Attributes \rightarrow Distinctness ($=, \neq$) \rightarrow (D, eye color, zip)
- ii) Ordinal Attributes \rightarrow Order ($<, >$) \rightarrow Ranking, Grade, floors
- iii) Interval Attributes \rightarrow meaningful diff ($+/-$) \rightarrow Calendar date, Temp (c/f)
- iv) Ratio ($\times, /, =$) \rightarrow length, count, time

v) Ratio

Q4-4,5,6,7 (+ CT-2)

- 1 ① How a tree based classification works? What are the advantages of tree classification.
- 2 Define the terms: True positive, and false negative & F-measure.
- 3 Write the procedures of tree based classification how does one can check the validity of a tree?
- 4 When do we stop constructing classification tree?
- 5 Among tree based classification & Rule-based ones, which one will you prefer, why?
- 6 When do we use learning curve and ROC curve?
- 7 Explain the Ensemble method of classification.
- 8 What is OLAP? Why do researchers use Slicing and Dicing in Data analysis.
- 9 Data may be affected by various kinds of reason - These reasons are may become in data quality problem.
- 10 Why EDA is necessary for Data Mining. Why do we use cross validation in classification.
- 11 Define the terms EDA, why EDA is necessary for data mining?
- 12 What is the main principle of Gene Endex, explain. Whether have to use Gene Endex in splitting.
- 13 Why do we use data aggregation and dimension reduction.
- 14 Give an example of stratified and write its purpose of use in data mining.
- 15 When do we need discretization and binarization of data.

Data Processing Chapter 3.

What do we use data aggregation?

Ans:- We use data aggregation for below purposes:-

- i) Data reduction - Reduce the number of entities or objects.
- ii) Change of Scale - Cities aggregated into regions, states, countries etc. Days aggregated into weeks, months, or years.
- iii) More Stable Data - Aggregated data tends to have less variability.

Types of Sampling:-

- A) Sample Random Sampling
 - i) Sampling without replacement.
 - ii) Sampling with replacement.
- B) Stratified Sampling.

Dimensionality Reductions-

Purposes:-

- i) Avoid curse of dimensionality
- ii) Reduce amount of time and memory required by data mining algorithm
- iii) Allow data to be more easily visualized
- iv) May help to eliminate irrelevant features or reduce noise.

Techniques:-

- i) Principle Component Analysis (PCA)
- ii) Singular value decomposition
- iii) Other: supervised and non-linear techniques.

Step for PCA Algorithm:-

- 1 Getting the dataset
- 2 Representing data into a structure
- 3 Standardizing the data
- 4 Calculating the Covariance matrix
- 5 Calculating the Eigenvalues and Eigenvectors
- 6 Sorting the Eigen Vectors
- 7 Calculating the new features or Principal components.

Chapler 9. Similarity and Dissimilarity Measures.

Similarity Measures:-

- Numerical measure of how alike two data objects are
- If higher then objects are more alike.
- Often falls in the range [0, 1]

Dissimilarity Measures:-

- Numerical measure of how different two data objects are.
- Lower when objects are more alike.
- Minimum dissimilarity is often 0.
- Upper limit varies.
- Proximity refers to a similarity or dissimilarity.

Similarity/Dissimilarity for Sample Attributes

The following table shows the similarity and dissimilarity between two objects, x and y , with respect to a single sample attribute $(x - y) = (b)$

Attribute Type	Dissimilarity	Similarity
Nominal $\rightarrow d = 0 \text{ if } x=y \text{ else } 1$	$S=1 \text{ if } x=y \text{ else } 0$	
Ordinal $\rightarrow d = x-y /(n-1)$ (values mapped to integers 0 to $n-1$)	$S=1-d$	
Interval or Ratio $\rightarrow d = x-y /l$	$S=1-d, S=1+d, S=1-d \text{ when } d \neq l$	

Dissimilarity - Euclidean Distance :-

$$d(x, y) = \sqrt{\sum_{k=1}^n (x_k - y_k)^2} = \sqrt{(x_{i1} - y_{i1})^2 + (x_{i2} - y_{i2})^2}$$

where n is the no. of dimensions (attribute) and x_k and y_k are, respectively, the k th attribute (components) of data objects x and y .

Dissimilarity - Minkowski Distance :-

$$d(x, y) = \left(\sum_{k=1}^n |x_k - y_k|^p \right)^{1/p} = (|x_{i1} - y_{i1}|^p + |x_{i2} - y_{i2}|^p + \dots)^{1/p}$$

where p is the order parameter, n is the no. of dimensions (attribute) and x_k and y_k are respectively, the k th attribute (components) of data objects x and y .

Manhattan Distance :- $d(x, y) = |x_{i1} - y_{i1}| + |x_{i2} - y_{i2}| + \dots + |x_{in} - y_{in}|$

Examples: - Give Distances

Given two objects represented by the tuples $(22, 1, 42, 10)$ and $(20, 0, 36, 8)$:

a) Compute the Euclidean distance between the objects

" Manhattan

" Minkowski

$$\text{Distance } p=3: \sqrt[3]{(x_1 - x_2)^3 + (x_2 - x_3)^3 + (x_3 - x_4)^3}$$

$$d(c, d) = \sqrt[3]{(22 - 20)^3 + (1 - 0)^3 + (42 - 36)^3 + (10 - 8)^3} = 0.71$$

b) Manhattan Distance

$$d(c, d) = |x_1 - x_2| + |x_2 - x_3| + |x_3 - x_4| + |x_4 - x_1|$$

$$= |22 - 20| + |1 - 0| + |42 - 36| + |10 - 8|$$

c) Minkowski Distance

$$d(c, d) = \sqrt[p]{(x_1 - x_2)^p + (x_2 - x_3)^p + (x_3 - x_4)^p + (x_4 - x_1)^p}$$

$$= \sqrt[p]{(22 - 20)^p + (1 - 0)^p + (42 - 36)^p + (10 - 8)^p}$$

- Example - Euclidean

$$d(p, k) = \sqrt{(p_1 - k_1)^2 + (p_2 - k_2)^2 + \dots + (p_n - k_n)^2}$$

Distance between p and k is the Euclidean Distance:

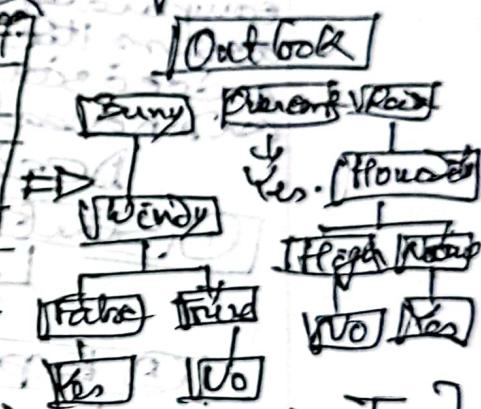
$$d(p, k) = \sqrt{(p_1 - k_1)^2 + (p_2 - k_2)^2 + \dots + (p_n - k_n)^2} = \sqrt{(p_1 - k_1)^2 + (p_2 - k_2)^2 + \dots + (p_n - k_n)^2}$$

Chapter-5.

What is decision tree Classification?

Ans:- Decision tree builds classification or regression models in the form of a tree structure. It breakdown a dataset into smaller and smaller subsets while at the same time an associated decision tree is incrementally developed. The final result is a tree with decision nodes and leaf nodes. Example: A decision node (e.g. Outlook) has two or more branches (e.g. Sunny, Overcast, Rain). Leaf node (e.g. Play) represents a classification or decision. The topmost decision node in a tree which corresponds to the best predictor called root node. Decision trees can handle both categorical and numerical data.

outlook	temp	Humidity	Windy	Play
Sunny	High	High	False	No
Sunny	High	Low	False	No
Sunny	Low	High	True	No
Rainy	High	High	False	Yes
Rainy	High	Low	False	Yes
Rainy	Low	High	True	Yes
Overcast	Low	Low	False	Yes



For what aspect, you may like decision tree classification? [Decision Tree]

Ans:- A decision tree could be used to help a company decide which city to move its headquarters or whether to open a satellite office. Decision trees are popular tools in machine learning, as they can be used to build predictive models.

2. D.P.

How does decision tree works?

Sol: A decision tree is a flow chart that starts with one main idea and then branches out based on the consequences of your decision. It's called a decision tree because the model typically looks like a tree with branches. These trees are used for decision tree analysis, which involves visually outlining the potential outcomes, cost and consequence of a complex decision. You can use a decision tree to calculate the expected value of each outcome based on the decision and consequence. If it led to another decision, you can compare the outcomes to one another, you can also compare the outcomes to one another, you can also assess the best course of action. You can also use a decision tree to solve problems, manage cost, and reveal opportunities.

Decision tree analysis work in five steps:-

① Start with your idea

② Add chance and decision nodes

③ Expand until you reach end point

④ Calculate Tree Value

⑤ Evaluate outcomes

How does the Decision tree algorithm work?

Ans:- In a decision tree, for predicting the class of the given dataset, the algorithm starts from the root node of the tree. This algorithm compares the values of root attribute with the record (real dataset) attribute and, based on the comparison, follows the branch and jumps to the next node. For the next node, the algorithm again compare the attribute value with the other sub-nodes and move further. It continues the process until it reaches the leaf node of the tree. The complete process can be better understood, using the below algorithm.

Step-1 - Begin the tree with the root node, which contains the complete dataset.

Step-2 - Find the best attribute in the dataset using Attribute Selection Measure (GSM).

Step-3 - Divide the S into subsets that contains possible values for the best attribute.

Step-4 - Generate the decision tree using the subset of the dataset node which contains the best attribute.

Step-5 - Recursively make new decision trees using the subset created in step-3.

Examples:-

Advantages of Decision Trees :-

- ① It is simple to understand & easy to learn.
- ② It can be very useful for solving decision-related problems.
- ③ It helps to think about all the possible outcomes for problem.
- ④ There is less requirement of data cleaning compared to other algorithms.

Attribute Selection measure: - There are popular technique for Attribute Selection Measure (ASM), which are based upon Information Gain.

a) Gini's Index

Information Gain: - Information measures the measurement of changes in entropy after the segmentation of a dataset based on an attribute.

It can be calculated by below formula:

$$\text{Information Gain} = \text{Entropy}(S) - (\text{P}_1 \text{Entropy}(P_1) + \text{P}_2 \text{Entropy}(P_2))$$

Entropy: - Entropy is a metric to measure the impurity on a given attribute. It specifies randomness in data. Entropy can be calculated as

$$\text{Entropy}(S) = - P_{\text{Yes}} \log_2 P_{\text{Yes}} - P_{\text{No}} \log_2 P_{\text{No}}$$

Gini Index: - Gini Index is a measure of impurity used while creating a decision tree in the CART (Classification and Regression Tree) algorithm.

$$\text{Gini Index} = 1 - \sum P_j^2$$

Stopping Criteria

- i) Number of entries in the node is less than some pre-specified limit.
- ii) Purity of the node is more than some specified limit.
- iii) Depth of the node is more than some specified limit.
- iv) Predictor value for all records are identical.

→ Classification :- Lec-05.

Definition :- Classification is the process of identifying and grouping objects or ideas into predetermined categories.



Classification Task :- Binary Classification can be categorised into one of two classes.

- i) Binary Classification & (No. of class is longer than 2).
- ii) Multiclass Classification

A classification model serves two important roles in data mining & must provide accurate prediction with a fast response time.
i) Predictive model (To classify previously unlabeled instances).
ii) Descriptive model (To identify the characteristics that distinguish instances from different classes).
It is particularly used for medical application, such as medical diagnosis, it help to determine characteristics.

Training Set :- A classifier model is created using a given set of instances, known as the training set.

Learning Algorithm :- The systematic approach for learning a classification model given a training set is known as a learning algorithm.

Induction :- The process of creating learning algorithm to build a classification model from the training data is known as induction.
Deduction :- The process of applying a classification model on unseen test instances to predict their class labels is known as deduction.

97	IT
WT	WT

Classification Techniques:- A classification technique refers to a general approach to classification:-

(A) Base Classification

- i) Decision Tree Based Methods
- ii) Rule Based Methods
- iii) Nearest Neighbor
- iv) Naïve Bayes and Bayesian Belief Network.
- v) Support Vector Machines.
- vi) Neural & Deep Neural Network.

(B) Ensemble Classifier

- i) Boosting
- ii) Bagging.
- iii) Random Forests.

Confusion Matrix:- The performance of a classifier model can be evaluated by comparing the predicted labels against the true labels of instances. This information can be summarized in a table called a confusion matrix. A confusion matrix is nothing but a table with two dimensions.

(i) Actual and

(ii) Predicted

And further more, both the dimension have "True Positive (TP)", True Negative (TN), False Positive (FP) and False Negative (FN).

	TP	FP
	FN	TN

TP - It is the case when both actual class & predicted class of data point is 1.

TN - It is the case when both actual class & predicted class of data point is 0.

FP - It is the case when actual class of data point is 0 & predicted class of data point is 1.

FN - It is the case when actual class of data point is 1 & predicted class of data point is 0.

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + FN}$$

$$\text{Error Rate} = \frac{\text{No. of wrong prediction}}{\text{Total No. of prediction}}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$F_1 \text{ Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

$$\text{Recall or Sensitivity} = \frac{TP}{TP + FN}$$

$$\text{Decision Tree Induction}$$

Neary Algorithms

(i) Hunt's Algorithm - (ii) H. vector

(i) CART

(ii) ID 3, C4.5

(iii) SLIQ, SPRPART (are in epoch 1)
(iv) C4.5 (is in epoch 2)

(i) Decision tree (ii)

(iii) naive bayes (iv)

Design Issues of Decision Tree Induction

- > how should training record be split?
 - method for expressing test condition
 - depending on attribute types,
 - measure for evaluating the goodness of a test of condition
 - > how should the splitting procedure stop?
 - Stop splitting if all the records belong to the same class or have identical attribute values
 - Early termination
- Test Condition for Nominal/ordinal/Attributes
- (A) Multway Split - Use many partitions and test values.
 - (B) Binary Split - Divides values into two subsets.

Measures of Node Impurity

- ① Gini Index
- ② Entropy $H(S) = \sum_{t=0}^{C-1} P_t \log P_t = -P_0 \log P_0 - P_1 \log P_1 - P_2 \log P_2$
- ③ Misclassification Error = $1 - \max_t P_t$
 - i) Entropy is zero (minimum) when all the examples same class.
 - ii) Entropy is 1 (maximum) when all are equal w.r.t. classes.
- ④ Average Entropy
- ⑤ Information Gain (IG)

What are the steps in ID3 algorithm?

The steps in ID3 algorithm are as follows:

1. Calculate entropy for dataset.
2. For each attribute/feature.
 - 2.1. Calculate entropy for all its categorical values.
 - 2.2. Calculate information gain for the feature.
3. Find the feature with maximum information gain.
4. Repeat it until we get the desired tree.

Use ID3 algorithm on a data We'll discuss it here mathematically and later see its implementation in Python. So, Let's take an example to make it more clear.

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

Here, dataset is of binary classes (yes and no), where 9 out of 14 are "yes" and 5 out of 14 are "no".

Complete entropy of dataset is:

$$H(S) = - p(\text{yes}) * \log_2(p(\text{yes})) - p(\text{no}) * \log_2(p(\text{no})) \\ = - (9/14) * \log_2(9/14) - (5/14) * \log_2(5/14) = - (-0.41) - (-0.53) = 0.94$$

For each attribute of the dataset, let's follow the step-2 of pseudocode : -First Attribute - Outlook.

Categorical values - sunny, overcast and rain

$$H(\text{Outlook}=\text{sunny}) = -(2/5) * \log_2(2/5) - (3/5) * \log_2(3/5) = 0.971$$

$$H(\text{Outlook}=\text{rain}) = -(3/5) * \log_2(3/5) - (2/5) * \log_2(2/5) = 0.971$$

$$H(\text{Outlook}=\text{overcast}) = -(4/4) * \log_2(4/4) = 0 = 0$$

Average Entropy Information for Outlook -

$$I(\text{Outlook}) = p(\text{sunny}) * H(\text{Outlook}=\text{sunny}) + p(\text{rain}) * H(\text{Outlook}=\text{rain}) + p(\text{overcast}) * H(\text{Outlook}=\text{overcast}) \\ = (5/14) * 0.971 + (5/14) * 0.971 + (4/14) * 0 = 0.693$$

$$\text{Information Gain} = H(S) - I(\text{Outlook}) = 0.94 - 0.693 = 0.247$$

I suggest you to do the same calculations for all the remaining attributes without scrolling down....

Second Attribute - Temperature, Categorical values - hot, mild, cool

$$H(\text{Temperature}=\text{hot}) = -(2/4) * \log_2(2/4) - (2/4) * \log_2(2/4) = 1$$

$$H(\text{Temperature}=\text{cool}) = -(3/4) * \log_2(3/4) - (1/4) * \log_2(1/4) = 0.811$$

$$H(\text{Temperature}=\text{mild}) = -(4/6) * \log_2(4/6) - (2/6) * \log_2(2/6) = 0.9179$$

Average Entropy Information for Temperature -

$$I(\text{Temperature}) = p(\text{hot}) * H(\text{Temperature}=\text{hot}) + p(\text{mild}) * H(\text{Temperature}=\text{mild}) + p(\text{cool}) * H(\text{Temperature}=\text{cool}) = (4/14) * 1 + (6/14) * 0.9179 + (4/14) * 0.811 = 0.9108$$

$$\text{Information Gain} = H(S) - I(\text{Temperature}) = 0.94 - 0.9108 = 0.0292$$

Third Attribute - Humidity, Categorical values - high, normal

$$H(\text{Humidity}=\text{high}) = -(3/7) * \log(3/7) - (4/7) * \log(4/7) = 0.983$$

$$H(\text{Humidity}=\text{normal}) = -(6/7) * \log(6/7) - (1/7) * \log(1/7) = 0.591$$

Average Entropy Information for Humidity -

$$I(\text{Humidity}) = p(\text{high}) * H(\text{Humidity}=\text{high}) + p(\text{normal}) * H(\text{Humidity}=\text{normal}) \\ = (7/14) * 0.983 + (7/14) * 0.591 = 0.787$$

$$\text{Information Gain} = H(S) - I(\text{Humidity}) = 0.94 - 0.787 = 0.153$$

Fourth Attribute - Wind

Categorical values - weak, strong

$$H(\text{Wind}=\text{weak}) = -(6/8) * \log(6/8) - (2/8) * \log(2/8) = 0.811$$

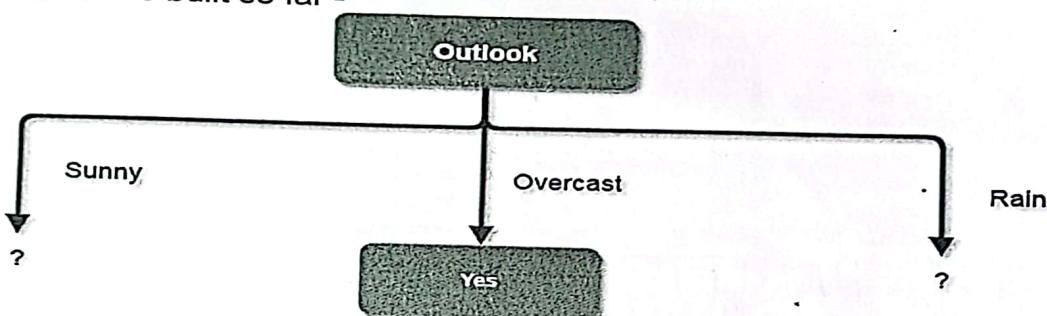
$$H(\text{Wind}=\text{strong}) = -(3/6) * \log(3/6) - (3/6) * \log(3/6) = 1$$

Average Entropy Information for Wind -

$$I(\text{Wind}) = p(\text{weak}) * H(\text{Wind}=\text{weak}) + p(\text{strong}) * H(\text{Wind}=\text{strong}) \\ = (8/14) * 0.811 + (6/14) * 1 = 0.892$$

$$\text{Information Gain} = H(S) - I(\text{Wind}) = 0.94 - 0.892 = 0.048$$

Here, the attribute with maximum information gain is Outlook. So, the decision tree built so far -



Here, when Outlook == overcast, it is of pure class(Yes).

Now, we have to repeat same procedure for the data with rows consist of Outlook value as Sunny and then for Outlook value as Rain.

I again recommend you to perform further calculations and then cross check by scrolling down...

Now, finding the best attribute for splitting the data with Outlook=Sunny values{ Dataset rows = [1, 2, 8, 9, 11]}.

Complete entropy of Sunny is -

$$H(S) = - p(\text{yes}) * \log_2(p(\text{yes})) - p(\text{no}) * \log_2(p(\text{no})) \\ = - (2/5) * \log_2(2/5) - (3/5) * \log_2(3/5) = 0.971$$

First Attribute - Temperature

Categorical values - hot, mild, cool

$$H(\text{Sunny}, \text{Temperature}= \text{hot}) = -0 - (2/2) * \log(2/2) = 0$$

$$H(\text{Sunny}, \text{Temperature}= \text{cool}) = -(1) * \log(1) - 0 = 0$$

$$H(\text{Sunny}, \text{Temperature}= \text{mild}) = -(1/2) * \log(1/2) - (1/2) * \log(1/2) = 1$$

Average Entropy Information for Temperature -

$$I(\text{Sunny}, \text{Temperature}) = p(\text{Sunny, hot}) * H(\text{Sunny, Temperature}= \text{hot}) + p(\text{Sunny, mild}) * H(\text{Sunny, Temperature}= \text{mild}) + p(\text{Sunny, cool}) * H(\text{Sunny, Temperature}= \text{cool}) = (2/5) * 0 + (1/5) * 0 + (2/5) * 1 = 0.4$$

$$\text{Information Gain} = H(\text{Sunny}) - I(\text{Sunny}, \text{Temperature}) = 0.971 - 0.4 = 0.571$$

Second Attribute - Humidity

Categorical values - high, normal

$$H(\text{Sunny, Humidity=high}) = -(3/3) * \log(3/3) = 0$$

$$H(\text{Sunny, Humidity=normal}) = -(2/2) * \log(2/2) = 0$$

Average Entropy Information for Humidity -

$$I(\text{Sunny, Humidity}) = p(\text{Sunny, high}) * H(\text{Sunny, Humidity=high}) + p(\text{Sunny, normal}) * H(\text{Sunny, Humidity=normal}) = (3/5) * 0 + (2/5) * 0 = 0$$

$$\text{Information Gain} = H(\text{Sunny}) - I(\text{Sunny, Humidity}) = 0.971 - 0 = 0.971$$

Third Attribute - Wind

Categorical values - weak, strong

$$H(\text{Sunny, Wind=weak}) = -(1/3) * \log(1/3) - (2/3) * \log(2/3) = 0.918$$

$$H(\text{Sunny, Wind=strong}) = -(1/2) * \log(1/2) - (1/2) * \log(1/2) = 1$$

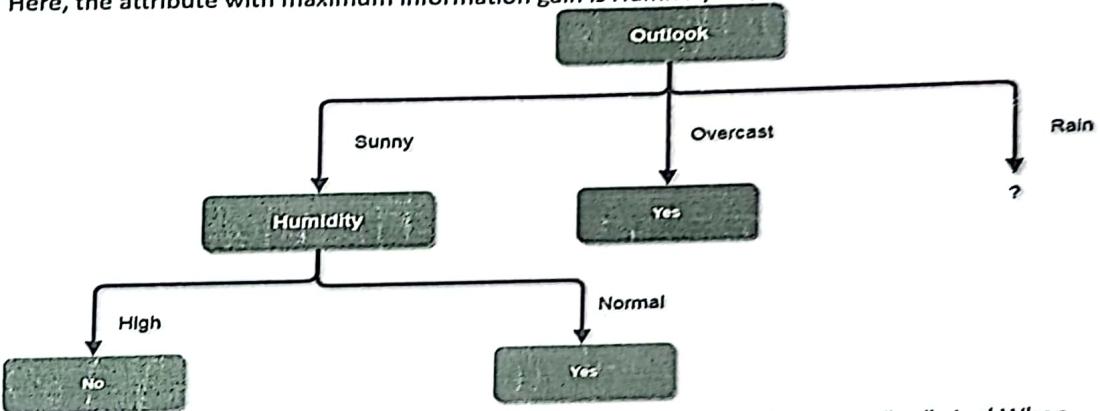
Average Entropy Information for Wind -

$$I(\text{Sunny, Wind}) = p(\text{Sunny, weak}) * H(\text{Sunny, Wind=weak}) + p(\text{Sunny, strong}) * H(\text{Sunny, Wind=strong})$$

$$= (3/5) * 0.918 + (2/5) * 1 = 0.9508$$

$$\text{Information Gain} = H(\text{Sunny}) - I(\text{Sunny, Wind}) = 0.971 - 0.9508 = 0.0202$$

Here, the attribute with maximum information gain is Humidity. So, the decision tree built so far -



Here, when Outlook = Sunny and Humidity = High, it is a pure class of category "no". And When Outlook = Sunny and Humidity = Normal, it is again a pure class of category "yes". Therefore, we don't need to do further calculations.

Now, finding the best attribute for splitting the data with Outlook=Sunny values {Dataset rows = [4, 5, 6, 10, 14]}.

Complete entropy of Rain is -

$$H(S) = - p(\text{yes}) * \log_2(p(\text{yes})) - p(\text{no}) * \log_2(p(\text{no}))$$

$$= - (3/5) * \log_2(3/5) - (2/5) * \log_2(2/5) = 0.971$$

First Attribute - Temperature

Categorical values - mild, cool

$$H(\text{Rain, Temperature=cool}) = -(1/2) * \log_2(1/2) - (1/2) * \log_2(1/2) = 1$$

$$H(\text{Rain, Temperature=mild}) = -(2/3) * \log_2(2/3) - (1/3) * \log_2(1/3) = 0.918$$

Average Entropy Information for Temperature -

$$I(\text{Rain, Temperature}) = p(\text{Rain, mild}) * H(\text{Rain, Temperature=mild}) + p(\text{Rain, cool}) * H(\text{Rain, Temperature=cool}) = (2/5) * 1 + (3/5) * 0.918 = 0.9508$$

Information Gain = $H(\text{Rain}) - I(\text{Rain}, \text{Temperature}) = 0.971 - 0.9508 = 0.0202$

Second Attribute - Wind

Categorical values - weak, strong

$$H(\text{Wind=weak}) = -(3/3) * \log(3/3) * 0 = 0$$

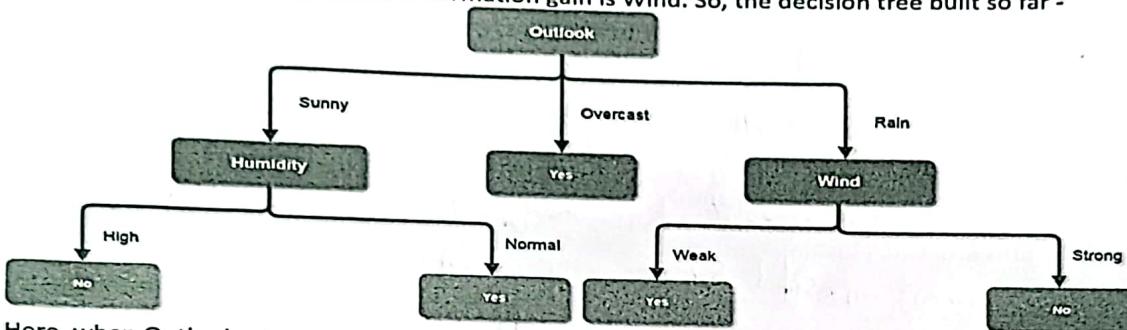
$$H(\text{Wind=strong}) = 0 - (2/2) * \log(2/2) = 0$$

Average Entropy Information for Wind -

$$I(\text{Wind}) = p(\text{Rain, weak}) * H(\text{Rain, Wind=weak}) + p(\text{Rain, strong}) * H(\text{Rain, Wind=strong}) = (3/5) * 0 + (2/5) * 0 = 0$$

Information Gain = $H(\text{Rain}) - I(\text{Rain, Wind}) = 0.971 - 0 = 0.971$

Here, the attribute with maximum information gain is Wind. So, the decision tree built so far -



Here, when Outlook = Rain and Wind = Strong, it is a pure class of category "no". And When Outlook = Rain and Wind = Weak, it is again a pure class of category "yes".

And this is our final desired tree for the given dataset.

What are the characteristics of ID3 algorithm?

Finally, I am concluding with Characteristics of ID3.

Characteristics of ID3 Algorithm are as follows:

ID3 uses a greedy approach that's why it does not guarantee an optimal solution; it can get stuck in local optimums.

ID3 can overfit to the training data (to avoid overfitting, smaller decision trees should be preferred over larger ones).

This algorithm usually produces small trees, but it does not always produce the smallest possible tree.

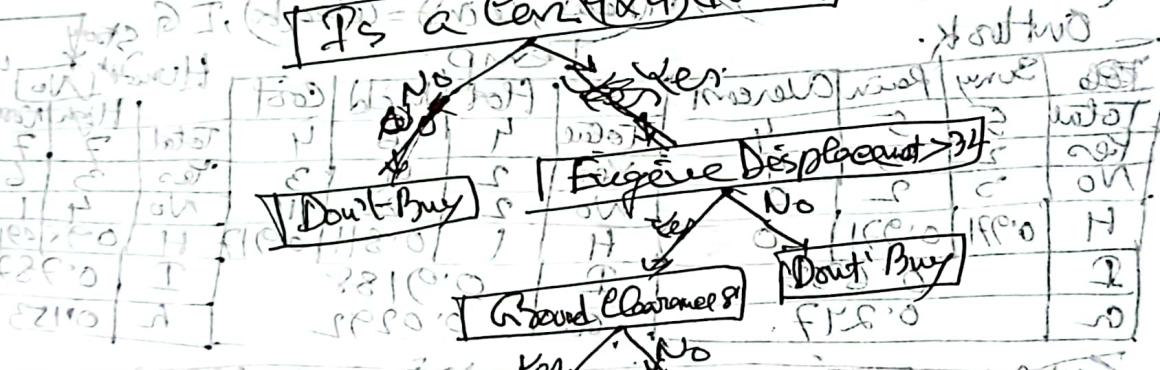
ID3 is harder to use on continuous data (if the values of any given attribute is continuous, then there are many more places to split the data on this attribute, and searching for the best value to split by can be time consuming).

With this technical article at OpenGenus, you must have a strong idea of ID3 Algorithm and how to use ID3 Algorithm to build a Decision Tree to predict the weather.

Decision Tree - Decision tree is a popular machine learning algorithm mainly used for classification. Concept of entropy and information gain are required to apply decision tree for non-linear classification and regression model.

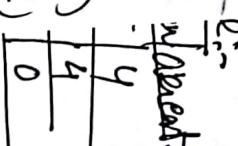
The main idea of a decision tree is to split datasets in many steps based on features until you reach a node which contains data points that fall under the same label.

Decision Tree for Classification and Prediction



Entropy - Entropy is a measure of the amount of uncertainty in the dataset. It is defined as:

Entropy is $H(S) = -\sum p_i \log_2 p_i$. Information Gain (IG) tells us how much uncertainty in S was reduced after splitting set S on attribute A . $IG(A, S) = H(S) - \sum_{t \in T} p(t) H(t)$.



$$H(S) = -\sum p_i \log_2 p_i$$

PD-3

(A) Complete Entropy of $H(3)$

1st AttrIBUTES $\rightarrow H(\text{outlook} = \text{Sunny}), H(\text{outlook} = \text{Rain}), H(\text{outlook} = \text{Cloudy}), P, G$
 2nd " $\rightarrow H(\text{Temp} = \text{Hot}), H(\text{Temp} = \text{Mild}), H(\text{Temp} = \text{Cool}), P, G$
 3rd " $\rightarrow H(\text{Humidity} = \text{High}), H(\text{Humidity} = \text{Normal}), P, G$
 4th " $\rightarrow H(\text{Wind} = \text{Weak}), H(\text{Wind} = \text{Strong}), P, G$ Outlook

(B) Splitting: Outlook (Sunny), Entropy $H(\text{Sunny}, \text{Temp} = \text{Hot}), H(\text{Sunny}, \text{Temp} = \text{Mild}), H(\text{Sunny}, \text{Temp} = \text{Cool}), P, G$ Sun, Overcast

b) $H(\text{Sunny}, \text{Humidity} = \text{High}), H(\text{Sunny}, \text{Humidity} = \text{Normal}), P, G$ Humidity
 c) $H(\text{Sunny}, \text{Wind} = \text{Weak}), H(\text{Sunny}, \text{Wind} = \text{Strong}), P, G$ Wind

(C) Splitting: Outlook (Rain), Entropy $H(\text{Rain}, \text{Temp} = \text{Hot}), H(\text{Rain}, \text{Temp} = \text{Mild}), H(\text{Rain}, \text{Temp} = \text{Cool}), P, G$ Rain

$H(\text{Rain}, \text{Humidity} = \text{High}), H(\text{Rain}, \text{Humidity} = \text{Normal}), P, G$ (Wind)

$H(\text{Rain}, \text{Wind} = \text{Strong}), H(\text{Rain}, \text{Wind} = \text{Weak}), P, G$ Wind

Outlook.

			Temp				Humidity			Wind		
Sunny Rain Cloudy			Hot	Mild	Cool	Total	High	Normal	Low	Yes	No	Strong
Total	5	5	4	Total	1	6	4	Total	7	7	Total	8
Yes	2	3	4	Yes	1	4	3	Yes	3	6	Yes	6
No	3	2	0	No	2	2	1	No	4	1	No	2
H	0.971	0.971	0	H	1	0.811	0.917	H	0.98091	0.0191	H	0.981
P	0.693	0.693	0	P	0.9188			P	0.787	0.212	P	0.787
G	0.247	0.247	0	G	0.0292			G	0.153	0.847	G	0.153

			Humidity				Wind		
Sunny Rain Cloudy			High	Normal	Low	Yes	No	Strong	Weak
Total	2	2	1	3	2	1	2	3	1
Yes	0	1	0	0	2	1	1	1	0
No	2	1	0	3	0	1	1	2	0

			Humidity				Wind		
Rain = Yes Temp = Cool Wind = Weak			High	Normal	Low	Yes	No	Strong	Weak
Total	0	3	2	2	3	2	3	2	3
Yes	0	2	1	1	2	1	1	1	0
No	0	1	1	1	0	1	1	1	0

Continue...

$$f(S) = P_{yes} \log_2 P_{yes} - P_{no} \log_2 P_{no}$$

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
1-D1	Sunny	Hot	High	Weak	No
2-D2	Sunny	Hot	High	Strong	No
3-D3	Overcast	Hot	High	Weak	Yes ✓
4-D4	Rain	Mild	High	Weak	Yes ✓
5-D5	Rain	Cool	Normal	Weak	Yes ✓
6-D6	Rain	Cool	Normal	Strong	No
7-D7	Overcast	Cool	Normal	Strong	Yes ✓
8-D8	Sunny	Mild	High	Weak	No
9-D9	Sunny	Cool	Normal	Weak	Yes ✓
10-D10	Rain	Mild	Normal	Weak	Yes ✓
11-D11	Sunny	Mild	Normal	Strong	Yes ✓
12-D12	Overcast	Mild	High	Strong	Yes ✓
13-D13	Overcast	Hot	Normal	Weak	Yes ✓
14-D14	Rain	Mild	High	Strong	No

Total = 14
Yes = 9.
No = 5

Categorical Value:-

Outlook = Sunny Rain Overcast	Total	5	4
T = 2	3	4	
N = 3	2	0	

Round - 1 :- Outlook: $H(\beta) = -\frac{1}{4} \log_2 \frac{1}{4} - \frac{3}{4} \log_2 \frac{3}{4} = 0.94$

outlook			Temperature			Humidity		Wind.		
Total	Sunny	Rain	overlook	Hot	Mild	Cool.	High	Normal	Weak	Strong
Total	5	5	4	Total	4	6	4	Total	7	7
Yes	2	3	4	Yes	2	4	3	Yes	3	6
No	3	2	0	No	2	2	1	No	4	1
H	0.971	0.971	0	H	1	0.811	0.979	H	0.983	0.591
P	0.693	P		P		0.9188	T	0.787	T	0.892
G.	0.247	a		a		0.0292	a	0.1153	a.	0.048

Round - 2 :- Sunny Value: $H(\beta) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.971$

Round - 3 :- Rain Value:

Round - 4 :- Wind Value:

11/18

Roll-232133.
Md. Monir Ullah

Ans. to the que No-1 -

Qn: Given that

$$TP = 6$$

$$TN = 2$$

$$FP = 1$$

$$FN = 1$$

$$\text{Accuracy} = \frac{TP+TN}{TP+TN+FP+FN}$$

$$= \frac{6+2}{6+2+1+1}$$

$$= 0.8$$

$$\text{Precision} = \frac{TP}{TP+FP} = \frac{6}{6+1} = \frac{6}{7} = 0.857$$

$$\text{Recall} = \frac{TP}{TP+FN} = \frac{6}{6+1} = \frac{6}{7} = 0.857$$

$$F_1 \text{ Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

$$= 2 \times \frac{0.857 \times 0.857}{0.857 + 0.857} = 2 \times \frac{0.7344}{1.714}$$

$$= \frac{1.4688}{1.714} = 0.856$$

To illustrate, suppose S is a collection of 14 examples of some boolean concept, includes 9 positive and 5 negative examples (we adopt notation [9+, 5-] to summarize such a sample of data). Then the entropy of S relative to this boolean classification is

Data Műveg.

08/09/23.

Association Analysis

Association Anayres:
- P. Dules Henley - Support

Association Régionale de la Confédération

Support Count() , Ego-(Friends, Boss, Doctor) = 2

Support, E.g S({Milk, Bread, Diaper}) = 2/5

minsup threshold

P. P. - English in Matrees:- \rightarrow supposed

Rule Evaluation → Confederal (c).

(26) Rule 2018 Accepted

Or ~~the~~ ~~the~~ ~~the~~ support

Confidem.

Baute-force approach:-

$$R = 3^d - 2^{d+1} + 1$$

↳ practical acceptable

Suggest for acceptance

At one 200, Support 80, 60

ଓମ୍ବାଲୁ-କୁ,

Approximate Algorithms \Leftarrow $C = \infty$

Strelk, Detlef = Borsig

$$\frac{G(\text{Mills, Deeper Bed})}{T} = \frac{2}{5} = 0.4$$

Aprioric Algoselection

threshold 50% min

Frequent Item

Mean Support = Total No. of Migrants
= 3

$$c = \alpha$$

Support based proving.

Min Support threshold = 60.

Q. Suppose threshold = 50%; Confidence = 60%.

Solution :-

Support threshold 50%.
 $= 6 \text{ items} \times 50\% = 3.$

Step-1 :- $k=1$.

I_5 deleted, because
 minimum threshold is 3.
 $\text{But } D_5 = 2.$

Step-2 :- $k=2$

Only I_1, I_2, I_3 is frequent.

Generate Association Rule :-

Rule :-

$$I_1, I_2 \Rightarrow I_3$$

$$I_1, I_3 \Rightarrow I_2$$

$$I_2, I_3 \Rightarrow I_1$$

$$\begin{array}{l} I_1, I_2 \\ I_1, I_3 \\ I_2, I_3 \end{array}$$

$$I_1 = \frac{3}{4} = 75\%$$

$$I_2 = \frac{3}{5} = 60\%$$

$$I_3 = \frac{3}{4} = 75\%$$

Table-1. Transaction

Item	Count
I_1	4
I_2	5
I_3	4
I_4	3
I_5	2

Item	Count
I_1, I_2	4
I_1, I_3	3
I_2, I_3	4
I_2, I_4	3
I_3, I_4	2

Item	Count
I_1, I_2, I_3	3
I_1, I_2, I_4	2
I_1, I_3, I_4	1
I_2, I_3, I_4	2

Item	Count
I_1, I_2, I_3, I_4	1

Item	Count
I_1, I_2, I_3	3
I_1, I_2, I_4	2
I_1, I_3, I_4	1
I_2, I_3, I_4	2

Item	Count
I_1, I_2, I_3, I_4	1

Item	Count
I_1, I_2, I_3	3
I_1, I_2, I_4	2
I_1, I_3, I_4	1
I_2, I_3, I_4	2

Item	Count
I_1, I_2, I_3, I_4	1

Item	Count
I_1, I_2, I_3	3
I_1, I_2, I_4	2
I_1, I_3, I_4	1
I_2, I_3, I_4	2

Item	Count
I_1, I_2, I_3, I_4	1

Item	Count
I_1, I_2, I_3	3
I_1, I_2, I_4	2
I_1, I_3, I_4	1
I_2, I_3, I_4	2

Item	Count
I_1, I_2, I_3, I_4	1

Item	Count
I_1, I_2, I_3	3
I_1, I_2, I_4	2
I_1, I_3, I_4	1
I_2, I_3, I_4	2

Item	Count
I_1, I_2, I_3, I_4	1

Item	Count
I_1, I_2, I_3	3
I_1, I_2, I_4	2
I_1, I_3, I_4	1
I_2, I_3, I_4	2

Item	Count
I_1, I_2, I_3, I_4	1

Item	Count
I_1, I_2, I_3	3
I_1, I_2, I_4	2
I_1, I_3, I_4	1
I_2, I_3, I_4	2

Item	Count
I_1, I_2, I_3, I_4	1

Item	Count
I_1, I_2, I_3	3
I_1, I_2, I_4	2
I_1, I_3, I_4	1
I_2, I_3, I_4	2

Item	Count
I_1, I_2, I_3, I_4	1

Item	Count
I_1, I_2, I_3	3
I_1, I_2, I_4	2
I_1, I_3, I_4	1
I_2, I_3, I_4	2

Item	Count
I_1, I_2, I_3, I_4	1

Item	Count
I_1, I_2, I_3	3
I_1, I_2, I_4	2
I_1, I_3, I_4	1
I_2, I_3, I_4	2

Item	Count
I_1, I_2, I_3, I_4	1

Item	Count
I_1, I_2, I_3	3
I_1, I_2, I_4	2
I_1, I_3, I_4	1
I_2, I_3, I_4	2

Item	Count
I_1, I_2, I_3, I_4	1

Item	Count
I_1, I_2, I_3	3
I_1, I_2, I_4	2
I_1, I_3, I_4	1
I_2, I_3, I_4	2

Item	Count
I_1, I_2, I_3, I_4	1

Item	Count
I_1, I_2, I_3	3
I_1, I_2, I_4	2
I_1, I_3, I_4	1
I_2, I_3, I_4	2

Item	Count
I_1, I_2, I_3, I_4	1

Item	Count

</tbl_r

Example - Q2

Milk, Bread, Butter = 3

Milk, Bread	Bread, C=60%
Bread, Butter	Milk, C=100%
Milk, Butter	Bread, C=60%

Item	Support (from example)
Milk	8
Bread	7
Sugar	5
Butter	7

Item ID	Items
T1	Milk, Bread
T2	Bread, Sugar
T3	Bread, Butter
T4	Milk, Bread, Sugar
T5	Milk, Bread, Butter
T6	Milk, Butter
T7	Milk, Sugar
T8	Milk, Sugar
T9	Sugar, Butter
T10	Milk, Sugar, Butter
T11	Milk, Bread, Butter

Based on Association rules mentioned in the above table, we can recommended products to the customer to optimize product placement in retail stores.

Example - Step-1

Product	Frequency
Rice (R)	9
Pulse (P)	5
Oil (O)	4
Milk (M)	4

Step-2: threshold = 50%

Item	Frequency
RP	4
RO	3
RM	2
PO	4
PM	3
OM	2

Item ID	Rice	Pulse	Oil	Butter	Apple
T1	1	1	1	1	0
T2	0	1	1	1	0
T3	0	0	0	0	1
T4	1	1	1	1	0
T5	1	1	1	1	1
T6	1	1	1	1	1

Step-3: we get
RP, RO, PO & PM.

Step-4:

RP, RO & RPO
PO, PM & POM.

Step-5:

Itemset	Frequency
RPO	4
POM	3

Example 04:-

Step - 1 :- K = 1.

K = 1 (Table 2).

Itemset	Sup. count
I ₁	6
I ₂	7
I ₃	6
I ₄	2
I ₅	2

Itemset	Sup. count
I ₁ , I ₂	4
I ₁ , I ₃	4
I ₁ , I ₅	2
I ₂ , I ₃	4
I ₂ , I ₄	2
I ₃ , I ₅	2
I ₂ , I ₅	2

Confidence :- We will show the rule generation:-

$$I_1, I_2 \Rightarrow I_3 \Rightarrow \frac{2}{7} \times 100 = 50\%$$

$$I_1, I_3 \Rightarrow I_2 \Rightarrow \frac{2}{4} \times 100 = 50\%$$

$$I_2, I_3 \Rightarrow I_1 \Rightarrow \frac{2}{4} \times 100 = 50\%$$

$$I_1 \Rightarrow I_2, I_3 \Rightarrow \frac{2}{6} \times 100 = 33\%$$

$$I_2 \Rightarrow I_1, I_3 \Rightarrow \frac{2}{7} \times 100 = 27\%$$

$$I_3 \Rightarrow I_1, I_2 \Rightarrow \frac{2}{6} \times 100 = 33\%$$

So if minimum confidence is 50%.

the first 3 rules can be considered as strong association.

Step - 2:-

K = 2 (Table 3).

Itemset	Sup. count
I ₁ , I ₂	4
I ₁ , I ₃	4
I ₁ , I ₄	1
I ₁ , I ₅	2
I ₂ , I ₃	4
I ₂ , I ₄	2
I ₂ , I ₅	2
I ₃ , I ₄	0
I ₃ , I ₅	1
I ₄ , I ₅	0

Step 3, K=3.

Itemset	Sup. count
I ₁ , I ₂ , I ₃	2
I ₁ , I ₂ , I ₅	2

We stop here, because no frequent itemset are found further.

Min Support Count = 2
Table 1
C = 60%

TID	Item
T ₁	I ₁ , I ₂ , I ₅
T ₂	I ₂ , I ₄
T ₃	I ₂ , I ₃
T ₄	I ₁ , I ₂ , I ₄
T ₅	I ₁ , I ₃
T ₆	I ₂ , I ₃
T ₇	I ₁ , I ₃
T ₈	I ₁ , I ₂ , I ₃ , I ₅
T ₉	I ₁ , I ₂ , I ₃

Example-05:-

Support threshold = 60% ~~(Table-1)~~, so item Support $= \frac{60}{100} \times 5 = 3$.
 $k=1$ (Table-3).

TID	Item
1	Bread, Milk
2	Beer, Bread, Diaper, Egg
3	Beer, Coke, Diaper, Milk
4	Beer, Bread, Diaper, Milk
5	Bread, Coke, Diaper, Milk

Item	Support Count
Bread	4
Coke	2
Milk	4
Beer	3
Diaper	4
Egg	1

Table-3 -

Item	Support Count
Bread	1
Milk	4
Beer	3
Diaper	4

Table-4 ($k=2$).

Itemset	Count
Bread, Milk	4
Bread, Beer	2
Bread, Diaper	3
Milk, Beer	2
Milk, Diaper	3
Beer, Diaper	3

Table-5.

Itemset	Count
Bread, Milk	4
Bread, Diaper	3
Milk, Diaper	3
Beer, Diaper	3

Table-6 ($k=3$).

Itemset	Count
Beer, Diaper, Milk	2
Beer, Bread, Diaper	2

Example-06:- Minimum Support is 60%, then support $= \frac{60}{100} \times 5 = 3$

TID	Itemset
T1	6, 7, 8, 5, 4, 10
T2	3, 8, 7, 5, 4, 10
T3	6, 1, 5, 4
T4	6, 9, 2, 5, 10
T5	2, 8, 5, 4

Item	Sup
4, 5, 8	3
4, 5, 6	2 ✗
4, 5, 10	2 ✗
5, 6, 8	1 ✗
5, 6, 10	2 ✗
5, 8, 10	2 ✗

Item	Sup
4, 5, 8	3

Item	Supp.
1	1 ✗
2	2 ✗
3	1 ✗
4	4
5	5
6	3
7	2 ✗
8	4
9	1 ✗
10	3

Item	Supp.
4, 5	4
5, 5	5
6, 3	3
8, 4	4
10, 3	3

Item	Sup
4, 5	4
4, 6	2 ✗
4, 8	3
5, 6	3
5, 8	3
4, 10	2 ✗
5, 6	3
5, 8	3
5, 10	3
6, 8	1 ✗
6, 10	2 ✗
8, 10	2 ✗

Support :- Support is an indicator of how frequently the items appear in the data.

Confidence:- Confidence indicates the number of times the if-then statements are found true.

Lift:- A third metric called lift can be used to compare confidence with expected confidence or how many times an if-then statement is expected to be found true.

$$\text{Support}(X) = \frac{\text{Number of transactions containing } (X)}{\text{Total Number of Transactions}}$$

$$\text{Confidence } (X \Rightarrow Y) = \frac{\text{Number of transactions containing } X \text{ and } Y}{\text{Number of transactions containing } X}$$

$$\text{Support}$$

① It is measure of the number of times an itemset appears in a dataset.

② Support is calculated by dividing the number of transactions containing at least one itemset by the total number of transactions.

③ Support is used to identify itemsets that occur frequently in the dataset.

④ Support is often used with a threshold to identify itemsets that occur frequently enough to be of interest.

(5)

① Confidence is a measure of the likelihood that an itemset will appear in another itemset.

② Confidence is calculated by dividing the number of transactions containing both itemsets by the number of transactions containing the first itemset.

③ Confidence is used to evaluate the strength of a rule.

④ Confidence is interpreted as the percentage of transactions in which the second itemset appears given that the first itemset appears.

(5)

X is the itemset