

December 2023

BCA(DS)- V SEMESTER

Machine Learning - 1 (BCA-DS-302)

Time: 3 Hours

Max. Marks:75

- Instructions:**
1. It is compulsory to answer all the questions (1.5 marks each) of Part -A in short.
 2. Answer any four questions from Part -B in detail.
 3. Different sub-parts of a question are to be attempted adjacent to each other.

PART -A

- Q1 (a) What is the role of preprocessing of data in machine learning? (1.5)
 (b) How can outlier values be treated? (1.5)
 (c) What is Ensemble Learning? (1.5)
 (d) What is the standard approach to supervised learning? (1.5)
 (e) What is overfitting? How it can affect model generalization? (1.5)
 (f) What is Model Selection in Machine Learning? (1.5)
 (g) Define Logistic Regression. (1.5)
 (h) Why is ML important? (1.5)
 (i) What is 'training Set' and 'test Set' in a Machine Learning Model? (1.5)
 (j) Write any two applications of Machine Learning. (1.5)

PART -B

- Q2 (a) How many types of Data can be represented in machine learning algorithms. (7)
 (b) Explain different types of learning in Neural Network with example. (8)
- Q3 (a) Explain K means clustering by taking suitable example. (8)
 (b) What are the parameters to evaluate different machine learning algorithm. (7)
- Q4 What are Support Vector Machines? Why do we need to use Support Vector Machines? How do Support Vector Machines Work? Describe the significance of Kernel functions in SVM. List any two kernel functions. (15)
- Q5 (a) What is Self Organizing Maps? Explain Unsupervised learning of clusters in detail with example. (10)
 (b) Distinguish between bagging and boosting. (5)
- Q6 (a) What is Regression ? Explain with its Types. Also explain by taking a suitable example. (10)
 (b) What are the benefits of pruning in decision tree induction? Explain different approaches to tree pruning? (5)
-
- Q7 How naïve Bayes is suitable for supervised learning, how does it works? explain by taking a suitable example What are merit and demerits of it? (15)

PYQ Solution

PART-A: Questions 1

What is the role of data preprocessing in ml ?

Data preprocessing in machine learning prepares raw data for model training by cleaning, transforming, and organizing it. This step ensures that the data is complete, consistent, and in a suitable format, improving the model's accuracy and performance.

How can outlier values be treated?

Outliers can be treated by removing them, transforming them, or replacing them. Common methods include:

1. **Removing Outliers:** Discarding values beyond a certain threshold, like 1.5 times the interquartile range (IQR).
2. **Capping/Flooring:** Setting outliers to a maximum or minimum limit.
3. **Transforming:** Using log or square root transformations to reduce the impact of extreme values.

What is ensemble learning?

Ensemble learning combines multiple models to improve prediction accuracy and reduce errors. Common methods include bagging, boosting, and stacking, which work together to create more reliable and robust results than individual models.

Why is ML important?

Answer: ML automates data analysis, helps identify patterns, and enables systems to learn from data and make decisions without explicit programming. It is crucial for tasks such as fraud detection, customer insights, and predictive analytics.

What is over fitting?

Over fitting is when a machine learning model learns the training data too well, capturing noise and details that don't generalize to new data. This results in high accuracy on training data but poor performance on unseen data, as the model becomes too specific to the training set.

Differentiate between supervised and unsupervised machine learning.

Answer: Supervised learning involves training a model using labeled data to make predictions (e.g., spam detection). Unsupervised learning uses unlabeled data to find hidden patterns (e.g., clustering customers).

What is training set and testing set in a ML model?

Answer: A training set is used to teach a model, allowing it to learn from data. A testing set evaluates the model's performance on new, unseen data.

Define Logistic regression.

Answer: Logistic regression is a statistical method for binary classification that predicts the probability of a binary outcome (0 or 1) based on input features.

Differentiate between classification and regression.

Answer: Classification involves predicting discrete labels (e.g., spam or not spam), while regression involves predicting continuous values (e.g., predicting house prices).

viii. What is Pruning in Decision Tree?

Answer: Pruning is the process of removing unnecessary branches from a decision tree to simplify the model and prevent overfitting.

ix. Discuss how linear regression is different from logistic regression.

Answer: Linear regression predicts continuous outcomes (e.g., temperature), while logistic regression predicts binary or categorical outcomes (e.g., yes or no).

x. Name any two applications of supervised learning.

Answer: Email spam detection, Credit card fraud detection.

Question: What parameters to evaluate the different machine learning algorithms are

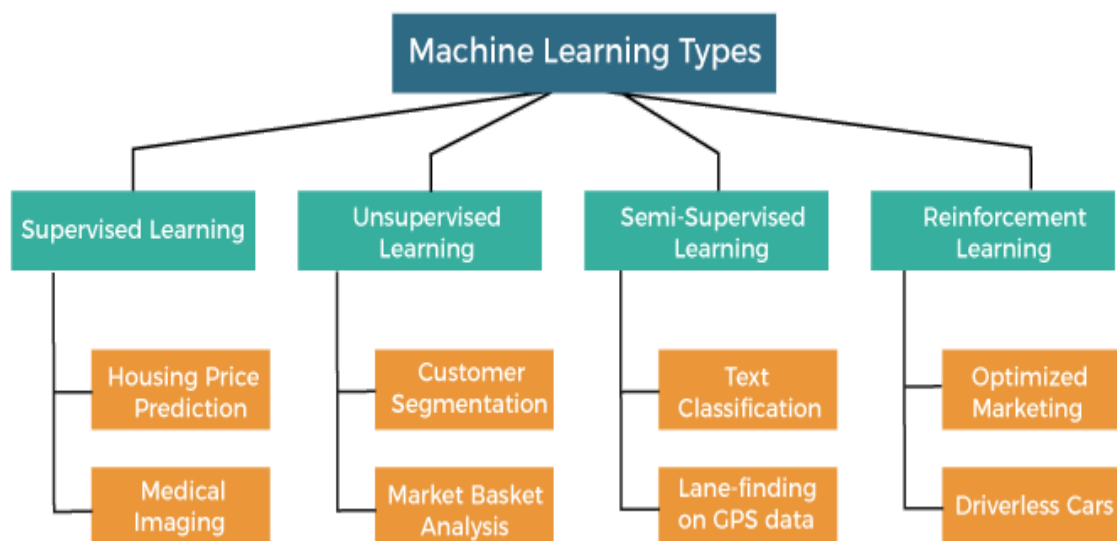
- Machine learning is a branch of artificial intelligence (AI) and computer science which focuses on the use of data and algorithms to imitate the way that humans learn, gradually improving its accuracy.
- Machine Learning Technique :

Machine learning is a subset of AI, which enables the machine to automatically learn from data, improve performance from past experiences, and make predictions. Machine learning contains a set of algorithms that work on a huge amount of data. Data is fed to these algorithms to train them, and on the basis of training, they build the model & perform a specific task.

These ML algorithms help to solve different business problems like Regression, Classification, Forecasting, Clustering, and Associations, etc.

Based on the methods and way of learning, machine learning is divided into mainly four types, which are:

1. Supervised Machine Learning
2. Unsupervised Machine Learning
3. Semi-Supervised Machine Learning
4. Reinforcement Learning



In this topic, we will provide a detailed description of the types of Machine Learning along with their respective algorithms:

1. Supervised Machine Learning

1. **The main goal of the supervised learning technique is to map the input variable(x) with the output variable(y).** Some real-world applications of supervised learning are **Risk Assessment, Fraud Detection, Spam filtering**, etc.

Categories of Supervised Machine Learning

2. Supervised machine learning can be classified into two types of problems, which are given below:
 - o **Classification**
 - o **Regression**

a) Classification

Classification algorithms are used to solve the classification problems in which the output variable is categorical, such as **"Yes" or No, Male or Female, Red or Blue, etc.** The classification algorithms predict the categories present in the dataset. Some real-world examples of classification algorithms are **Spam Detection, Email filtering, etc.**

Some popular classification algorithms are given below:

- o **Random Forest Algorithm**
- o **Decision Tree Algorithm**
- o **Logistic Regression Algorithm**
- o **Support Vector Machine Algorithm**

b) Regression

Regression algorithms are used to solve regression problems in which there is a linear relationship between input and output variables. These are used to predict continuous output variables, such as market trends, weather prediction, etc.

Some popular Regression algorithms are given below:

- o **Simple Linear Regression Algorithm**
- o **Multivariate Regression Algorithm**
- o **Decision Tree Algorithm**
- o **Lasso Regression**

Advantages and Disadvantages of Supervised Learning

Advantages:

- o Since supervised learning work with the labelled dataset so we can have an exact idea about the classes of objects.
- o These algorithms are helpful in predicting the output on the basis of prior experience.

Disadvantages:

- o These algorithms are not able to solve complex tasks.
- o It may predict the wrong output if the test data is different from the training data.
- o It requires lots of computational time to train the algorithm.

Applications of Supervised Learning

Some common applications of Supervised Learning are given below:

- o **Image Segmentation:**
Supervised Learning algorithms are used in image segmentation. In this process, image classification is performed on different image data with pre-defined labels.
- o **Medical Diagnosis:**
Supervised algorithms are also used in the medical field for diagnosis purposes. It is done by using medical images and past labelled data with labels for disease conditions. With such a process, the machine can identify a disease for the new patients.

- o **Fraud Detection** - Supervised Learning classification algorithms are used for identifying fraud transactions, fraud customers, etc. It is done by using historic data to identify the patterns that can lead to possible fraud.
- o **Spam detection** - In spam detection & filtering, classification algorithms are used. These algorithms classify an email as spam or not spam. The spam emails are sent to the spam folder.
- o **Speech Recognition** - Supervised learning algorithms are also used in speech recognition. The algorithm is trained with voice data, and various identifications can be done using the same, such as voice-activated passwords, voice commands, etc.

2. Unsupervised Machine Learning

Unsupervised learning is different from the Supervised learning technique; as its name suggests, there is no need for supervision. It means, in unsupervised machine learning, the machine is trained using the unlabeled dataset, and the machine predicts the output without any supervision.

In unsupervised learning, the models are trained with the data that is neither classified nor labelled, and the model acts on that data without any supervision.

The main aim of the unsupervised learning algorithm is to group or categories the unsorted dataset according to the similarities, patterns, and differences. Machines are instructed to find the hidden patterns from the input dataset.

Let's take an example to understand it more precisely; suppose there is a basket of fruit images, and we input it into the machine learning model. The images are totally unknown to the model, and the task of the machine is to find the patterns and categories of the objects.

So, now the machine will discover its patterns and differences, such as colour difference, shape difference, and predict the output when it is tested with the test dataset.

Categories of Unsupervised Machine Learning

Unsupervised Learning can be further classified into two types, which are given below:

- o **Clustering**
- o **Association**

1) Clustering

The clustering technique is used when we want to find the inherent groups from the data. It is a way to group the objects into a cluster such that the objects with the most similarities remain in one group and have fewer or no similarities with the objects of other groups. An example of the clustering algorithm is grouping the customers by their purchasing behaviour.

Some of the popular clustering algorithms are given below:

- o **K-Means Clustering algorithm**
- o **Mean-shift algorithm**

- o **DBSCAN Algorithm**
- o **Principal Component Analysis**
- o **Independent Component Analysis**

2) Association

Association rule learning is an unsupervised learning technique, which finds interesting relations among variables within a large dataset. The main aim of this learning algorithm is to find the dependency of one data item on another data item and map those variables accordingly so that it can generate maximum profit. This algorithm is mainly applied in **Market Basket analysis, Web usage mining, continuous production**, etc.

Some popular algorithms of Association rule learning are **Apriori Algorithm, Eclat, FP-growth algorithm**.

Advantages and Disadvantages of Unsupervised Learning Algorithm

Advantages:

- o These algorithms can be used for complicated tasks compared to the supervised ones because these algorithms work on the unlabeled dataset.
- o Unsupervised algorithms are preferable for various tasks as getting the unlabeled dataset is easier as compared to the labelled dataset.

Disadvantages:

- o The output of an unsupervised algorithm can be less accurate as the dataset is not labelled, and algorithms are not trained with the exact output in prior.
- o Working with Unsupervised learning is more difficult as it works with the unlabelled dataset that does not map with the output.

Applications of Unsupervised Learning

- o **Network Analysis:** Unsupervised learning is used for identifying plagiarism and copyright in document network analysis of text data for scholarly articles.
- o **Recommendation Systems:** Recommendation systems widely use unsupervised learning techniques for building recommendation applications for different web applications and e-commerce websites.
- o **Anomaly Detection:** Anomaly detection is a popular application of unsupervised learning, which can identify unusual data points within the dataset. It is used to discover fraudulent transactions.
- o **Singular Value Decomposition:** Singular Value Decomposition or SVD is used to extract particular information from the database. For example, extracting information of each user located at a particular location.

3. Semi-Supervised Learning

Semi-Supervised learning is a type of Machine Learning algorithm that lies between Supervised and Unsupervised machine learning. It represents the intermediate ground between Supervised (With Labelled training data) and Unsupervised learning (with no

labelled training data) algorithms and uses the combination of labelled and unlabeled datasets during the training period.

Although Semi-supervised learning is the middle ground between supervised and unsupervised learning and operates on the data that consists of a few labels, it mostly consists of unlabeled data. As labels are costly, but for corporate purposes, they may have few labels. It is completely different from supervised and unsupervised learning as they are based on the presence & absence of labels.

To overcome the drawbacks of supervised learning and unsupervised learning algorithms, the concept of Semi-supervised learning is introduced. The main aim of [semi-supervised learning](#) is to effectively use all the available data, rather than only labelled data like in supervised learning. Initially, similar data is clustered along with an unsupervised learning algorithm, and further, it helps to label the unlabeled data into labelled data. It is because labelled data is a comparatively more expensive acquisition than unlabeled data.

We can imagine these algorithms with an example. Supervised learning is where a student is under the supervision of an instructor at home and college. Further, if that student is self-analyzing the same concept without any help from the instructor, it comes under unsupervised learning. Under semi-supervised learning, the student has to revise himself after analyzing the same concept under the guidance of an instructor at college.

Advantages and disadvantages of Semi-supervised Learning

Advantages:

- o It is simple and easy to understand the algorithm.
- o It is highly efficient.
- o It is used to solve drawbacks of Supervised and Unsupervised Learning algorithms.

Disadvantages:

- o Iterations results may not be stable.
- o We cannot apply these algorithms to network-level data.
- o Accuracy is low.

Question: KNN Algorithm:

K-Nearest Neighbor (KNN) Algorithm for Machine Learning

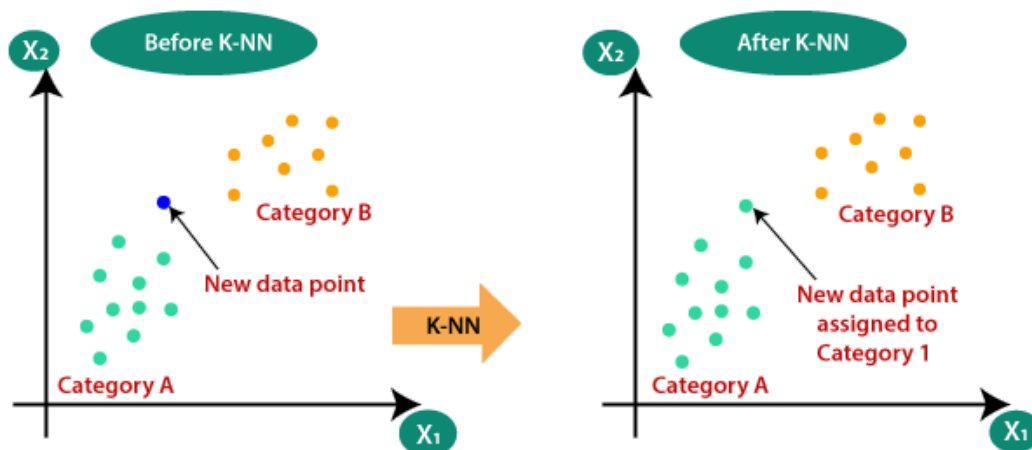
- o K-Nearest Neighbour is one of the simplest Machine Learning algorithms based on Supervised Learning technique.
- o K-NN algorithm assumes the similarity between the new case/data and available cases and put the new case into the category that is most similar to the available categories.

- o K-NN algorithm stores all the available data and classifies a new data point based on the similarity. This means when new data appears then it can be easily classified into a well suite category by using K- NN algorithm.
- o K-NN algorithm can be used for Regression as well as for Classification but mostly it is used for the Classification problems.
- o K-NN is a **non-parametric algorithm**, which means it does not make any assumption on underlying data.
- o It is also called a **lazy learner algorithm** because it does not learn from the training set immediately instead it stores the dataset and at the time of classification, it performs an action on the dataset.
- o KNN algorithm at the training phase just stores the dataset and when it gets new data, then it classifies that data into a category that is much similar to the new data.
- o **Example:** Suppose, we have an image of a creature that looks similar to cat and dog, but we want to know either it is a cat or dog. So for this identification, we can use the KNN algorithm, as it works on a similarity measure. Our KNN model will find the similar features of the new data set to the cats and dogs images and based on the most similar features it will put it in either cat or dog category.



Why do we need a K-NN Algorithm?

Suppose there are two categories, i.e., Category A and Category B, and we have a new data point x_1 , so this data point will lie in which of these categories. To solve this type of problem, we need a K-NN algorithm. With the help of K-NN, we can easily identify the category or class of a particular dataset. Consider the below diagram:

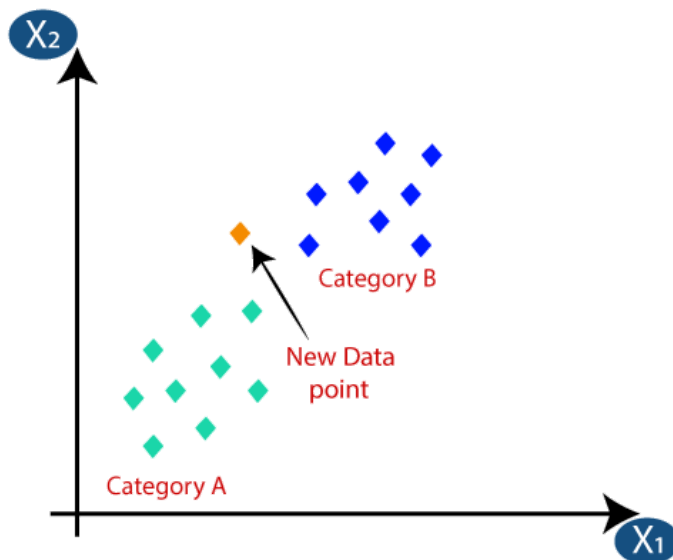


How does K-NN work?

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

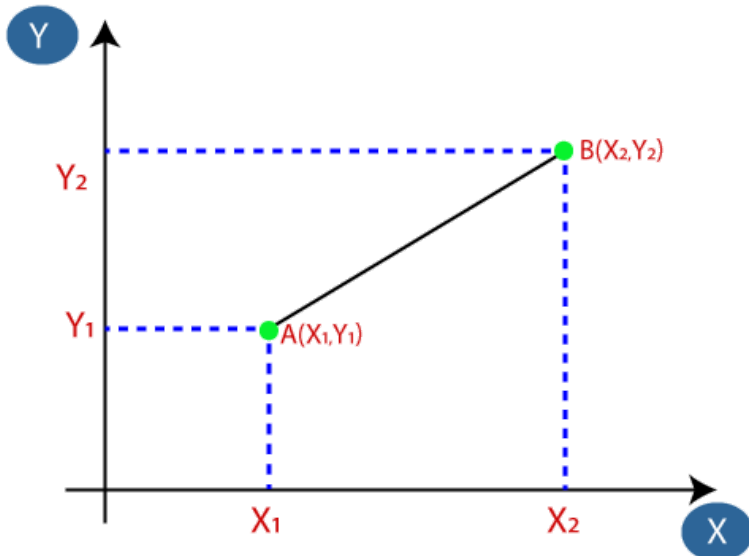
- o **Step-1:** Select the number K of the neighbors
- o **Step-2:** Calculate the Euclidean distance of **K number of neighbors**
- o **Step-3:** Take the K nearest neighbors as per the calculated Euclidean distance.
- o **Step-4:** Among these k neighbors, count the number of the data points in each category.
- o **Step-5:** Assign the new data points to that category for which the number of the neighbor is maximum.
- o **Step-6:** Our model is ready.

Suppose we have a new data point and we need to put it in the required category. Consider the below image:



- o Firstly, we will choose the number of neighbors, so we will choose the $k=5$.

- o Next, we will calculate the **Euclidean distance** between the data points. The Euclidean distance is the distance between two points, which we have already studied in geometry. It can be calculated as:



$$\text{Euclidean Distance between } A_1 \text{ and } B_2 = \sqrt{(X_2 - X_1)^2 + (Y_2 - Y_1)^2}$$

- o By calculating the Euclidean distance we got the nearest neighbors, as three nearest neighbors in category A and two nearest neighbors in category B. Consider the below image:



- o As we can see the 3 nearest neighbors are from category A, hence this new data point must belong to category A.

How to select the value of K in the K-NN Algorithm?

Below are some points to remember while selecting the value of K in the K-NN algorithm:

- o There is no particular way to determine the best value for "K", so we need to try some values to find the best out of them. The most preferred value for K is 5.
- o A very low value for K such as K=1 or K=2, can be noisy and lead to the effects of outliers in the model.
- o Large values for K are good, but it may find some difficulties.

Advantages of KNN Algorithm:

- o It is simple to implement.
- o It is robust to the noisy training data
- o It can be more effective if the training data is large.

Disadvantages of KNN Algorithm:

- o Always needs to determine the value of K which may be complex some time.
- o The computation cost is high because of calculating the distance between the data points for all the training samples.

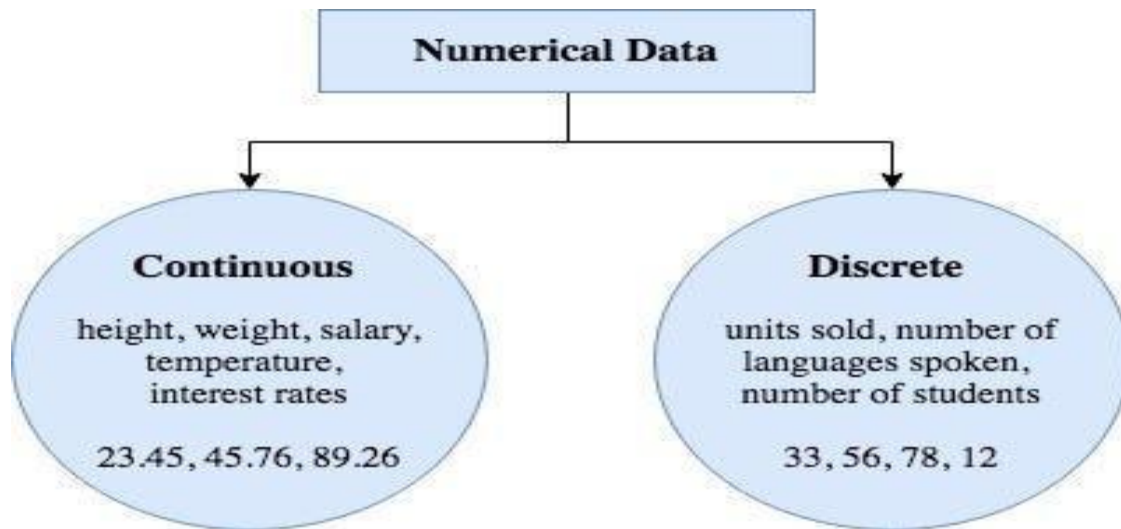
Question

Data Represtation:

1. The word data refers to constituting people, things, events, ideas. It can be a title, an integer, or any cast. After collecting data the investigator has to condense them in tabular form to study their salient features. Such an arrangement is known as the presentation of data.

It refers to the process of condensing the collected data in a tabular form, Numerical or graphically. This arrangement of data is known as Data Representation.

Numerical Representation: Numerical data is any data where data points are exact numbers. Statisticians also might call numerical data, quantitative data. This data has meaning as a **measurement** such as house prices or as a count, such as a number of residential properties in Los Angeles or how many houses sold in the past year. Numerical data can be characterized by continuous or discrete data. Continuous data can assume any value within a range whereas discrete data has distinct values.



Graphical Representation of Data

Graphical Representation of Data,” where numbers and facts become lively pictures and colorful diagrams. Instead of staring at boring lists of numbers, we use fun charts, cool graphs, and interesting visuals to understand information better. In this exciting concept of data visualization, we’ll learn about different kinds of graphs, charts, and pictures that help us see patterns and stories hidden in data.

There is an entire branch in mathematics dedicated to dealing with collecting, analyzing, interpreting, and presenting numerical data in visual form in such a way that it becomes easy to understand and the data becomes easy to compare as well, the branch is known as Statistics.

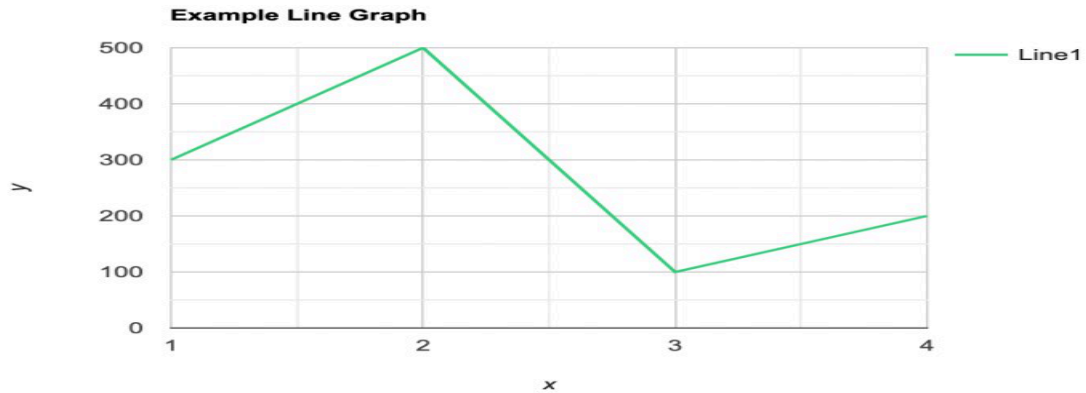
The branch is widely spread and has a plethora of real-life applications such as Business Analytics, demography, Astor statistics, and so on. In this article, we have provided everything about the graphical representation of data, including its types, rules, advantages, etc.

Types of Graphical Representations

2. Comparison between different items is best shown with graphs, it becomes easier to compare the crux of the data about different items. Let’s look at all the different types of graphical representations briefly:

Line Graphs

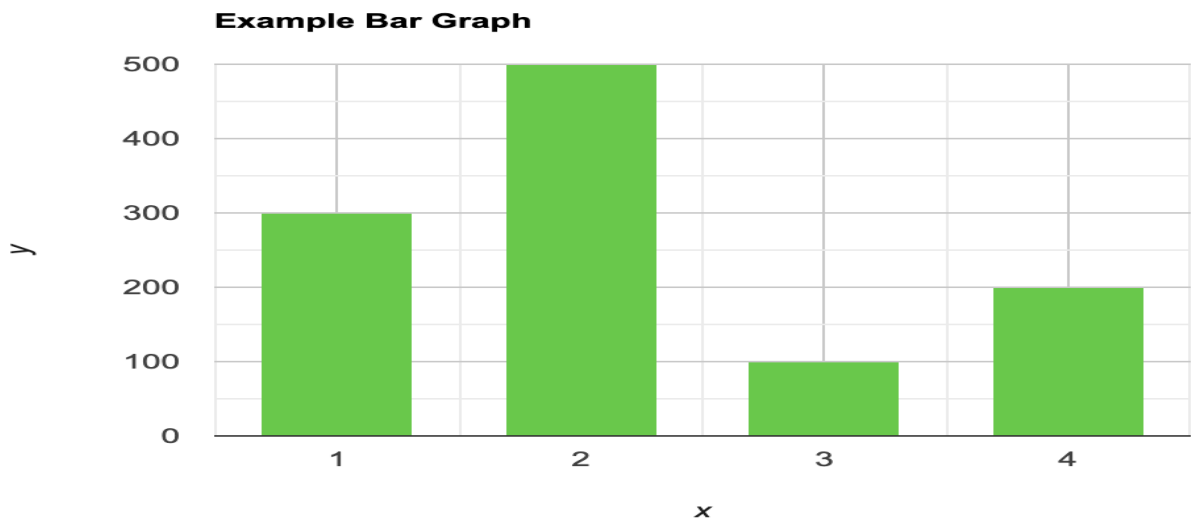
3. A **line graph** is used to show how the value of a particular variable changes with time. We plot this graph by connecting the points at different values of the variable. **It can be useful for analyzing the trends in the data and predicting further trends.**



4.

Bar Graphs

5. A bar graph is a type of graphical representation of the data in which bars of uniform width are drawn with equal spacing between them on one axis (x-axis usually), depicting the variable. **The values of the variables are represented by the height of the bars.**

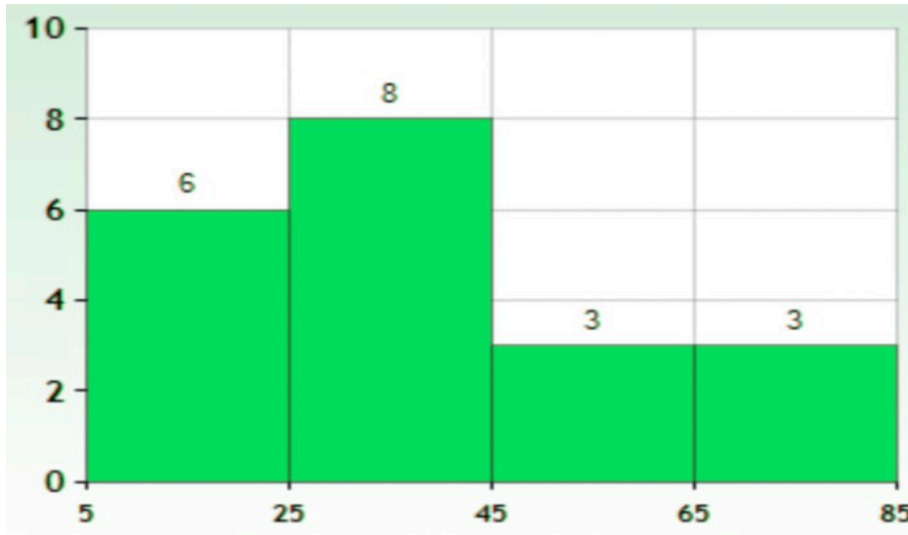


6.

Histograms

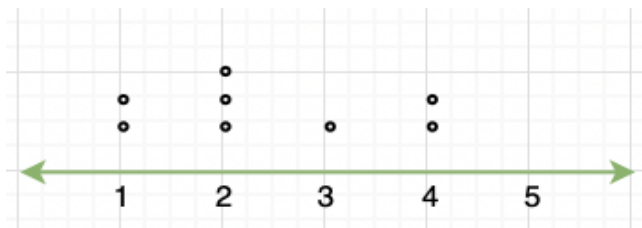
7. This is similar to [bar graphs](#), but it is based frequency of numerical values rather than their actual values. The data is organized into intervals and the bars represent the frequency of the values in that range. That is, it counts how many values of the data lie in a particular range.

8.



Line Plot

It is a plot that displays data as points and checkmarks above a number line, showing the frequency of the point.



Question 3(B)

Naive Bayes:

- o Naïve Bayes algorithm is a supervised learning algorithm, which is based on Bayes theorem and used for solving classification problems.
- o It is mainly used in text classification that includes a high-dimensional training dataset.
- o Naïve Bayes Classifier is one of the simple and most effective Classification algorithms which helps in building the fast machine learning models that can make quick predictions.
- o It is a probabilistic classifier, which means it predicts on the basis of the probability of an object.

- o Some popular examples of Naïve Bayes Algorithm are spam filtration, Sentimental analysis, and classifying articles.

Why is it called Naïve Bayes?

The Naïve Bayes algorithm is comprised of two words Naïve and Bayes, Which can be described as:

- o Naïve: It is called Naïve because it assumes that the occurrence of a certain feature is independent of the occurrence of other features. Such as if the fruit is identified on the bases of color, shape, and taste, then red, spherical, and sweet fruit is recognized as an apple. Hence each feature individually contributes to identify that it is an apple without depending on each other.
- o Bayes: It is called Bayes because it depends on the principle of Bayes' Theorem.

Bayes' Theorem:

- o Bayes' theorem is also known as Bayes' Rule or Bayes' law, which is used to determine the probability of a hypothesis with prior knowledge. It depends on the conditional probability.
- o The formula for Bayes' theorem is given as:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Where,

$P(A|B)$ is Posterior probability: Probability of hypothesis A on the observed event B.

$P(B|A)$ is Likelihood probability: Probability of the evidence given that the probability of a hypothesis is true.

$P(A)$ is Prior Probability: Probability of hypothesis before observing the evidence.

$P(B)$ is Marginal Probability: Probability of Evidence.

Working of Naïve Bayes' Classifier:

Working of Naïve Bayes' Classifier can be understood with the help of the below example:

Suppose we have a dataset of weather conditions and corresponding target variable "Play". So using this dataset we need to decide that whether we should play or not on a particular day according to the weather conditions. So to solve this problem, we need to follow the below steps:

1. Convert the given dataset into frequency tables.
2. Generate Likelihood table by finding the probabilities of given features.
3. Now, use Bayes theorem to calculate the posterior probability.

Problem: If the weather is sunny, then the Player should play or not?

Solution: To solve this, first consider the below dataset:

S.No.	Outlook	Play
0	Rainy	Yes
1	Sunny	Yes
2	Overcast	Yes
3	Overcast	Yes
4	Sunny	No
5	Rainy	Yes
6	Sunny	Yes
7	Overcast	Yes
8	Rainy	No
9	Sunny	No
10	Sunny	Yes
11	Rainy	No
12	Overcast	Yes
13	Overcast	Yes

Frequency table for the Weather Conditions:

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

Likelihood table weather condition:

Weather	No	Yes	
Overcast	0	5	$5/14 = 0.35$
Rainy	2	2	$4/14 = 0.29$
Sunny	2	3	$5/14 = 0.35$
All	$4/14 = 0.29$	$10/14 = 0.71$	

Applying Bayes theorem:

$$P(\text{Yes} | \text{Sunny}) = P(\text{Sunny} | \text{Yes}) * P(\text{Yes}) / P(\text{Sunny})$$

$$P(\text{Sunny} | \text{Yes}) = 3/10 = 0.3$$

$$P(\text{Sunny}) = 0.35$$

$$P(\text{Yes}) = 0.71$$

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

$$P(\text{No} | \text{Sunny}) = P(\text{Sunny} | \text{No}) * P(\text{No}) / P(\text{Sunny})$$

$P(\text{Sunny}|\text{NO}) = 2/4 = 0.5$

$P(\text{No}) = 0.29$

$P(\text{Sunny}) = 0.35$

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

So as we can see from the above calculation that $P(\text{Yes}|\text{Sunny}) > P(\text{No}|\text{Sunny})$

Hence on a Sunny day, Player can play the game.

Advantages of Naïve Bayes Classifier:

- o Naïve Bayes is one of the fast and easy ML algorithms to predict a class of datasets.
- o It can be used for Binary as well as Multi-class Classifications.
- o It performs well in Multi-class predictions as compared to the other Algorithms.
- o It is the most popular choice for text classification problems.

Disadvantages of Naïve Bayes Classifier:

- o Naive Bayes assumes that all features are independent or unrelated, so it cannot learn the relationship between features.

Applications of Naïve Bayes Classifier:

- o It is used for Credit Scoring.
- o It is used in medical data classification.
- o It can be used in real-time predictions because Naïve Bayes Classifier is an eager learner.
- o It is used in Text classification such as Spam filtering and Sentiment analysis.

Types of Naïve Bayes Model:

There are three types of Naive Bayes Model, which are given below:

- o Gaussian: The Gaussian model assumes that features follow a normal distribution. This means if predictors take continuous values instead of discrete, then the model assumes that these values are sampled from the Gaussian distribution.
- o Multinomial: The Multinomial Naïve Bayes classifier is used when the data is multinomial distributed. It is primarily used for document classification problems, it means a particular document belongs to which category such as Sports, Politics, education, etc.

The classifier uses the frequency of words for the predictors.

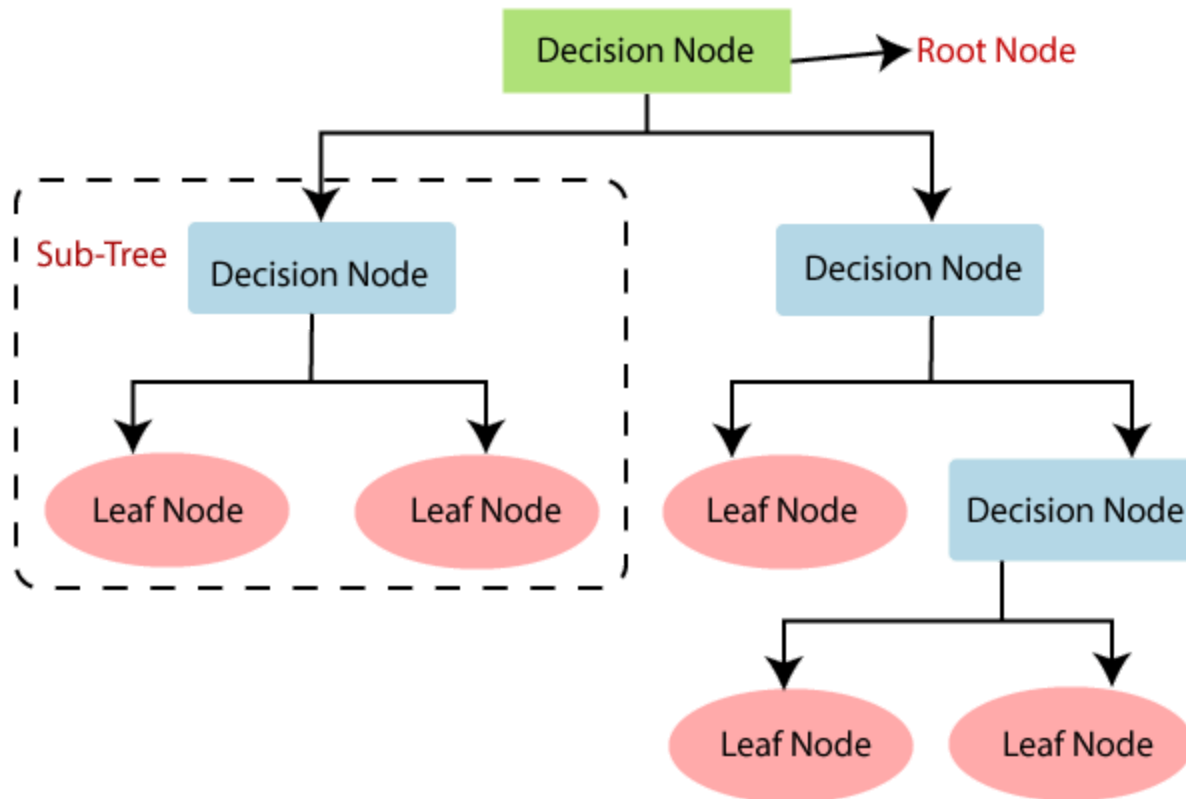
- o Bernoulli: The Bernoulli classifier works similar to the Multinomial classifier, but the predictor variables are the independent Booleans variables. Such as if a particular word is present or not in a document. This model is also famous for document classification tasks.

Question 4

Decision Tree:

- o Decision Tree is a **Supervised learning technique** that can be used for both classification and Regression problems, but mostly it is preferred for solving Classification problems. It is a tree-structured classifier, where **internal nodes represent the features of a dataset, branches represent the decision rules** and **each leaf node represents the outcome**.
- o In a Decision tree, there are two nodes, which are the **Decision Node** and **Leaf Node**. Decision nodes are used to make any decision and have multiple branches, whereas Leaf nodes are the output of those decisions and do not contain any further branches.
- o The decisions or the test are performed on the basis of features of the given dataset.
- o ***It is a graphical representation for getting all the possible solutions to a problem/decision based on given conditions.***
- o It is called a decision tree because, similar to a tree, it starts with the root node, which expands on further branches and constructs a tree-like structure.
- o In order to build a tree, we use the **CART algorithm**, which stands for **Classification and Regression Tree algorithm**.
- o A decision tree simply asks a question, and based on the answer (Yes/No), it further split the tree into subtrees.
- o Below diagram explains the general structure of a decision tree:

Note: A decision tree can contain categorical data (YES/NO) as well as numeric data.



Why use Decision Trees?

There are various algorithms in Machine learning, so choosing the best algorithm for the given dataset and problem is the main point to remember while creating a machine learning model. Below are the two reasons for using the Decision tree:

- o Decision Trees usually mimic human thinking ability while making a decision, so it is easy to understand.
- o The logic behind the decision tree can be easily understood because it shows a tree-like structure.

Decision Tree Terminologies

Root Node: Root node is from where the decision tree starts. It represents the entire dataset, which further gets divided into two or more homogeneous sets.

Leaf Node: Leaf nodes are the final output node, and the tree cannot be segregated further after getting a leaf node.

Splitting: Splitting is the process of dividing the decision node/root node into sub-nodes according to the given conditions.

Branch/Sub Tree: A tree formed by splitting the tree.

Pruning: Pruning is the process of removing the unwanted branches from the tree.

Parent/Child node: The root node of the tree is called the parent node, and other nodes are called the child nodes.

