



Practical 1

Aim :- To introduce R-tools, install the R programming environment & basic packages.

Theory

Introduction to R:

R is a programming language & free software environment for statistical computing & graphics supported by the R Foundation for statistical computing. It is widely used among statisticians & data miners for data analysis & developing statistical software. R provides a wide variety of statistical & graphical techniques, including linear & non-linear modeling, classical statistical tests, time-series analysis, classification, clustering & others.

Features:

- 1] Open source: R is an open-source programming language which means it is free to use, modify & distribute.
- 2] Comprehensive statistical analysis: R includes a vast array of tools for statistical analysis and data mining.
- 3] Visualization: R excels at creating detailed & advanced plots & graphs offering extensive visualization capabilities.



- 4) Community support: R has large & large & active community that contributes packages & provide support to users.
- 5) Extensible :- R can be extended easily via packages. There are thousands of packages available on CRAN (Comprehensive R Archive Network) for various applications.
- 6) Cross-platform - R is available on multiple platforms including windows, Mac OS & Linux.
- 7) Data handling: R can handle & manipulate large datasets efficiently.
- 8) Interfacing: R can integrate with other languages like C/C++, Python & can be used with database management systems.



R

Python

- R is a language & environment for statistical programming which includes statistical computing & graphics.
 - It has many features which are useful for statistical analysis & representation.
 - It has many easy to use packages for performing tasks.
 - Various popular R IDE are Rstudio, RKward, R Commander etc.
 - There are many packages & libraries like ggplot 2, lapply etc.
 - It is mainly used for complex data analysis in data science.
- Python is a general purpose programming language for data analysis & scientific computing.
 - It can be used to develop GUI applications & web applications as well as with embedded systems.
 - It can easily perform matrix computation as well as optimization.
 - Various popular Python IDEs are Spyder, Eclipse + Pydev, Atom.
 - Some essential packages & libraries are pandas, numpy, scipy.
 - It takes a more streamlined approach for data science projects.



Steps to install R & Basic packages:-

1. Install R:

- Go to the R project website
- click on "CRAN" under "download" to find a CRAN mirror
- Select a CRAN mirror near you.
- Download the R installer for your OS.

2. Install R studio:- Download the R studio Desktop installer for your OS. Run the installer & follow the on-screen instructions to complete the installation

3. Installing basic packages:-

Commands:

~~install.packages("data.table")~~

~~install.packages("dplyr")~~

~~install.packages("ggplot2")~~

~~install.packages("caret")~~

~~- data.table :- It is highly efficient for handling large datasets, making operations like filtering, aggregating & joining data much faster compared to base R or other packages.~~



- dplyr:- It simplifies data manipulation tasks & makes code more readable & expressive, especially for those coming from a SQL background.
- ggplot2:- It is powerful, flexible & makes it easy to create ~~as~~ aesthetically pleasing & informative visualizations.
- caret:- It simplifies the process of training & evaluating machine learning models providing a unified interface to a wide range of model types & tool for model comparison & tuning.

Conclusion:

In this practical we introduced the R programming environment & its features. We also covered the difference betⁿ R & Python.

~~Mr
25/7/2024~~



Practical 2:

Aim: Implement different data types & data structures in R, including vectors, lists, matrices, arrays, factors & data frames.

Theory:

- 1] **Vector :-** A vector is a basic data structure in R that contains elements of the same type. It can be numeric, integer, character, logical or complex.

Ex:

numeric_vector \leftarrow c(1.1, 2.2, 3.3)

print(numeric_vector)

char_vector \leftarrow c("apple", "banana")

print(char_vector)

logical_vector \leftarrow c(TRUE, FALSE, TRUE)

print(logical_vector)

- 2] **List :-** A list is an R data structure that can hold elements of different types number, string, vector. It is similar to a vector but more flexible.



Ex:

```
my_list <- list(name = "john", age = 25,  
score = c(35, 40, 39))
```

```
print(my_list)
```

```
print(my_list$name)
```

```
print(my_list$age)
```

```
print(my_list$scores)
```

3)

Matrices :- A matrix is a two dimensional array that contains elements of the same type. It is used for mathematical computations involving rows & columns.

Ex:

```
my_matrix <- matrix(1:9, nrow = 3, ncol = 3)
```

```
print(my_matrix)
```

```
print(my_matrix[1, 2])
```

```
print(my_matrix[, 3])
```

```
print(my_matrix[2, ])
```

4)

Array :- An array is a multi-dimensional data structure that can hold elements of the same type. It can be considered as a generalized form of a matrix.



Ex:

```
my_array <- array(1:12, dim = c(3, 2, 2))  
print(my_array)  
print(my_array[1, 2, 1])  
print(my_array[, , 2])
```

- 5) **Factor:** A factor is a data structure used for fields that take only pre-defined, finite no. of values [categorical data]. Factors are useful in statistical modeling.

Ex:

```
my_factor <- factor(c("low", "medium", "high"))  
print(my_factor)  
print(levels(my_factor))  
my_vector <- c("yes", "no", "yes")  
factor_vector <- as.factor(my_vector)  
print(factor_vector)
```

- 6) **Data Frame:** A data frame is a table or a 2-dimensional array like structure in which each column contains values of one variable & each row contains one set of values from each column. It's similar to a matrix but can contain different types of data.

Ex:

```
my_df <- data.frame(  
  name = c("Alice", "Bob", "Charlie"),  
  age = c(20, 21, 22),
```



Score = c(90, 91, 92)

Print c(my_df)

print c(my_df\$name)

print c(my_df[1,])

print c(my_df[, "age"])

Conclusion:-

In this practical, we explored & implemented different data types & data structures in R.

Including vectors, lists, matrices, arrays & data frame. Understanding these structures is fundamental for performing data analysis & manipulation tasks in R.

WZTHOM
25TH Nov

Practical 3

Aim: To implement a program in R that calculates basic statistics including mean, median, mode, variance, Std deviation, range, summary, skewness & kurtosis

Theory:

Descriptive Statistics

It is used to summarize & describe the main features of a dataset. They provide simple summaries about the sample & the measures. Together with simple graphics analysis, they form the basis of virtually every quantitative analysis of data.

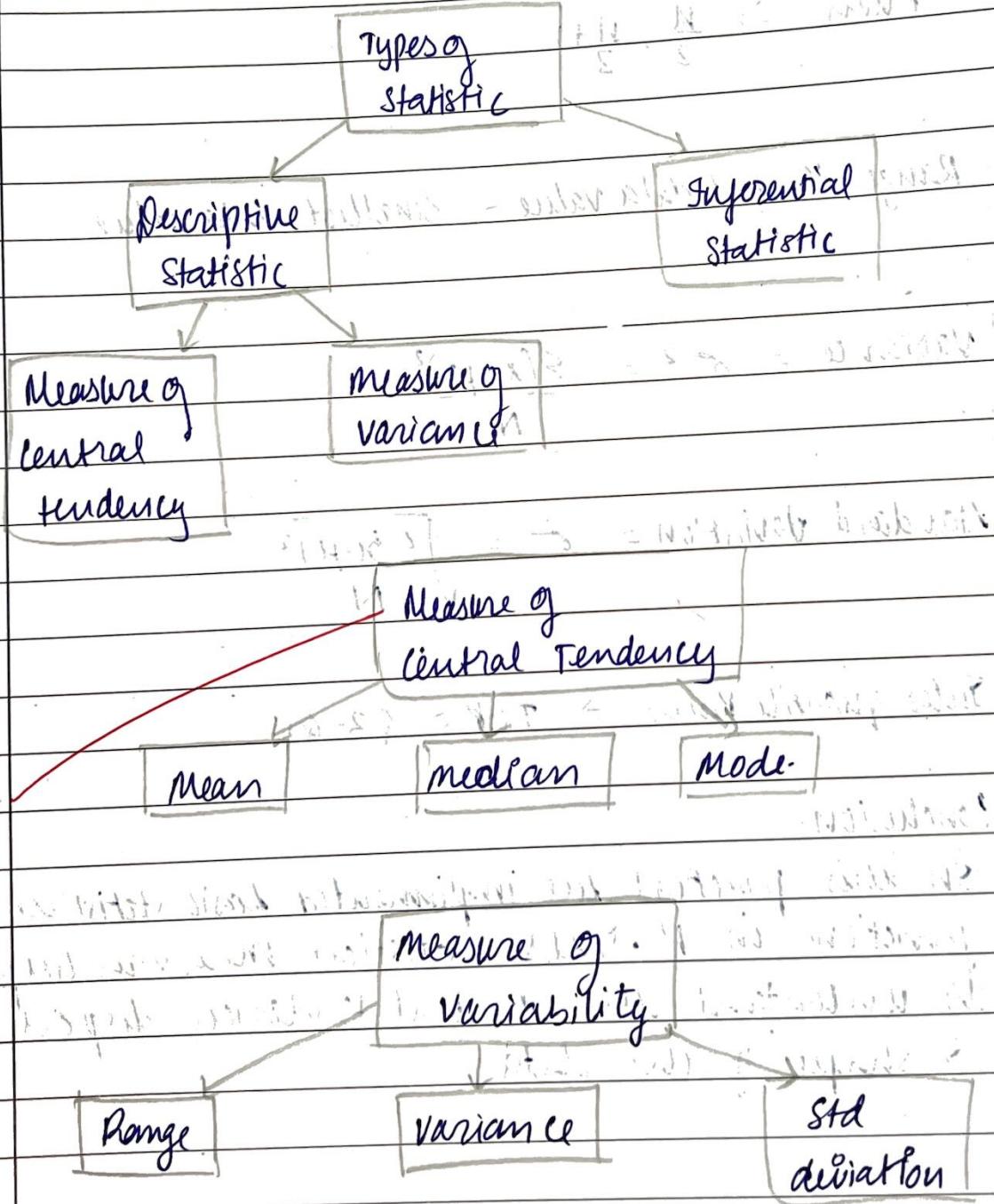
1. Mean: The mean is the sum of all the numbers in a dataset divided by the count of numbers in that dataset.
2. Median: The median is the middle value in a data set when the numbers are arranged in ascending or descending order. If the dataset has an even number of observations, the median is taking the two middle numbers.



3. Mode: the mode is the value that appears most frequently in a data set.
4. Variance: variance measures how far a set of numbers are spread out from their average values. It is the average of the squared differences from the mean.
5. Standard deviation: it is the square root of the variance. It indicates the extent to which the values in a dataset deviate from the mean.
6. Range: the range is the difference between the highest & lowest values in a dataset.
7. summary: the summary function in R provides a summary of the data, including the min, 1st quartile, median, mean, 3rd quartile & maximum.
8. Skewness: it measures the asymmetry of the data distribution. Positive skewness indicates a distribution with a long right tail, while negative skewness indicates a distribution with a long left tail.



q) Kurtosis: It measures the tailedness of the data distribution. High kurtosis means more of the variance is due to infrequent extreme deviations, as opposed to frequent modestly sized deviations.





Mean $\bar{x} = \frac{1}{n} \sum x_i$ is a measure of location for

individual data or Ranks involving measurement without units.

Median: measure of central tendency based on

$$\text{odd } \Rightarrow \frac{n+1}{2}$$

$$\text{Even } \Rightarrow \frac{n}{2}, \frac{n+1}{2}$$

Range = Largest data value - Smallest data value

$$\text{Variance} = \sigma^2 = \frac{\sum (x - \bar{x})^2}{N}$$

$$\text{Standard deviation} = \sigma = \sqrt{\frac{\sum (x - \bar{x})^2}{N}}$$

$$\text{Inter quartile Range} \Rightarrow IQR = Q_3 - Q_1$$

Conclusion:

In this practical we implemented basic statistical functions in R. These statistical measures help to understand the central tendency, dispersion & shape of the data.



Practical 4

Aim: Write a program in R for data visualization using histogram etc.

Theory

Data Visualization is the graphical representation of information & data. By using visual elements like charts, graphs & maps. Data visualization tools provide an accessible way to see & understand trends, outliers & patterns in data.

1. Histogram

It is a graphical representation that organizes a group of data point into user-specified ranges.

Similar in appearance to a bar graph, the histogram condenses a data series into bins by taking many points & grouping them into logical ranges or bins.

2. Scatter plot

It is a type of data visualization that displays values for typically two recent variable for a set of data.

The data is displayed as a collection of points, each having the value of one variable determining the position



on the horizontal axis & the value of the other variable determining the position on the vertical axis.

3. Dot plot:

A dot plot is a simple statistical chart that consists of data points plotted as dots on a simple scale, used to show frequency distributions or patterns in data.

4. Line chart:

A line chart is a type of chart that displays information as a series of data points called "markers" connected by straight line segments. It is similar to a scatter plot except that the measurement points are ordered and joined with straight line segments.

5. Pie chart:

A pie chart is a circular graphic which is divided into slices to illustrate numerical proportions. In a pie chart, the area of each slice is proportional to the quantity it represents.



6. Box plot:-

It is a standardized way of displaying the distribution of data based on five number summary: min, first quartile (Q_1), median, third quartile (Q_3) & maximum.

7. Density plot:-

It is a smoothed continuous version of a histogram estimated from the data. It is used to visualize the distribution of the data.

8. Normal QQ Plot:-

It is a graphical tool to help assess if a set of data plausibly came from some theoretical distribution such as a Normal or exponential. The data are plotted against a theoretical distribution in such a way that the points should form an approximate straight line.

Conclusion:

In this lab we successfully created various types of data visualizations using the base R functions.

These visualizations help us to understand the underlying pattern, distribution & relationship of the dataset.

These plots help to obtain different perspective making it easier to derive meaningful insights from the data.



Practical 5

Aim: Implementation of data processing and EDA in R.

Theory: - Data mining, machine learning, and deep learning

Data mining is a process of extracting useful information from large amounts of data.

1. Data processing:

Data processing is a critical step in the data analysis pipeline that involves preparing raw data for analysis. Data preprocessing includes several tasks such as data cleaning, data transformation, and feature scaling.

2. Data cleaning:

Data cleaning involves detecting & correcting errors & inconsistencies in the data to improve its quality.

It is an essential step to ensure that the dataset is reliable & accurate for analysis.

Key Steps:

- Detecting missing values: Identify missing values in the dataset, which can occur due to various reasons such as data entry error or missing information.

- Handling missing values: Decide on a strategy to handle missing values. Common methods include:-



1) Removal: Deleting rows or columns with missing values if they are few & not critical to the analysis.

2) Imputation: Filling in missing values using various methods such as mean, median, mode or more advanced techniques like regression or ML models.

Detecting Outliers: Identify outliers in the data which can skew the analysis. Outliers can be detected using statistical methods or visualization techniques.

- Handling Outliers: Outliers can be handled by removing them, transforming them or using robust statistical methods that are less sensitive to outliers.
- Correcting Data Entry Errors: Detect & correct errors such as typos, incorrect values or inconsistent format.

• Standardizing data format: Ensure consistency in data formats, such as date formats, numeric precision & categorical labels.



3. Data Transformation

Data transformation involves converting data into a suitable format or structure for analysis. This process helps in enhancing the quality & usability of the data.

Steps:-

- Encoding Categorical Variables: Convert categorical variables into numerical format for analysis.
Common techniques:
 - 1) Label Encoding: Assigning a unique integer to each category.
 - 2) One-hot Encoding: Creating binary columns for each category.
- Creating New feature: Derive new features from existing data to capture additional information.
- Scaling numerical Features:
Standardize the range of numerical features to ensure that they contribute equally to the analysis.
 - D) Min-Max Scaling: Rescaling the features to a fixed range, usually $[0,1]$



2) Standardization (Z-score scaling):
Transforming the feature to have a mean of 0 & a std deviation of 1

- Normalization: Adjusting the data to a common scale without distorting differences in the ranges of values

4. Feature Scaling:
Feature scaling is the process of normalizing the range of independent variables or features of data. It is crucial for alg algorithms that are sensitive to the scale of data, such as gradient descent - based methods & distance-based algo.

Conclusion:

In this practical we have performed data pre-processing & EDA, that are fundamental steps in the data analysis process



Practical 6:

Aim: To implement analysis of variance (ANOVA) on a given dataset & calculate the difference of variances b/w the groups.

Theory:

ANOVA is a statistical method used to compare the means of three or more samples to understand if at least one sample mean is significantly different from the others. It helps in determining whether the observed variations among sample means are due to actual differences in the population or merely due to random chance.

Key Concepts:-

1. Null Hypothesis (H_0): - Assumes that all group means are equal.
2. Alternative hypothesis (H_1): - Assumes that at least one group mean is different.
3. Between Group Variance: - The variation due to the interaction b/w the different groups.
4. Within Group Variance: - The variation within each group.
5. F-statistic: - A ratio of b/w-group variance to within group variance. A higher F-statistic indicates a greater disparity b/w grp means.



Steps in ANOVA:

1. Calculate the group means: Determine the mean of each group in the dataset.
2. Calculate the overall mean: Determine the mean of all data points combined.
3. Compute the between-group sum of squares (SSB): Measure the variance between the group means & the overall mean.
4. Compute the within-group sum of squares (SSW): Measure the variance within each group.
5. Calculate the mean square between (MSB): Divide SSB by the degrees of freedom between the group.
6. Calculate the mean square within (MSW): Divide SSW by the degree of freedom within the groups.
7. Compute the F-statistic: Divide MSB by MSW to obtain the F-value.
8. Compare with critical value: Compare the F-value with the critical value from the F-distribution table to determine the significance.

Conclusion:-

ANOVA allows us to statistically test if there are any significant differences in the variances between groups, aiding in understanding the impact of different factors on the observed data.



Practical 7

Aim: - To implement simple linear regression on a given dataset

Theory:

Simple linear regression is a statistical method used to model the "relationship bet" a dependent variable & an independent variable by fitting a linear eqⁿ to the observed data. The objective is to predict the value of the dependent variable based on the value of the independent variable.

Key concept:

1) Dependent Variable (Y) :- The variable we are trying to predict or explain.

2) Independent Variable (X) :- The variable we are using to make predictions about the dependent variables.

3) Linear Relationship :- The relationship betⁿ x & y is assumed to be linear, meaning it can be represented by a st. line.



- u] Regression line: The line that best fits the data points, described by the eqn $[y = b_0 + b_1 x]$ where b_0 is the intercept & b_1 is the slope
- s] Least square method:- A method used to determine the best fitting line by minimizing the sum of the squares of the vertical distances of the points from the line.

Steps in Simple Linear Regression:

- 1] Plot the data: visualize the data points on a scatter plot to observe the relationship betⁿ x & y.
- 2] Calculate the Mean of x & y: Determine the average values of the independent & dependent variables.
- 3] Determine the slope (b_1): Use the formula
$$b_1 = \frac{\sum_{i=1}^n (x - \bar{x})(y - \bar{y})}{\sum_{i=1}^n (x - \bar{x})^2}$$
 to calculate the slope of the regression line.
- u] Determine the intercept (b_0): Use the formula, to calculate the intercept of regression line.

$$b_0 = \bar{y} - b_1 \bar{x}$$



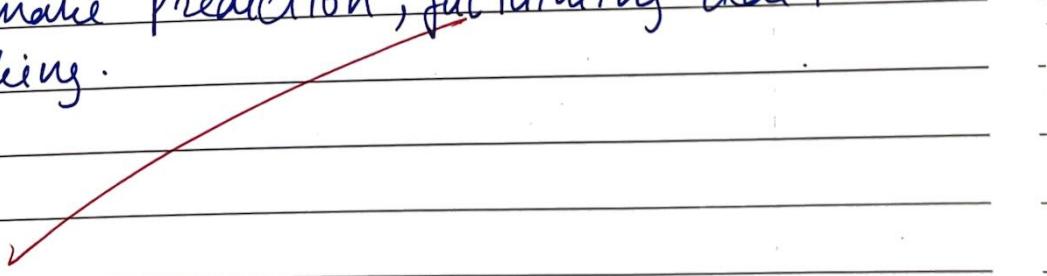
- 5] Form the regression eqn: Combine the slope & intercept to form the regression eqn

$$y = b_0 + b_1 x$$

- 6] Make predictions: Use the regression equation to predict the value of the dependent variable based on new values of the independent variable.
- 7] Evaluate the model: Assess the goodness-of-fit using metrics such as the coefficient of determination (R^2) & the residuals (differences betⁿ observed & predicted values)

Conclusion:-

Simple linear regression provides a straightforward method to model the relationship betⁿ two variables & make prediction, facilitating data driven decision-making.





Practical 8:

Aim: To implement simple logistic regression on a given dataset.

Theory:

Simple logistic regression is a statistical method used to model the relationship between a binary dependent variable & one independent variable. Unlike linear regression, which predicts a continuous outcome, logistic regression predicts the probability of a binary outcome [eg: success / failure, yes / no].

Key Concepts:

1. Dependent Variable (Y): The binary outcome we are trying to predict, often coded as 0 or 1.
2. Independent Variable (X): The variable used to predict the binary outcome.
3. Logit Function: The natural logarithm of the odds of the dependent variable being 1. The logistic regression model is based on the logit function.
$$\text{logit}(Y) = \ln(Y/(1-Y)) = b_0 + b_1 X$$

where p is the probability that the outcome equals 1.



4. Sigmoid function :-

The function that converts the logit back to a probability given by

$$P = \frac{1}{1 + e^{-x}}$$

5. Maximum likelihood Estimation [MLE]:

A method used to estimate the parameters [coefficients] of the logistic regression model by maximizing the likelihood function.

Steps in simple Logistic Regression:-

1. Plot the data:

visualize the data points on a scatter plot to observe the relationship betn x & the binary outcome y.

2. Initialize the model: define the logistic regression model with the logit function.

3. Estimate Parameters (b_0 & b_1): Use maximum likelihood estimation to find the coefficients that best fit the data.

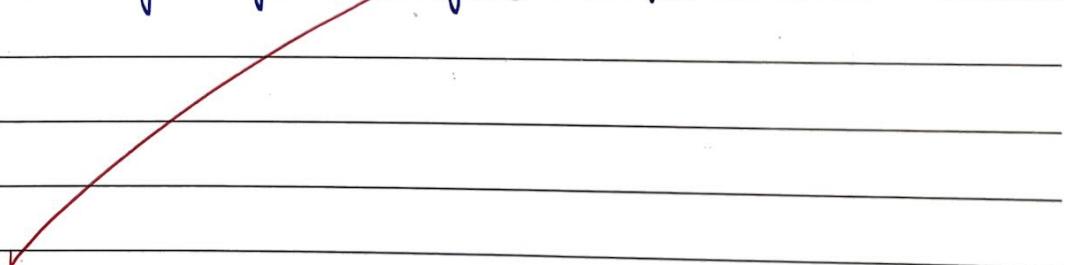
4. Form the logistic regression eqn:- Combine the coefficients to form the logistic regression eqn.



5. Make Predictions: use the logistic regression equation to predict the probability of the binary outcome for new values of the independent variable.
6. Classify Outcomes: Determine the binary outcome by applying a threshold [eg: if $P > 0.5$, predict 1 else 0]
7. Evaluate the model : Assess the model's performance using metrics such as accuracy, precision, recall & the area under the ROC curve [AUC]

Conclusion:-

Simple logistic regression allows us to model the probability of a binary outcome based on one predictor variable making it useful for classification tasks.





Practical 9

Aim: Building a classification Model using K-nearest neighbours (KNN)

Theory:-

K-Nearest Neighbors [KNN] is a simple yet powerful classification & regression algorithm that is a non-parametric & instance-based. It belongs to a category of algorithms called Lazy-Learners because, unlike many algo. that explicitly build models during training, KNN does not learn a model. Instead, it memorizes the training data & performs classification or regression only when it receives a query instance. The fundamental assumption of KNN is that similar data points are likely to belong to the same class, making it a strong tool for classification tasks.

Key concepts:-

1] Instance based learning:- It means that it stores the training instances rather than constructing an explicit model. For each new data point, it finds its ~~to~~ to those stored instances. When a prediction is made, KNN uses the training instances that are closest to the new data points to classify it.



2. Distance Metric: -

KNN relies on a distance metric to measure the similarity b/w data points.

i) Euclidean distance: - The dist. b/w 2 points in space [st. line]

ii) Manhattan distance: - The sum of the absolute differences b/w the co-ordinates of two points.

3. Majority Voting: -

In KNN classification the class label of a new data point is determined by a vote voting mechⁿ. Once the 'k' nearest neighbors are identified their class labels are observed & the new point is assigned to the class that is most frequent among these 'k' neighbors. This is the essence of how KNN makes predictions.

4. Parameter 'k': -

The parameter 'k' determines how many neighbors are considered for the majority voting process.

- if 'k' is small makes the algo. more sensitive to local variations in the data & lead to overfitting.

- if 'k' is large, it smoothes the decision boundary, reducing the risk of overfitting but it can potentially ignore subtle but imp. structures in the data.



Steps in building a KNN classification model :

1. Data preprocessing:-

Before applying KNN, the data must be preprocessed to ensure that the distance metric works correctly.

Apply normalization or standardization to scale each feature to a common range. Handle the missing values by either imputing the missing values or by removing affected rows & columns.

2. Selecting the value of 'k':-

Choosing the right 'k' is critical. A small 'k' leads to more flexible decision boundaries but can make the model sensitive to noise. A large 'k' produces smoother, more generalized boundaries but may miss imp. local patterns.

3. Computing distance:-

Once the model receives a new data point to classify, it computes the distance from this point to all the data points in the training set using the chosen distance metric.

4. Identifying nearest neighbors:-

After calculating the distances the model selects the 'k' train instances with the smallest



distances to the new data point.

5. Majority Voting:

The model then examines the class labels of the identified ' k ' nearest neighbors. The class with the most votes is assigned as the predicted label for the new data points.

6. Model Evaluation:

Once the KNN model has made predictions, it needs to be evaluated to assess its performance.

Conclusion:-

The ~~K-means~~ KNN builds a classification model by assigning class labels based on the majority votes among the nearest neighbors, making it a straightforward and versatile classifier.