

Name Supriya S

Std.: 6th Sem

Telephone No.

### Blood Group.

Div. 'E' Sec

Sub. M1 lab

Roll No.

Roll No. IBM22CS350

E-mail ID:

Birth Day.

Sr.No.	Title	Page No.	Sign./Remarks
1.	Different ways of importing datasets	1	9
2.	Different data preprocessing techniques.	3	Cpt 1.3 -63m
3.	TD3 implementation on weather dataset	5	Implementation
4.	Linear Regression	8	8
5.	Logistic Regression	13	8
6.	KNN Regression	15	8
7.	Random Forest	17	8
8.	Ada Boost	22	8
9.	K Means Clustering	24	8
10.	PCA	26	#1st 2025

## Lab - 0

1. To Do Exercise ways of importing datasets  
 i) Four different ways of initializing values directly into DataFrame
- ```
df = pd.DataFrame({ 'USN': ['450', '451', '452',
                             '453', '454'],
                      'Name': ['Sita', 'Rita', 'Meena', 'Manya',
                               'Keerthi'],
                      'Marks': [94, 89, 87, 86, 84] })
```

print(df)

- 2) Importing datasets from sklearn.datasets  
 from sklearn.datasets import load\_diabetes  
 diabetes = load\_diabetes()

```
df = pd.DataFrame(diabetes.data, columns=
                   diabetes.feature_names)
```

df['target'] = diabetes.target  
 print(df)

- 3) Importing datasets from a specific .csv file

file = 'sample sales data.csv'

df = pd.read\_csv(file)

print(df.head())

- 4) Downloading datasets from existing dataset repositories like Kaggle, UCI, etc

file = 'Diabetes.csv'

df = pd.read\_csv(file)

print(df.head())

2. import yfinance as yf

import pandas as pd

import matplotlib.pyplot as plt

```
# tickers = ["HDFCBANK.NS", "ICICIBANK.NS",  
           "KOTAKBANK.NS"]
```

```
data = yf.download(tickers, start="2024-01-01",  
                   end="2024-12-31",  
                   groupby='ticker')
```

```
hdfc_data = data['HDFCBANK.NS']
```

```
hdfc_data['Daily Return'] = hdfc_data['Close'].pct_change()
```

```
icici_data = data['ICICIBANK.NS']
```

```
icici_data['Daily Return'] = icici_data['Close'].pct_change()
```

```
kotak_data = data['KOTAKBANK.NS']
```

```
kotak_data['Daily Return'] = kotak_data['Close'].pct_change()
```

```
plt.figure(figsize=(12,6))
```

```
plt.subplot(2, 1, 1)
```

```
hdfc_data['Close'].plot(title="Closing Price")
```

```
icici_data['Close'].plot(title="Closing Price")
```

```
kotak_data['Close'].plot(title="Closing Price")
```

```
plt.subplot(2, 1, 2)
```

```
hdfc_data['Daily Return'].plot(title="Daily Returns")
```

```
icici_data['Daily Return'].plot(title="Daily Returns")
```

```
kotak_data['Daily Return'].plot(title="Daily Returns")
```

```
plt.tight_layout()
```

```
plt.show()
```

## Lab - 1

1. housing.csv : Demo of various data pre-import pandas as pd processing technique  
df = pd.read\_csv("housing.csv")  
print("Information of all columns: ")  
print(df.info())  
print("InStatistical information of all numerical columns: ")  
print(df.describe())  
print(df['Ocean Proximity'].value\_counts)  
missingvalues = df.isnull().sum()  
col = missingvalues[missingvalues > 0]  
print(col)

### 2. Diabetes dataset

1. Which columns in the dataset had missing values? How did you handle them?

None of the columns had missing values but if were present then numerical column's NaN can be replaced with median and categorical column's null can be replaced with mode.

2. Which categorical columns did you identify in the dataset? How did you encode them?

The diabetes dataset had gender, and class as categorical columns and adult dataset had workclass

12-3

education, marital-status, occupation, relationship, race, gender, native-country and income as categorical column.

Using Ordinal Encoder we can encode categorical columns.

3. What is the difference b/w Min - Max Scaling and Standardization? When would you use one over the other?  
Min - Max Scaling :- transforms data to fit with a specific range (usually 0 to 1)

Standardization scales data by subtracting mean and dividing by S.D., with mean as 0 and S.D as 1

*Get rid of*

Use Min-Max scaling when you need to preserve the relative relationships b/w values within a feature and use standardization when your algorithm is sensitive to outliers or assumes a normal distribution.

## Lab No. 2

Implement ID3 (Iterative Dichotomiser 3 algorithm on Tennis Weather dataset

```
import pandas as pd
```

```
data = pd.read_csv('tennis.csv')
```

```
print(data)
```

```
from sklearn.preprocessing import LabelEncoder
```

```
outlook = LabelEncoder()
```

```
temp = LabelEncoder()
```

```
humidity = LabelEncoder()
```

```
wind = LabelEncoder()
```

```
play = LabelEncoder()
```

```
data['outlook'] = outlook.fit_transform
```

```
(data['outlook'])
```

```
data['temp'] = temp.fit_transform
```

```
(data['temp'])
```

```
data['humidity'] = humidity.fit_transform
```

```
(data['humidity'])
```

```
data['windy'] = wind.fit_transform
```

```
(data['windy'])
```

```
data['play'] = play.fit_transform
```

```
(data['play'])
```

```
print(data)
```

features cols = ['outlook', 'temp', 'humidity', 'windy']

```
x = data[features_cols]
```

```
y = data.play
```

```
print(x)
```

```
print(y)
```

from sklearn.model\_selection import  
train\_test\_split

x\_train, x\_test, y\_train, y\_test =  
train\_test\_split(x, y, test\_size=0.3)

from sklearn.tree import

DecisionTreeClassifier

classifier = DecisionTreeClassifier(criterion  
= 'entropy')

classifier.fit(x\_train, y\_train)

classifier.predict(x\_test)

classifier.score(x\_test, y\_test)

clf = DecisionTreeClassifier(criterion  
= 'entropy')

clf.fit(x, y)

importances = clf.feature\_importances\_

feature\_importance\_df = pd.DataFrame

({'Feature': x.columns,}

'Information Gain': importances}

})

feature\_importance\_df = feature

importance\_df.sort\_values(by=

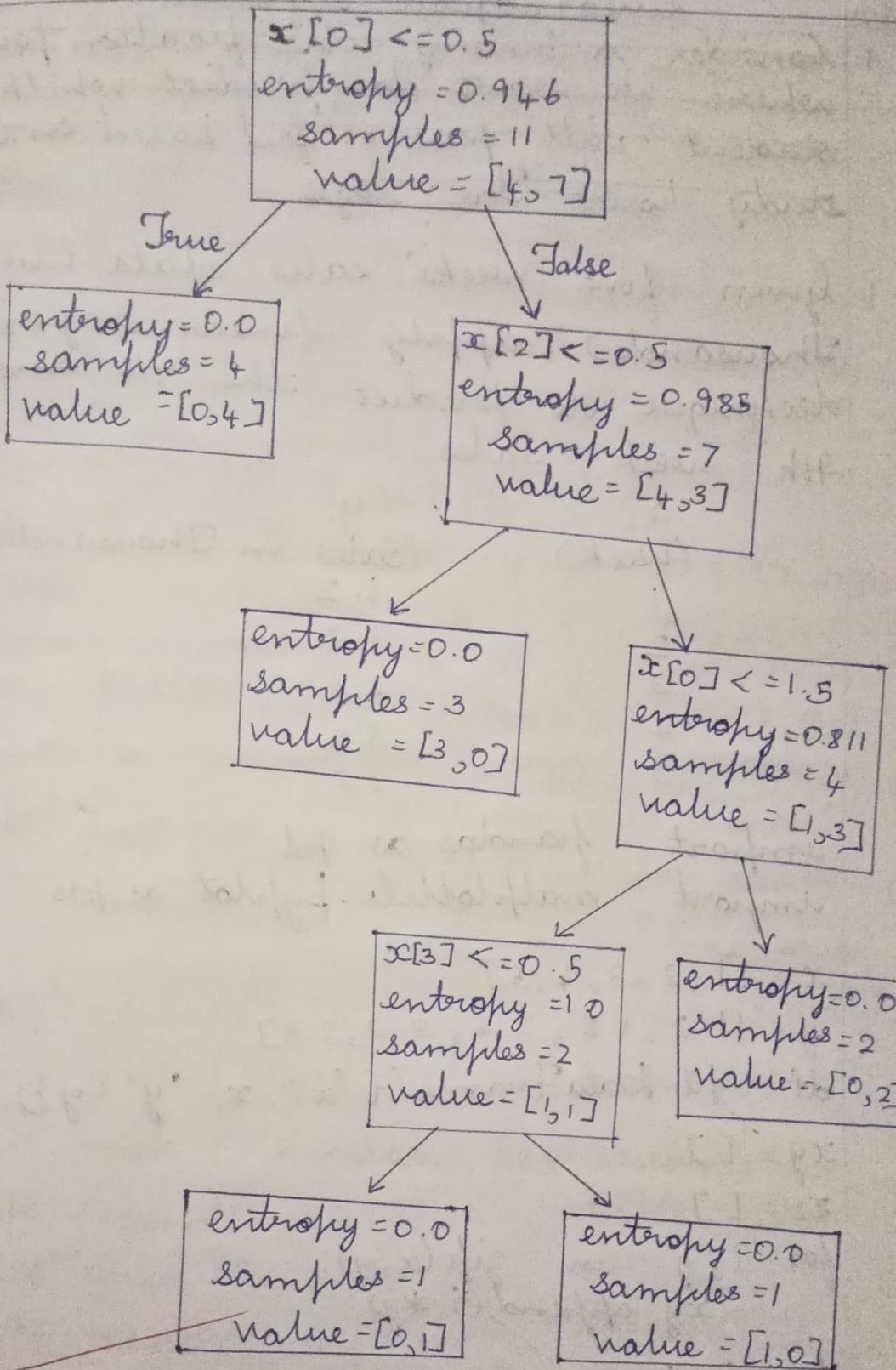
'Information Gain':

ascending=False)

print(feature\_importance\_df)

from sklearn import tree

tree.plot\_tree(classifier)



|   | Feature  | Information Gain |
|---|----------|------------------|
| 0 | outlook  | 0.362629         |
| 3 | windy    | 0.274205         |
| 2 | humidity | 0.211237         |
|   | temp     | 0.151929         |

## Lab - 3

19-3-25

## Linear Logistic Regression

+ Consider a binary classification problem where we want to predict whether a student will pass or fail based on their study hours. The logic

- Given five weeks' sales data (in Thousands), apply linear regression technique to predict the 7th and 9th week sales.

| $x_i$<br>(Week) | $y_i$<br>(Sales in Thousands) |
|-----------------|-------------------------------|
| 1               | 1.2                           |
| 2               | 1.8                           |
| 3               | 2.6                           |
| 4               | 3.2                           |
| 5               | 3.8                           |

import pandas as pd

import matplotlib.pyplot as plt

$x = [1, 2, 3, 4, 5]$

$y = [1.2, 1.8, 2.6, 3.2, 3.8]$

df = pd.DataFrame({ 'x': x, 'y': y })

$xy = []$

$x2 = []$

for i, j in zip(x, y):  
 $xy.append(i * j)$

for i in x:

~~$x2.append(i ** 2)$~~

$xmean = sum(x) / len(x)$

$ymean = sum(y) / len(y)$

$xy\_mean = sum(xy) / len(xy)$

$x2\_mean = sum(x2) / len(x2)$

$$a_1 = (x_{\text{y-mean}} - x_{\text{mean}} * y_{\text{mean}}) / (x_{\text{2-mean}} - x_{\text{mean}} * x_{\text{mean}})$$

$$a_0 = y_{\text{mean}} - a_1 * x_{\text{mean}}$$

```
print(f "a0 = {a0}, a1 = {a1}")
```

```
x_input = int(input("Enter the  
value of x:"))
```

$$y_{\text{pred}} = a_0 + a_1 * x_{\text{input}}$$

```
print(f "Predicted y for x = {x_input}  
is : {y_pred}")
```

```
plt.scatter(x, y, color='blue', label  
= 'Data Points')
```

```
plt.plot(x, [a0 + a1 * xi for xi in x],  
color='red', label='Regression  
line')
```

```
plt.xlabel('x')
```

```
plt.ylabel('y')
```

```
plt.title('Linear Regression')
```

```
plt.legend()
```

```
plt.grid(True)
```

```
plt.show()
```

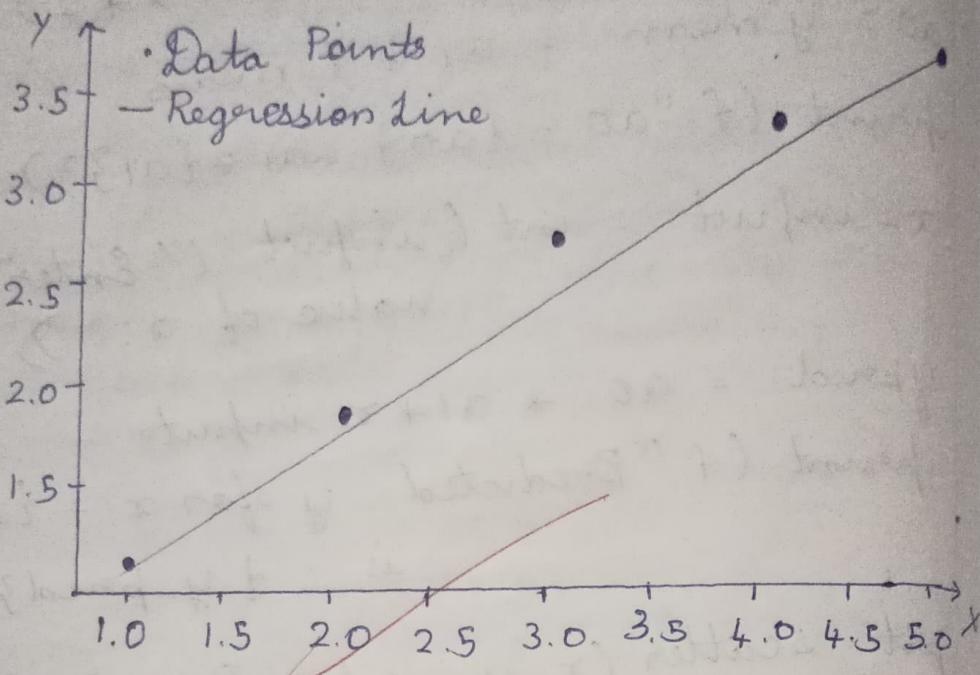
O/P:-

$$a_0 = 0.54, a_1 = 0.68$$

Enter the value of x: 7

Predicted y for x = 7 is : 5.16

## Linear Regression



2. Find linear regression of the data of week & product sales. Use matrix approach for finding Linear Regression

| $X_i$<br>(week) | $Y_i$<br>(sales, in thousand) |
|-----------------|-------------------------------|
| 1               | 1                             |
| 2               | 3                             |
| 3               | 4                             |
| 4               | 8                             |

```
import numpy as np
import matplotlib.pyplot as plt
def matrix(x, y, n):
    ax = np.ones((n, 2))
```

~~ax[:, 1] = x~~

~~ay = np.array(y).reshape(-1, 1)~~

~~ay = np.array(y).reshape(-1, 1)~~

```

 $xt = ax.T$ 
 $A = np.dot(xt, ax)$ 
 $\det A = np.linalg.det(A)$ 
if  $\det A \neq 0$ :
     $A_{inv} = np.linalg.inv(A)$ 
     $a = np.dot(A_{inv}, np.dot$ 
         $(xt, ay))$ 
    return a.flatten()
else:
    return "Inverse doesn't exist"

```

$n = \text{int(input())}$   
 $x = [\text{int(input('Enter')) for i in range(n)}]$   
 $y = [\text{int(input()) for i in range(n)}]$   
 $res = \text{matrix}(x, y, n)$   
if  $\text{isinstance(res, str)}$ :  
 print(res)
else:  
 $a_0, a_1 = res$   
 print(f" Slope ( $a_1$ ) = {a\_1},  
 Intercept ( $a_0$ ) = {a\_0}")

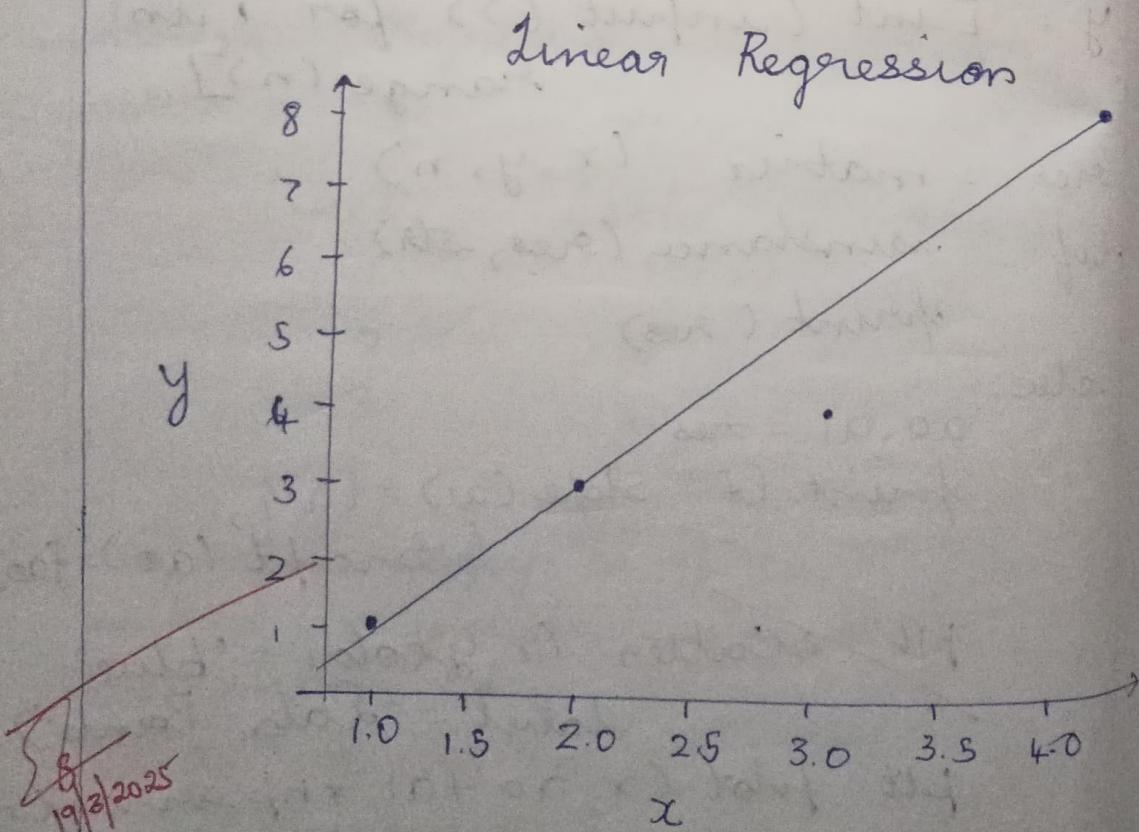
~~plt.scatter(x, y, color = 'blue',  
 label = 'Data Points')~~  
 $plt.plot(x, a_0 + a_1 * np.array(x),$   
 color = 'red', label  
 = 'Regression Line?')

```
plt.xlabel('x')  
plt.ylabel('y')  
plt.title('Linear Regression')  
plt.legend()  
plt.grid(True)  
plt.show()
```

O/P:-

4  
1  
2  
3  
4  
1  
3  
4  
8

Slope ( $a_1$ ) = 2.2 , Intercept ( $a_0$ ) = 1



## Lab - 5

## Logistic Regression

1. Consider a binary classification problem where we want to predict whether a student will pass or fail based on their study hours. The logistic regression model has been trained, and the learned parameters are  $a_0 = -5$  (intercept) and  $a_1 = 0.8$  (coefficient for study hours).

a. Write the logistic regression equation for this problem.

$$P(y=1|x) = \frac{1}{1+e^{-(a_0+a_1x)}} = \frac{1}{1+e^{-(-5+0.8x)}}$$

b. Calculate the probability that a student who studies for 7 hours will pass

$$P(y=1|7) = \frac{1}{1+e^{-(65+0.8(7))}} = \frac{1}{1+e^{-0.6}}$$

$$\approx 0.6457 \text{ or } 64.57\%$$

c. Determine the predicted class (pass or fail) for this student based on a threshold of 0.5

Since  $P(y=1|7) = 64.57\%$ .

alone 0.5 or 50%.  
Predicted class is pass.

2. Consider  $z = [2, 1, 0]$  for three classes. Apply SoftMax function to find the probability values of three classes.

$$P(y=1|z) = \frac{e^{z_1}}{\sum_j e^{z_j}}, \quad z = [2, 1, 0]$$

is the SoftMax function

$$z_1 = 2 \quad z_2 = 1 \quad z_3 = 0$$

$$P(y=1|z) = \frac{e^{z_1}}{e^{z_3} + e^{z_1} + e^{z_2}} = \frac{e^2}{e^2 + e^1 + e^0} \approx 0.665$$

$$P(y=2|z) = \frac{e^{z_2}}{e^{z_3} + e^{z_1} + e^{z_2}} = \frac{e^1}{e^0 + e^2 + e^1} = \frac{1}{1.107} \approx 0.248$$

~~$$P(y=3|z) = \frac{e^{z_3}}{e^{z_3} + e^{z_1} + e^{z_2}} = \frac{e^0}{e^0 + e^2 + e^1} = \frac{1}{1.107} \approx 0.090$$~~

Lab - 6

Build KNN Classification model for diabetes a given dataset

- Consider the following dataset, for  $K = 3$  and test data  $(X, 35, 100)$  as  $(\text{Person}, \text{Age}, \text{Salary}_K)$  solve using Knn classifier model and predict the target.

| Person | Age | Salary <sub>K</sub> | Target | Euclidean Distance                     |
|--------|-----|---------------------|--------|----------------------------------------|
| A      | 18  | 50                  | N      | $\sqrt{(35-18)^2 + (100-50)^2} = 52.8$ |
| B      | 23  | 55                  | N      | 46.6                                   |
| C      | 24  | 70                  | N      | 31.9                                   |
| D      | 41  | 60                  | Y      | 40.5                                   |
| E      | 43  | 70                  | Y      | 31.1                                   |
| F      | 38  | 40                  | Y      | 60.1                                   |
| X      | 35  | 100                 | ?      |                                        |

$$d = \sqrt{\frac{(\text{Age}_1 - \text{Age}_2)^2 + (\text{Salary}_K - \text{Salary}_K)^2}{3}} \text{ as } K=3$$

$(41, 60) \quad 31.1 \quad Y$   
 $(24, 70) \quad 31.9 \quad N$   
 $(41, 60) \quad 40.5 \quad Y$

Since majority is Y  
 $(35, 100)$  Target is Y

- For this dataset How to choose the k value? Demonstrate using accuracy rate and error rate
  - Split the data into training and testing sets.
  - Train the model with various values of  $K$  (e.g. 1, 3, 5, 7, 9)

3. Evaluate accuracy and error rate for each  $k$ .

4. Plot accuracy =  $\frac{\text{Number of correct predictions}}{\text{Total predictions}}$

$$\text{error rate} = 1 - \text{accuracy}$$

4. Plot accuracy vs  $k$  or error rate vs  $k$  and select  $k$  with the highest accuracy or lowest error rate.

For diabetes dataset

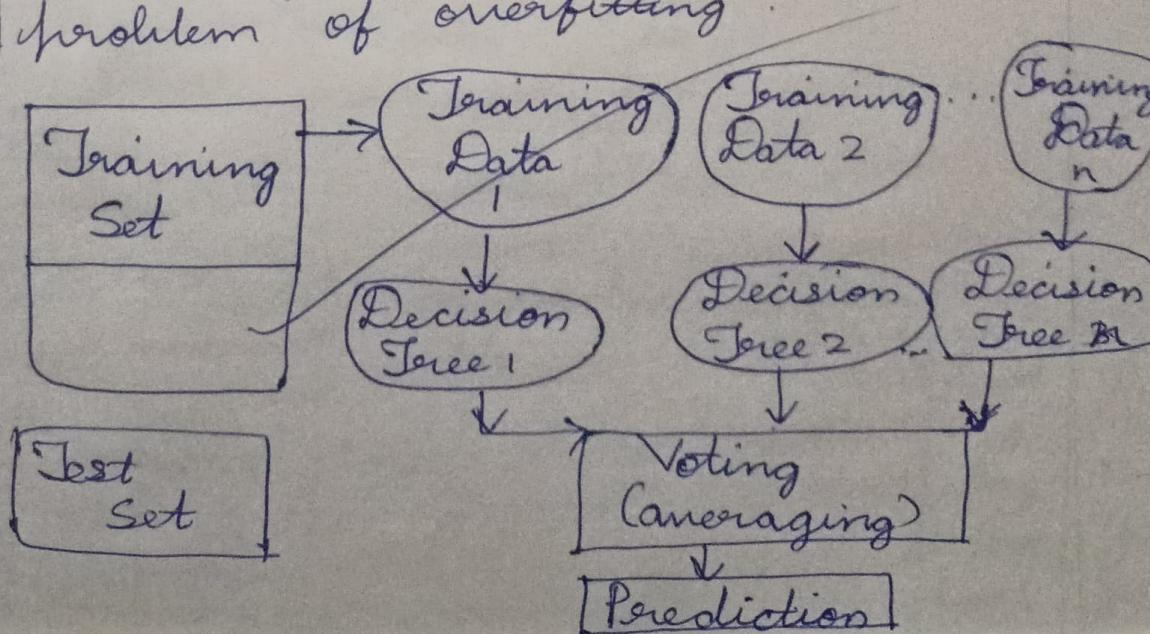
What is the purpose of feature scaling? How to perform it?

KNN uses distance, so features with larger values dominate. Scaling ensures fairness. Using standardization (z-score normalization).

8  
21/4/2025

## Random Forest

- Random Forest is a popular machine learning algorithm that belongs to the supervised learning technique. Can be used for both Classification and Regression problems in ML.
- Based on the concept of ensemble learning, which is a process of combining multiple classifiers to solve a complex problem and to improve the performance of the model.
- Random Forest is a classifier that contains a number of decision trees on various subsets of the given dataset and takes the average to improve the predictive accuracy of that dataset.
- The greater number of trees in the forest leads to higher accuracy and prevents the problem of overfitting.



## Random Forest algorithm

1. Select random  $K$  data points from the training set.
2. Build the decision trees associated with the selected data points (Subsets)
3. Choose the number  $N$  for decision trees that you want to build.
4. Repeat 1 & 2
5. For new data points, find the predictions of each decision tree and assign the new data points to the category that wins the majority votes.

✓ Gov  
16.0

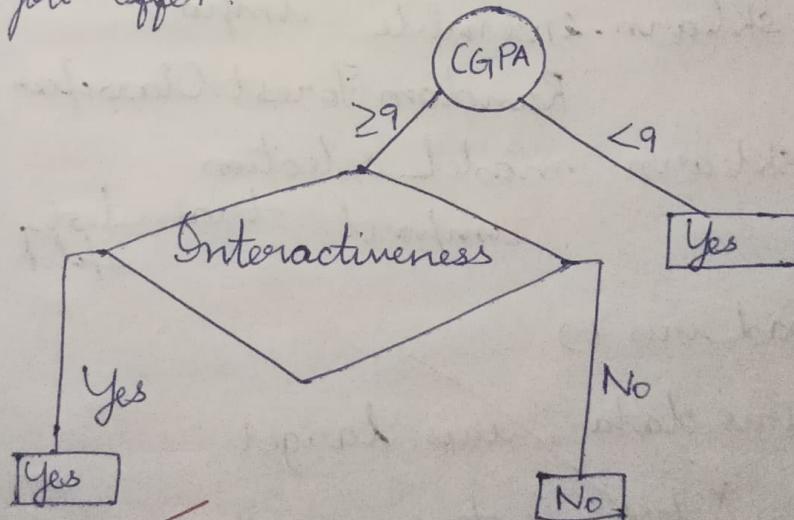
Implement Random Forest ensemble method

Draw the Decision tree considering CGPA as root node

Sample - S<sub>1</sub>

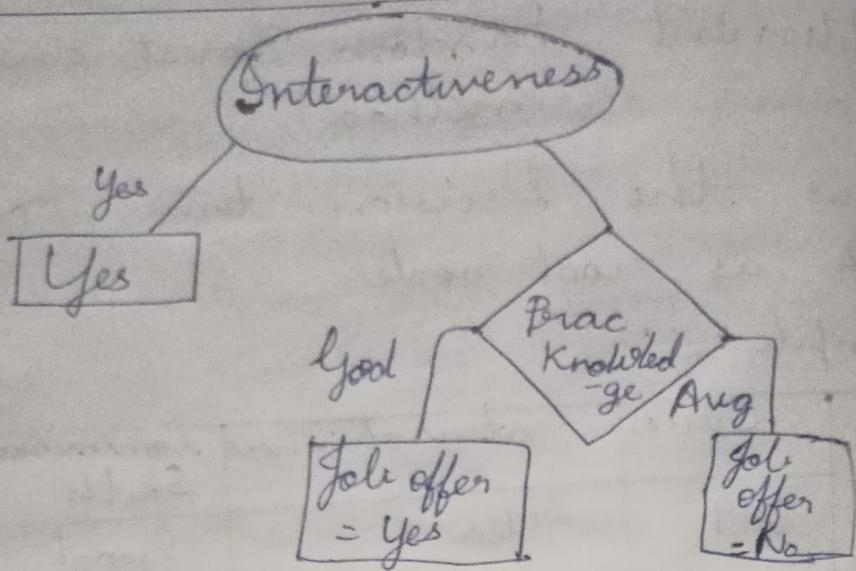
| SNo. | CGPA     | Interactiveness | Communication Skills | Prac Know | Job off |
|------|----------|-----------------|----------------------|-----------|---------|
| 1.   | $\geq 9$ | Yes             | Good                 | Good      | Yes     |
| 2.   | < 9      | No              | Moderate             | Good      | Yes     |
| 3.   | $\geq 9$ | No              | Moderate             | Aug       | No      |
| 4.   | $\geq 9$ | No              | Moderate             | Aug       | No      |
| 5.   | $\geq 9$ | Yes             | Moderate             | Good      | Yes     |

Class Entropy for the target class  
(Job Offer).



Sample S<sub>2</sub>

| S.No. | CGPA     | Interactiveness | Communication Skills | Prac Know | Job off |
|-------|----------|-----------------|----------------------|-----------|---------|
| 2.    | < 9      | No              | Moderate             | Good      | Yes     |
| 3.    | $\geq 9$ | No              | Moderate             | Aug       | No      |
| 3.    | $\geq 9$ | No              | Moderate             | Aug       | No      |
| 5.    | $\geq 9$ | Yes             | "                    | Good      | Yes     |
| 5.    | $\geq 9$ | Yes             | "                    | Good      | Yes     |



Write Python Code to implement the following using "iris.csv" dataset  
 Build a RF classifier to classify IRIS flower dataset

```

import pandas as pd
from sklearn.ensemble import
    RandomForestClassifier
from sklearn.model_selection
    import train_test_
split
iris = load_iris()
X, y = iris.data, iris.target
Xtrain, Xtest, ytrain, ytest =
    train_test_split(X, y, test_size=0.3,
                     random_state=42)
  
```

```

rf = RandomForestClassifier(n_estimators=10)
rf.fit(Xtrain, ytrain)
  
```

point ("Accuracy with 10 trees .")

O/P:-

Accuracy with 10 trees : 1

Best accuracy : 1

Best accuracy Score : 1

Confusion Matrix

|    |    |
|----|----|
| TP | FP |
| FN | TN |

Lecture-8

Implement Boosting ensemble method.

Using AdaBoost Algorithm show the decision stump calculation steps for the attribute CGPA.

| CGPA     | Interactivity | Prof<br>Know | Comm<br>Skill | Job<br>Prof |
|----------|---------------|--------------|---------------|-------------|
| $\geq 9$ | Yes           | Good         | Good          | Yes         |
| $< 9$    | No            | Good         | Moderate      | Yes         |
| $\geq 9$ | No            | Avg          | "             | No          |
| $< 9$    | "             | "            | Good          | "           |
| $\geq 9$ | Yes           | Good         | Moderate      | Yes         |
| $\geq 9$ | "             | "            | "             | "           |

$\Sigma$   
if initial weight =  $1/6$

$$\epsilon_i = \sum_{j=1}^6 H_i(d_j) W_t(d_j)$$

$$\epsilon_{CGPA} = 2 \times \frac{1}{6} = 0.33$$

$$\Delta_{CGPA} = \frac{1}{2} \frac{\ln(1 - \epsilon_{CGPA})}{\epsilon_{CGPA}}$$

$$= 0.347$$

$$Z_{LGP\Delta} = \frac{1}{6} \times 4 \times e^{-0.347} + \frac{1}{6} \times 2 \times e^{+0.347}$$

$$= 0.9428$$

$$w + (d_j)_{i+1} = \frac{\frac{1}{6} \times e^{-0.347}}{0.9428} = 0.1249$$

for incorrect instance

$$= \frac{\frac{1}{6} \times e^{+0.347}}{0.9428} = 0.2601$$

| CGPA | Initial weight | Updated weight |
|------|----------------|----------------|
| >=9  | $\frac{1}{6}$  | 0.1249         |
| <9   | $\frac{1}{6}$  | 0.2601         |
| >=9  | $\frac{1}{6}$  | 0.1249         |
| <9   | $\frac{1}{6}$  | 0.2601         |
| >=9  | $\frac{1}{6}$  | 0.1249         |
| >=9  | $\frac{1}{6}$  | 0.1249         |

Best accuracy score = ,

~~Confusion matrix~~

TN FP

FN TP

Lab - 9

## K Means Algorithm

Compute two clusters using K-means algorithm for clustering where initial cluster centers are  $(1.0, 1.0)$  and  $(5.0, 7.0)$ . Execute for 2 iterations.

Iteration 1

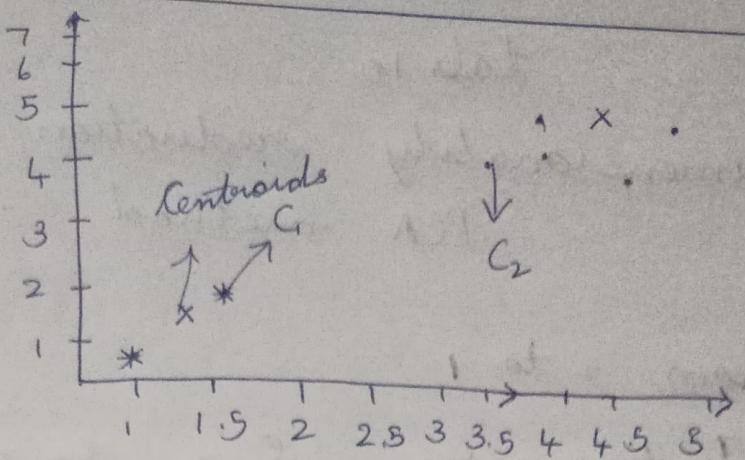
| Record No.     | A   | B   | Cluster<br>center<br>$(1.0, 1.0)$<br>$(5.0, 7.0)$ |
|----------------|-----|-----|---------------------------------------------------|
| R <sub>1</sub> | 1.0 | 1.0 | C <sub>1</sub>                                    |
| R <sub>2</sub> | 1.5 | 2.0 | C <sub>1</sub>                                    |
| R <sub>3</sub> | 3.0 | 4.0 | C <sub>1</sub>                                    |
| R <sub>4</sub> | 5.0 | 7.0 | C <sub>2</sub>                                    |
| R <sub>5</sub> | 3.5 | 5.0 | C <sub>2</sub>                                    |
| R <sub>6</sub> | 4.5 | 5.0 | C <sub>2</sub>                                    |
| R <sub>7</sub> | 3.5 | 4.5 | C <sub>2</sub>                                    |

New centroids  $(1.83, 2.33)$

$(4.128, 5.375)$

Iteration 2

| Record No.     | A   | B   |                |
|----------------|-----|-----|----------------|
| R <sub>1</sub> | 1.0 | 1.0 | C <sub>1</sub> |
| R <sub>2</sub> | 1.5 | 2.0 | C <sub>1</sub> |
| R <sub>3</sub> | 3.0 | 4.0 | C <sub>2</sub> |
| R <sub>4</sub> | 5.0 | 7.0 | C <sub>2</sub> |
| R <sub>5</sub> | 3.5 | 5.0 | C <sub>2</sub> |
| R <sub>6</sub> | 4.5 | 5.0 | C <sub>2</sub> |
| R <sub>7</sub> | 3.5 | 4.5 | C <sub>2</sub> |



New centroids

$$[1.25, 1.5] \quad [3.9, 5.1]$$

Elbow Method for Optimal  $k$

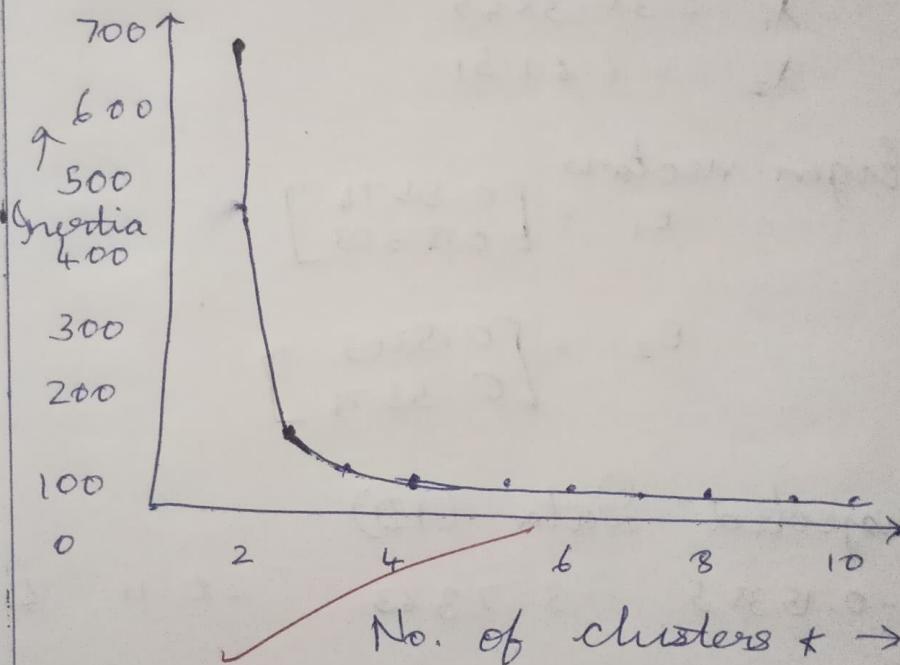


Table 10

## Dimensionality reduction using PCA method

From 2 to 1

| Feature | Ex 1 | Ex 2 | Ex 3 | Ex 4 |
|---------|------|------|------|------|
| $x_1$   | 4    | 8    | 13   | 7    |
| $x_2$   | 11   | 4    | 5    | 14   |

Eigen Values

$$\lambda_1 = 30.3849$$

$$\lambda_2 = 6.6151$$

Eigen vectors

$$e_1 = \begin{bmatrix} 0.5574 \\ 0.8303 \end{bmatrix}$$

$$e_2 = \begin{bmatrix} 0.8303 \\ 0.5574 \end{bmatrix}$$

Projected Data (1D)

$$[-0.15385 \quad -3.7363 \quad -0.11 \quad 4.0]$$

heart.csv

Accuracy before PCA : 0.85

Accuracy after PCA : 0.83

DS  
21/5/2025