

**VISVESVARAYA TECHNOLOGICAL
UNIVERSITY**
“JnanaSangama”, Belgaum -590014, Karnataka.



**LAB REPORT
on**

Machine Learning (23CS6PCMAL)

Submitted by

Navya Billalar (1BM22CS175)

in partial fulfillment for the award of the degree of

**BACHELOR OF ENGINEERING
in
COMPUTER SCIENCE AND ENGINEERING**



**B.M.S. COLLEGE OF ENGINEERING
(Autonomous Institution under VTU)
BENGALURU-560019
Sep-2024 to Jan-2025**

B.M.S. College of Engineering,
Bull Temple Road, Bangalore 560019
(Affiliated To Visvesvaraya Technological University, Belgaum)

Department of Computer Science and Engineering



CERTIFICATE

This is to certify that the Lab work entitled “Machine Learning (23CS6PCMAL)” carried out by **Navya Billalar (1BM22CS175)**, who is bonafide student of **B.M.S. College of Engineering**. It is in partial fulfillment for the award of **Bachelor of Engineering in Computer Science and Engineering** of the Visvesvaraya Technological University, Belgaum. The Lab report has been approved as it satisfies the academic requirements in respect of an Machine Learning (23CS6PCMAL) work prescribed for the said degree.

Lab Faculty Incharge	
Name: Ms. Saritha A N Assistant Professor Department of CSE, BMSCE	Dr. Kavitha Sooda Professor & HOD Department of CSE, BMSCE

Index

Sl. No.	Date	Experiment Title	Page No.
1	21-2-2025	Write a python program to import and export data using Pandas library functions	3
2	3-3-2025	Demonstrate various data pre-processing techniques for a given dataset	9
3	10-3-2025	Implement Linear and Multi-Linear Regression algorithm using appropriate dataset	13
4	17-3-2025	Build Logistic Regression Model for a given dataset	16
5	24-3-2025	Use an appropriate data set for building the decision tree (ID3) and apply this knowledge to classify a new sample	18
6	7-4-2025	Build KNN Classification model for a given dataset	22
7	21-4-2025	Build Support vector machine model for a given dataset	26
8	5-5-2025	Implement Random forest ensemble method on a given dataset	29
9	5-5-2025	Implement Boosting ensemble method on a given dataset	32
10	12-5-2025	Build k-Means algorithm to cluster a set of data stored in a .CSV file	35
11	12-5-2025	Implement Dimensionality reduction using Principal Component Analysis (PCA) method	38

Github Link:

https://github.com/navya-08/ML_LAB/tree/main

Program 1

Write a python program to import and export data using Pandas library functions

Observation:

Date 3/3/25
Page _____

LAB-1

```
import pandas as pd
from sklearn.datasets import load_
```

Method - 1

```
import pandas as pd
data = {
    'IBN': ['IBM2QCS175', 'IBM2QCS176', 'IBM2QCS177', 'IBM2QCS178'],
    'Name': ['Navya', 'Nehal', 'Nidhi', 'Nihal', 'Nikhilesh'],
    'Marks': [80, 80, 80, 80, 80]
}
df = pd.DataFrame(data)
print(df)
```

Method - 2

```
from sklearn.datasets import load_diabetes
diabetes = load_diabetes()
df = pd.DataFrame(diabetes.data, columns = diabetes.feature_names)
df['target'] = diabetes.target
print(df.head())
```

Method - 3

```
file_path = 'sample_sales_data.csv'
df = pd.read_csv(file_path)
print("Sample data : ")
print(df.head())
```

Method - 24 :

```
df = pd.read_csv("Mobiles-Dataset.csv")
print("Sample data : ")
print(df.head())
```

II

To Do:

1. import yfinance as yf
import pandas as pd

```
Ticker = ["HDFCBANK.NS", "ICICIBANK.NS",
          "KOTAKBANK.NS"]
```

For Ticker in Ticker:

```
data = yf.download(Ticker, start =
                    "2024-01-01", end = "2024-12-31")
print(f"First 5 rows of {Ticker} dataset:")
print(data.head())
data.to_csv(f"{Ticker}-stock-data.csv")
print(f"{Ticker}-stock-data saved to {Ticker}-stock-data.csv")
```

2.

Closing price and daily returns.

```
import yfinance as yf
import matplotlib.pyplot as plt
plt.figure(figsize=(12, 5))
plt.subplot(2, 1, 1)
```

Date _____
Page _____

```

hdfc = 'HDFC'
hdfc = data['HDFC.BANK.NS']
icici = data['ICICIBANK.NS']
kotak = data['KOTAK.NS']
hdfc['Daily Return'] = hdfc['Close'].pct_change()
icici['Daily Return'] = icici['Close'].pct_change()
kotak['Daily Return'] = kotak['Close'].pct_change()

hdfc['Close'].plot(label='HDFC')
icici['Close'].plot(label='ICICI')
kotak['Close'].plot(label='Kotak')

hdfc['Daily Return'].plot(label='HDFC')
icici['Daily Return'].plot(label='ICICI')
kotak['Daily Return'].plot(label='Kotak')
plt.show()

```

Code:

```

import pandas as pd

# Reading data from a CSV file

data = {

'Name': ['Alice', 'Bob', 'Charlie', 'David', 'Evangline'],
'USN': ['1BM22CS025', '1BM22CS030', '1BM22CS035', '1BM22CS040', '1BM22CS045'],
'Marks': [25, 30, 35, 40, 45]

}

df = pd.DataFrame(data)

print("Sample data:")

print(df.head())

```

Sample data:

	Name	USN	Marks
0	Alice	1BM22CS025	25
1	Bob	1BM22CS030	30
2	Charlie	1BM22CS035	35
3	David	1BM22CS040	40
4	Evangline	1BM22CS045	45

```
from sklearn.datasets import load_diabetes

dia = load_diabetes()

df = pd.DataFrame(dia.data, columns=dia.feature_names)

df['target'] = dia.target

print("Sample data:")

print(df.head())
```

Output:

Sample data:

```
    age      sex      bmi      bp      s1      s2      s3  \
0  0.038076  0.050680  0.061696  0.021872 -0.044223 -0.034821 -0.043401
1 -0.001882 -0.044642 -0.051474 -0.026328 -0.008449 -0.019163  0.074412
2  0.085299  0.050680  0.044451 -0.005670 -0.045599 -0.034194 -0.032356
3 -0.089063 -0.044642 -0.011595 -0.036656  0.012191  0.024991 -0.036038
4  0.005383 -0.044642 -0.036385  0.021872  0.003935  0.015596  0.008142

      s4      s5      s6  target
0 -0.002592  0.019907 -0.017646  151.0
1 -0.039493 -0.068332 -0.092204    75.0
2 -0.002592  0.002861 -0.025930  141.0
3  0.034309  0.022688 -0.009362  206.0
4 -0.002592 -0.031988 -0.046641  135.0
```

```
import yfinance as yf
```

```
import pandas as pd
```

```
import matplotlib.pyplot as plt
```

```
print("\nShape of the dataset:")
```

```

print(data.shape)

# Check column names

print("\nColumn names:")

print(data.columns)

# Summary statistics for a specific stock (e.g., Reliance)

reliance_data = data['RELIANCE.NS']

print("\nSummary statistics for Reliance Industries:")

print(reliance_data.describe())

# Calculate daily returns

# Create a copy of the Reliance data to avoid modifying a slice of the original dataframe

reliance_data = data['RELIANCE.NS'].copy()

# Now, apply the calculation

reliance_data['Daily Return'] = reliance_data['Close'].pct_change()

```

Output:

```

Shape of the dataset:
(247, 15)

Column names:
MultiIndex([('RELIANCE.NS',      'Open'),
            ('RELIANCE.NS',      'High'),
            ('RELIANCE.NS',      'Low'),
            ('RELIANCE.NS',      'Close'),
            ('RELIANCE.NS',      'Volume'),
            ('TCS.NS',          'Open'),
            ('TCS.NS',          'High'),
            ('TCS.NS',          'Low'),
            ('TCS.NS',          'Close'),
            ('TCS.NS',          'Volume'),
            ('INFY.NS',          'Open'),
            ('INFY.NS',          'High'),
            ('INFY.NS',          'Low'),
            ('INFY.NS',          'Close'),
            ('INFY.NS',          'Volume')],
           names=['Ticker', 'Price'])

Summary statistics for Reliance Industries:
   Price      Open      High      Low     Close    Volume
count  247.000000  247.000000  247.000000  247.000000  2.470000e+02
mean   1155.033899  1163.758985  1144.612976  1154.002433  1.316652e+07
std    65.890843   66.876907   65.755901   66.726021   6.754099e+06
min    1015.178443  1017.470038  999.137216  1008.876526  3.370033e+06
25%   1106.532938  1111.081861  1092.347974  1104.997559  8.717141e+06
50%   1155.424265  1163.078198  1146.716157  1155.240967  1.158959e+07
75%   1202.667031  1209.102783  1193.235594  1201.447937  1.530302e+07
max   1297.045129  1308.961472  1281.920577  1302.476196  5.708188e+07

```

```

# Plot the closing price and daily returns

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

plt.subplot(2, 1, 1)
reliance_data['Close'].plot(title="Reliance Industries - Closing Price")

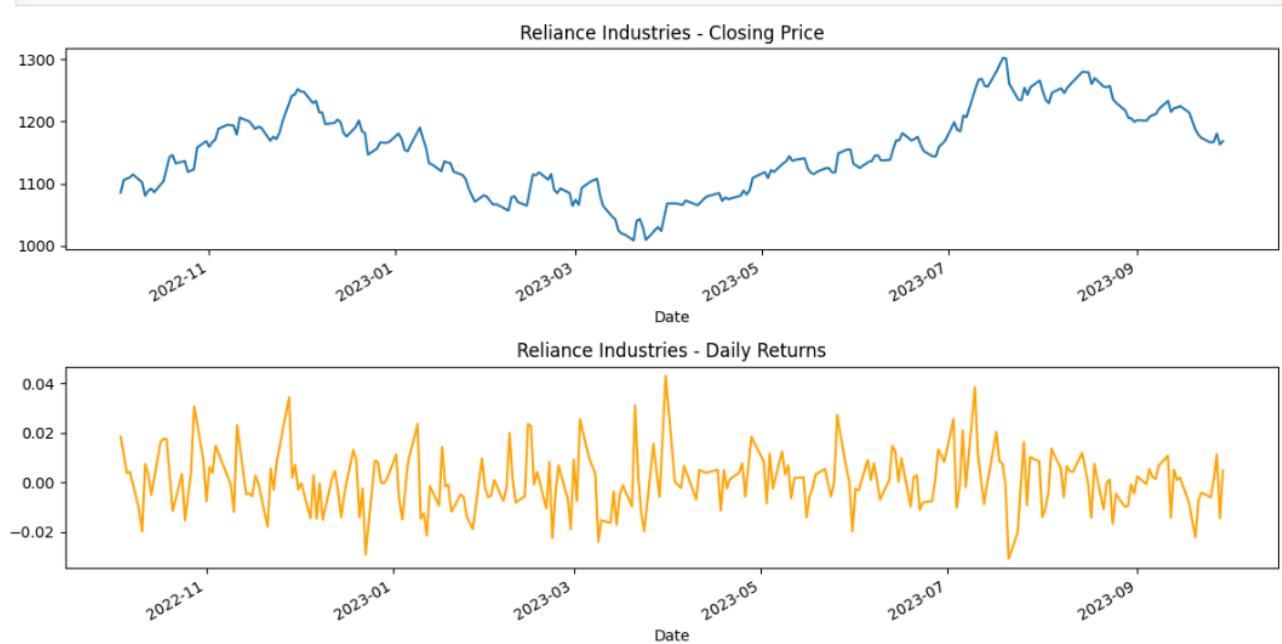
plt.subplot(2, 1, 2)
reliance_data['Daily Return'].plot(title="Reliance Industries - Daily Returns", color='orange')

plt.tight_layout()

plt.show()

```

Output:



Program 2

Observation:

Date / /
Page / /

LAB - 1

I. i. import pandas as pd
file_path = "housing.csv"
df = pd.read_csv(file_path)

ii print("Column Information: ")
print(df.columns())

iii print("Statistical information")
print(df.describe())

iv print(df['Ocean Proximity'].value_counts())

v missing_values = df.isnull().sum()
missing_att = missing_values[missing_values > 0]
print(missing_att)
Suv...,""

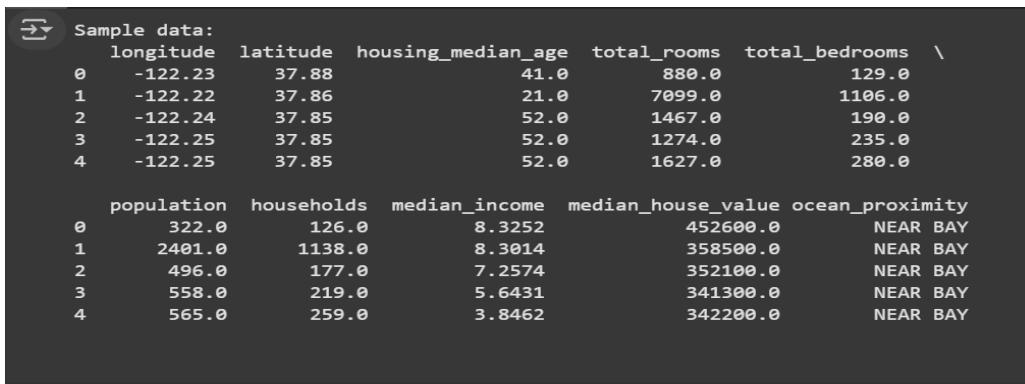
MAP

Code:

i. To load .csv file into the data frame

```
import pandas as pd  
  
file_path = '/content/housing.csv'  
  
df = pd.read_csv(file_path)  
  
print("Sample data:")  
  
print(df.head())  
  
print("\n")
```

output:



The screenshot shows a Jupyter Notebook cell with the title "Sample data:" followed by a table of housing data. The table has 5 rows and 7 columns. The columns are labeled: longitude, latitude, housing_median_age, total_rooms, total_bedrooms, population, households, median_income, median_house_value, and ocean_proximity. The data is as follows:

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_income	median_house_value	ocean_proximity
0	-122.23	37.88	41.0	880.0	129.0	322.0	126.0	8.3252	452600.0	NEAR BAY
1	-122.22	37.86	21.0	7099.0	1106.0	2401.0	1138.0	8.3014	358500.0	NEAR BAY
2	-122.24	37.85	52.0	1467.0	190.0	496.0	177.0	7.2574	352100.0	NEAR BAY
3	-122.25	37.85	52.0	1274.0	235.0	558.0	219.0	5.6431	341300.0	NEAR BAY
4	-122.25	37.85	52.0	1627.0	280.0	565.0	259.0	3.8462	342200.0	NEAR BAY

ii. To display information of all columns

```
print(df.info())
```

Output:

```

→ <bound method DataFrame.info of
   0    -122.23    37.88      41.0    880.0     129.0
   1    -122.22    37.86      21.0    7099.0     1106.0
   2    -122.24    37.85      52.0    1467.0      190.0
   3    -122.25    37.85      52.0    1274.0      235.0
   4    -122.25    37.85      52.0    1627.0      280.0
...
   ..    ...
20635   -121.09    39.48      25.0    1665.0      374.0
20636   -121.21    39.49      18.0     697.0      150.0
20637   -121.22    39.43      17.0    2254.0      485.0
20638   -121.32    39.43      18.0    1860.0      409.0
20639   -121.24    39.37      16.0    2785.0      616.0

population  households  median_income  median_house_value \
   0        322.0       126.0      8.3252    452600.0
   1       2401.0       1138.0      8.3014    358500.0
   2        496.0       177.0       7.2574    352100.0
   3       558.0       219.0       5.6431    341300.0
   4       565.0       259.0       3.8462    342200.0
...
   ..    ...
20635   845.0       330.0      1.5603    78100.0
20636   356.0       114.0       2.5568    77100.0
20637   1007.0      433.0       1.7000    92300.0
20638   741.0       349.0       1.8672    84700.0
20639   1387.0      530.0       2.3886    89400.0

ocean_proximity
   0      NEAR BAY
   1      NEAR BAY
   2      NEAR BAY
   3      NEAR BAY
   4      NEAR BAY
...
   ..    ...
20635   INLAND
20636   INLAND
20637   INLAND
20638   INLAND
20639   INLAND

[20640 rows x 10 columns]

```

iii. To display the count of unique labels for “Ocean Proximity” column

```
print(df['ocean_proximity'].value_counts())
```

Output:

```

→ ocean_proximity
<1H OCEAN      9136
INLAND         6551
NEAR OCEAN     2658
NEAR BAY       2290
ISLAND          5
Name: count, dtype: int64

```

iv. To display which attributes (columns) in a dataset have missing values count greater than zero

```
print(df.isnull().sum())
```

```

→ longitude      0
latitude        0
housing_median_age  0
total_rooms      0
total_bedrooms   207
population       0
households       0
median_income     0
median_house_value 0
ocean_proximity   0
dtype: int64

```

V.To display statistical information of all numerical

```
print(df.describe())
```

Output:

```
    longitude      latitude  housing_median_age  total_rooms \
count  20640.000000  20640.000000  20640.000000  20640.000000 \
mean   -119.569704   35.631861    28.639486  2635.763081
std     2.003532    2.135952   12.585558  2181.615252
min    -124.350000   32.540000    1.000000   2.000000
25%    -121.800000   33.930000   18.000000  1447.750000
50%    -118.490000   34.260000   29.000000  2127.000000
75%    -118.010000   37.710000   37.000000  3148.000000
max    -114.310000   41.950000   52.000000  39320.000000

      total_bedrooms  population  households  median_income \
count  20433.000000  20640.000000  20640.000000  20640.000000 \
mean   537.870553  1425.476744   499.539680   3.870671
std    421.385070  1132.462122   382.329753   1.899822
min     1.000000    3.000000    1.000000   0.499900
25%    296.000000   787.000000   280.000000   2.563400
50%    435.000000  1166.000000   409.000000   3.534800
75%    647.000000  1725.000000   605.000000   4.743250
max    6445.000000  35682.000000  6082.000000  15.000100

      median_house_value
count      20640.000000
mean     206855.816909
std      115395.615874
min     14999.000000
25%    119600.000000
50%    179700.000000
75%    264725.000000
max    500001.000000
```

Program 3

Observation:

24.3.25 LAB-4 Date / / Page /

Linear Regression: [univariate]

$\text{x} \leftarrow$ array of independent values
 $y \leftarrow$
 $n \leftarrow \text{len}(x)$
 $x_mean \leftarrow \text{mean}(x)$
 $y_mean \leftarrow \text{mean}(y)$
Compute numerator = $\sum (x_i - x_mean) * (y_i - y_mean)$
compute denominator = $\sum (x_i - x_mean)^2$
Slope m = numerator / denominator

Intercept b $\leftarrow y_mean - (m * x_mean)$

Print "Slope (m) : ", m
Print "Intercept (b) : ", b

→ Multivariate

$n \leftarrow \text{len}(y)$
 $x = \text{np. c_}[\text{np. ones}(n), x]$
~~theta = np. linalg. inv(x.T * x) * x.T @ y~~
 $b = \text{theta}[0]$
 $m = \text{theta}[1:7]$

Code:

Linear Regression:

```
import numpy as np
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.datasets import fetch_california_housing

california_housing = fetch_california_housing()

X = pd.DataFrame(california_housing.data, columns=california_housing.feature_names)
y = pd.Series(california_housing.target)

X = X[['MedInc', 'AveRooms']]

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
model = LinearRegression()
model.fit(X_train, y_train)

y_pred = model.predict(X_test)

results = pd.DataFrame({'Actual': y_test, 'Predicted': y_pred})
print("Linear Regression Results:")
print(results.head())
```

Output :

```
Linear Regression Results:
      Actual   Predicted
20046  0.47700  1.162302
3024   0.45800  1.499135
15663  5.00001  1.955731
20484  2.18600  2.852755
9814   2.78000  2.001677
```

2. Multiple regression

Code:

```
import numpy as np
import pandas as pd
```

```

import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.datasets import fetch_california_housing
california_housing = fetch_california_housing()
X = pd.DataFrame(california_housing.data, columns=california_housing.feature_names)
y = pd.Series(california_housing.target)
X = X[['MedInc', 'AveRooms']]
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
model = LinearRegression()
model.fit(X_train, y_train)
y_pred = model.predict(X_test)
# Print the actual vs predicted values
results = pd.DataFrame({'Actual': y_test, 'Predicted': y_pred})
print(results.head())

```

Output:

	Actual	Predicted
20046	0.47700	1.162302
3024	0.45800	1.499135
15663	5.00001	1.955731
20484	2.18600	2.852755
9814	2.78000	2.001677

Program 4

Logistic regression

Observation:

Date _____
Page _____

Logistic Regression

Pseudocode:

```
function Sigmoid(x)
    return 1/(1+exp(-x))

m ← no. of samples of x
n ← no. of features of x
w ← zero vector of size n
b ← 0

for epoch in range (1 to epochs)
    p ← sigmoid (x * w + b)
    dw ← (1/m) * x^T * (p - y)
    db ← (1/m) * sum (p - y)
    w ← w - l_rate * dw
    b ← b - l_rate * db
    return w, b

function predict (x, w, b)
    return sigmoid (x * w + b) > 0.5

function DecisionFn(x, w, b)
    return Sigmoid (x * w + b)

// prints the output of a sample
print predict (DecisionFn (x, w, b), w, b)
```

(Handwritten notes below pseudocode)

function Sigmoid(x)
return 1/(1+exp(-x))

m ← no. of samples of x
n ← no. of features of x
w ← zero vector of size n
b ← 0

for epoch in range (1 to epochs)
 p ← sigmoid (x * w + b)
 dw ← (1/m) * x^T * (p - y)
 db ← (1/m) * sum (p - y)
 w ← w - l_rate * dw
 b ← b - l_rate * db
 return w, b

function predict (x, w, b)
return sigmoid (x * w + b) > 0.5

function DecisionFn(x, w, b)
return Sigmoid (x * w + b)

print predict (DecisionFn (x, w, b), w, b)

Code:

```
import numpy as np
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.datasets import load_iris
iris = load_iris()
X = pd.DataFrame(iris.data, columns=iris.feature_names)
y = pd.Series(iris.target)
X = X[y.isin([0, 1])] # Select only classes 0 and 1
y = y[y.isin([0, 1])]
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
model = LogisticRegression()
model.fit(X_train, y_train)
y_pred = model.predict(X_test)
results = pd.DataFrame({'Actual': y_test, 'Predicted': y_pred})
print("Logistic Regression Results:")
print(results.head())
```

output:

Logistic Regression Results:		
	Actual	Predicted
83	1	1
53	1	1
70	1	1
45	0	0
44	0	0

Program 5

Observation:

The image shows handwritten Python code on lined paper. At the top left, there is a large handwritten '9'. To the right of the date and page fields, it says 'Date / /' and 'Page / /'. Below the header, the code is organized into several functions:

- entropy**:
Imports pandas and numpy.
Reads a CSV file named 'weather.csv'.
Creates a DataFrame df from the data.
Defines entropy as the negative sum of probabilities times the log base 2 of probabilities.
- information_gain**:
Takes a DataFrame df, a feature column, and a target column.
Calculates total entropy of the target column.
Finds unique values in the feature column.
Initializes weighted_entropy to 0.
Iterates over each unique value in the feature:
 - Creates a subset of the DataFrame where the feature value matches the current iteration.
 - Calculates the weighted entropy for this subset by multiplying its length by the entropy of the target column for this subset.
Returns the total entropy minus the weighted entropy.
- ID3**:
Takes a DataFrame df and a list of features.
If the target column has only one unique value, returns the index of that value.
Otherwise, initializes gains to an empty list.
Iterates over each feature:
 - Creates a list of information gains for each feature.
 - Chooses the feature with the maximum gain and adds its index to the gains list.
Returns the index of the feature with the maximum gain.

Date _____
Page _____

```

best_feature = max(gains, key=gains.get)
tree = {} best_feature : 2 3 3
for value in df[best_feature].unique():
    subset = df[df[best_feature] == value]
    tree[best_feature][value] = id3(subset, features)
    for f in features if f != best_feature]:
        target)
return tree
features = ['Outlook', 'Temperature', 'Humidity', 'Wind']
target = 'PlayTennis'
decision_tree = id3(df, features, target)
print (decision_tree)

Output:
{'Outlook': {'Sunny': {'Humidity': {'high': 'No',
'low': 'Yes'}}}, 'Overcast': 'Yes', 'Rain': {'Wind': {'weak': 'Yes', 'Strong': 'No'}}}

```

Code:

```

import numpy as np
import pandas as pd
from graphviz import Digraph

def entropy(dataset):
    class_counts = dataset.iloc[:, -1].value_counts()
    prob = class_counts / len(dataset)
    return -np.sum(prob * np.log2(prob))

def information_gain(dataset, feature):
    total_entropy = entropy(dataset)
    feature_values = dataset[feature].value_counts()
    weighted_entropy = 0
    for value, count in feature_values.items():
        subset = dataset[dataset[feature] == value]
        weighted_entropy += (count / len(dataset)) * entropy(subset)
    return total_entropy - weighted_entropy

def best_feature(dataset):
    features = dataset.columns[:-1]
    best_info_gain = -1
    best_feature = None
    for feature in features:
        info_gain = information_gain(dataset, feature)

```

```

if info_gain > best_info_gain:
    best_info_gain = info_gain
    best_feature = feature
return best_feature

def id3(dataset, max_depth=None, depth=0):
    if len(dataset.iloc[:, -1].unique()) == 1:
        return dataset.iloc[0, -1]
    if len(dataset.columns) == 1:
        return dataset.iloc[:, -1].mode()[0]
    if max_depth is not None and depth >= max_depth:
        return dataset.iloc[:, -1].mode()[0]
    best = best_feature(dataset)
    tree = {best: {}}
    for value in dataset[best].unique():
        subset = dataset[dataset[best] == value]
        tree[best][value] = id3(subset.drop(columns=[best]), max_depth=max_depth, depth=depth+1)
    return tree

def create_tree_diagram(tree, dot=None, parent_name="Root", parent_value ""):
    if dot is None:
        dot = Digraph(format="png", engine="dot")

    # Recursively add nodes to the graph
    if isinstance(tree, dict):
        for feature, branches in tree.items():
            feature_name = f"{parent_name}_{feature}"
            dot.node(feature_name, feature)
            dot.edge(parent_name, feature_name, label=parent_value)

            for value, subtree in branches.items():
                value_name = f"{feature_name}_{value}"
                dot.node(value_name, f"{feature}: {value}")
                dot.edge(feature_name, value_name, label=str(value))

                # Recurse for each subtree
                create_tree_diagram(subtree, dot, value_name, str(value))
    else:
        dot.node(parent_name + "_class", f"Class: {tree}")
        dot.edge(parent_name, parent_name + "_class", label="Leaf")

    return dot

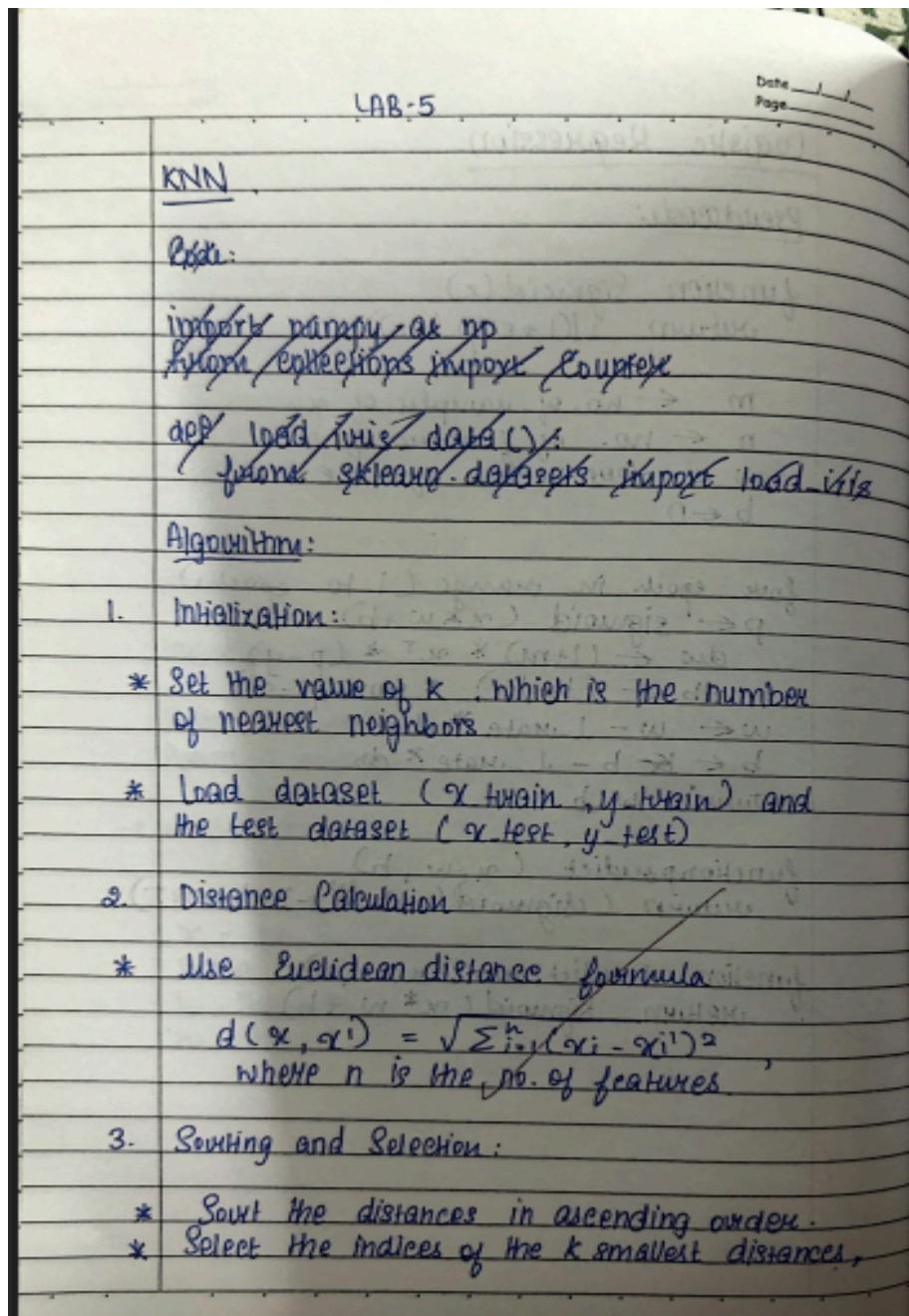
data = {
    'Outlook': ['Sunny', 'Sunny', 'Overcast', 'Rainy', 'Rainy', 'Rainy', 'Overcast', 'Sunny', 'Sunny', 'Rainy', 'Sunny', 'Overcast', 'Overcast', 'Overcast', 'Rainy'],
    'Temperature': ['Hot', 'Hot', 'Hot', 'Mild', 'Cool', 'Cool', 'Cool', 'Mild', 'Cool', 'Mild', 'Mild', 'Mild', 'Mild', 'Hot', 'Mild'],
    'Humidity': ['High', 'High', 'High', 'High', 'Normal', 'Normal', 'Normal', 'High', 'Normal', 'Normal', 'Normal', 'Normal', 'High', 'Normal', 'High'],
}

```

```
'Windy': ['Weak', 'Strong', 'Weak', 'Weak', 'Weak', 'Strong', 'Strong', 'Weak', 'Weak', 'Weak', 'Strong', 'Strong',  
'Weak', 'Strong'],  
    'PlayTennis': ['No', 'No', 'Yes', 'Yes', 'Yes', 'No', 'Yes', 'No', 'Yes', 'Yes', 'Yes', 'Yes', 'Yes', 'Yes', 'Yes', 'No']  
}  
df = pd.DataFrame(data)  
  
tree = id3(df, max_depth=3)  
  
dot = create_tree_diagram(tree)  
dot.render("decision_tree", view=True)
```

Program 6

Observation:



4. Label Determination:

- * Retrieve the label of the k nearest neighbors from y-train.
 - * Determine the most common label among these K neighbors.

5. Prediction:

- * Assign the most common label found in step H as the predicted label for the sample x .

6. Evaluation:

- * Compare the predicted labels from all samples in X-test with their actual labels in y-test.

* Calculate accuracy

accuracy = $\frac{\text{no. of correct prediction}}{\text{total no. of test samples}}$

Date _____
Page _____

Pseudocode:

```

def euclidean_distance(a,b):
    return np.sqrt(np.sum((a-b)**2))

def knn_predict(x_train, y_train, x_test,
                k=3):
    d = []
    for i in range(len(x_train)):
        dist = euclidean_distance(x_test, x_train[i])
        d.append((dist, y_train[i]))
    d.sort(key=lambda x: x[0])

    k_nearest_labels = [label for _, label in
                        distances[:k]]
    return Counter(k_nearest_labels).most_common(1)[0][0]

```

Code:

```

# Importing necessary libraries
import numpy as np
from sklearn import datasets
from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score

```

```

# Load the Iris dataset
iris = datasets.load_iris()
X = iris.data # Features

```

```
y = iris.target # Target labels

# Split the dataset into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

# Create the KNN classifier (k=3 for this example)
knn = KNeighborsClassifier(n_neighbors=3)

# Fit the model
knn.fit(X_train, y_train)

# Predict using the test set
y_pred_knn = knn.predict(X_test)

# Evaluate the model's accuracy
accuracy_knn = accuracy_score(y_test, y_pred_knn)
print(f"KNN Model Accuracy: {accuracy_knn * 100:.2f}%")
```

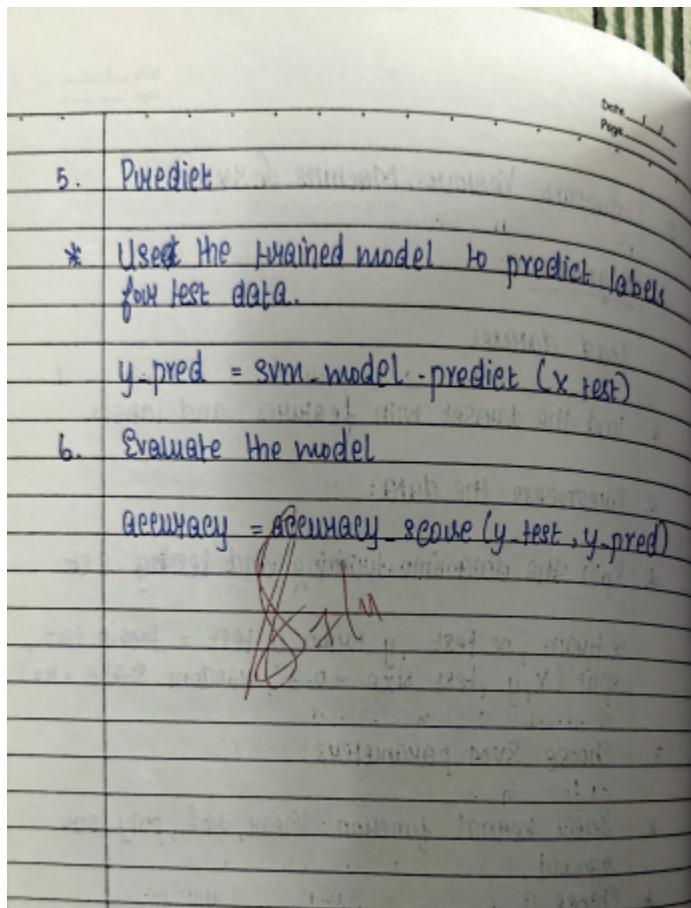
Output:

KNN Model Accuracy: 100.00%

Program 7

Observation:

Date _____	Page _____
2. Support Vector Machine (SVM)	
<u>Algorithm:</u>	
1. Load dataset	
* load the dataset with features and labels.	
2. Preprocess the data:	
* Split the data into training and testing sets	
x_train, x-test, y-train, y-test = train-test-split (X, y, test_size = 0.3, random_state = 42)	
3. Choose SVM parameters:	
* Select kernel function: linear, rbf, poly or sigmoid.	
* Choose C	
* Choose gamma	
svm_model = SVC (kernel = 'rbf', C = 1.0, gamma = 'scale')	
4. Train the SVM model	
* Fit the SVM model using the training data	
svm_model.fit (x-train, y-train)	



Code:

```
import numpy as np
from sklearn import datasets
from sklearn.model_selection import train_test_split
from sklearn.svm import SVC
from sklearn.metrics import accuracy_score

# Load the Iris dataset
iris = datasets.load_iris()
X = iris.data # Features
y = iris.target # Target labels

# Split the dataset into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

# Create the SVM classifier (using a linear kernel for this example)
svm = SVC(kernel='linear')

# Fit the model
svm.fit(X_train, y_train)
```

```
# Predict using the test set  
y_pred_svm = svm.predict(X_test)  
  
# Evaluate the model's accuracy  
accuracy_svm = accuracy_score(y_test, y_pred_svm)  
print(f"SVM Model Accuracy: {accuracy_svm * 100:.2f}%")
```

Output:

SVM Model Accuracy: 100.00%

Program 8

Observation:

LAB - 6	
Date	/ /
Page	/ /
1.	Implement random forest ensemble method on a given dataset.
	<u>Algorithm:</u>
1.	1. load dataset D with features X and target Y.
2.	a. Preprocess the data: b. Handle the missing value c. encode categorical values
	c. Split dataset into training set D_train and test set D_test.
3.	3. Initialize RandomForest model with parameters: a. n_estimators (no of trees) b. max_depth (max depth of trees) c. random state
4.	4. For each tree t = 1 to n_estimators: a. Randomly select a bootstrap sample of data points from D_train. b. For each node in the tree: i. randomly select a subset of features. ii. Gini impurity : $P_{\text{ini}} = 1 - \sum (p_{i,t})^2$ - Entropy : Entropy (S) = $\sum (p_{i,t} * \log_2(p_{i,t}))$
5.	b. For each test instance x-test, aggregate predictions from all trees using majority voting. Prediction = mode (predictions from all trees)

Date	/ /
Page	/ /
6.	6. Calculate accuracy or other metrics such as precision, recall, F1-score.
7.	7. Return model performance.

Code:

```

from sklearn.datasets import load_iris
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score

# Load Iris dataset
iris = load_iris()
X = iris.data
y = iris.target

# Split dataset into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

# Initialize Random Forest Classifier
clf = RandomForestClassifier(n_estimators=100, random_state=42)

# Train the model
clf.fit(X_train, y_train)

# Predict on the test set
y_pred = clf.predict(X_test)

# Calculate accuracy
accuracy = accuracy_score(y_test, y_pred)
print(f'Random Forest Classifier Accuracy: {accuracy:.4f}')

```

output:

Random Forest Classifier Accuracy: 1.0000

Random Forest on a Kaggle Dataset (Visualization)

Code:

```

import pandas as pd
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score

# Load Titanic dataset (Replace URL with the Kaggle Titanic dataset URL or local file)
url = "https://raw.githubusercontent.com/datasciencedojo/datasets/master/titanic.csv"
data = pd.read_csv(url)

# Preprocessing: Select features and target variable
data = data.dropna(subset=['Survived', 'Pclass', 'Sex', 'Age', 'Fare']) # Drop missing values
X = data[['Pclass', 'Sex', 'Age', 'Fare']]
X['Sex'] = X['Sex'].map({'male': 0, 'female': 1}) # Encode 'Sex' column as binary
y = data['Survived']

```

```

# Split the data into training and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

# Initialize Random Forest Classifier
clf = RandomForestClassifier(n_estimators=100, random_state=42)

# Train the model
clf.fit(X_train, y_train)

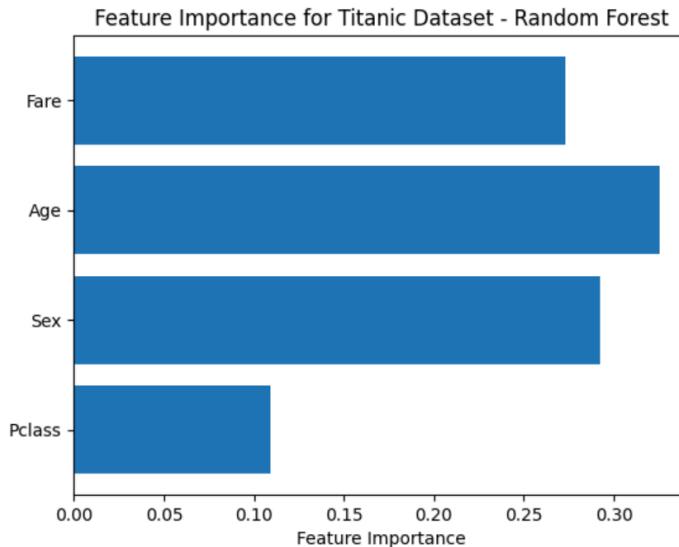
# Make predictions and calculate accuracy
y_pred = clf.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print(f'Random Forest Classifier Accuracy on Titanic Dataset: {accuracy:.4f}')

# Feature importance visualization
importances = clf.feature_importances_
features = X.columns

# Plot feature importance
plt.barh(features, importances)
plt.xlabel('Feature Importance')
plt.title('Feature Importance for Titanic Dataset - Random Forest')
plt.show()

```

Output:



Program 9

Observation:

II	Implement Boosting ensemble method on a given dataset
	<u>Algorithm:</u>
1.	Load dataset with features X and target y
2.	Preprocess the data.
3.	Initialize the weights of all training samples: $w_1 = w_2 = \dots = w_N = 1/N$
4.	For $t = 1$ to n estimators:
a.	Calculate the error of the classifier: $\epsilon_t = (\sum w_i * I(h_t(x_i) \neq y_i)) / \sum w_i$
b.	Compute the classifier weight (α_t): $\alpha_t = 0.5 * \log((1 - \epsilon_t) / \epsilon_t)$
c.	Update the weights: $w_i = w_i * \exp(\alpha_t * I(h_t(x_i) \neq y_i))$
d.	Normalize the weights: $w_i = w_i / \sum w_i$

Date _____	Page _____
5.	Final prediction for test point $x = \text{sign}(\sum \alpha_t * h_t(x))$
6.	Evaluate the model
7.	Return model performance.

Code:

```

from sklearn.datasets import load_iris
from sklearn.ensemble import GradientBoostingClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score

# Load Iris dataset
iris = load_iris()

```

```

X = iris.data
y = iris.target

# Split dataset into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

# Initialize Gradient Boosting Classifier
gb_clf = GradientBoostingClassifier(n_estimators=100, random_state=42)

# Train the model
gb_clf.fit(X_train, y_train)

# Make predictions
y_pred = gb_clf.predict(X_test)

# Calculate accuracy
accuracy = accuracy_score(y_test, y_pred)
print(f'Gradient Boosting Classifier Accuracy on Iris Dataset: {accuracy:.4f}')

```

Output:

Gradient Boosting Classifier Accuracy on Iris Dataset: 1.0000

Boosting Ensemble (XGBoost) on Iris Dataset (Visualization)

```

import matplotlib.pyplot as plt
import xgboost as xgb
from sklearn.datasets import load_iris
from sklearn.model_selection import train_test_split

# Load Iris dataset
iris = load_iris()
X = iris.data
y = iris.target

# Split dataset into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

# Initialize XGBoost Classifier
xgb_clf = xgb.XGBClassifier(n_estimators=100, random_state=42)

# Train the model
xgb_clf.fit(X_train, y_train)

# Make predictions
y_pred = xgb_clf.predict(X_test)

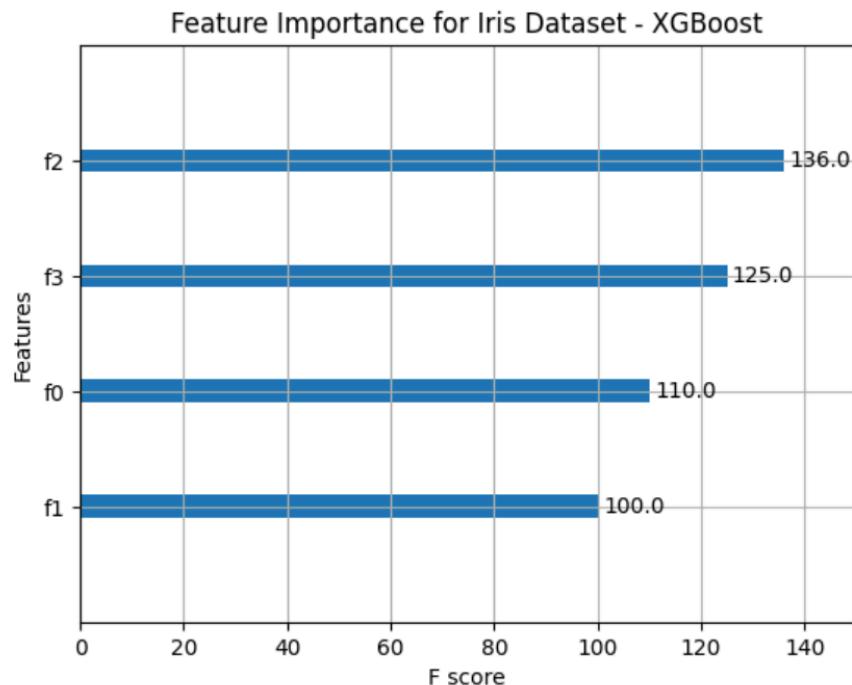
# Calculate accuracy
accuracy = (y_pred == y_test).mean()

```

```
print(f'XGBoost Classifier Accuracy on Iris Dataset: {accuracy:.4f}')
```

```
# Feature importance visualization
xgb.plot_importance(xgb_clf, importance_type='weight', max_num_features=10)
plt.title('Feature Importance for Iris Dataset - XGBoost')
plt.show()
```

Output:



Program 10

Observation:

III	Build K-Means algorithm a for cluster a set of data stored in a .csv file.
1.	load dataset D with data points x_i , where $x_i \in \mathbb{R}^d$
2.	Preprocess the data.
3.	Initialize k centroids randomly: c_1, c_2, \dots, c_k
4.	For each data point x_i :
i.	$\text{distance}(x_i, c_j) = \ x_i - c_j\ _2$
ii.	cluster(x_i) = argmin j distance(x_i, c_j)
5.	If centroids do not change between iterations, stop: convergence.
6.	Output:
i.	Final centroids c_1, c_2, \dots, c_k
ii.	Cluster assignment for each data point.

Code:

```

from sklearn.datasets import load_iris
from sklearn.cluster import KMeans
from sklearn.metrics import accuracy_score
import numpy as np
from sklearn.model_selection import train_test_split

# Load Iris dataset
iris = load_iris()
X = iris.data
y = iris.target

# Perform K-means clustering (3 clusters for 3 species in Iris)
kmeans = KMeans(n_clusters=3, random_state=42)
y_kmeans = kmeans.fit_predict(X)

# Since K-means is an unsupervised algorithm, the cluster labels don't directly match the target labels
# To calculate accuracy, we map cluster labels to actual species labels
# Here we assume that the true labels are known and compare them with the clusters' output

# We use the mode (most frequent) of the true class labels for each cluster to match them
from scipy.stats import mode

```

```

cluster_labels = np.zeros_like(y_kmeans)
for i in range(3): # 3 clusters
    mask = (y_kmeans == i)
    cluster_labels[mask] = mode(y[mask])[0]

# Calculate accuracy by comparing the mapped cluster labels with true labels
accuracy = accuracy_score(y, cluster_labels)
print(f'K-means Clustering Accuracy on Iris Dataset: {accuracy:.4f}')

```

output:

K-means Clustering Accuracy on Iris Dataset: 0.8867

K-means on a Kaggle Dataset (Visualization)

```

import pandas as pd
import matplotlib.pyplot as plt
from sklearn.cluster import KMeans
from sklearn.decomposition import PCA

# Load Titanic dataset (You can replace the URL with your own Kaggle dataset link)
url = "https://raw.githubusercontent.com/datasciencedojo/datasets/master/titanic.csv"
data = pd.read_csv(url)

# Preprocessing: Selecting relevant features and dropping missing values
data = data.dropna(subset=['Pclass', 'Age', 'Fare'])
X = data[['Pclass', 'Age', 'Fare']]

# Normalize the data (optional, but helps with clustering)
X = (X - X.mean()) / X.std()

# Apply K-means clustering (let's assume we use 3 clusters for this example)
kmeans = KMeans(n_clusters=3, random_state=42)
y_kmeans = kmeans.fit_predict(X)

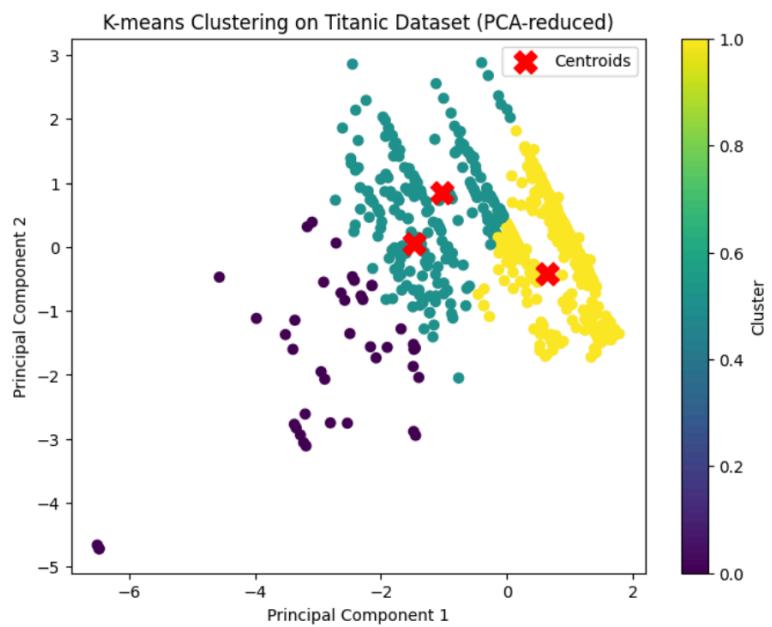
# Apply PCA for 2D visualization (reduce to 2 components)
pca = PCA(n_components=2)
X_pca = pca.fit_transform(X)

# Plot the clusters in 2D
plt.figure(figsize=(8, 6))
plt.scatter(X_pca[:, 0], X_pca[:, 1], c=y_kmeans, cmap='viridis')
plt.scatter(kmeans.cluster_centers_[:, 0], kmeans.cluster_centers_[:, 1], s=200, c='red', marker='X',
label='Centroids')
plt.title('K-means Clustering on Titanic Dataset (PCA-reduced)')
plt.xlabel('Principal Component 1')
plt.ylabel('Principal Component 2')
plt.colorbar(label='Cluster')

```

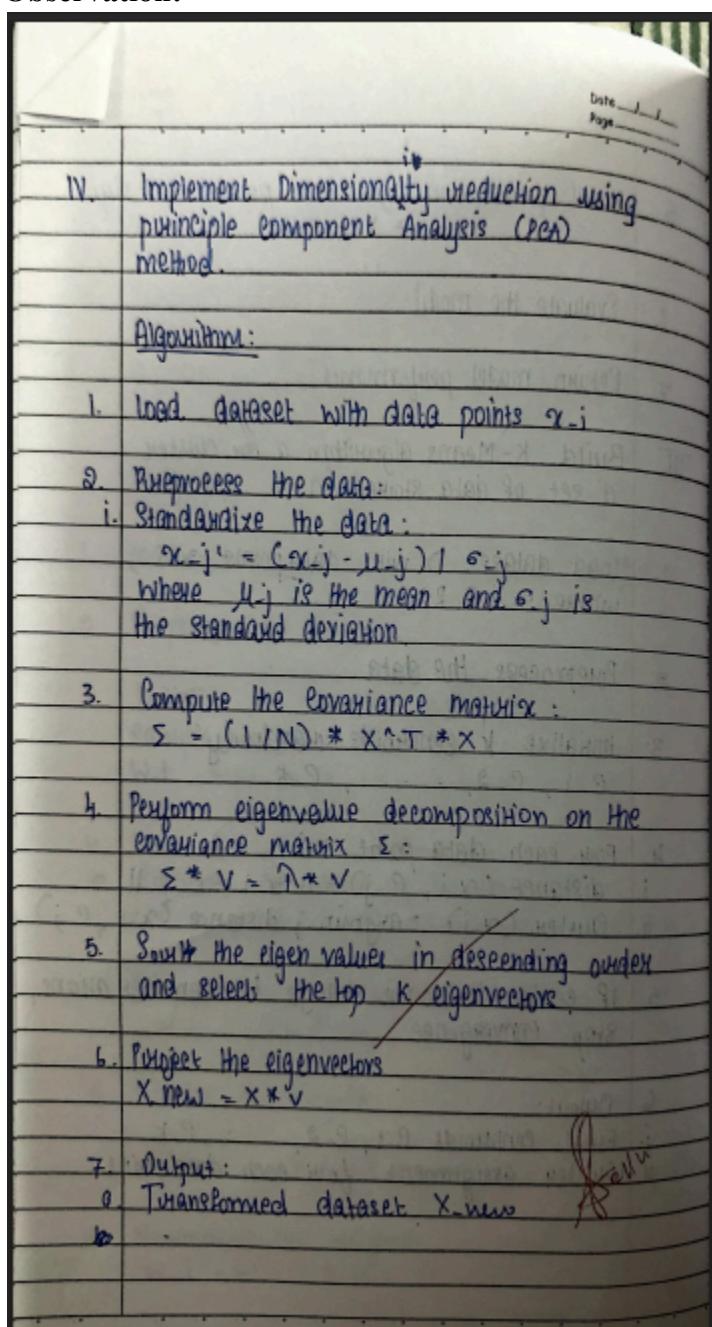
```
plt.legend()  
plt.show()
```

Output:



Program 11

Observation:



Code:

```
from sklearn.datasets import load_iris
from sklearn.decomposition import PCA
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score
```

```
# Load Iris dataset
iris = load_iris()
X = iris.data
```

```

y = iris.target

# Apply PCA to reduce data to 2D
pca = PCA(n_components=2)
X_pca = pca.fit_transform(X)

# Split dataset into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X_pca, y, test_size=0.3, random_state=42)

# Train a Logistic Regression model
clf = LogisticRegression(max_iter=200)
clf.fit(X_train, y_train)

# Make predictions
y_pred = clf.predict(X_test)

# Calculate accuracy
accuracy = accuracy_score(y_test, y_pred)
print(f'Accuracy after PCA on Iris Dataset: {accuracy:.4f}')

```

Output:

Accuracy after PCA on Iris Dataset: 1.0000

PCA on Iris Dataset (Visualization)

```

import matplotlib.pyplot as plt
from sklearn.datasets import load_iris
from sklearn.decomposition import PCA

# Load Iris dataset
iris = load_iris()
X = iris.data
y = iris.target

# Apply PCA to reduce data to 2D
pca = PCA(n_components=2)
X_pca = pca.fit_transform(X)

# Plot the PCA-reduced data
plt.figure(figsize=(8, 6))
plt.scatter(X_pca[:, 0], X_pca[:, 1], c=y, cmap='viridis')
plt.title('PCA on Iris Dataset')
plt.xlabel('Principal Component 1')
plt.ylabel('Principal Component 2')
plt.colorbar(label='Target Class')
plt.show()

```

Output:

