# VISVESVARAYA TECHNOLOGICAL UNIVERSITY

**"JnanaSangama", Belgaum -590014, Karnataka.**

**LAB REPORT**
**on**

# Machine Learning (23CS6PCMAL)

***Submitted by***

**Abhishek Gouda (1BM22CS006)**

*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 **Abhishek Gouda (1BM22CS006),** 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.

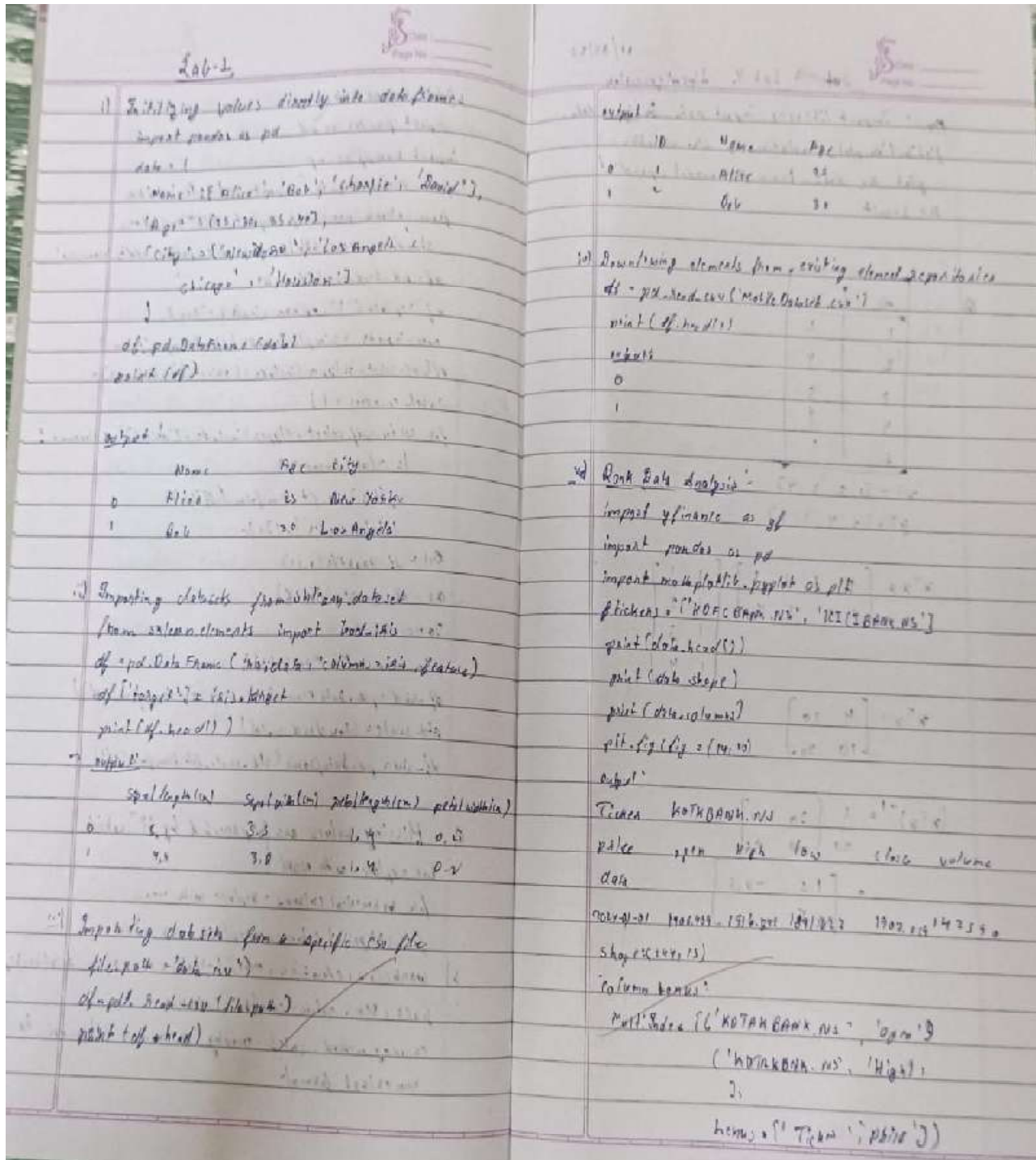| Dr. Seema Patil<br>Assistant Professor<br>Department of CSE, BMSCE | Dr. Kavitha Sooda<br>Professor & HOD<br>Department of CSE, BMSCE |
|---|---|

# Index

Github Link: https://github.com/Abhii2404/6thSem-ML-Lab

**Program 1**

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

Screenshot

Lab-1

1] Initializing values directly into data frame.
import pandas as pd
data = {
  'Name': ['Alice', 'Bob', 'Charlie', 'David'],
  'Age': [25, 30, 35, 40],
  'City': ['New York', 'Los Angeles', 'Chicago', 'Houston']
}
df = pd.DataFrame(data)
print(df)

output:
        Name   Age   City
0      Alice   25    New York
1      Bob    30    Los Angeles

2] Importing datasets from sklearn dataset
from sklearn.datasets import load_iris
df = pd.DataFrame (iris.data, columns = iris.feature)
df['target'] = iris.target
print(df.head())

output:-
   sepal length(cm)  sepal width(cm)  petal length(cm)  petal width(cm)
0       5.1              3.5              1.4              0.2
1       4.9              3.0              1.4              0.2

Importing datasets from a specific csv file
filepath = 'data.csv'
df = pd.read_csv(filepath)
print(df.head())

output:
     ID    Name    Age   City
0     1    Alice    25
1          Bob    30

10] Downloading elements from existing element export to disc
df = pd.read_csv('Mobile Data.csv')
print(df.head())

output:
0
1

4] Bank Data Analysis:
import yfinance as yf
import pandas as pd
import matplotlib.pyplot as plt
tickers = ['HDFCBANK.NS', 'ICICIBANK.NS']
print(data.head())
print(data.shape)
print(data.columns)
plt.fig(fig = (14, 10))

output:
Ticker    KOTAKBANK.NS
Price    Open  High  Low  Close  Volume
Date
2024-01-01  1902.999  1916.201  1841.022  1902.013  14359...
shape: (244, 75)
column banks:
MultiIndex [('KOTAKBANK.NS', 'Open'),
            ('KOTAKBANK.NS', 'High'),
            ],
names = ['Ticker', 'Price'])

5

Code:

```python
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="2024-12-30", group_by='ticker')

print("First 5 rows of the dataset:")

print(data.head())

print("\nShape of the dataset:")

print(data.shape)

print("\nColumn names:")

print(data.columns)

hdfc_data = data['HDFCBANK.NS']

print("\nSummary statistics for HDFC Bank:")

print(hdfc_data.describe())

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

icici_data = data['ICICIBANK.NS']

print("\nSummary statistics for ICICI Bank:")

print(icici_data.describe())

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

kotak_data = data['KOTAKBANK.NS']

print("\nSummary statistics for Kotak Mahindra Bank:")

print(kotak_data.describe())
```

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

plt.figure(figsize=(14, 10))

plt.subplot(3, 2, 1)

hdfc_data['Close'].plot(title="HDFC Bank - Closing Price")

plt.subplot(3, 2, 2)

hdfc_data['Daily Return'].plot(title="HDFC Bank - Daily Returns", color='orange')

plt.subplot(3, 2, 3)

icici_data['Close'].plot(title="ICICI Bank - Closing Price")

plt.subplot(3, 2, 4)

icici_data['Daily Return'].plot(title="ICICI Bank - Daily Returns", color='orange')

plt.subplot(3, 2, 5)

kotak_data['Close'].plot(title="Kotak Mahindra Bank - Closing Price")

plt.subplot(3, 2, 6)

kotak_data['Daily Return'].plot(title="Kotak Mahindra Bank - Daily Returns", color='orange')

plt.tight_layout()

plt.show()


hdfc_data.to_csv('hdfc_bank_data.csv')

icici_data.to_csv('icici_bank_data.csv')

kotak_data.to_csv('kotak_bank_data.csv')


print("\nHDFC Bank data saved to 'hdfc_bank_data.csv'.")

print("ICICI Bank data saved to 'icici_bank_data.csv'.")

print("Kotak Bank data saved to 'kotak_bank_data.csv'.")
```
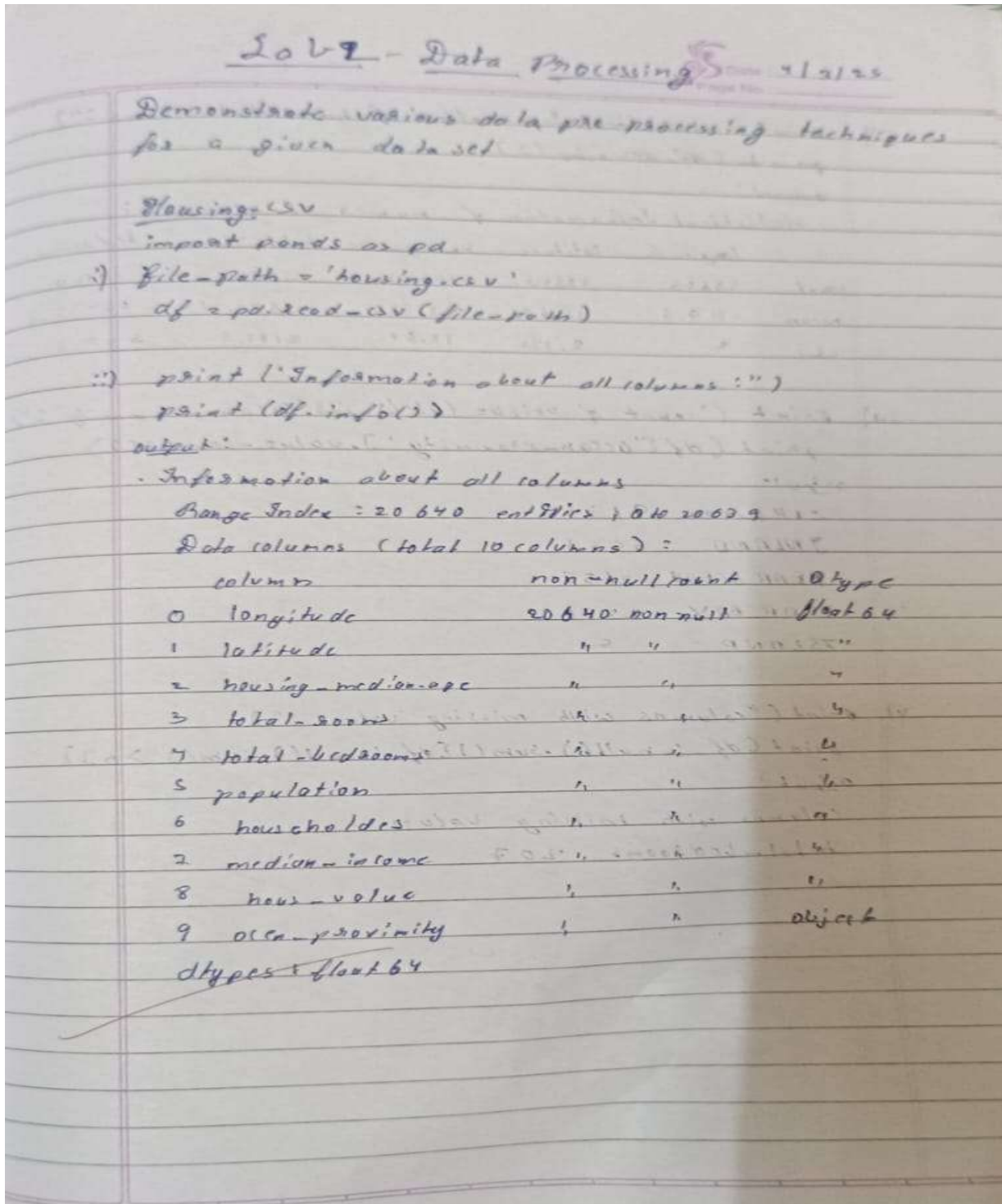
**Program 2**

Demonstrate various data pre-processing techniques for a given dataset

Screenshot



2.02.1 - Data Processing                                    1/2/25

Demonstrate various data pre processing techniques
for a given dataset

Housing.csv
import pands as pd
i) File-path = 'housing.csv'
    df = pd.read-csv (file-path)

ii) print ("Information about all columns :")
    print (df.info())
output:
. Information about all columns
    Range Index : 20 640 entries } 0 to 20639
    Data columns (total 10 columns):
    column                      non-null count        Dtype
    0  longitude                20640 non null        float64
    1  latitude                   "        "            "
    2  housing-median-age         "        "            "
    3  total-rooms                "        "            "
    4  total-bedrooms             "        "            "
    5  population                 "        "            "
    6  householders               "        "            "
    7  median-income              "        "            "
    8  hous-value                 "        "            "
    9  occu-proximity             "        "           object
    dtypes : float64

iii) print ("statistical Information of numerical columns :")
print (df.describe ())

output:
statistical Information of numerical columns ::

| | longitude | latitude | housing_age | total_rooms | bedrooms |
|---|---|---|---|---|---|
| count | 20640 | 20640 | 20640 | 20640 | 20640 |
| mean | -119.5 | 35.631 | 28.639 | 2635.76 | 26493.1 |
| std | e | 2.173 | 12.58 | 2181.6 | 537.8 |

iv) Print ("count of unique labels for 'ocean proximity':")
print (df ['ocean-proximity'].value-counts())

output:
<1HOcean      9136
INLAND        6551
NEAR OCEAN    2658
NEAR BAY      2290
ISLAND        5

v) print ("columns with missing values :")
print (df.isnull().sum ()[df.isnull().sum() >0])

output:
columns with missing values :
total-bedrooms   207

## Diabetes and Adult Income data sets

```python
import pandas as pd
import numpy as np
diabetes_df = pd.read_csv('diabetes.csv')
adult_df = pd.read_csv('adult.csv')
```

1) Missing values

```python
print(diabetes_df.isnull().sum())
print(adult_df.isnull().sum())
```

output

Missing values in diabetes dataset:

ID        0
No_pation  0

Missing values in Adult income dataset:

age         0
workclass   0

```python
from sklearn.processing import MinMaxScaler, standard
import matplotlib.pyplot as plt
file_path = 'diabetes.csv'
df = pd.read_csv(file_path)
df_num = df.select_dtypes(include=['number']).copy()
input = SimpleImputer(strategy='mean')
df_num.iloc[:, :] = input.fit_transform(df_numeric)
df[df_num.columns] = df.numeric
Q1 = df_num.quantile(0.25)
Q3 = df_num.quantile(0.75)
IQR = Q3 - Q1
min_max_scalar = MinMaxScaler(df)
```

Left column:

```
... std_scaler = StandardScaler()
df_standard = pd.DataFrame (std_scaler.fit_transform (...))
print ("The Raw data Dataset (min-max scaled):")
print (df_min_max_scaled())
print ("The Standard Dataset (standard scaler):")
print (df_standard_head())
```

→ Which categorical columns did you scale by?
   None of the categorical ones.

→ What categorical columns did you identifying encoding?

→ What is the diff b/w minmax scaling and standardization?
   when would you use one over the other?

→ plain these scaling transform able to find range [0,1]

$$x' = \frac{x - x_{min}}{x_{max} - x_{min}}$$

used when :- data does not follow a normal dist.
          - features have diff ranges and need to be bound

Standardization transforms data to have zero means
and unit variance.

$$x' = \frac{x - \mu}{\sigma}$$

used when : data set follows a Gaussian dist.
          - many ML algo assume normality

Right column:

```
Adult_income.csv
import pandas as pd
import numpy as np
from sklearn.impute import SimpleImputer
from sklearn.preprocessing import LabelEncoder,
   StandardScaler, MinMaxScaler   file_path = "adult_income"
df = pd.read_csv (file_path)
df.replace (' ?', np.nan, inplace = True)
num_imputer = SimpleImputer (strategy = 'mean')
df [df.select_dtypes (include = ['object']).columns] =
label_encoders = {}
for col in df.select_dtypes (include = ['object']).columns:
   le = LabelEncoder()
   df [col] = le.fit_transform (df [col])
   label_encoders [col] = le
Q1 = df.quantile (0.25)
Q3 = df.quantile (0.75)
IQR = Q3 - Q1
min_max = MinMaxScaler()
df_std = pd.DataFrame (std_scale.fit_transform (...))
std_scale = StandardScaler()
df_std = pd.DataFrame (std_scale.fit_transform (df))
```

a) Missing values are represented by " ?" which
   we replace with nan
   for numerical columns → replace with mean
   for categorical columns → replace with mode

b) workclass, education, marital-status, occupation, relationship
   race, sex, native-country, income
   through method label encoding be converted categorical to
   numerical format

Code:

```python
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import MinMaxScaler, StandardScaler
from sklearn.impute import SimpleImputer
from sklearn.preprocessing import LabelEncoder
import seaborn as sns
import matplotlib.pyplot as plt

diabetes_data = pd.read_csv('/content/Dataset of Diabetes .csv')
adult_income_data = pd.read_csv('/content/adult.csv')

print("Diabetes Dataset:")
print(diabetes_data.head())

print("\nAdult Income Dataset:")
print(adult_income_data.head())

diabetes_numerical_cols = diabetes_data.select_dtypes(include=[np.number]).columns
diabetes_categorical_cols = diabetes_data.select_dtypes(include=[object]).columns

diabetes_imputer_num = SimpleImputer(strategy='median')
diabetes_data[diabetes_numerical_cols] =
diabetes_imputer_num.fit_transform(diabetes_data[diabetes_numerical_cols])

diabetes_imputer_cat = SimpleImputer(strategy='most_frequent')
diabetes_data[diabetes_categorical_cols] =
diabetes_imputer_cat.fit_transform(diabetes_data[diabetes_categorical_cols])

adult_income_numerical_cols = adult_income_data.select_dtypes(include=[np.number]).columns
adult_income_categorical_cols = adult_income_data.select_dtypes(include=[object]).columns

adult_income_imputer_num = SimpleImputer(strategy='median')
adult_income_data[adult_income_numerical_cols] =
adult_income_imputer_num.fit_transform(adult_income_data[adult_income_numerical_cols])

adult_income_imputer_cat = SimpleImputer(strategy='most_frequent')
adult_income_data[adult_income_categorical_cols] =
adult_income_imputer_cat.fit_transform(adult_income_data[adult_income_categorical_cols])

categorical_columns_adult = adult_income_data.select_dtypes(include=['object']).columns
label_encoder = LabelEncoder()

for col in categorical_columns_adult:
    adult_income_data[col] = label_encoder.fit_transform(adult_income_data[col])
```

```python
def detect_and_remove_outliers(df):
    numerical_df = df.select_dtypes(include=[np.number])
    Q1 = numerical_df.quantile(0.25)
    Q3 = numerical_df.quantile(0.75)
    IQR = Q3 - Q1
    return df[~((numerical_df < (Q1 - 1.5 * IQR)) | (numerical_df > (Q3 + 1.5 * IQR))).any(axis=1)]

diabetes_data_cleaned = detect_and_remove_outliers(diabetes_data)
adult_income_data_cleaned = detect_and_remove_outliers(adult_income_data)

min_max_scaler = MinMaxScaler()

diabetes_numerical_cols = diabetes_data_cleaned.select_dtypes(include=[np.number]).columns
diabetes_data_normalized = diabetes_data_cleaned.copy()

diabetes_data_normalized[diabetes_numerical_cols] =
min_max_scaler.fit_transform(diabetes_data_cleaned[diabetes_numerical_cols])

adult_income_numerical_cols =
adult_income_data_cleaned.select_dtypes(include=[np.number]).columns
adult_income_data_normalized = adult_income_data_cleaned.copy()

adult_income_data_normalized[adult_income_numerical_cols] =
min_max_scaler.fit_transform(adult_income_data_cleaned[adult_income_numerical_cols])

standard_scaler = StandardScaler()

diabetes_data_standardized = diabetes_data_cleaned.copy()
diabetes_data_standardized[diabetes_numerical_cols] =
standard_scaler.fit_transform(diabetes_data_cleaned[diabetes_numerical_cols])

adult_income_data_standardized = adult_income_data_cleaned.copy()
adult_income_data_standardized[adult_income_numerical_cols] =
standard_scaler.fit_transform(adult_income_data_cleaned[adult_income_numerical_cols])
```

## Program 3

Use an appropriate data set for building the decision tree (ID3) and apply this knowledge to classify a new sample.

Screenshot:

Lab-5.

| Instance | $A_2$ | $A_3$ | Classification |
|---|---|---|---|
| 1 | Hot | High | No |
| 2 | Hot | High | No |
| 6 | cool | High | No |
| 7 | Hot | High | No |
| 8 | Hot | Normal | Yes |

$$Entropy = -\frac{4}{5} \log\left(\frac{4}{5}\right) - \frac{1}{5} \log\left(\frac{1}{5}\right)$$

$$= 0.7219$$

For $a_2$,

$$S_{hot}\ [1+, 3-] = -\frac{1}{4} \log\left(\frac{1}{4}\right) - \frac{3}{4} \log\left(\frac{3}{4}\right)$$

$$= 0.8113$$

$S_{cool} = [0+, 1-] = 0$

$$Gain\ (s, a_2) = 0.7219 - \frac{4}{5} \times 0.8113 - 0 = 0.07286$$

For $a_3$

$S_{high}\ [0+, 4-] = 0$

$S_{normal}\ [0-, 1+] = 0$

$Gain\ (s, a_3) = 0.7219$

∴ $a_3$ has highest gain value, it is taken as root

```
dtc = Decision Tree Classifier()
dtc.fit(X_train, y_train)
y_pred = dtc.predict(X_test)

X_petal = petal.iloc[:,:-1]
y_petal = petal.iloc[:,-1]
X_train, y_train, y_test, X_train =
train_test_split(X_petal, Y_petal, test_size=0.2,
                 random_state=42)

dtr = Decision Tree Regressor()
dtc.fit(X_train, y_train)
y_pred = dtr.predict(X_test)

print(mean_absolute_error(y_test, y_pred))
```

Output)
Decision tree classification for Isc
Accuracy : 1.0
confusion Matrix :
[[10  0  0]
 [ 0  0  0]
 [ 0  0  10]]
classification report

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| setosa | 1.00 | 1.00 | 1.00 | 10 |
| versicolor | 1.00 | 1.00 | 1 | 9 |
| virginica | 1.00 | 1.00 | 1 | 11 |
| accuracy |  |  | 1 | 30 |
| macro avg |  |  |  | 30 |
| weighted avg |  |  |  | 30 |

Decision tree for drug :-
Accuracy : 1.0
Confusion matrix :

$$\begin{bmatrix} [6 & 0 & 0 & 0 & 0] \\ [0 & 3 & 0 & 0 & 0] \\ [0 & 0 & 5 & 0 & 0] \\ [0 & 0 & 0 & 11 & 0] \\ [0 & 0 & 0 & 0 & 15] \end{bmatrix}$$

Classification report :

|        | precision | recall | f1 score | support |
|--------|-----------|--------|----------|---------|
| Drug A |           |        |          | 6       |
| B      |           |        |          | 3       |
| C      |           |        |          | 5       |
| X      |           |        |          | 11      |
| Y      |           |        |          | 15      |
|        |           |        |          |         |
| Accuracy |         |        |          | 40      |
| macro avg |        |        |          | 40      |
| weighted avg |     |        |          | 40      |

Decision Tree for Petrol consumption
Mean Absolute error : 98.9
Mean square Error : 17957.3
Root Mean square Squared Error : 133.6312089

Code:

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier, DecisionTreeRegressor, plot_tree
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report,
mean_absolute_error, mean_squared_error
from sklearn.preprocessing import LabelEncoder

iris = pd.read_csv("/content/iris (4).csv")
drug = pd.read_csv("/content/drug.csv")
petrol = pd.read_csv("/content/petrol_consumption.csv")

X_iris = iris.iloc[:, :-1]
y_iris = iris.iloc[:, -1]
X_train, X_test, y_train, y_test = train_test_split(X_iris, y_iris, test_size=0.2, random_state=42)

dtc = DecisionTreeClassifier()
dtc.fit(X_train, y_train)
y_pred = dtc.predict(X_test)

print("Decision Tree Classification for IRIS Dataset:")
print("Accuracy:", accuracy_score(y_test, y_pred))
print("Confusion Matrix:\n", confusion_matrix(y_test, y_pred))
print("Classification Report:\n", classification_report(y_test, y_pred))

X_drug = drug.iloc[:, :-1]
y_drug = drug.iloc[:, -1]

le = LabelEncoder()

for col in X_drug.select_dtypes(include=['object']).columns:
    X_drug[col] = le.fit_transform(X_drug[col])

X_train, X_test, y_train, y_test = train_test_split(X_drug, y_drug, test_size=0.2, random_state=42)

dtc = DecisionTreeClassifier()
dtc.fit(X_train, y_train)
y_pred = dtc.predict(X_test)

print("\nDecision Tree Classification for Drug Dataset:")
print("Accuracy:", accuracy_score(y_test, y_pred))
print("Confusion Matrix:\n", confusion_matrix(y_test, y_pred))
print("Classification Report:\n", classification_report(y_test, y_pred))

X_petrol = petrol.iloc[:, :-1]
```

```python
y_petrol = petrol.iloc[:, -1]
X_train, X_test, y_train, y_test = train_test_split(X_petrol, y_petrol, test_size=0.2, random_state=42)

dtr = DecisionTreeRegressor()
dtr.fit(X_train, y_train)
y_pred = dtr.predict(X_test)

print("\nDecision Tree Regression for Petrol Consumption:")
print("Mean Absolute Error:", mean_absolute_error(y_test, y_pred))
print("Mean Squared Error:", mean_squared_error(y_test, y_pred))
print("Root Mean Squared Error:", np.sqrt(mean_squared_error(y_test, y_pred)))
```

## Program 4

Implement Linear and Multi-Linear Regression algorithm using appropriate dataset

Screenshot

Code:

Linear Regression

Code:
```
import pandas as pd
import numpy as np
from sklearn import linear_model
import matplotlib.pyplot as plt

df = pd.read_csv('having-per-price.csv')
plt.xlabel('area')
plt.ylabel('price')
plt.scatter(df.area, df.price, color='red', marker='+')

new_df = df.drop('price', axis='columns')
price = df.price

reg = linear_model.LinearRegression()
reg.fit(new_df, price)
reg.predict([[3300]])
reg.coef_
reg.intercept_

y = m*x + b

predict_val = reg.predict([[3000]])
```

Output:

1/ predicted per capita income for Canada in 2020 : 41,027.69...

predicted salary for 1st year of experience
=> 140337.54253139565.

2/ Mean Absolute Error (Canada) = 3240.91...

Mean Absolute Error (Salary) = 4519.18...

## Multiple

```
import pandas as pd
import numpy as np
from sklearn import linear_model
df = pd.read_csv('housing price Multiple.csv')
df
df.drop(columns)
df.columns = df[bedrooms, fill na (df.bedrooms.
median())]
df

reg = multiple_model LinearRegression()
reg.fit(df.drop('Price', axis = 'columns'),
    df.price)

reg.coef_
reg.intercept

reg.predict([[3000, 3, 40]])
```

1] predicted salary for candidate 1 (2 yr experience, test score, 6 interview score)

   6.54560.509006661

2] predicted salary for candidate 1 (12 yr experience, 10 test score, 10 interview score) = 90563.49999

3] Mean-Absolute Error for salary prediction 1667
   Predicted profit for the given inputs
   (Florida State) 155449.00.6929852
   Mean absolute for profit prediction:
   1278.869449333337

Code:

```python
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.preprocessing import LabelEncoder
from sklearn.metrics import mean_absolute_error
import matplotlib.pyplot as plt

hiring_data = pd.read_csv('hiring.csv')
print(hiring_data.head())
hiring_data = hiring_data.dropna()

experience_mapping = {
    'one': 1, 'two': 2, 'three': 3, 'four': 4, 'five': 5, 'six': 6, 'seven': 7, 'eight': 8,
    'nine': 9, 'ten': 10, 'eleven': 11, 'twelve': 12, 'thirteen': 13, 'fourteen': 14,
}

hiring_data['experience'] = hiring_data['experience'].replace(experience_mapping)
hiring_data['experience'] = pd.to_numeric(hiring_data['experience'], errors='coerce')

if hiring_data['experience'].isnull().any():
    print("Warning: There are still non-numeric values in the 'experience' column.")
    hiring_data = hiring_data.dropna(subset=['experience'])

X_hiring = hiring_data[['experience', 'test_score(out of 10)', 'interview_score(out of 10)']]
y_hiring = hiring_data['salary($)']

X_train_hiring, X_test_hiring, y_train_hiring, y_test_hiring = train_test_split(X_hiring, y_hiring, test_size=0.2, random_state=42)

regressor_hiring = LinearRegression()
regressor_hiring.fit(X_train_hiring, y_train_hiring)

candidate_1 = np.array([[2, 9, 6]])
candidate_2 = np.array([[12, 10, 10]])

salary_1 = regressor_hiring.predict(candidate_1)
salary_2 = regressor_hiring.predict(candidate_2)

print(f"Predicted salary for candidate 1 (2 yr experience, 9 test score, 6 interview score): {salary_1[0]}")
print(f"Predicted salary for candidate 2 (12 yr experience, 10 test score, 10 interview score): {salary_2[0]}")
```

```python
companies_data = pd.read_csv('/content/1000_Companies.csv')
print(companies_data.head())
companies_data = companies_data.dropna()

label_encoder = LabelEncoder()
companies_data['State'] = label_encoder.fit_transform(companies_data['State'])

X_companies = companies_data[['R&D Spend', 'Administration', 'Marketing Spend', 'State']]
y_companies = companies_data['Profit']

X_train_companies, X_test_companies, y_train_companies, y_test_companies = train_test_split(X_companies, y_companies, test_size=0.2, random_state=42)

regressor_companies = LinearRegression()
regressor_companies.fit(X_train_companies, y_train_companies)

input_data = np.array([[91694.48, 515841.3, 11931.24, label_encoder.transform(['Florida'])[0]]])
predicted_profit = regressor_companies.predict(input_data)

print(f"Predicted profit for the given inputs (Florida State): {predicted_profit[0]}")

y_pred_hiring = regressor_hiring.predict(X_test_hiring)
mae_hiring = mean_absolute_error(y_test_hiring, y_pred_hiring)
print(f"Mean Absolute Error for Salary Prediction: {mae_hiring}")

y_pred_companies = regressor_companies.predict(X_test_companies)
mae_companies = mean_absolute_error(y_test_companies, y_pred_companies)
print(f"Mean Absolute Error for Profit Prediction: {mae_companies}")
```

## Program 5

Build Logistic Regression Model for a given dataset

Screenshot

Code:

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, confusion_matrix

file_path = 'HR_comma_sep.csv'
data = pd.read_csv(file_path)

print(data.info())

print(data.head())

print(data.describe())

plt.figure(figsize=(8, 5))
sns.countplot(x='salary', hue='left', data=data)
plt.title('Impact of Salary on Employee Retention')
plt.xlabel('Salary')
plt.ylabel('Count')
plt.legend(title='Employee Retention', labels=['Stayed', 'Left'])
plt.show()

plt.figure(figsize=(10, 6))
sns.countplot(x='Department', hue='left', data=data)
plt.title('Impact of Department on Employee Retention')
plt.xlabel('Department')
plt.ylabel('Count')
plt.legend(title='Employee Retention', labels=['Stayed', 'Left'])
plt.xticks(rotation=45)
plt.show()

data_encoded = pd.get_dummies(data, columns=['salary', 'Department'], drop_first=True)

print(data_encoded.info())

X = data_encoded.drop('left', axis=1)
y = data_encoded['left']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)
```

```
logreg = LogisticRegression(max_iter=1000)

logreg.fit(X_train_scaled, y_train)

y_pred = logreg.predict(X_test_scaled)

accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy of the Logistic Regression Model: {accuracy * 100:.2f}%")

cm = confusion_matrix(y_test, y_pred)
plt.figure(figsize=(6, 5))
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', cbar=False, xticklabels=['Stayed', 'Left'],
yticklabels=['Stayed', 'Left'])
plt.title('Confusion Matrix')
plt.xlabel('Predicted')
plt.ylabel('Actual')
plt.show()
```

## Program 6

Build KNN Classification model for a given dataset.

Screenshot



Lab-6  K-NN   Date: 1/4/2

| Person | Age | Salary | Target |
|--------|-----|--------|--------|
| A | 18 | 50 | N |
| B | 23 | 53 | N |
| C | 24 | 70 | N |
| D | 41 | 60 | Y |
| E | 43 | 70 | Y |
| F | 38 | 40 | Y |
| X | 35 | 100 | ? |

$(x_2, x_1) = (35, 100)$

for $(18, 50)$  $d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$

$= \sqrt{(35-18)^2 + (100-50)^2}$

$= 52.81$

for $(23, 53)$  $= \sqrt{(35-23)^2 + (100-53)^2}$

$= 46.57$

for $(24, 70)$  $= \sqrt{(35-24)^2 + (100-70)^2}$

$= 31.95$

for $(41, 60)$  $= \sqrt{(35-41)^2 + (100-60)^2}$

$= 40.44$

for $(43, 70)$  $= \sqrt{(35-43)^2 + (100-70)^2}$

$= 31.04$

for $(38, 40)$  $= \sqrt{(35-38)^2 + (100-40)^2}$

$= 60.07$

| Person | Age | salary | Donget | distance | Rank |
|--------|-----|--------|--------|----------|------|
| A | 18 | 30 | N | 52.81 | 5 |
| B | 23 | 55 | N | 76.51 | 4 |
| C | 24 | 70 | N | 31.93 | 2 |
| D | 41 | 60 | Y | 40.44 | 3 |
| E | 43 | 70 | X | 31.04 | 1 |
| F | 38 | 40 | Y | 60.07 | 6 |
| X | 35 | 100 | Y | | |

K=3, 2 → Y   1 → N

K=1          → Y   (35, 100) → (x, y)

k=2          → N

k=3          → (x, y)   → (35, 100)

So X(35, 100) is Y.

Sol-1  KNN

1) Import pandas as pd
Import numpy as np
from sklearn.model_selection import test_train_split
import sklearn.

head_off = pd.read_csv("heart.csv")
print(head_off.head())
x = head_off.drop(columns=['target'])
y = head_off['target']
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2, random_state=42)
scaler = StandardScaler()
x_train_scaled = scaler.fit_transform(x_train)
k_values = range(1, 21)
best_k = 1
best_score = 0
for k in k_values:
    knn = KNeighborsClassifier(n_neighbors=k)
    knn.fit(x_train_scaled, y_train)
    y_pred = knn.predict(x_test_scaled)
    if score > best_score:
        best_score = score
        best_k = k
accuracy = accuracy_score(y_test, y_pred_best)
print("accuracy on test data is ", accuracy, "90-96")

Output:
Accuracy on test data : 91.80%

Output:
• accuracy rate of this test data 100.00%
• accuracy on diabetes test data 69.70%

1) By plotting accuracy and error rate various k
values select the k with highest accuracy and
lowest error rate. This value is considered optimal
for dataset.

2) Feature scaling ensures that all features contribute
equally to the model especially for distance based
algorithms like KNN. It prevents features with large
values from dominating the training process.
How to perform it:
  Standard scaling : Transform feature to have
  mean = 0 and standard deviation = 1
  min-max scaling : scales features to a
  fixed range usually [0,1].

  These scaling techniques help improve
model performance and convergence speed.

Code:

```python
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
import matplotlib.pyplot as plt
import seaborn as sns

iris_df = pd.read_csv('/content/iris (3).csv')

print(iris_df.head())

X_iris = iris_df.drop(columns=['species'])
y_iris = iris_df['species']

X_train_iris, X_test_iris, y_train_iris, y_test_iris = train_test_split(X_iris, y_iris, test_size=0.2,
random_state=42)

scaler = StandardScaler()
X_train_iris = scaler.fit_transform(X_train_iris)
X_test_iris = scaler.transform(X_test_iris)

knn_iris = KNeighborsClassifier(n_neighbors=3)

knn_iris.fit(X_train_iris, y_train_iris)

y_pred_iris = knn_iris.predict(X_test_iris)

accuracy_iris = accuracy_score(y_test_iris, y_pred_iris)
print(f"Accuracy on Iris test data: {accuracy_iris * 100:.2f}%")

cm_iris = confusion_matrix(y_test_iris, y_pred_iris)
sns.heatmap(cm_iris, annot=True, fmt="d", cmap="Blues", xticklabels=knn_iris.classes_,
yticklabels=knn_iris.classes_)
plt.title("Confusion Matrix for Iris Dataset")
plt.xlabel('Predicted')
plt.ylabel('True')
plt.show()

print("Classification Report for Iris Dataset:")
print(classification_report(y_test_iris, y_pred_iris))

diabetes_df = pd.read_csv('diabetes.csv')
print(diabetes_df.head())
```
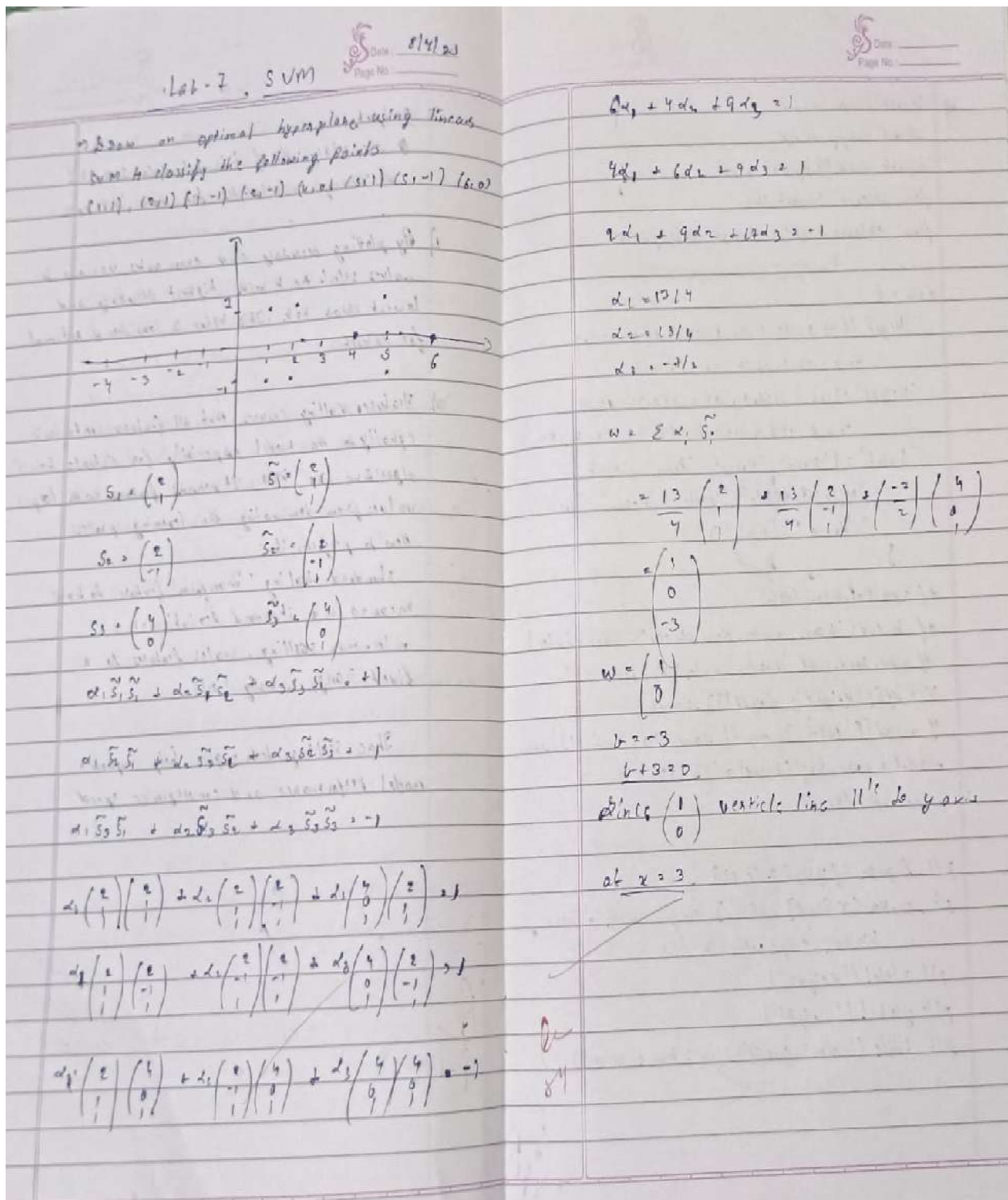
```python
X_diabetes = diabetes_df.drop(columns=['Outcome'])
y_diabetes = diabetes_df['Outcome']

X_train_diabetes, X_test_diabetes, y_train_diabetes, y_test_diabetes = train_test_split(X_diabetes, y_diabetes, test_size=0.2, random_state=42)

scaler = StandardScaler()
X_train_diabetes = scaler.fit_transform(X_train_diabetes)
X_test_diabetes = scaler.transform(X_test_diabetes)

knn_diabetes = KNeighborsClassifier(n_neighbors=5)

knn_diabetes.fit(X_train_diabetes, y_train_diabetes)

y_pred_diabetes = knn_diabetes.predict(X_test_diabetes)

accuracy_diabetes = accuracy_score(y_test_diabetes, y_pred_diabetes)
print(f"Accuracy on Diabetes test data: {accuracy_diabetes * 100:.2f}%")

cm_diabetes = confusion_matrix(y_test_diabetes, y_pred_diabetes)
sns.heatmap(cm_diabetes, annot=True, fmt="d", cmap="Blues", xticklabels=knn_diabetes.classes_,
yticklabels=knn_diabetes.classes_)
plt.title("Confusion Matrix for Diabetes Dataset")
plt.xlabel('Predicted')
plt.ylabel('True')
plt.show()

print("Classification Report for Diabetes Dataset:")
print(classification_report(y_test_diabetes, y_pred_diabetes))
```

## Program 7

Build Support vector machine model for a given dataset

Screenshot

Code:

```python
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.svm import SVC
from sklearn.metrics import accuracy_score, confusion_matrix, roc_auc_score, roc_curve
from sklearn.preprocessing import LabelEncoder, label_binarize
import matplotlib.pyplot as plt
import numpy as np
import seaborn as sns

df = pd.read_csv("/content/letter-recognition.csv")

top_classes = df['letter'].value_counts().head(5).index.tolist()
df = df[df['letter'].isin(top_classes)]

X = df.iloc[:, 1:]
y = df.iloc[:, 0]

label_encoder = LabelEncoder()
y_encoded = label_encoder.fit_transform(y)

y_bin = label_binarize(y_encoded, classes=np.unique(y_encoded))
n_classes = y_bin.shape[1]

X_train, X_test, y_train, y_test_bin = train_test_split(X, y_bin, test_size=0.2, random_state=42)

svm_model = SVC(kernel='linear', probability=True)
svm_model.fit(X_train, y_train.argmax(axis=1))
y_score = svm_model.predict_proba(X_test)

y_pred = svm_model.predict(X_test)
y_true = y_test_bin.argmax(axis=1)

print("Accuracy:", accuracy_score(y_true, y_pred))
print("Confusion Matrix:\n", confusion_matrix(y_true, y_pred))

plt.figure()
for i in range(n_classes):
    fpr, tpr, _ = roc_curve(y_test_bin[:, i], y_score[:, i])
    auc = roc_auc_score(y_test_bin[:, i], y_score[:, i])
    plt.plot(fpr, tpr, label=f"{label_encoder.inverse_transform([i])[0]} AUC={auc:.2f}")
plt.plot([0, 1], [0, 1], 'k--')
plt.title("ROC Curve (Top 5 Classes)")
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.legend(loc="lower right")
plt.tight_layout()
```

```
plt.show()

macro_auc = roc_auc_score(y_test_bin, y_score, average="macro")
print("Macro AUC Score:", macro_auc)
```

## Program 8

Implement Random forest ensemble method on a given dataset.

Screenshot



Lab-8

Date: 13/04/25

Q] Get the difference in decision tree and random forest classifier.

→

| Decision tree | Random forest |
|---|---|
| - A decision tree is a single tree structure where decisions are made by splitting the dataset at each node | • It is an ensemble learning method that builds multiple decision trees and combines their results |
| • These are highly prone to overfitting, especially with complex datasets | • These reduce overfitting by averaging multiple decision trees, which generally improves performance |
| - A model has high variance and low bias. | - A models tend to have low variance and moderate bias as they combine the predictions from many trees |

Q. Discuss all the parameters of random forest classifier
- n-estimators: The no of trees in the forest
- criterion: The functions to measure the quality of a split
- max-depth: The max depth of the tree
- min-sample-split: The min no of sample required to split an internal node
- min-samples-leaf: The min no of samples required to be featured in leaf node
- max-features: The no of features to consider when looking for the best split
- bootstrap: whether to use bootstrapping when creating trees
- oob-score: whether to use out-of-bag samples to estimate the generalization accuracy.
- n-jobs: The no of jobs to run in parallel for both fit() and predict()
- random-state: Controls the randomness of the estimator

35

3) Algorithm of Random Forest

→ Step 1: The data set is divided into input features (x) and the target labels (y)

Step 2: For each tree, create a bootstrap sample from the dataset

Step 3: For each tree, randomly select a subset of features. It is determined by max-features parameter

Step 4: Build a decision tree on the bootstrapped dataset using the selected subset of features

Step 5: Once all trees are built, make predictions by aggregating the predictions of all the trees

Step 6: During training, the data points that were not models generalization performance

Step 7: The final model is evaluated using metrics like accuracy, confusion matrix and AUC score, depending on the task at hand.

Code:

```python
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score, confusion_matrix
from sklearn import preprocessing

df = pd.read_csv('/content/train.csv')

X = df.iloc[:, :-1]
y = df.iloc[:, -1]

for column in X.columns:
    if X[column].dtype == 'object':
      le = preprocessing.LabelEncoder()
      X[column] = le.fit_transform(X[column])

if y.dtype == 'object':
  le = preprocessing.LabelEncoder()
  y = le.fit_transform(y)

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

rf_classifier = RandomForestClassifier(random_state=42)
rf_classifier.fit(X_train, y_train)

y_pred = rf_classifier.predict(X_test)

accuracy = accuracy_score(y_test, y_pred)
conf_matrix = confusion_matrix(y_test, y_pred)

print(f"Accuracy: {accuracy}")
print(f"Confusion Matrix:\n{conf_matrix}")
```

**Program 9**

Implement Boosting ensemble method on a given dataset.

Screenshot

Code:

```python
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.datasets import load_iris
from sklearn.ensemble import AdaBoostClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score

iris = load_iris()
X = iris.data
y = iris.target

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

results = []

n_estimators_list = [10, 50, 100]
learning_rates = [0.01, 0.1, 1]

for n in n_estimators_list:
    for lr in learning_rates:
        tree_base = DecisionTreeClassifier(max_depth=1)
        model = AdaBoostClassifier(estimator=tree_base, n_estimators=n, learning_rate=lr,
random_state=42)
        model.fit(X_train, y_train)
        y_pred = model.predict(X_test)
        acc = accuracy_score(y_test, y_pred)
        results.append({
            'Base': 'DecisionTree',
            'n_estimators': n,
            'learning_rate': lr,
            'Accuracy': acc
        })

for n in n_estimators_list:
    for lr in learning_rates:
        log_reg_base = LogisticRegression(max_iter=1000)
        model = AdaBoostClassifier(estimator=log_reg_base, n_estimators=n, learning_rate=lr,
random_state=42)
        model.fit(X_train, y_train)
        y_pred = model.predict(X_test)
        acc = accuracy_score(y_test, y_pred)
        results.append({
            'Base': 'LogisticRegression',
```

```
        'n_estimators': n,
        'learning_rate': lr,
        'Accuracy': acc
    })

results_df = pd.DataFrame(results)
print(results_df)

import seaborn as sns
plt.figure(figsize=(12, 6))
sns.barplot(x='n_estimators', y='Accuracy', hue='Base', data=results_df, ci=None)
plt.title('AdaBoost Accuracy with Different Estimators and n_estimators')
plt.show()
```

## Program 10

Build k-Means algorithm to cluster a set of data stored in a .CSV file.

Screenshot



Lab-10    K-means

→ Work on Algorithm of K-means?

→ 1) Select the number K to decide the no of clusters
   2) Select random k points or centroids
   3) Assign each data points to their closest centroid which will form the predefined k clusters
   4) Calculate the variance and place a new centroid of each cluster
   5) Repeat the third steps which means reassign each data point to new closest centroid of each cluster
   6) If any reassignment occurs then go to step 4 else go to finish
   7) The model is ready

→ How to determine no of clusters?
   Elbow method
   Silhouette score
   gap statistic

Q) What is the formula of sum of squared errors? Plot SSE vs no of clusters.

→ $SSE = \sum_{k=1}^{K} \sum_{i \in C_k} ||x_i - u_k||^2$

The plot of SSE vs no of clusters
1) Run the k-means algorithm for a range of clusters values
2) Compute the SSE for each no of clusters
3) Plot SSE vs no of clusters
4) Calculate the variance and place a new centroid of each cluster.
5) Repeat the step 3 which means reassign each data point to the new closest centroid of each cluster
6) If any reassignment occurs, then go to step 4 else go to Finish.

Q Describe Elbow technique?

i) Run k-means for range of values of k
ii) Compute the within cluster sum of squares for each values of k
iii) Plot the wcss against the no of cluster (k)
iv) look for the elbow point - the value of k at which the rate of decrease in wcss slows significantly

(3) Discuss all parameters used in kmeans()

3 i) n_clusters : no of clusters to form
ii) n-init : no of times the k-means algorithm will be run without different centroid seeds
iii) max-iter : max no of iterations of the k-means.
iv) tol : Relative tolerance with regards to inertia to declare convergence.
vi) Random-state : controls the regards to inertia to declare convergence
vii) Algorithm : k-means implementation
viii) verbose : Enable verbose output for debugging purposs

Q) cluster eight point with (x,y) respresenting locations into 3 cluster A₁(2,10), A₂(2,5), A₃(8,4), A₄(5,8), A₅(7,5), A₆(6,4), A₇(1,2), A₈(4,9)
Initial cluster: A₁(2,10), A₄(5,8), A₇(1,2)

→ Calculation of distance b/w points A₁(2,10) and C₁(2,10)
$P_1(P_1, C_1) = |x_2 - x_1| + |y_2 - y_1| = |2-2| + |10-10| = 0$
calculating dis b/w A₂(2,10) and C₁(5,8):
$P(A_1, C_2) = |5-2| + |8-10| = 3 - 2 = 3$
calculating distance b/w A₁(2,10) and C₃(1,2):
$P(A_1, C_3) = |2-1| + |10-2| = 9 + 1$
T^like we calculate the distance of other points from each of the centre of 3 clusters.

for cluster-01
→ we have only one point A₁(2,10) in cluster-01
→ so cluster centre remains the same

Iteration(01) Given Pts | Distance from C₁ | Distance from C₂ | Distance from C₃ | Point belongs to
| A/P₁ | | | | |
| P₁(2,10) | 0 | 5 | 9 | C₁ |
| P₂(3,12) | 5 | 6 | 2 | C₃ |
| P₃(1,4) | 12 | 7 | 9 | C₂ |
| P₄(5,8) | 5 | 0 | 10 | C₂ |
| P₅(7,5) | 10 | 5 | 9 | C₂ |
| P₆(6,3) | 10 | 3 | 4 | C₂ |
| P₇(1,2) | 9 | 10 | 0 | C₃ |
| P₈(4,9) | 5 | 2 | 10 | C₂ |

For cluster 02
New centre = C₂ $(5 + 8 + 7 + 6)/4$, $(8 + 8 + 5 + 4)/4$
$4 + 6 + 5 + 4/4 = (6, 6)$

For cluster = C₃ $(1 + 10)/2$, $(5, 2)/2$
$= (1.5, 3.5)$  This is completion of Iteration-01

Iteration-02

| Given Pts | Dist from C₁ | Dist from C₂ | Dist from C₃ | Point belongs to cluster |
|---|---|---|---|---|
| A₁(2,10) | 0 | 8 | 7 | C₁ |
| A₂(2,5) | 5 | 3 | 2 | C₃ |
| A₃(8,4) | 12 | 4 | 7 | C₂ |
| A₄(5,8) | 5 | 3 | 8 | C₂ |
| A₅(2,5) | 10 | 2 | 7 | C₂ |
| A₆(5,4) | 15 | 2 | 5 | C₂ |
| A₇(1,2) | 9 | 9 | 2 | C₃ |
| A₈(4,9) | 3 | 5 | 8 | C₁ |

First  centre = $(2 + 4)/2$, $(10 + 9)/2 = (3, 9.5)$
Second  Centre = $((8 + 5 + 7 + 6)/4$, $(4 + 8 + 5 + 4)/4) = (6.5, 5.25)$
Third  centre = $(2 + 1)/2$, $(5 + 2)/2 = (1.5, 2.5)$

Iteration-03

| Given pts | Distance from $c_1$ | Dist from $c_2$ | Dist from $c_3$ | Point belongs to cluster |
|---|---|---|---|---|
| $A_1(2,10)$ | 1.5 | 9.23 | 9 | $c_1$ |
| $A_2(2,5)$ | 5.3 | 4.75 | 2 | $c_3$ |
| $A_3(8,4)$ | 10.5 | 2.23 | 2 | $c_2$ |
| $A_4(5,8)$ | 3.3 | 4.75 | 8 | $c_1$ |
| $A_5(7,5)$ | 8.8 | 0.78 | 7 | $c_2$ |
| $A_6(6,4)$ | 8.5 | 1.75 | 5 | $c_2$ |
| $A_7(1,2)$ | 9.5 | 8.75 | 2 | $c_3$ |
| $A_8(4,9)$ | 1.5 | 6.25 | 8 | $c_1$ |

$c_1$ ⟶ $((2+5+4)/3 = 3.66 ; (10+8+9)/3 = 9$

$(8+7+6)/3 = 7 , (4+5+4)/3 = 4.33$

$(2+1)/2 = 1.5 , 3.5$

| Given point | Dist from $c_1$ | Dist from $c_2$ | Dist from $c_3$ | Point |
|---|---|---|---|---|
| $A_1(2,10)$ | 1.66 | 10.63 | 8 | $c_1$ |
| $A_2(2,5)$ | 5.66 | 5.67 | 2 | $c_3$ |
| $A_3(8,4)$ | 4.34 | 1.33 | 7 | $c_2$ |
| $A_4(5,8)$ | 2.34 | 5.67 | 8 | $c_1$ |
| $A_5(7,5)$ | 7.34 | 0.67 | 7 | $c_2$ |
| $A_6(6,4)$ | 7.34 | 1.33 | 6 | $c_2$ |
| $A_7(1,2)$ | 9.66 | 8.33 | 2 | $c_3$ |
| $A_8(4,9)$ | 0.34 | 7.66 | 8 | $c_1$ |

⟶ Centre of clusters are $(3.66, 9) (7, 4.33) (1.5, 3.5)$

44

Code:

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.cluster import KMeans
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import accuracy_score

data = {
    'Name': [f'Person_{i+1}' for i in range(50)],
    'Age': np.random.randint(18, 70, size=50),
    'Income': np.random.randint(20000, 120000, size=50)
}

df = pd.DataFrame(data)

df.to_csv('income.csv', index=False)

df = pd.read_csv('income.csv')

X = df[['Age', 'Income']]

scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

X_train, X_test = train_test_split(X_scaled, test_size=0.2, random_state=42)

sse = []
k_range = range(1, 11)
for k in k_range:
    kmeans = KMeans(n_clusters=k, random_state=42)
    kmeans.fit(X_train)
    sse.append(kmeans.inertia_)

plt.plot(k_range, sse, marker='o')
plt.title('SSE vs Number of Clusters')
plt.xlabel('Number of Clusters')
plt.ylabel('Sum of Squared Errors (SSE)')
plt.show()

optimal_k = 3
kmeans = KMeans(n_clusters=optimal_k, random_state=42)
kmeans.fit(X_train)
y_pred = kmeans.predict(X_test)

print(f'Predicted Clusters for Test Data: {y_pred}')
```

# Program 11

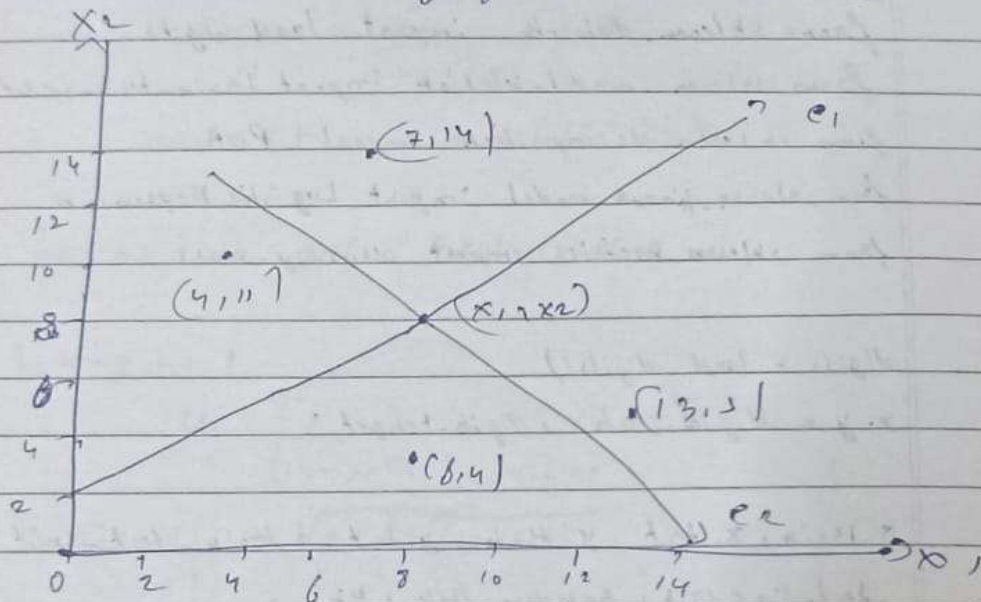Implement Dimensionality reduction using Principal Component Analysis (PCA) method.

Screenshot

Q) Define PCA ?

→ It is a dimensionality reduction and ml method used to a simplify a large data set into a smaller set while still maintaining significant patterns and trends

step: Geometrical meaning of first principal component

Code:

```python
import pandas as pd
import numpy as np
from sklearn.preprocessing import StandardScaler, LabelEncoder
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier
from sklearn.svm import SVC
from sklearn.metrics import accuracy_score
from sklearn.decomposition import PCA
from scipy import stats

df = pd.read_csv('heart (2).csv')

z_scores = np.abs(stats.zscore(df.select_dtypes(include=[np.number])))
df_no_outliers = df[(z_scores < 3).all(axis=1)]

df_cleaned = df_no_outliers.copy()
for col in df_cleaned.select_dtypes(include='object').columns:
    df_cleaned[col] = LabelEncoder().fit_transform(df_cleaned[col])

X = df_cleaned.drop('HeartDisease', axis=1)
y = df_cleaned['HeartDisease']

scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=0.2, random_state=42, stratify=y)

models = {
    "Logistic Regression": LogisticRegression(max_iter=1000),
    "Random Forest": RandomForestClassifier(),
    "SVM": SVC()
}

print("Accuracy without PCA:")
for name, model in models.items():
    model.fit(X_train, y_train)
    y_pred = model.predict(X_test)
    acc = accuracy_score(y_test, y_pred)
    print(f"{name}: {acc:.4f}")

pca = PCA(n_components=5)
X_pca = pca.fit_transform(X_scaled)
X_train_pca, X_test_pca, y_train, y_test = train_test_split(X_pca, y, test_size=0.2, random_state=42, stratify=y)
```

```python
print("\nAccuracy with PCA:")
for name, model in models.items():
    model.fit(X_train_pca, y_train)
    y_pred = model.predict(X_test_pca)
    acc = accuracy_score(y_test, y_pred)
    print(f"{name}: {acc:.4f}")
```