

VISVESVARAYATECHNOLOGICAL UNIVERSITY

“JnanaSangama”, Belgaum -590014, Karnataka.



LAB REPORT
on

Machine Learning(23CS6PCMAL)

Submitted by

Aditya(1BM22CS015)

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 **Aditya (1BM22CS015)**, 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 anMachine Learning (23CS6PCMAL) work prescribed for the said degree.

Dr. Seema Patil Assistant Professor DepartmentofCSE,BMSCE	Dr.Kavitha Sooda Professor & HOD Department of CSE, BMSCE
---	---

Index

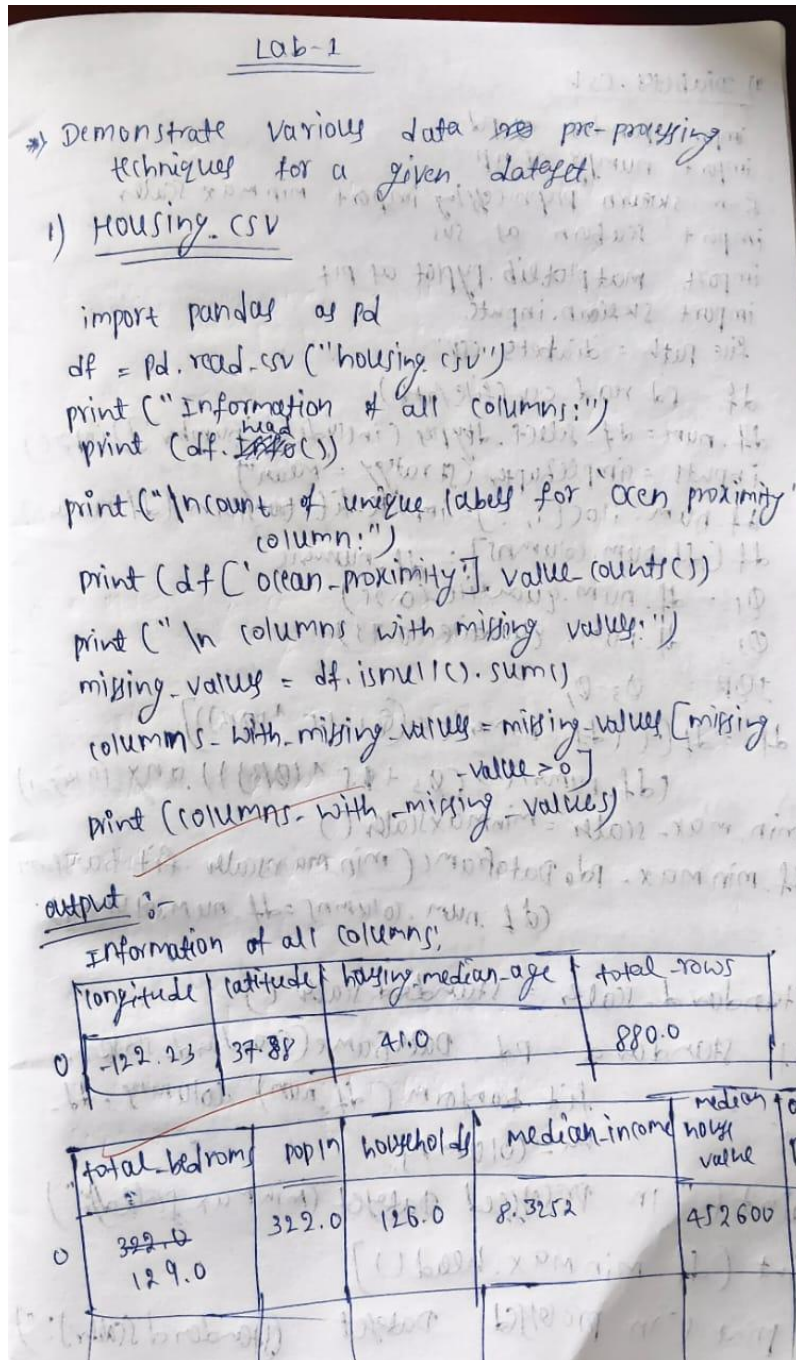
Sl. No.	Date	Experiment Title	Page No.
1	4-3-2025	Write a python program to import and export data using Pandas library functions	5-8
2	11-3-2025	Demonstrate various data pre-processing techniques for a given dataset	9-13
3	18-3-2025	Use an appropriate data set for building the decision tree (ID3) and apply this knowledge to classify a new sample.	14-16
4	1-4-2025	Implement Linear and Multi-Linear Regression algorithm using appropriate dataset	17-20
5	8-4-2025	Build Logistic Regression Model for a given dataset	21-23
6	15-4-2025	Build KNN Classification model for a given dataset.	24-27
7	15-4-2025	Build Support vector machine model for a given dataset	28-30
8	22-4-2025	Implement Random forest ensemble method on a given dataset.	31-32
9	22-4-2025	Implement Boosting ensemble method on a given dataset.	33-35
10	29-4-2025	Build k-Means algorithm to cluster a set of data stored in a .CSV file.	36-39
11	29-4-2025	Implement Dimensionality reduction using Principal Component Analysis (PCA) method.	40-44

GithubLink: <https://github.com/AdityaMK15/machine-learning-lab>

Program 1

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

Screenshot



2) Diabetes.csv

```
import pandas as pd
import numpy as np
from sklearn.preprocessing import minmax_scaler
import sklearn as sk
import matplotlib.pyplot as plt
import sklearn.inputs

file_path = "diabetes.csv"
df = pd.read_csv(file_path)
df_num = df.select_dtypes(include=["number"]).copy()
inputs = SimpleImputer(strategy="mean")
df_num.iloc[:, :] = inputs.fit_transform(df_num)
df[df_num.columns] = df_num
Q1 = df_num.quantile(0.25)
Q3 = df_num.quantile(0.75)
IQR = Q3 - Q1
df = df[(df_num < (Q1 + 1.5 * IQR)) && (df_num > (Q3 - 1.5 * IQR))]
min_max_scaler = minmax_scaler()
df_min_max = pd.DataFrame(min_max_scaler.fit_transform(df_num), columns=df_num.columns)

standard_scaler = standard_scaler()
df_standard = pd.DataFrame(standard_scaler.fit_transform(df_min_max), columns=df_min_max.columns)
```

```
print(df_standard.head())
```

Diabetes Dataset

Missing values are present in numerical column is present, which are replaced by the mean of the respective column.

2) No categorical columns, hence no encoding.

3) min-max scaling transform data to a fixed range (0, 1) using

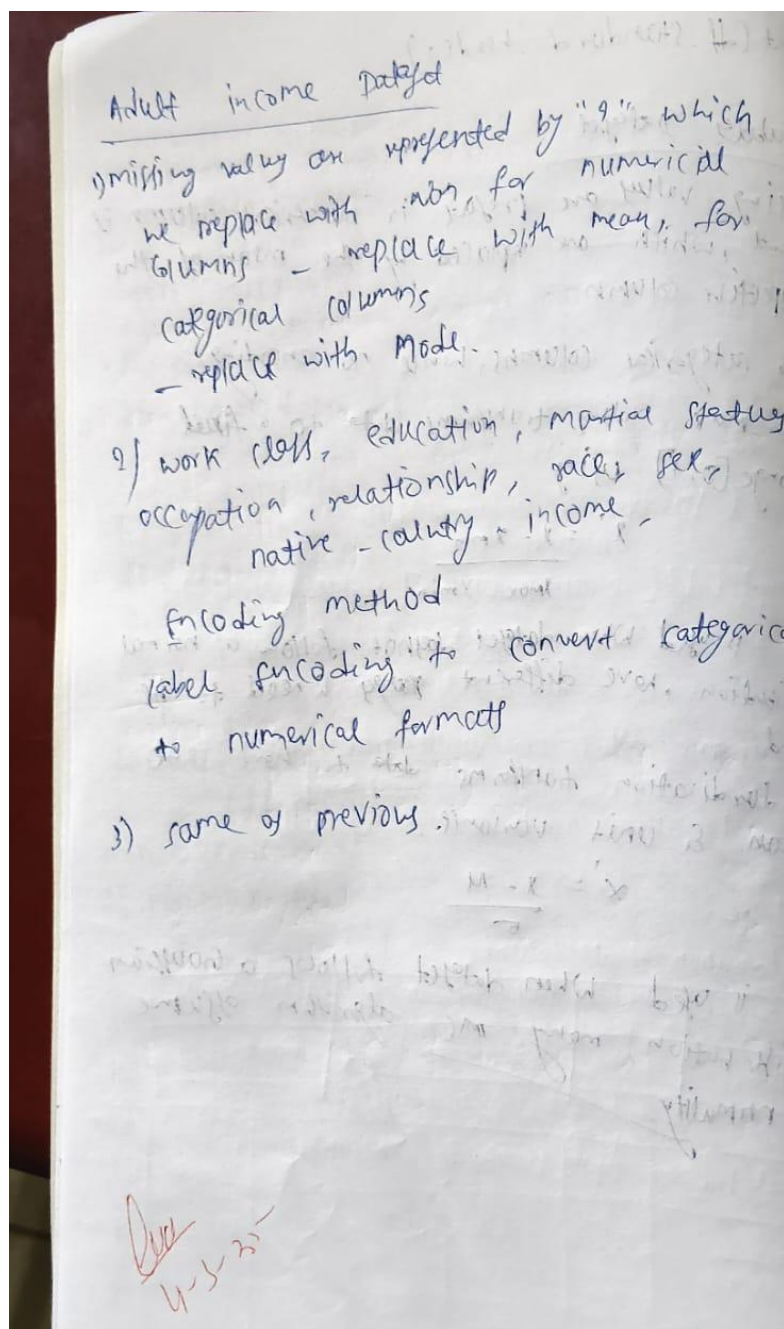
$$x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

It is used when dataset does not follow a normal distribution, have different ranges & need to be bound.

Standardization transforms data to have zero mean & unit variance.

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

It is used when dataset follows a Gaussian distribution, many ML algorithms assume normality.



Code:

```
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())

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

```
LAB-2

1) import pandas as pd
data = {
    'id': ['001', '002', '003', '004'],
    'name': ['A', 'B', 'C', 'D'],
    'marks': [30, 40, 50, 60]
}

df = pd.DataFrame(data)
print("sample data:")
print(df.head(1))

output: sample data
```

	id	name	marks
0	001	A	30
1	002	B	40
2	003	C	50
3	004	D	60

```
2) from sklearn.datasets import load_diabetes
diabetes = load_diabetes()
df = pd.DataFrame(diabetes.data, columns =
    diabetes.feature_names)

df["target"] = diabetes.target
print("sample data")
print(df.head())
```

sample data

age	sex	ppv	bp	sh
0.038076	0.05064	0.061696	0.01172	-0.044223

S2	S3	S4	S5	S6	target
-0.034121	-0.0634	-0.0025	0.0999	-0.0176	11.0

3) file-path = 'data.csv'
 df = pd.read_csv('data.csv')
 print("sample data:")
 print(df.head())
 print("\n")

output:

sample data

	product	quantity	price	value	Region
0	laptop	5	1000	5000	North
1	mouse	15	20	300	West
2	keyboard	10	50	500	East
3	monitor	8	200	1600	South

4) df = pd.read_csv('data.csv')
 print(df)
 print("first few rows of the data set")
 print(df.head())

output:-

product	quantity	price	value	Region
0 laptop	5	1000	5000	North

1) Bank Data Analysis

using the code of stock market data analysis consider the following

① HDFC Bank Ltd, ICICI Bank Ltd, Kotak Bank Ltd, Mahindra Bank Ltd.

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

② start date = 2024-01-01, End Date = 2024-12-30

1) plot the closing price & daily returns for all the 3 banks mentioned.

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)


```

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])

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 13 Logistic Regression

① $a_0 = -5$
 $a_1 = 0.8$

i) $f(x) = \frac{1}{1 + \exp(-(-5 + 0.8x))}$

ii) $f(7) = \frac{1}{1 + \exp(-(-5 + 0.8 \times 7))}$
 $= 0.6417$

iii) If
 $f(x) < 0.5$ then the student is fail,
 $f(x) > 0.5$ then the student is pass.

② $z = [2, 1, 0]$

$\text{softmax}(z_k) = \frac{e^{z_k}}{\sum_{j=1}^n e^{z_j}}$

$\text{softmax}(z_1) = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}} = \frac{e^2}{e^2 + e^1 + e^0} \approx 0.6652$

$\text{softmax}(z_2) = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}} \approx \frac{e^1}{e^2 + e^1 + e^0}$

$$\text{softmax}(z_3) = \frac{e^0}{e^1 + e^2 + e^0} \approx 0.091$$

18-3

Code:

```
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]
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

LAB - 4

Q)

Week	Sales
1	2
2	4
3	5
4	9
5	9

soln: $X^T = \begin{bmatrix} 1 & 2 & 3 & 4 \end{bmatrix}$ $Y = \begin{bmatrix} 2 \\ 4 \\ 5 \\ 9 \\ 9 \end{bmatrix}$

$X = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \\ 1 & 4 & 9 & 16 \\ 1 & 9 & 16 & 25 \end{bmatrix}$ $Y = \begin{bmatrix} 2 \\ 4 \\ 5 \\ 9 \\ 9 \end{bmatrix}$

$\beta = [(X^T X)^{-1} X^T] Y$

compute $X^T X = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ 1 & 4 \end{bmatrix} = \begin{bmatrix} 4 & 10 \\ 10 & 30 \end{bmatrix}$

$(X^T X)^{-1} = \begin{bmatrix} 1.5 & -0.5 \\ -0.5 & 0.2 \end{bmatrix}$

$[(X^T X)^{-1} X^T] = \begin{bmatrix} 1 & 0.5 & 0 & -0.5 \\ -0.3 & -0.1 & 0.1 & 0.3 \end{bmatrix}$

$[(X^T X)^{-1} X^T] Y = \begin{bmatrix} 1 & 0.5 & 0 & -0.5 \\ -0.3 & -0.1 & 0.1 & 0.3 \end{bmatrix} \begin{bmatrix} 2 \\ 4 \\ 5 \\ 9 \end{bmatrix}$

14-21A2

Year	Age
1	1
2	2
3	3
4	4
5	5

$$p_1 = 9.2$$

$$\begin{bmatrix} 1 & 0 & 0 & 1 \end{bmatrix} = r_4$$

$$\begin{bmatrix} 7.0 & 2.1 \\ 2.0 & 2.0 \end{bmatrix} = 1.7 \times 10^3$$

```
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]}")

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

lab-5

Instance	A_1	A_2	classification
1	hot	high	no
2	hot	high	no
6	cool	high	no
7	hot	high	no
8	hot	normal	yes

Entropy = $-\frac{4}{5} \log \frac{4}{5} - \frac{1}{5} \log \frac{1}{5}$
 $= 0.8219$

for a_1 ,
 $Split(1+, 3-) = \frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4}$
 $= 0.8113$

$Split(0+, 1-) = 0$
 $Gain(0, 2) = 0.8219 - \frac{4}{5} \times 0.8113 - 0 = 0.0771$

for a_2 ,
 $Split(0+, 4-) = 0$
 $Split(1-, 1-) = 0$
 $Gain(1, 0) = 0.8219$

$\therefore a_2$ has highest gain value it is taken as root

Diagram of a decision tree:
 root node a_2 splits into 'high' and 'normal'.
 'high' leads to a leaf node 'no' with instances 1, 2, 6.
 'normal' leads to a leaf node 'yes' with instance 8.

Code:

```
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

person	Age	Salary	Target
A	18	50	N
B	23	55	N
C	29	70	Y
D	41	70	Y
E	43	70	Y
F	38	40	Y
X	35	100	Y

$(x_L, y_L) = (35, 100)$

for $(18, 30)$ d = $\sqrt{(x_L - x_1)^2 + (y_L - y_1)^2}$
 $= \sqrt{(35 - 18)^2 + (100 - 50)^2}$
 $= 82.81$

for $(23, 55)$ = $\sqrt{(35 - 23)^2 + (100 - 55)^2}$
 $= 46.57$

for $(29, 70)$ = $\sqrt{(35 - 29)^2 + (100 - 70)^2}$
 $= 31.95$

for $(41, 60)$ = $\sqrt{(35 - 41)^2 + (100 - 60)^2}$
 $= 40.44$

for $(43, 70)$ = $\sqrt{(35 - 43)^2 + (100 - 70)^2}$
 $= 51.04$

for $(38, 40)$ = $\sqrt{(35 - 38)^2 + (100 - 40)^2}$
 $= 60.07$

person	Age	Salary	Target	Distance	Rank
A	18	50	N	82.81	5
B	23	55	N	46.57	4
C	29	70	Y	31.95	2
D	41	60	Y	40.44	3
E	43	70	Y	51.04	1
F	38	40	Y	60.07	6
X	35	100	Y	-	-

K=3 $\boxed{Y-Y}$ $\boxed{1-N}$

K=1 -Y

K=2 -N

K=3 -Y

So, $x(35, 100)$ is Y.

Code:

```
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())
```

```

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()

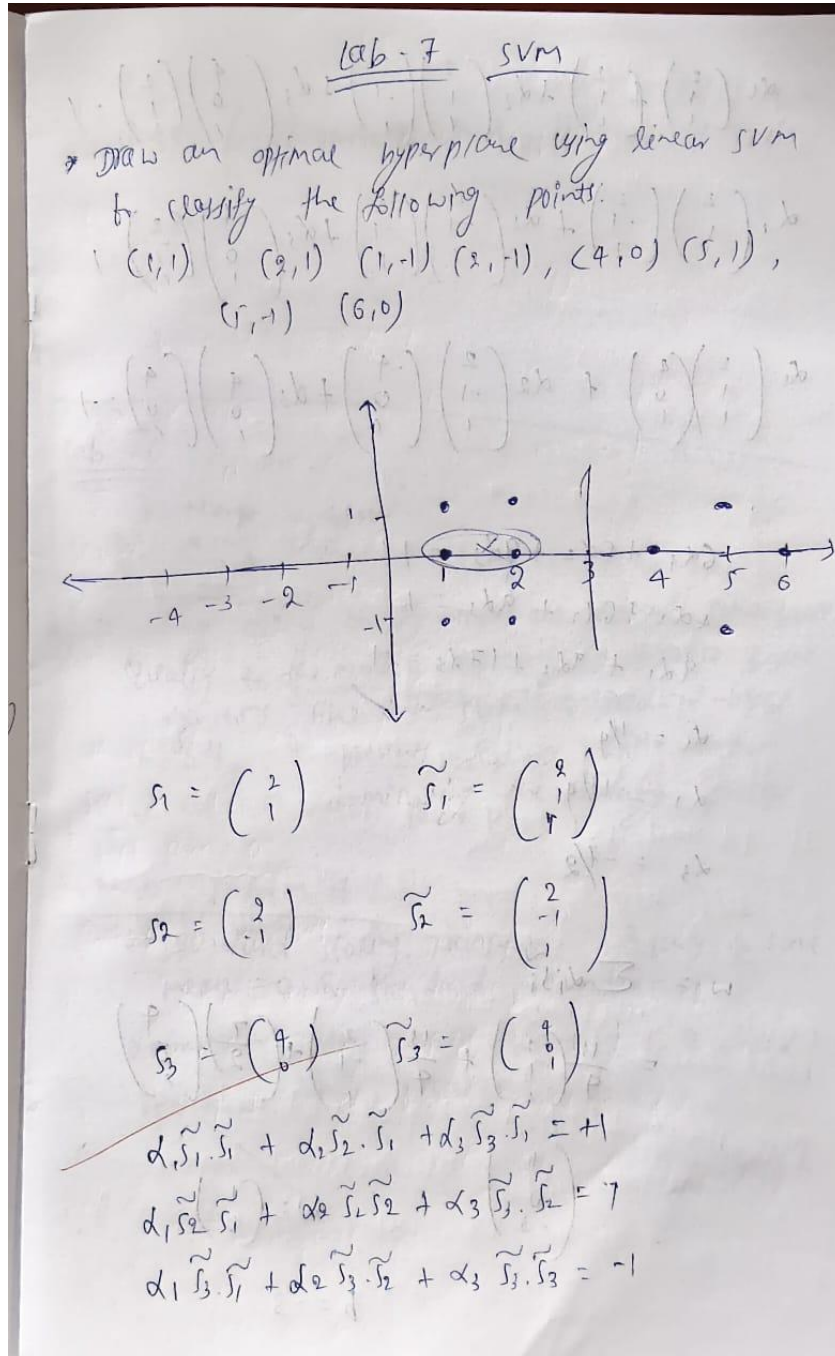
print("Classification Report for Diabetes Dataset:")
print(classification_report(y_test_diabetes, y_pred_diabetes))

```

Program 7

Build Support vector machine model for a given dataset

Screenshot



$$\alpha_1 \begin{pmatrix} 2 \\ 1 \\ 1 \end{pmatrix} + \alpha_2 \begin{pmatrix} 2 \\ 1 \\ 1 \end{pmatrix} + \alpha_3 \begin{pmatrix} 4 \\ 0 \\ 1 \end{pmatrix} = 1$$

$$\alpha_1 \begin{pmatrix} 2 \\ 1 \\ 1 \end{pmatrix} + \alpha_2 \begin{pmatrix} 2 \\ 1 \\ 1 \end{pmatrix} + \alpha_3 \begin{pmatrix} 4 \\ 0 \\ 1 \end{pmatrix} = 1$$

$$\alpha_1 \begin{pmatrix} 2 \\ 1 \\ 1 \end{pmatrix} + \alpha_2 \begin{pmatrix} 2 \\ 1 \\ 1 \end{pmatrix} + \alpha_3 \begin{pmatrix} 4 \\ 0 \\ 1 \end{pmatrix} = 1$$

$$6\alpha_1 + 4\alpha_2 + 9\alpha_3 = 1$$

$$4\alpha_1 + 6\alpha_2 + 9\alpha_3 = 1$$

$$9\alpha_1 + 9\alpha_2 + 17\alpha_3 = -1$$

$$\alpha_1 = 1/9$$

$$\alpha_2 = 13/9$$

$$\alpha_3 = -7/2$$

$$w = \sum \alpha_i \tilde{x}_i$$

$$= \frac{13}{9} \begin{pmatrix} 2 \\ 1 \\ 1 \end{pmatrix} + \frac{13}{9} \begin{pmatrix} 2 \\ 1 \\ 1 \end{pmatrix} + \left(-\frac{7}{2}\right) \begin{pmatrix} 4 \\ 0 \\ 1 \end{pmatrix}$$

$$= \begin{pmatrix} 1 \\ 0 \\ -3 \end{pmatrix} \quad w = \begin{pmatrix} 1 \\ 0 \\ -3 \end{pmatrix}$$

$$b = -3$$

since $\begin{pmatrix} 1 \\ 0 \end{pmatrix}$ vertical line, \parallel to y axis

at $x=3$

$$b+3=0$$

Lab - 6

k-nearest neighbour

Question-2

2) Feature scaling ensures that all features contribute equally to the model, especially for distance-based algorithms like k-NN, SVM & gradient-based models, it prevents features with larger values from dominating the learning process.

How to perform it

→ standard scaling: Transforms features to have mean = 0 & standard deviation = 1.

→ min-max scaling: scales features to a fixed range, usually [0, 1].

Code:

```
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

Date: 15/04/2025
Page No.: _____

Lab - 8 (Random Forest)

- 1) Decision Tree vs Random Forest:
 - single tree vs Multiple.
 - less accurate vs More accurate
 - Fast to train vs slower to train
 - single prediction vs majority vote.
- 2) Parameters of Random Forest Classification:
 - n-estimate: no. of trees in forest
 - criterion: to nearest quality of split
 - max-depth: min depth of tree
 - min-sample-split: min sample req to split node
 - Max feature: no. of feature to look for best fit
- 3) Algorithm:
 - 1) training dataset = X, y
 - 2) for n times:
 - randomly select sample with replacement
 - grow a decision tree
 - at each split choose random subset of features
 - split nodes using best feature
 - 3) Aggregate prediction:
 - classification - majority vote
 - regression - avg

Ⓐ O/p prediction

Code:

```
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)

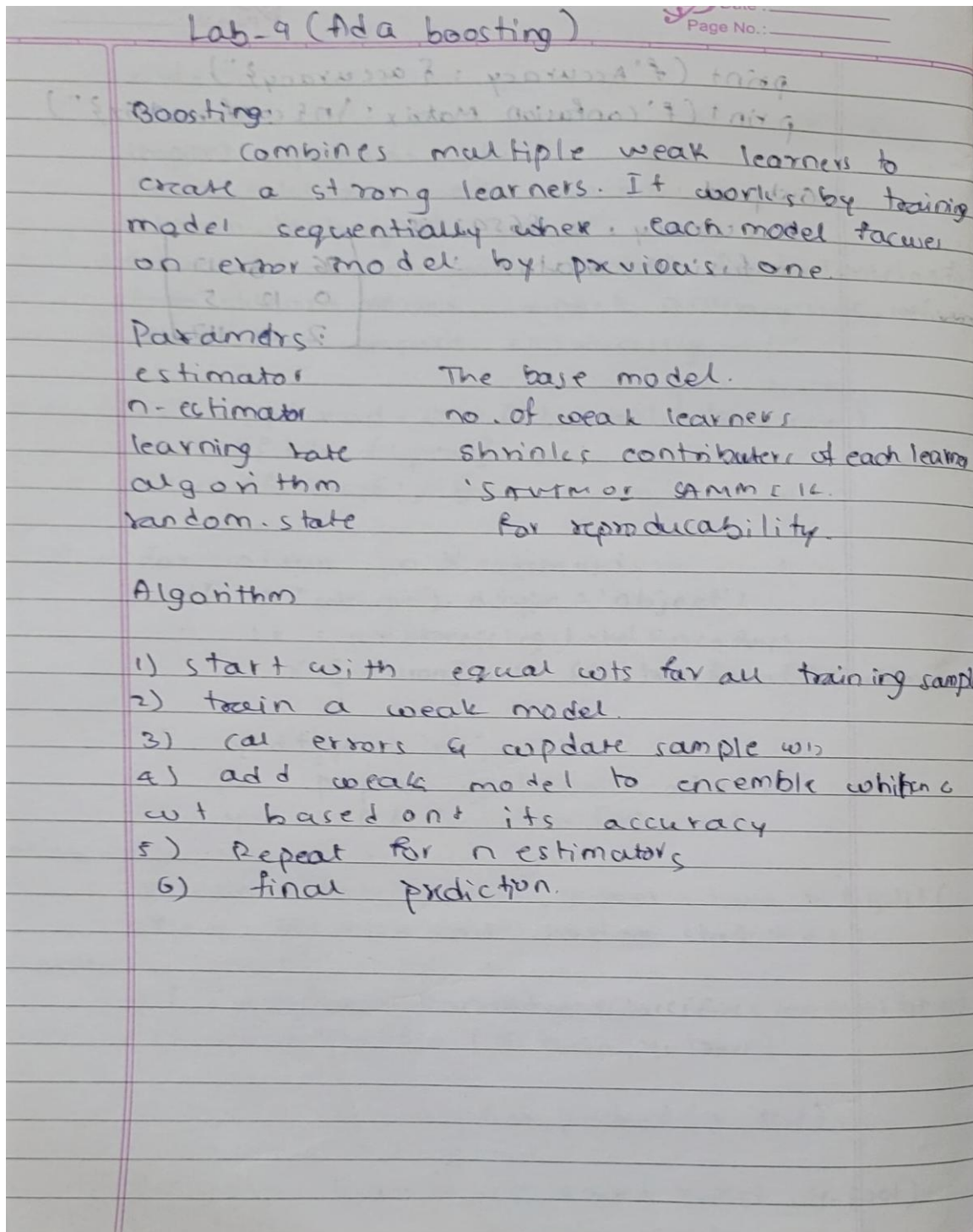
accuracy = accuracy_score(y_test, y_pred)
conf_matrix = confusion_matrix(y_test, y_pred)

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

Program 9

Implement Boosting ensemble method on a given dataset.

Screenshot



Code:

```
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 clustering

2) Cluster eight points (with (x,y) representing location) into 3 clusters
 $A_1(2,10), A_2(2,5), A_3(6,4), A_4(5,8), A_5(7,5), A_6(6,9), A_7(1,2), A_8(4,2)$
 initial clusters: $A_1(2,10), A_4(5,8), A_7(1,2)$
 1st step: calc distance A_1 & each centre of 3 clusters, calc dist b/w A_1 & C_1, C_2, C_3

$P(A_1, C_1)$
 $= |x_1 - x_1| + |y_1 - y_1|$
 $= |2 - 2| + |10 - 10| = 0$
 $P(A_1, C_2) = |5 - 2| + |8 - 10| = 3 + 2 = 5$
 $P(A_1, C_3) = |1 - 2| + |2 - 10| = 1 + 8 = 9$
 1st calc dist of other pts from each centre.

Iteration - 1

Point	Dist to $C_1(2,10)$	Dist to $C_2(5,8)$	Dist to $C_3(1,2)$	Point belongs to cluster
$A_1(2,10)$	0	5	9	C_1
A_2	5	6	4	C_3
A_3	12	7	9	C_2
A_4	5	0	10	C_2
A_5	10	5	9	C_2
A_6	10	5	7	C_2
A_7	9	10	0	C_3
A_8	3	2	10	C_2

For C_1 : only A_1 in C_1 so cluster centre remains same.
 For C_2 : $\left(\frac{8+5+7+6+5}{5}, \frac{4+8+4+9}{5}\right) = (6, 6)$
 For C_3 : $\left(\frac{2+1}{2}, \frac{10+2}{2}\right) = (1.5, 6)$

Iteration 2

Row \ Point	$C_1 (0, 10)$	$C_2 (6, 6)$	$C_3 (14, 3.5)$	
A1 (2, 10)	0	8	7	C_1
A2	5	5	2	C_3
A3	12	4	7	C_2
A4	5	3	8	C_2
A5	10	2	7	C_2
A6	10	2	5	C_2
A7	9	9	2	C_3
A8	3	8	8	C_1

$$C_1 = \left(\frac{2+4}{2+2}, \frac{20+9}{20} \right) = (3, 9.5)$$

$$C_2 = \left(\frac{8+5+12}{4}, \frac{4+8+7+4}{4} \right) = (6.5, 8.25)$$

$$C_3 = \left(\frac{5+10}{2}, \frac{5+2}{2} \right) = (7.5, 3.5)$$

Iteration 3

Row \ Point	$C_1 (3, 9.5)$	$C_2 (6.5, 8.25)$	$C_3 (7.5, 3.5)$	
A1 (2, 10)	9.5	9.25	7	C_1
A2	5	4.25	2	C_3
A3	9.5	2.25	7	C_2
A4	3.5	9.25	7	C_1
A5	8.5	0.25	8	C_2
A6	8.5	1.75	7	C_2
A7	9.5	8.25	5	C_2
A8	1.5	0.25	8	C_3

$$C_1 = \left(\frac{2+9+4}{3}, \frac{10+9+9}{3} \right) = (6.67, 9.67)$$

$$C_2 = \left(\frac{8+7+6}{3}, \frac{4+11+4}{3} \right) = (7, 9.33)$$

$$C_3 = \left(\frac{5+1}{2}, \frac{5+2}{2} \right) = (3, 3.5)$$

Iteration: 4

Cluster	1 (356, 9)	2 (7, 43)	3 (18, 35)	4 (1, 10)
A1 (3, 10)	7.66	1.66	8	9
A2	5	8.66	7	3
A3	9.33	1.33	8	2
A4	2.33	5.67	7	1
A5	7.34	0.67	6	2
A6	7.34	1.33	2	2
A7	9.66	8.33	8	3
A8	0.33	7.66	8	1

same cluster pop

3 clusters: 1 (356, 9), 2 (7, 43), 3 (18, 35)

K-means code:

important random by pd

from sklearn.model_selection

names = [f'Person_{i+1}' for i in range(50)]

ages = np.random.randint(20, 60, size=50)

incomes = np.random.randint(20000, 120000, size=50)

data = pd.DataFrame({'name': names, 'age': ages, 'income': incomes})

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

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

x = data[['age', 'income']]

x_train, x_test = train_test_split(x, test_size=0.2, random_state=42)

scaler = StandardScaler()

x_train_scaled = scaler.fit_transform(x_train)

x_test_scaled = scaler.transform(x_test)

k = 3

cluster = range(k)

for n in cluster_range:

kmeans = KMeans(n_clusters=k, random_state=42, n_init=10)

kmeans.fit(x_train_scaled)

score = kmeans.score(x_test_scaled)

accuracy = kmeans.score(x_test_scaled)

score = accuracy_score(y_test, y_pred)

print('Accuracy: %f' % accuracy)

0/f

accuracy = 0.6

Code:

```
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)

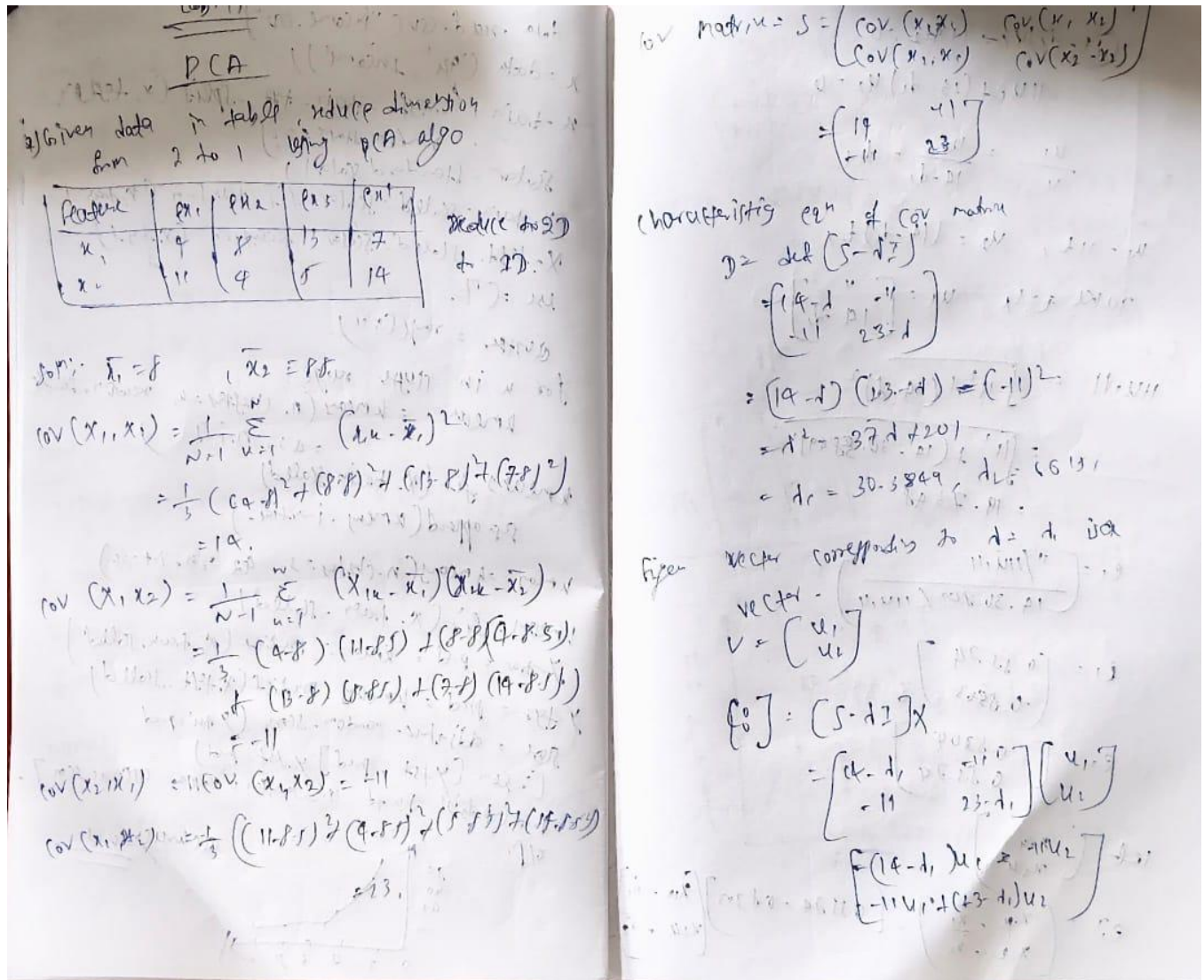
print(f'Predicted Clusters for Test Data: {y_pred}')

```

Program 11

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

Screenshot



$$(14-d_1)u_1 + 11u_2 = 0$$

$$-11u_1 + (14-d_1)u_2 = 0$$

$$\frac{u_1}{11} = \frac{u_2}{14-d_1} = t_1$$

$$u_1 = 11t_1, \quad u_2 = (14-d_1)t_1$$

$$\text{Take } t_1 = 1, \quad u_1 = \begin{bmatrix} 11 \\ 14-d_1 \end{bmatrix}$$

$$\|u_1\| = \sqrt{11^2 + (14-d_1)^2}$$

$$= \sqrt{11^2 + (14-30.8849)^2}$$

$$= 19.7398$$

$$e_1 = \left(\frac{11}{19.7398}, \frac{14-30.8849}{19.7398} \right)$$

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

$$e_2 = \begin{bmatrix} 0.8304 \\ 0.5574 \end{bmatrix}$$

$$\text{Let } \begin{pmatrix} x_{1k} \\ x_{2k} \end{pmatrix}$$

$$e_1^T = \begin{pmatrix} x_{1k} - \bar{x}_1 \\ x_{2k} - \bar{x}_2 \end{pmatrix} \begin{bmatrix} 0.5574 & -0.8303 \end{bmatrix} \begin{pmatrix} x_{1k} - \bar{x}_1 \\ x_{2k} - \bar{x}_2 \end{pmatrix}$$

$$= 0.5574(x_{1k} - \bar{x}_1) - 0.8303(x_{2k} - \bar{x}_2)$$

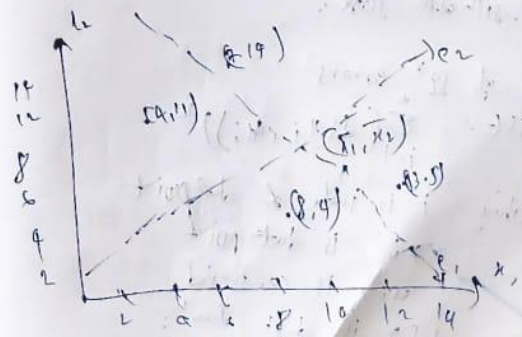
$$\begin{pmatrix} u_1 \\ u_2 \end{pmatrix} = \begin{pmatrix} 11 \\ 11 \end{pmatrix}$$

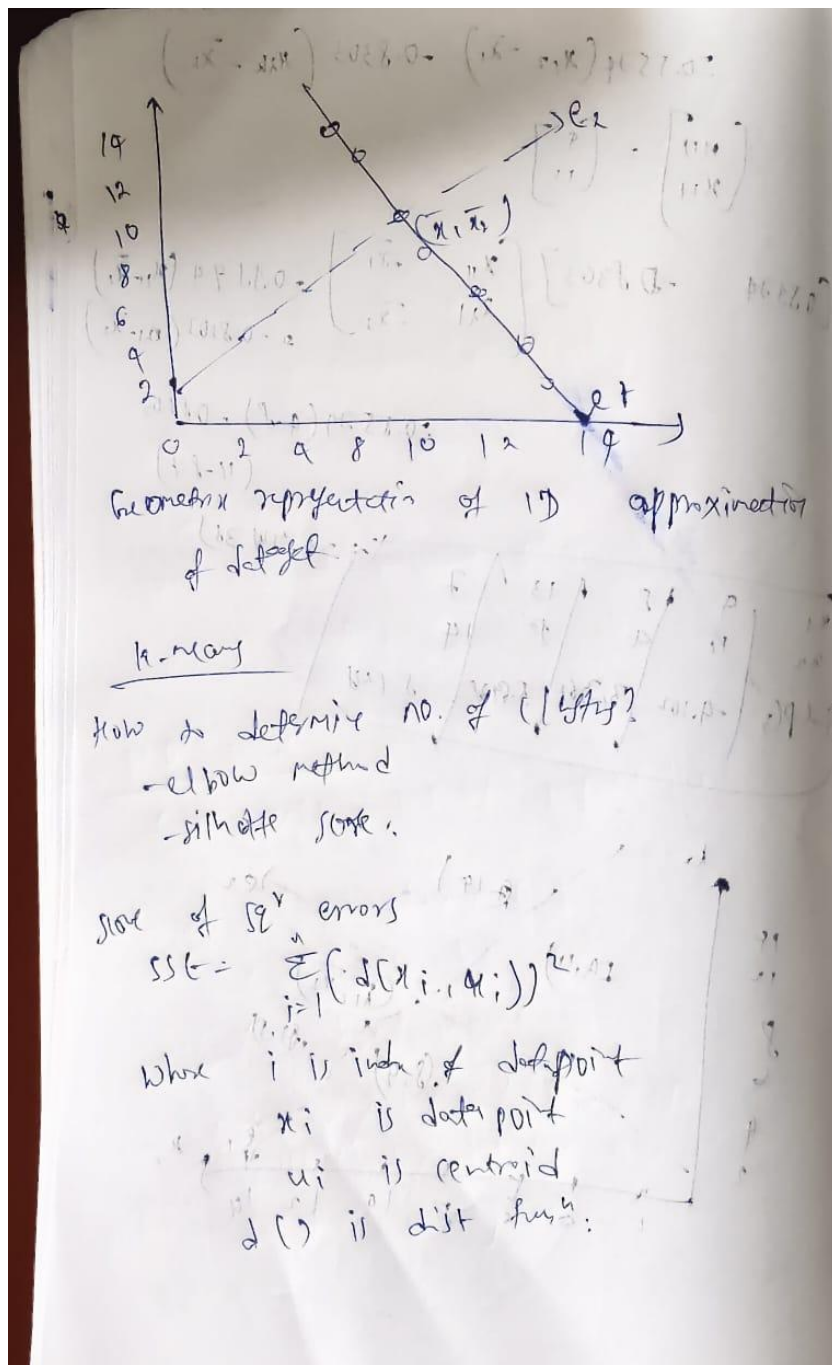
$$\begin{bmatrix} 0.5574 & -0.8303 \end{bmatrix} \begin{pmatrix} x_{11} - \bar{x}_1 \\ x_{21} - \bar{x}_2 \end{pmatrix} = 0.5574(x_{11} - \bar{x}_1) - 0.8303(x_{21} - \bar{x}_2)$$

$$= 0.5574(4-8) - 0.8303(11-8)$$

$$= -4.2030$$

x_1	9	8	13	7
x_2	11	4	5	14
PCs	-4.302	3.236	5.928	-8.128





Code:

```

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)

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}")

```

