

VISVESVARAYA TECHNOLOGICAL UNIVERSITY

“JnanaSangama”, Belgaum- 590014, Karnataka.



LAB REPORT on Machine Learning (23CS6PCMAL)

Submitted by

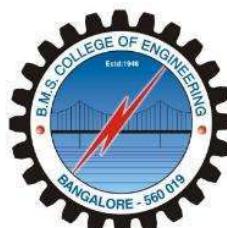
ARBAJ WADAGERA (1BM22CS051)

in partial fulfilment 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
February 2025 – June 25

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 **Arbaj Wadagera (1BM22CS051)**, who is Bonafide student of **B.M.S. College of Engineering**. It is in partial fulfilment for the award of **Bachelor of Engineering in Computer Science and Engineering** of the Visvesvaraya Technological University, Belgaum. The Laboratory report has been approved as it satisfies the academic requirements in respect of a Machine Learning (23CS6PCMAL) work prescribed for the said degree.

Lab Faculty InCharge's:	
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	21-02-2025	Write a python program to import and export data using pandas library functions.	1
2	03-03-2025	Demonstrate various data pre-processing techniques for a given dataset.	7
3	10-03-2025	Use an appropriate data set for building the decision tree (ID3) and apply this knowledge to classify a new sample.	14
4	17-03-2025	Implement Linear and Multi-Linear Regression algorithm using appropriate dataset.	26
5	24-03-2025	Build Logistic Regression Model for a given dataset.	36
6	07-04-2025	Build KNN Classification model for a given dataset.	41
7	07-04-2025	Build Support vector machine model for a given dataset.	46
8	21-04-2025	Implement Random Forest ensemble method on a given dataset.	53
9	05-05-2025	Implement Boosting ensemble method on a given dataset.	56
10	05-05-2025	Build k-Means algorithm to cluster a set of data stored in a .CSV file.	59
11	12-05-2025	Implement Dimensionality reduction using Principal Component Analysis (PCA) method.	64

GitHub Link: <https://github.com/Arbaj-Wadagera/6thSem-ML-Lab>

Program 1

Write a python program to import and export data using Panda's library functions

Screenshot:

Lab - 1

① Initializing values directly into dataframes

```
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	Charlie	35	Chicago
3	David	40	Houston

② Importing datasets from sklearn.datasets

```
from sklearn.datasets import load_iris
iris = load_iris()
df = pd.DataFrame(iris.data, columns=iris.feature_names)
df['target'] = iris.target
print('sample data')
print(df.head(1))
```

Output

Sample data

	Sepal length (cm)	Sepal width (cm)	Petal length (cm)	Petal width (cm)	Target
0	5.1	3.5	1.4	0.2	Iris-setosa
1	4.9	5.4	2.0	0.3	Iris-versicolor

③ Importing datasets from a specific .CSV file

```
file_path = data.csv
```

```
df = pd.read_csv(file_path)
```

```
print(df.head())
```

Output

	Id	Name	Age
0	1	Alice	25
1	2	Bob	29
2	3	Charlie	21

④ Downloading datasets from existing dataset repositories like kaggle

```
df = pd.read_csv("mobileDataset.csv")
```

```
print(df.head(1))
```

Output

comp. name	modelname	weight	Ram	Front Camera
Apple	iphone 16	144g	6GB	12mp
Back Cam	Processor	Battery	Source Size	
48mp	A17 Bionic	3600 mah	6-inches	

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

③ 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",
end='2024-12-30', groupby=tickers)

print(df.head(1))

print(data.shape)

Output

Tickers KOTAKBANK.NS

Price	Open	High	Low	Close	Volume
2024-01-01	1906	1916	1891	1903.0	14259

Code:

```
from sklearn.datasets import load_iris

import pandas as pd

iris = load_iris()

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

df.head()

df['target'] = iris.target

df

import kagglehub

# Download latest version

path = kagglehub.dataset_download("abdulmalik1518/mobiles-dataset-2025")

print("Path to dataset files:", path)

df = pd.read_csv("/content/Mobiles_Dataset_(2025).csv", encoding='latin-1') # or 'ISO-8859-1', or
'cp1252'

df.head()

df['Company Name']

data = {"USN": [1, 2, 3], "Name": ["A", "B", "C"]}

df = pd.DataFrame(data)

df

from sklearn.datasets import load_diabetes
```

```
diabetes = load_diabetes()

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

df.head()

df.columns

df = pd.read_csv("/content/Dataset_of_Diabetes.csv")

df.head()

import yfinance as yf import pandas as pd

import matplotlib.pyplot as plt

tickers = ["RELIANCE.NS", "TCS.NS", "INFY.NS"]

# Fetch historical data for the last 1 year

data = yf.download(tickers, start="2022-10-01", end="2023-10-01", group_by='ticker')

# Display the first 5 rows of the dataset

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

print(data.head())

print("\nShape of the dataset:")

print(data.shape)
```

```
# 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

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

# 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()
```

Program 2

Demonstrate various data pre-processing techniques for a given dataset

Screenshot:

LAB-II

writing python code , considering filename as "housing.csv"

i) To load .csv file into the data frame

```
from google.colab import files  
uploaded = files.upload()
```

ii) To display info. of all columns

```
print("Information of all columns:")  
print(housing_data.info())
```

iii) To display statistical information of all numerical

```
print(housing_data.describe())
```

iv) To display the count of unique labels for "Ocean Proximity" column

```
print(housing_data['ocean-proximity'].value_counts())
```

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

```
print("In columns with missing values:")  
missing_values = housing_data.isnull().sum()  
missing_columns = missing_values[missing_values > 0]  
print(missing_columns)
```

Output

IT - 2 P.M.

i) → uploaded images, also writing practice

ii)	Column	Not-Null Count	¹⁹²³ present
-----	--------	----------------	-------------------------

longitude 20640 (non-null). head of

latitude 23°17' trigrad 20640 non-null

housing-median-age 20640 non-null

iii) ~~Geographical~~ No. to latitude of place housing medium

longitude latitude housing medium
(0°N, 0°E) - age

count 20640.0000 20640.000000 20640.00000

mean -119.569704 35,631861 28.639486

Min 9-003531 Date 9-13-98 15-08-008

2.00 3332 - 11.82 15.58558
11.82 15.58558

iv) median-house-value

count : 20640-000000 (1936) trying

mean ~~206885.816909~~

~~1928 tab as mi 296~~ 115395-615824

S + d

min ✓ 199999.00000 ✓ 199999.00000

~~14/15~~ 1/3 taking

(3) There, where you are.

~~pizzim720px~~ 10 3 29 slow-pizzim

(ii) Diabetes.csv

```
import pandas as pd
import numpy as np
from sklearn.preprocessing import MinMaxScaler, StandardScaler
import seaborn as sns
import matplotlib.pyplot as plt
import sklearn.impute
file_path = "diabetes.csv"
df = pd.read_csv(file_path)
df.num = df.select_dtypes(include=['number']).copy()
imputer = SimpleImputer(strategy="mean")
df_mean = imputer.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))].any(axis=1)
min_max_scaler = MinMaxScaler()
df_minmax = pd.DataFrame(min_max_scaler.fit_transform(df_num), columns=df_num.columns)
StandardScaler = StandardScaler()
df_Standard = pd.DataFrame(StandardScaler.fit_transform(df_num), columns=df_num.columns)
print("In Processed Dataset (minmax scaled):")
print(df_minmax.head())
print("In Provided Dataset (standard scalar):")
print(df_Standard.head())
```

a) which columns in the dataset had missing values.
How to handle them

Ans:- missing values are present in numerical columns
if present, which are replaced by the mean of the
respective column.

b) which categorical columns did you identify? How
did you encode them?

Ans:- No categorical columns, hence no encoding.

c) what is the difference b/w min-max scaling &
standardization? when would you use one
over the other?

Ans:- Min max scaling transforms data to a final
range [0, 1] using

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

used when:- dataset does not follow a normal
dist.

- features have diff. ranges & used
to be bound

Standardization transforms data to have
zero mean & unit variance

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

used when:- dataset follows a Gaussian dist.
- many ML algos. assume
normality.

(will be covered in detail later on)

((background, 16) training
(feature extraction) feature behavior in ((training

((background, 16) training

Adult Income Dataset

- Missing values are represented by "?" which we replace with "nan" for numerical columns
— replace with mean, for categorical columns
— replace with mode.
- work class, education, marital status, occupation, relationship, race, sex, native country, income.

Encoding methods: drop rows, one-hot encoding

- Label Encoding to convert categorical to numerical
- format address (e.g., 0.0)
 - same as previous, marital status, by = 7b
(b) married
 - same as previous, native country, by = 2b
(b) United States

age	work class	education	marital status	relationship	native country
28	Private	HS grad	Married	夫の親戚	United States
29	Private	HS grad	Married	夫の親戚	United States
25	Private	HS grad	Married	夫の親戚	United States
30	Private	HS grad	Married	夫の親戚	United States

④ Downloading dataset from UCI Machine Learning Repository

⑤ Importing dataset from UCI Machine Learning Repository

adult.csv (adult.csv) → read_csv("adult.csv")

→ df = pd.read_csv("adult.csv")

→ df.info()

→ df.describe()

→ df["age"]

→ df["work class"]

→ df["education"]

→ df["marital status"]

→ df["relationship"]

→ df["native country"]

→ df["income"]

Code:

```
import pandas as pd  
  
import numpy as np  
  
# Load dataset  
  
df=pd.read_csv("data.csv")  
  
print(df.head())  
  
# Check missing values  
  
print(df.isnull().sum())  
  
  
# Drop rows with missing values  
  
df_cleaned = df.dropna()  
  
# Or fill missing values with mean/median  
  
df['Age'].fillna(df['Age'].mean(), inplace=True)  
  
df['Salary'].fillna(df['Salary'].median(), inplace=True)  
  
  
# For nominal categories  
  
df = pd.get_dummies(df, columns=['Gender', 'Country'], drop_first=True)  
# For ordinal categories  
  
from sklearn.preprocessing import OrdinalEncoder  
  
encoder = OrdinalEncoder()  
  
df[['Education_Level']] = encoder.fit_transform(df[['Education_Level']])  
from sklearn.preprocessing import StandardScaler, MinMaxScaler  
  
# Standardization (Z-score)  
  
scaler = StandardScaler()
```

```

df[['Age', 'Salary']] = scaler.fit_transform(df[['Age', 'Salary']])

# Min-Max Normalization

minmax = MinMaxScaler()

df[['Age', 'Salary']] = minmax.fit_transform(df[['Age', 'Salary']])

# Using IQR method

Q1 = df['Salary'].quantile(0.25)

Q3 = df['Salary'].quantile(0.75)

IQR = Q3 - Q1

df = df[(df['Salary'] >= Q1 - 1.5*IQR) & (df['Salary'] <= Q3 + 1.5*IQR)]

df['Age_Salary_Ratio'] = df['Age'] / df['Salary']

# Drop irrelevant columns

df.drop(['User_ID', 'Name'], axis=1, inplace=True)

# Correlation-based filtering

correlation_matrix = df.corr()

print(correlation_matrix)

from sklearn.model_selection import train_test_split
X = df.drop('Purchased', axis=1)
y = df['Purchased']

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

```

Program 3

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

Screenshot:

Date: 18/03/2025 LAB - III Probability Distribution

i) $a_0 = -5 \quad a_1 = 0.8$

(ii) Logistic regression equation

$$\rightarrow P(n) = \frac{1}{1 + e^{-(a_0 + a_1 n)}} = \frac{1}{1 + e^{(-5 + 0.8n)}}$$

(iii) Calculate the probability that a student who studies for 7 hours will pass

$$\rightarrow n = 7 \quad P(n) = \frac{1}{1 + e^{(-5 + 0.8(7))}} = 0.6457$$

(iv) Determine the predicted class (Pass / Fail) for the student based on threshold of 0.5

$$\begin{aligned} \rightarrow P(n) &= 0.6457 \\ P(n) &= 0.5 \end{aligned}$$

$y = \begin{cases} 1 & \text{if } P(n) > 0.5 \\ 0 & \text{otherwise} \end{cases}$

Thus $y = 1$ (Pass)

2) Consider $\mathbf{z} = [2, 1, 0]$ for three classes. Apply Software function to find probability values of 3 classes

$$\rightarrow \text{Softmax}(\mathbf{z}_k) = \frac{e^{z_k}}{\sum_{j=1}^3 e^{z_j}}$$

$$\text{Softmax}(z_1) = \frac{e^2}{e^2 + e^1 + e^0} = \frac{7.38905}{11.107337} = 0.6652$$

$$\text{Softmax}(z_2) = \frac{e^1}{e^2 + e^1 + e^0} = \frac{2.71828}{11.107337} = 0.244$$

$$\text{Softmax}(z_3) = \frac{e^0}{e^2 + e^1 + e^0} = \frac{1.091}{11.107337} = 0.091$$

The probabilities of the 3 classes are approximately 66.5%, 24.4%, 9.1%.

Scikit-learn (train_X, test_X), fit, labor

Logistic Regression (lab-executed codes)

*Binary.py

```
import pandas as pd
from matplotlib import pyplot as plt
df = pd.read_csv("insurance-data.csv")
df.head()
plt.scatter(df.age, df.bought_insurance, marker='+',
            color='red')
from sklearn.model_selection import train_test_
split
X_train, X_test, y_train, y_test = train_test_
split(df[['age']], df.bought_insurance, train_size=
0.9, random_state=10)
X_train.shape
```

X-test

```
from sklearn.linear_model import LogisticRegression  
model = LogisticRegression()  
model.fit(X_train, y_train)
```

X-test

y-test

y-predicted = model.predict(X-test)

y-predicted

model.score(X-test, y-test)

model.predict_proba(X-test)

y-predicted = model.predict([[60]])

y-predicted

model.coef_

model.intercept_

import math

def sigmoid(x):

return 1 / (1 + math.exp(-x))

def prediction_function(age):

z = 0.127 * age - 4.973

y = sigmoid(z)

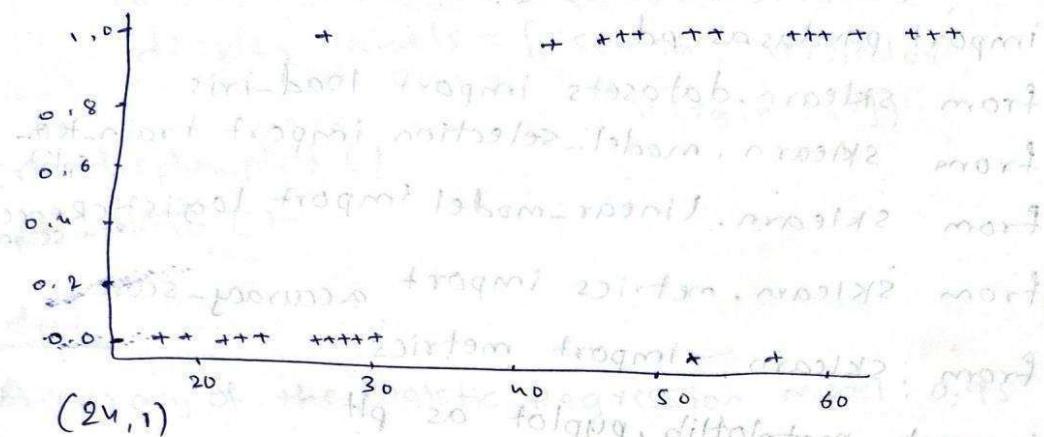
return y

age = 35

prediction_function(age)

DP

Age vs Insurance Purchase



(24, 1)

age

7 60

5 56

18 19

7 1

5 1

18 0

Name: bought_insurance

model Accuracy: 1.0

Predicted Probabilities: [0.06430655 0.93529345]

[0.10327333 0.89672667]

[0.92775258 0.07224742]

Prediction for age 60: [1 - 0.12740655 + 0.12740655 * 60 = 1.0]

model coefficient: [0.12740655]

model intercept: [-4.97339194]

Prediction for age 35: 0.3309834769552775

Worst predicted ton insurance but got 10 previous fit fairly

(18 - 0.12740655 * 35 + 0.12740655 * 35 = 0.12740655)

gives -0.12740655, which is not a valid value
(less than 0 or greater than 1)

*Multi-class.py

```
import pandas as pd
from sklearn.datasets import load_iris
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score
from sklearn import metrics
import matplotlib.pyplot as plt

iris = pd.read_csv("iris.csv")
iris.head()

X = iris.drop('species', axis='columns')
y = iris.species

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

model = LogisticRegression(multi_class='multinomial')
model.fit(X_train, y_train)

y_pred = model.predict(X_test)

accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy of the Multinomial Logistic Regression model on the test set : {accuracy:.2f}")

confusion_matrix = metrics.confusion_matrix(
    y_test, y_pred)
```

`cm_display = metrics.ConfusionMatrixDisplay(`

`confusion_matrix=confusion_matrix,`

`display_labels=["setosa", "versicolor",`

`"virginica"]])`

`cm_display.plot()`

`plt.show()`

Output

Accuracy of the Logistic Regression model: 0.95

Confusion matrix

	mammal	bird	Reptile	fish	Amphibian	Bug	Invertebrate	
mammal	12	0	0	0	0	0	0	12
bird	0	2	0	0	0	0	0	10
Reptile	0	0	0	0	0	0	1	8
fish	0	0	0	2	0	0	0	6
Amphibian	0	0	0	0	12	3	0	4
Bug	0	0	0	0	0	1	0	2
Invertebrate	0	0	0	0	0	0	0	0

Mammal Bird Reptile Fish Amphibian Bug Invertebrate

Predicted class

Get

Answering the Questions

- i) For dataset file "HR-comma-sep.csv"
 - i) which variable did you identify as having a direct and clear impact on employee retention? why?

→ The variables with a direct & clear impact on employee retention are

- Satisfaction level
- Average monthly hrs
- Time spent at the company
- Work accident
- Promotion in the last 5 years

- ii) what was the accuracy of your logistic regression model? Do you think this is a good accuracy? why or why not?

→ The accuracy of L-R model is 75.9%.

- This is reasonable accuracy but not perfect
- The model shows a relatively high no. of false negatives, indicating that some employees who left were incorrectly predicted to stay.

Q) For zoo dataset

- i) Did you perform any data preprocessing steps? If yes, what were they, & why were they necessary?

→ Preprocessing steps performed.

- Merged the zoo datasets using class-type as the key.
- Removed non-numeric columns before standardizing the features.
- Scaled the features using Standard Scalar for better model performance.
- Split the data into training & testing sets

(ii) Were there any missing or inconsistent values in the dataset? How did you handle them?

→ No missing or inconsistent values were found, so no special handling was necessary.

(iii) What does the confusion matrix tell you about the performance of your model?

→ The model achieved 93.15% accuracy, showing good performance.

(iv) Which class type were most frequently misclassified, why do you think this happened?

→ The class is Class 3, where all samples were incorrectly classified.

• This may have occurred due to the small number of samples, making it difficult for the model to learn its characteristics.

• Class 7 was also partially misclassified, likely due to overlapping features with Class 3.

$$P(86.0 = 0) = P(12.0 = 1) = P(86.0 = 1) = 0.2$$

$$\Theta = [e+o] \quad \text{mm}^2$$

$$\Theta = [e, o] \quad \text{mm}^2$$

$$P(86.0 = 0) = P(12.0 = 0) = P(86.0 = 1) = 0.1$$

Code:

```
import pandas as pd

import numpy as np

from graphviz import Digraph

# Calculate Entropy

def entropy(data):

    class_probabilities = data.iloc[:, -1].value_counts(normalize=True)

    return -np.sum(class_probabilities * np.log2(class_probabilities))

# Calculate Information Gain

def information_gain(data, feature):

    total_entropy = entropy(data)

    feature_values = data[feature].unique()

    weighted_entropy = 0

    for value in feature_values:

        subset = data[data[feature] == value]

        weighted_entropy += (len(subset) / len(data)) * entropy(subset)

    return total_entropy - weighted_entropy

# Find the best feature to split the data

def best_feature(data):

    features = data.columns[:-1] # Exclude the target column

    gains = {feature: information_gain(data, feature) for feature in features}

    return max(gains, key=gains.get)

# Create the decision tree

def id3(data, features=None):

    if len(data.iloc[:, -1].unique()) == 1: # All data points belong to the same class
```

```

return data.iloc[:, -1].iloc[0]

if len(features) == 0: # No more features to split on

    return data.iloc[:, -1].mode()[0]

best = best_feature(data)

tree = {best: {}}

new_features = features.copy()

new_features.remove(best)

for value in data[best].unique():

    subset = data[data[best] == value]

    tree[best][value] = id3(subset, new_features)

return tree

# Function to classify new examples based on the decision tree

def classify(tree, example):

    if not isinstance(tree, dict):

        return tree

    feature = list(tree.keys())[0]

    value = example[feature]

    return classify(tree[feature][value], example)

# Function to visualize the decision tree using Graphviz

def create_tree_diagram(tree, dot=None, parent_name="Root", parent_value=""):

    if dot is None:

        dot = Digraph(format="png", engine="dot")

```

```

if isinstance(tree, dict): # Tree node

    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: # Leaf node

    dot.node(parent_name + "_class", f"Class: {tree}")

    dot.edge(parent_name, parent_name + "_class", label="Leaf")

return dot

# Example usage

data = pd.DataFrame({

    'Outlook': ['Sunny', 'Sunny', 'Overcast', 'Rain', 'Rain', 'Rain', 'Overcast', 'Sunny', 'Sunny', 'Rain',
    'Sunny', 'Overcast', 'Overcast', 'Rain'],

    'Temperature': ['Hot', 'Hot', 'Hot', 'Mild', 'Cool', 'Cool', 'Mild', 'Cool', 'Mild', 'Mild', 'Mild', 'Mild', 'Hot',
    'Mild'],

    'Humidity': ['High', 'High', 'High', 'High', 'High', 'Low', 'Low', 'High', 'Low', 'Low', 'Low', 'High', 'Low',
    'High'],

    'Wind': ['Weak', 'Strong', 'Weak', 'Weak', 'Weak', 'Weak', 'Strong', 'Weak', 'Weak', 'Strong', 'Strong',
    'Weak', 'Strong', 'Weak'],

    'PlayTennis': ['No', 'No', 'Yes', 'Yes', 'Yes', 'No', 'Yes', 'No', 'Yes', 'Yes', 'Yes', 'Yes', 'Yes', 'Yes', 'No']
})

```

```
)  
  
# Train the decision tree  
  
tree = id3(data, features=list(data.columns[:-1]))  
  
print("Decision Tree:", tree)  
  
# Classify a new example  
  
example = {'Outlook': 'Sunny', 'Temperature': 'Cool', 'Humidity': 'Low', 'Wind': 'Strong'}  
  
prediction = classify(tree, example)  
  
print("Prediction for the example:", prediction)  
  
  
# Visualize the decision tree  
  
dot = create_tree_diagram(tree)  
  
dot.render("decision_tree", view=True) # This will generate and open the tree diagram
```

Program 4

Implement Linear and Multi-Linear Regression algorithm for appropriate dataset

Screenshot:

LAB-4

* Linear Regression

predict Canada's Pay co-operative program (1)

week	Sales
1	2
2	4
3	5
4	9

$X^T = [1 \ 2 \ 3 \ 4]$

$y^T = [2 \ 4 \ 5 \ 9]$

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

$y = \begin{bmatrix} 2 \\ 4 \\ 5 \\ 9 \end{bmatrix}$

$\beta = (X^T X)^{-1} X^T y$

compute $X^T X = \begin{bmatrix} 1 & 2 & 3 & 4 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 2 & 4 \\ 3 & 5 \\ 4 & 9 \end{bmatrix} = \begin{bmatrix} 30 & 10 \\ 10 & 30 \end{bmatrix}$

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

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

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

$\beta = \begin{bmatrix} 0.5 \\ 2.2 \end{bmatrix}$ Intercept
slope

$y = \beta_0 + \beta_1 x + \epsilon$

$y = -0.5 + 2.2x$

At $x = 5$
 $y = -0.5 + 2.2(5)$
 $\Rightarrow y = 10.5$

Steps of code (algo)

- 1) import libraries
- 2) import data-distribution
- 3) Analyze data-distribution
- 4) Distribution flat-visualization
- 5) Relationship b/w variables
- 6) split the data
- 7) Train the data model
- 8) Predict the result
- 9) visualize prediction
- 10) check value of co-efficients & intercept

Linear Regression

- * Predict Canada's Per capita income in year 2020.

```
import pandas as pd
import numpy as np
from sklearn import LinearModel
import matplotlib.pyplot as plt
df = pd.read_csv('canada-per-capita-income.csv')
print("Column Names:", df.columns)
df.rename(columns = {'per capita income (US$)': 'Income'}, inplace = True)
df.dropna(inplace = True)
plt.xlabel('Year')
plt.ylabel('Income')
plt.scatter(df['Year'], df['Income'], color='red',
            marker='+')
x = df[['Year']]
y = df['Income']
reg = LinearModel.LinearRegression()
reg.fit(x, y)
year_2020 = reg.predict([[2020]])[0]
year_2027 = reg.predict([[2027]])[0]
print(f'Year = 2020, Predicted Income = ${year_2020:.2f}')
print(f'Year = 2027, Predicted Income = ${year_2027:.2f}')
print("Coefficient (m): {reg.coef_[0]:.2f}")
print("-----")
```

```
print(f'Intercept(b) : {reg.intercept_}, {2f y}')  
print(" - _____")  
plt.plot(df['year'], reg.predict(x), color='blue')  
plt.show()
```

output

Column Names : Index(['year', 'per capita Income (US \$)'), dtype='object')

Year = 2020, Predicted Income = \$41,288.69

Year = 2027, Predicted Income = \$67,087.95

Coefficient (m) : 828.47

Intercept (b) : -1,632,210.76

Multiple Linear Regression

Hiring

```
import pandas as pd  
import numpy as np  
from sklearn import linear_model  
df = pd.read_csv('hiring.csv')  
print(df.isnull().sum())  
df['test-score(out of 10)'].fillna(df['test-score(out of 10)'].median(), inplace=True)  
df['experience'] = df['experience'].astype(str)
```

```

df_encoded = pd.get_dummies(df, columns=['experience'], drop_first=True)

x = df_encoded.drop('salary($)', axis='columns')
y = df_encoded['salary($)']

reg = linear_model.LinearRegression()
reg.fit(x, y)

def predict_salary(experience, test_score, interview_score):
    exp_col = 'experience' if experience != 'exp' else None
    input_data = {col: 0 for col in x.columns}
    if exp_col in input_data:
        input_data[exp_col] = 1
    input_data['test-score (out of 10)'] = test_score
    input_data['interview-score (out of 10)'] = interview_score
    input_df = pd.DataFrame([input_data])
    predicted_salary = reg.predict(input_df)[0]
    return f'Predicted Salary: ${predicted_salary:.2f}'

print(predict_salary('twelve', 10, 10))
print(predict_salary('two', 9.6))

```

Output

experience	2
test-score (out of 10)	1
interview-score (out of 10)	0
salary (\$)	0

Predicted Salary: \$ 67,803.92

Predicted Salary: \$ 63,509.80

Questions

* Linear_regression_housing_arcanprice.py
Linear regression, KNN predict good/bad_16 = 8

Output

saving_housing-area-price, CSV

predicted price for 3300 sq ft: \$628,715.75342465

$$\text{Coefficient (m)} = 135.78767123287673$$

Intercept (b) : 180616.43835616432

Manual price calculation for 733000 sq ft: 628715.
(example: $733000 \times 0.008 = 58640$) 7334246575

Predicted Price for 500 sq ft is £859554.794820

~~multiple-1r-homeprices.py~~

out put

Coefficients : [112.06244194 23388.88002779]

Wetlands & photo 2 habitat - 3231-71790863)

Intercept: 221323.00186540396

Predicted price for a home with 3000 sqr ft

~~3 bedrooms, 90 years old; 498408~~

2815803067

Manual Price calculation: wages

251574024

(8) *provis*

Linear-Regression → (Salary.csv)

Output

YearsExperience	Salary
0	39891
1	43525
2	37731
3	46205
4	39363

YearsExperience = 12, Predicted Income = \$139,980.89

$$\text{Coefficient } (m) = 9,398.64$$

$$\text{Intercept } (b) = 27,197.20$$

multiple-Regression (1000-companies.csv)

Output

Predicted Profit: \$510,570.99

Questions (revised)

- 1) Considering 3 data .csv files, did you perform any data preprocessing steps? If yes, why?
 - Canada-per-capita-income.csv → no encoding or scaling was needed. Checked for missing values & handled if any were found.
 - Salary.csv → verified for missing values, handled if present.
no encoding was needed.
 - Hiring → no scaling or encoding was needed.
- 2) For Canada per capita income did you visualize the regression line along with the data points? What does the plot tell you about the relationship below you & per capita income (as against time data points).
The plot indicated a strong linear relationship.
- 3) For hiring, what is the predicted salary for a candidate with 12 yrs of exp, 10 test score & 10 interview score?
→ \$82000.00
- 4) For 1000 companies, did you encode categorical variables? If yes, how? Did you scale the features? If yes, why?

- categorical encoding → The 'state' variable was one-hot encoded to convert categorical to numerical
- feature scaling → If features had varying magnitudes, scaling was applied to improve model efficiency.

Code:

Linear Regression

```
import pandas as pd  
  
df=pd.read_csv("/content/tvmarketing.csv")  
  
df  
  
# Visualise the relationship between the features and the response using scatterplots  
df.plot(x='TV',y='Sales',kind='scatter')  
  
from sklearn.model_selection import train_test_split  
  
x_train, x_test, y_train, y_test = train_test_split(df['TV'], df['Sales'], test_size=0.2, random_state=42)  
  
from sklearn.linear_model import LinearRegression model = LinearRegression()  
model.fit(x_train.values.reshape(-1, 1), y_train) y_train  
model.coef_  
  
model.intercept_
```

Multilinear Regression

```
import pandas as pd  
  
# Step 2 : import data  
  
house = pd.read_csv('https://github.com/YBIFoundation/Dataset/raw/main/Boston.csv')  
  
# display first 5 rows house.head()  
  
y = house['MEDV']
```

```
X = house.drop(['MEDV'],axis=1)

# Step 4 : train test split
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X,y, train_size=0.7, random_state=2529)

# Step 5 : select model
from sklearn.linear_model import LinearRegression
model = LinearRegression()

# Step 6 : train or fit model
model.fit(X_train,y_train)
model.intercept_

model.coef_
```

Program 5

Build Logistic Regression Model for a given dataset

Screenshot:

Lab - 5 decision Tree

instance a_1, a_2 new classification

- 1 Hot High N
- 2 Hot High N
- 3 Cool High N
- 4 Hot High N
- 5 Hot Normal Y

$$Empty(S) = -\frac{4}{6} \log_2 \left(\frac{4}{6} \right) - \frac{2}{3} \log_2 \left(\frac{1}{3} \right)$$

$$= 0.7219$$

$$S[1, 3] = -\frac{1}{4} \log_2 \left(\frac{1}{4} \right) - \frac{3}{4} \log_2 \left(\frac{3}{4} \right)$$

$$= 0.5113$$

$$S[0+1] = 0$$

$$cool$$

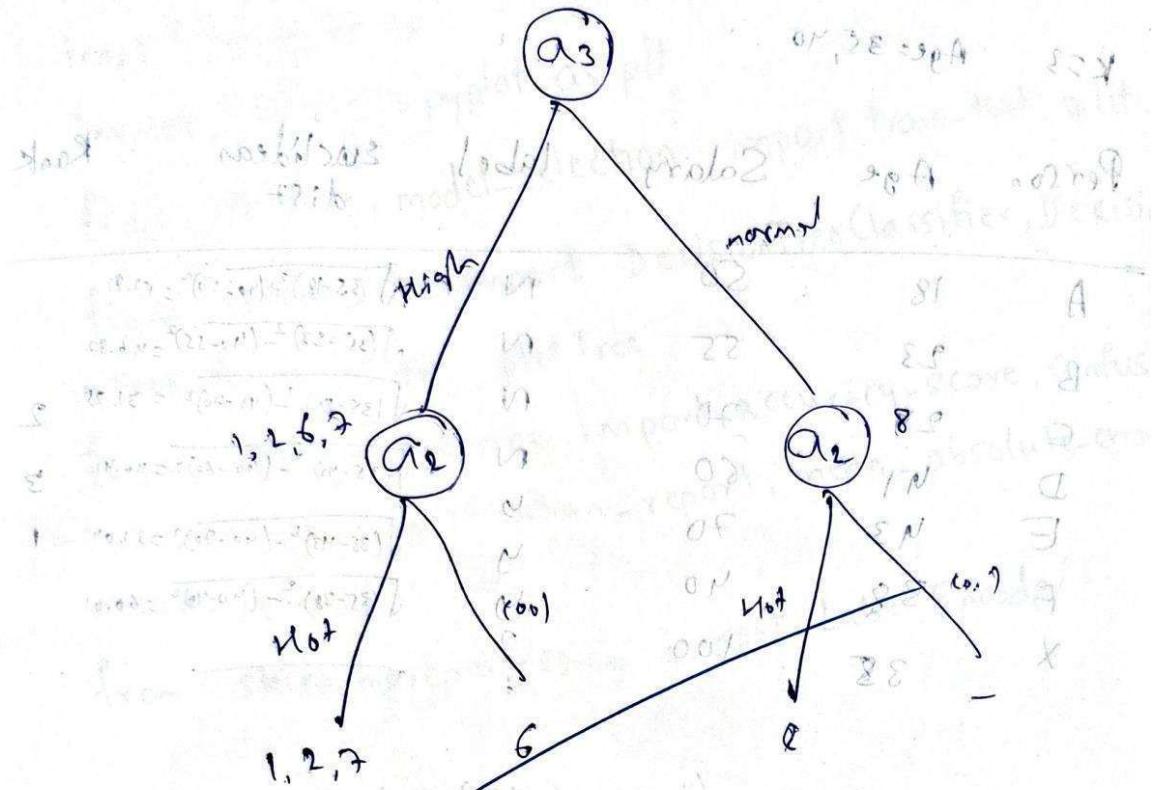
$$Gain(S, a.) = 0.7219 - \frac{4}{6} \cdot 0.5113 - 0 = 0.3252$$

$$\cancel{a_1} S_{all} [0+4] = 0$$

$$S_{normal} [1, 0] = 0$$

$$Gain(1, a) = 0.7219 - 0 = 0.7219$$

(i) a_3 is root node \Rightarrow since $l_p(s_0, a_3)$ is higher



Tab-5

```
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, plotTree
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("iris.csv")
drug = pd.read_csv("read.csv")
petrol = pd.read_csv("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("Accuracy", accuracy_score(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:

 print(col)

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("Accuracy", accuracy_score(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)

print("Mean Absolute Error", Mean_Absolute_Error)

(y_test, y_pred)

print("Mean Squared Error", Mean_Squared_Error)

(y_test, y_pred))

Output

Iris dataset

Accuracy : 1.0

Drug dataset

Accuracy : 1.0

Petrol dataset

Mean Absolute Error: 1.6

Mean Squared Error: 17609.8

Root Mean Squared Error: 13270.9

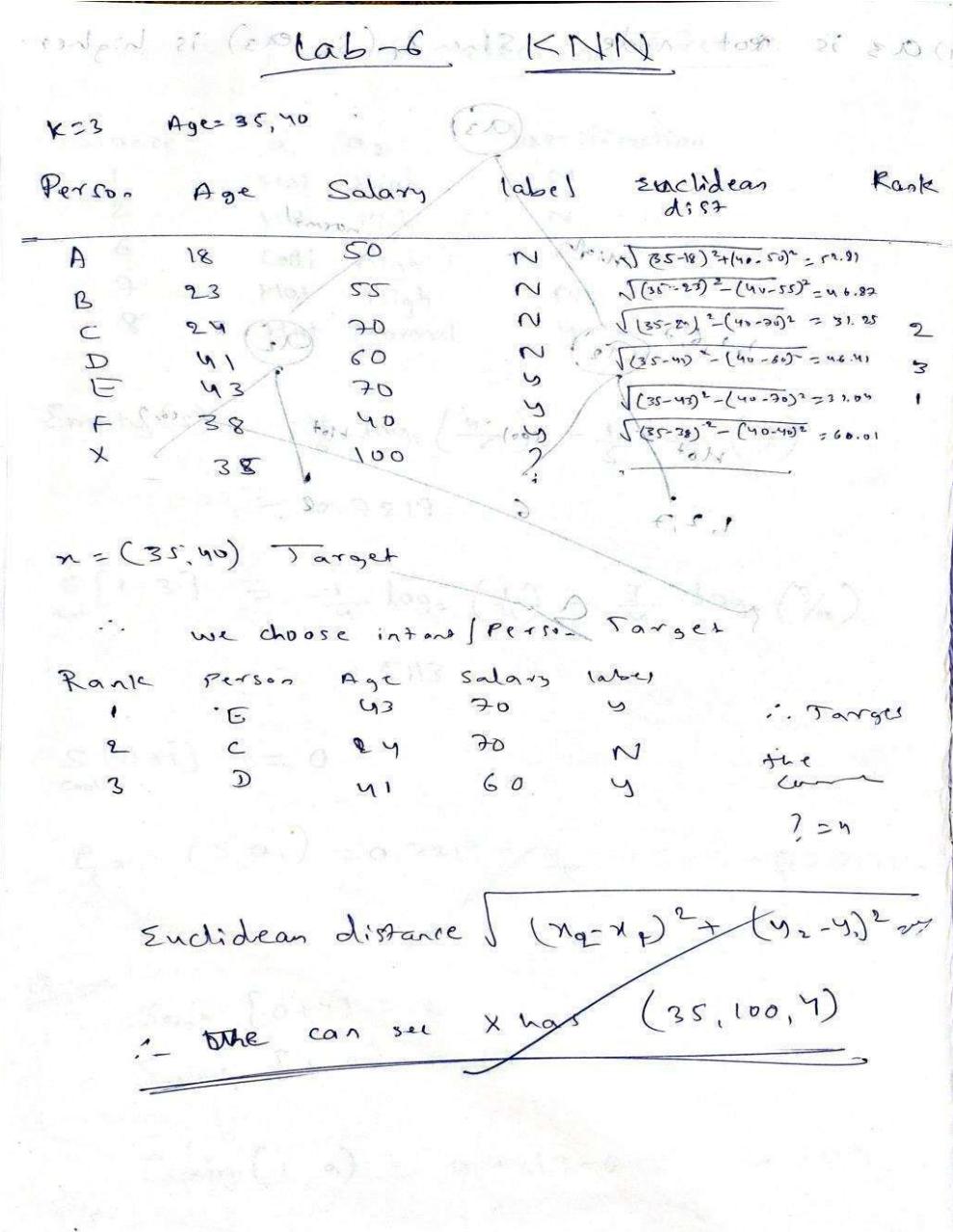
Code:

```
from sklearn.linear_model import LogisticRegression  
from sklearn.datasets import load_iris  
from sklearn.model_selection import train_test_split  
from sklearn.metrics import accuracy_score  
  
# Load sample dataset (binary classification - Iris with only 2 classes)  
iris = load_iris()  
X = iris.data[iris.target != 2]  
y = iris.target[iris.target != 2]  
  
# Train/Test split  
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)  
  
# Logistic Regression model  
model = LogisticRegression()  
model.fit(X_train, y_train)  
  
# Predict and evaluate  
y_pred = model.predict(X_test)  
print("Accuracy:", accuracy_score(y_test, y_pred))
```

Program 6

Build KNN Classification model for a given dataset

Screenshot:



#!/usr/bin/python3
 # Lab 6 - KNN
 # imports are same as previous program

```

heart_df = pd.read_csv('heart.csv')
print(heart_df.head())
x = heart_df.drop(['target'])
y = heart_df['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)
x_test_scaled = scaler.transform(x_test)

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)
    score = accuracy_score(y_test, y_pred)
    if score > best_score:
        best_score = score
        best_k = k

print(f"Best k value: {best_k} with accuracy: {best_score * 100 :.2f} %")
  
```

```

best_knn = KNeighborsClassifier(n_neighbors=best_k)
best_knn.fit(x_train_scaled, y_train)

y_pred_best = best_knn.predict(x_test_scaled)

accuracy = accuracy_score(y_test, y_pred_best)

print(f"Accuracy on test data: {accuracy * 100:.2f} %")

fig, ax = plt.subplots(figsize=(6, 4))

cm = confusion_matrix(y_test, y_pred_best)

sns.heatmap(cm, annot=True, fmt="d", cmap=plt.cm.Blues,
            xticklabels=best_knn.classes_,
            yticklabels=best_knn.classes_)

plt.show()

```

Output

best k value 7 with accuracy 91.80%
 Accuracy on test data: 91.80%

Actual	0	1
0	61032	1291
1	1291	61032

Confusion Matrix: $\begin{pmatrix} 61032 & 1291 \\ 1291 & 61032 \end{pmatrix}$

Classification Report:
 precision: 0.9180
 recall: 0.9180
 f1-score: 0.9180
 support: 20000

Support Vector Machine (SVM) Model

Coefficients (w): array([0.00000000e+00, -2.98747188e-01])

Code:

KNN

```
import numpy as np
from collections import Counter

class KNN:
    def __init__(self, k=3): self.k = k

    def fit(self, X, y):
        self.X_train = np.array(X)
        self.y_train = np.array(y)

    def euclidean_distance(self, x1, x2):
        return np.sqrt(np.sum((x1 - x2) ** 2))

    def predict(self, X):
        predictions = [self._predict(x) for x in X]
        return np.array(predictions)

    def _predict(self, x):
        # Compute distances to all training points
        distances = [self.euclidean_distance(x, x_train) for x_train in self.X_train]

        # Get indices of k nearest neighbors
        k_indices = np.argsort(distances)[:self.k]

        # Get the labels of those neighbors
        k_nearest_labels = [self.y_train[i] for i in k_indices]
```

```
# Return the most common label

most_common = Counter(k_nearest_labels).most_common(1)

return most_common[0][0]

# Sample dataset (like a mini version of Iris)

X_train = [[1, 2], [2, 3], [3, 1], [6, 5], [7, 7], [8, 6]]

y_train = [0, 0, 0, 1, 1, 1]

# Test data

X_test = [[5, 5], [1, 1]]

# Using the KNN model

knn = KNN(k=3)

knn.fit(X_train, y_train)

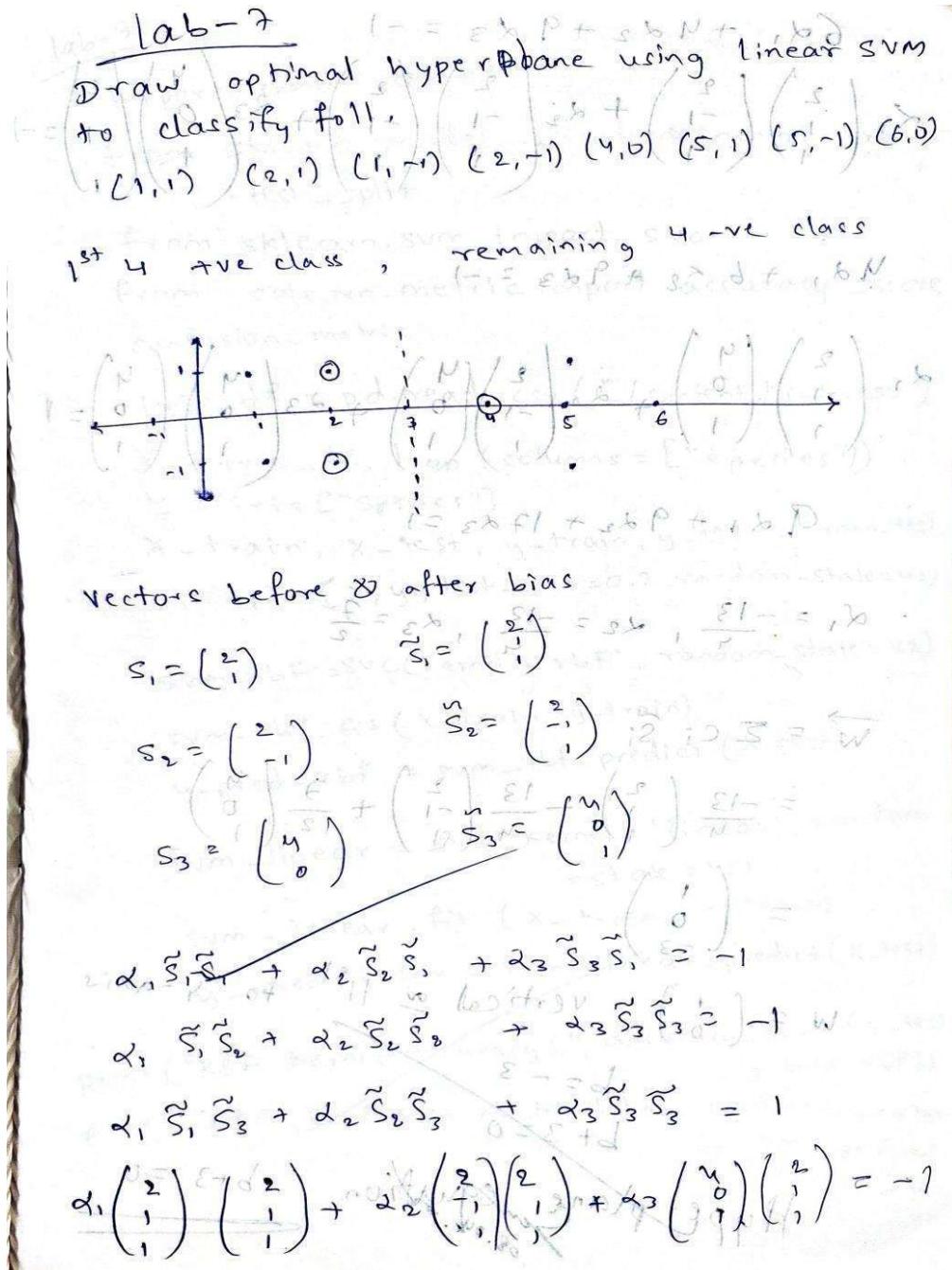
predictions = knn.predict(X_test)

print("Predictions:", predictions)
```

Program 7

Build Support vector machine model for a given dataset

Screenshot:



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

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

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

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

$$9\alpha_1 + 6\alpha_2 + 17\alpha_3 = 1$$

$$\alpha_1 = -\frac{13}{4}, \alpha_2 = \frac{-13}{4}, \alpha_3 = \frac{7}{2}$$

$$\vec{w} = \sum c_i \vec{s}_i$$

$$= -\frac{13}{4} \begin{pmatrix} 2 \\ -1 \\ 1 \end{pmatrix} + \frac{-13}{4} \begin{pmatrix} 2 \\ 0 \\ 1 \end{pmatrix} + \frac{7}{2} \begin{pmatrix} 4 \\ 0 \\ 1 \end{pmatrix}$$

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

$$w = \begin{pmatrix} 1 \\ 0 \\ -3 \end{pmatrix} \text{ vertical } \approx 11 \text{ to } y\text{-axis}$$

$$b = -3$$

$$b+3=0$$

~~Hyperplane equation~~

$$\cancel{\begin{pmatrix} 1 \\ 0 \\ -3 \end{pmatrix}} \cancel{b} = 0$$

lab-7 - ~~process A~~ (unidirectional) taking

```

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
iris_df = pd.read_csv("A\\content\\iris.csv")
X = iris_df.drop(columns=[ "species"])
y = iris["Species"]
X_train, X_test, y_train, y_test = train-test-
split(X, y, test_size=0.2, random_state=42)
SVM-RBF = SVC(kernel='rbf', random_state=42)
SVM-RBF.fit(X_train, y_train)
y_pred_rbf = SVM-RBF.predict(X_test)
SVM-Linear = SVC(kernel='linear', random-
state=42)
SVM-Linear.fit(X_train, y_train)
y_pred_linear = SVM-Linear.predict(X_test)
print("RBF Kernel Accuracy:", accuracy_score(y_test,
y_pred_rbf))
print("RBF Confusion matrix:\n", confusion_matrix(
y_test, y_pred_rbf))

```

print("Linear kernel Accuracy:", accuracy_score(
y-test, y-pred-linear))

print("Linear Confusion Matrix:\n", confusion_matrix(y-test, y-pred-linear))

0/8
RBF Kernel Accuracy: 0.6926511739153448

RBF Confusion Matrix: 0.6926511739153448

$\begin{bmatrix} [10 & 0 & 0] \\ [0 & 9 & 0] \\ [0 & 0 & 1] \end{bmatrix}$

(89+0)/89 = 0.988675134594683

Linear Kernel Accuracy: 1.0

Linear Confusion Matrix:

$\begin{bmatrix} [10 & 0 & 0] \\ [0 & 9 & 0] \\ [0 & 0 & 1] \end{bmatrix}$

(89+0)/89 = 1.0

$\begin{bmatrix} [10 & 0 & 0] \\ [0 & 9 & 0] \\ [0 & 0 & 1] \end{bmatrix}$

(89+0)/89 = 1.0

$\begin{bmatrix} [10 & 0 & 0] \\ [0 & 9 & 0] \\ [0 & 0 & 1] \end{bmatrix}$

$\begin{bmatrix} [10 & 0 & 0] \\ [0 & 9 & 0] \\ [0 & 0 & 1] \end{bmatrix}$

Q Answer the following questions:

* for "Iris.csv" dataset

1) what is the accuracy score of the classifier using

the linear kernel & RBF Kernel.

• Linear $\rightarrow 1.0$ {Accuracy}

• RBF $\rightarrow 1.0$

Both performed equally well with

100% accuracy! This is likely b/c the Iris dataset is relatively small & well separated,

making it easy for both kernels to classify

correctly.

2) which kernel (RBF or linear) gave better performance?
why do you think this is the case?

* for "letter-recognition.csv" dataset

1) Present & Interpret the confusion matrix. Are there any specific letters that are frequently confused with others?

2) what is the AUC (Area under the Curve) score, and how does it reflect the model performance?

3) How does the performance of the SVM model of this dataset compare to its performance on the IRIS dataset?

Ans

1) Accuracy : 0.9315 (93.15%)

Top 5 most confused letters are

$P \rightarrow S$ $B \rightarrow S$ $M \rightarrow U$ $D \rightarrow U$ $V \rightarrow U$

These letters might have 11th features, leading to confusion.

② ACU Score $\rightarrow 0.99817858$ almost perfect classification performance

③ Comparison to Iris dataset

\rightarrow SVM performed flawlessly on Iris, but had

minor misclassifications on the letter dataset

\rightarrow The letter dataset is much larger & has most complex patterns, making classification

harder than the simple 3-class Iris dataset.

\rightarrow Geometric view of SVM decision boundary

(SVM) \rightarrow max margin classifier

Max margin classifier \rightarrow max distance from classes

Margin = $\frac{2}{\|w\|}$

Max margin classifier \rightarrow max $\frac{2}{\|w\|}$

(SVM) \rightarrow max margin classifier

Max margin classifier \rightarrow max $\frac{2}{\|w\|}$

Max margin classifier \rightarrow max $\frac{2}{\|w\|}$

Code:

```
from sklearn import datasets  
  
from sklearn.model_selection import train_test_split  
  
from sklearn.svm import SVC  
  
import matplotlib.pyplot as plt  
  
from sklearn.decomposition import PCA  
  
  
# Load dataset  
  
iris = datasets.load_iris()  
  
X = iris.data  
  
y = iris.target  
  
# For binary classification (class 0 vs 1)  
  
X = X[y != 2]  
  
y = y[y != 2]  
  
# Train-test split  
  
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)  
  
  
# Train SVM  
  
clf = SVC(kernel='linear') # Try 'rbf', 'poly', etc.  
  
clf.fit(X_train, y_train)  
  
  
# Accuracy  
  
print("Test Accuracy:", clf.score(X_test, y_test))
```

Program 8

Implement Random Forest ensemble method on a given dataset

Screenshot:

- of lab-8 To understand working of ensemble learning
Difference between decision tree & random forest classifier
- 1) Difference between decision tree & random forest classifier
- Decision tree requires overfitting the training data → Random forest is a collection of decision trees trained on different subsets of data.
 - A simple tree model that splits data into multiple partitioning branches. → It is a collection of decision trees, usually trained on different subsets of data.
 - Easy to understand and visualize → More complex and harder to interpret.
 - It is faster to train → slower to train
 - Prone to overfitting → generally higher accuracy
 - mainly on noisy data
 - Low variance, hence more robust
 - low bias, high variance
 - Low bias, high variance
 - Low bias, high variance
 - Low bias, high variance
- 2) Random forest classifier parameters
- Some key parameters include:
1. n_estimators: Number of decision trees in the forest.
 2. max_depth: Maximum depth of each decision tree.
 3. min_samples_split: Minimum samples required to be at a split node.
 4. min_samples_leaf: Minimum samples required to be at a leaf node.

5. max_features : Maximum number of features to consider at each split.

③ Random Forest Algorithm

Here's a simplified algorithm:

1. Bootstrap Sampling : Create multiple subsets of the training data.

2. Decision tree training : Train a decision tree on each subset

3. Prediction : Each tree predicts the target variable.

4. Voting : Combine predictions from all trees

* to produce the final output

lab-9

① Boosting : It is an ensemble technique that combines weak models to form a strong model by focusing on correcting the errors of previous models

② Discuss all the parameters of AdaBoostClassifier()

→ Parameters :

1. base_estimator : The weak model (default is a decision tree)

2. n_estimators : No. of models to train (default is 50)

3. learning_rate : Controls model contribution (default is 1.0)

4. algorithm : Boosting method ('SAMME' or 'SAMME.R') default is SAMME.R

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 sample dataset  
iris = load_iris()  
X, y = iris.data, iris.target  
  
  
# Train/test split  
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)  
  
# Initialize Random Forest  
rf = RandomForestClassifier(n_estimators=100, random_state=42)  
rf.fit(X_train, y_train)  
  
# Predict and evaluate  
y_pred = rf.predict(X_test)  
print("Accuracy:", accuracy_score(y_test, y_pred))
```

Program 9

Implement Boosting ensemble method on a given dataset

Screenshot:

Lab-9

① Boosting: It is an ensemble technique that combines weak models to form a strong model by focusing on correcting the errors of previous models

② Discuss all the parameters of AdaBoostClassifier()

→ Parameters :

1. base-estimator: The weak model (default is a decision tree)
2. n_estimators: No. of models to train (default is 50)
3. learning_rate: Controls model contribution (default is 1.0)
4. algorithm: Boosting method ('SAMME' or 'SAMME.R')
default is SAMME

5. `random_state`: ensures reproducibility (optional)

AdaBoost Algorithm: an ensemble of weak learners

1. Assign equal wts. to all data pts.
2. Train a weak learner on weighted data
3. Increase wts. of misclassified pts.
4. Repeat for several iterations
5. Combine models into a final strong learner

topic final - 2023V03/3 2023/01/01

Code:

```
from sklearn.ensemble import AdaBoostClassifier  
  
from sklearn.datasets import load_iris  
  
from sklearn.model_selection import train_test_split  
  
from sklearn.metrics import accuracy_score  
  
  
# Load Iris dataset  
  
iris = load_iris()  
  
X, y = iris.data, iris.target  
  
  
# For AdaBoost, we'll use binary classification  
  
# Convert to binary (setosa vs. not-setosa)  
  
y = (y == 0).astype(int)  
  
  
# Split data  
  
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)  
  
# Train AdaBoost  
  
model = AdaBoostClassifier(n_estimators=50, learning_rate=1.0, random_state=42)  
model.fit(X_train, y_train)  
  
  
# Predict and evaluate  
  
y_pred = model.predict(X_test)  
  
print("AdaBoost Accuracy (sklearn):", accuracy_score(y_test, y_pred))
```

Program 10

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

Screenshot:

Lab 10

① Algorithm of K-Means

- choose the no. of clusters: K
- Randomly initialize K centroids
- Repeat until convergence:
 - Assign each pt. to the nearest centroid
 - Recalculate centroids as mean of assigned pts.

② How to Determine no. of Clusters?

- Use the elbow method, which plots SSE (Sum of squared Errors) vs no. of clusters. The "elbow" pt. suggests the optimal K .

③ Formula of SSE (Sum of squared Errors)

$$\rightarrow \text{SSE} = \sum_{i=1}^k \sum_{x \in C_i} \|x - \mu_i\|^2$$

where, μ_i is the centroid of cluster C_i

④ Elbow Technique

- plot SSE vs K
→ look for the "elbow" (sharp bend),
where adding clusters doesn't reduce SSE significantly
→ The K is considered optimal

5). Parameters of KMean()

- n-clusters : no. of clusters
- init : initialization method (e.g., 'k-means++')
- n-init : no. of initializations
- max_iter : max. iterations
- random_state : reproducibility
- algo. : opt. algo. ('lloyd', 'elkan')

Lab-10 Code

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.cluster import KMeans
# from sklearn.preprocessing import StandardScaler
from sklearn.datasets import load_iris
iris = load_iris()
iris_df = pd.DataFrame(iris.data, columns=iris.feature_names)[['petal length(cm)', 'petal width(cm)']]
iris_scaled = StandardScaler().fit_transform(iris_df)
sse = [KMeans(n_clusters=k, random_state=42).fit(iris_scaled).inertia_ for k in range(1, 11)]
```

```
plt.plot(range(1,11), sse, marker = 'o')
```

```
plt.title("Elbow method : SSE")
```

Curv
29.0 "n"

Code:

```
import pandas as pd

from sklearn.cluster import KMeans
import matplotlib.pyplot as plt

from sklearn.datasets import load_iris # Import load_iris

# Step 1: Load the Iris dataset directly
iris = load_iris()

# Create a DataFrame from the data and target
data = pd.DataFrame(data=iris.data, columns=iris.feature_names)

# Add the target column for potential reference, though not used for clustering
data['target'] = iris.target

# Step 2: Extract only numeric columns (or select required features)
# All features in the Iris dataset are numeric
X = data[iris.feature_names].values # Use the feature names to select columns

# Step 3: Apply KMeans
# Adjust n_clusters based on the expected number of clusters in your data (3 for Iris)
kmeans = KMeans(n_clusters=3, random_state=42, n_init=10) # Added n_init to suppress future
# warnings
data['Cluster'] = kmeans.fit_predict(X)

# Step 4: Plot clusters (for 2D data)
# Iris data has 4 features. We will plot the first two features for visualization.

if X.shape[1] >= 2:
```

```
plt.scatter(X[:, 0], X[:, 1], c=data['Cluster'], cmap='viridis')

plt.scatter(kmeans.cluster_centers_[:, 0], kmeans.cluster_centers_[:, 1], color='red', marker='x', s=200)

plt.title("K-Means Clustering of Iris Dataset")

plt.xlabel(iris.feature_names[0]) # Label with actual feature name

plt.ylabel(iris.feature_names[1]) # Label with actual feature name

plt.show()

else:

    print("Cannot plot clustering results directly for data with less than 2 features.")
```

Program 11

Implement Dimensionality reduction using Principal Component Analysis (PCA) method

Screenshot:

~~Table 11~~ PCA is a dimensionality reduction technique that transforms high-dimensional data into fewer dimensions while retaining most variance in the dataset.
Steps for Principal Component Analysis (PCA) algo.

1. Calculate mean
2. calculation of covariance matrix
3. Eigenvalues of the covariance matrix
4. Computation of the Eigen vectors - Unit eigenvectors
5. Computation of first principal components
6. Geometrical meaning of first principal components

Given the ^{data} table, reduce the dimension from 2 to 1 using the PCA.

feature	ex. 1	2	3	4
X ₁	4	8	13	7
X ₂	11	4	5	14

Step 1 \Rightarrow calculate mean

Fig. shows the sum

$$\bar{X}_1 = \frac{1}{4} (4+8+13+7) = 8,$$

$$\bar{X}_2 = \frac{1}{4} (11+4+5+14) = 8.5.$$

Step 2: calculate covariance matrix

$$\begin{aligned}\text{Cov}(x_1, x_1) &= \frac{1}{N-1} \sum_{k=1}^N (x_{1k} - \bar{x}_1)^2 \\ &= \frac{1}{3} (4-8)^2 + (8-8)^2 + (13-8)^2 \\ &= 14\end{aligned}$$

$$\begin{aligned}\text{cov}(x_1, x_2) &= \frac{1}{N-1} \sum_{k=1}^N (x_{1k} - \bar{x}_1)(x_{2k} - \bar{x}_2) \\ &= \frac{1}{3} [(4-8)(11-8.5) + (8-8)(4-8.5) \\ &\quad + (12-8)(5-8.5) + (7-8)(14-8.5)] \\ &= -11\end{aligned}$$

$$\text{cov}(x_2, x_1) = \text{cov}(x_1, x_2)$$

$$\begin{aligned}\text{cov}(x_2, x_2) &= \frac{1}{N-1} \sum_{k=1}^N (x_{2k} - \bar{x}_2)^2 \\ &= \frac{1}{3} [(11-8.5)^2 + (4-8.5)^2 + (5-8.5)^2 + (14-8.5)^2] \\ &= 23\end{aligned}$$

Eigenvalues of the covariance matrix

The cov. matrix is

$$\begin{bmatrix} \text{cov}(x_1, x_1) & \text{cov}(x_1, x_2) \\ \text{cov}(x_2, x_1) & \text{cov}(x_2, x_2) \end{bmatrix} = \begin{bmatrix} 14 & -11 \\ -11 & 23 \end{bmatrix}$$

Step 3: Eigen values of cov. matrix

The char. eqn² of cov. matrix is

$$\left[\det \begin{pmatrix} 14-\lambda & -11 \\ -11 & 23-\lambda \end{pmatrix} = 0 \right]$$

$$0 = \det((S - \lambda I)) \xrightarrow{\text{Simplifying}} 0 = ((x, x) v_0)$$

$$(14-\lambda) \begin{vmatrix} 1 & -1 \\ -1 & 23-\lambda \end{vmatrix} =$$

$$(14-\lambda)(23-\lambda) - (-1) \times (-11) =$$

$$= \lambda^2 - 37\lambda + 201$$

$$((x, x) v_0) = ((x, x) v_0)$$

giving the char. eqn we get,

$$\lambda = \frac{1}{2}(37 \pm \sqrt{565})$$

$$= 30.3849, 6.6151$$

$$+ (2.8-7) + (2.8-11) = (\lambda_1, \lambda_2) (S v_0)$$

Step 4: Computation of eigenvectors

→ the eigenvector corr. to $\lambda_1 = \lambda_1$ is a vector

$$v = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}$$

satisfying the foll. eqns:

$$\begin{bmatrix} 0 \\ 0 \end{bmatrix} = (S - \lambda_1 I) x$$

$$= \begin{bmatrix} 14 - \lambda_1 & -11 \\ -11 & 23 - \lambda_1 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}$$

$$\leftarrow \begin{cases} (14 - \lambda_1) u_1 - 11 u_2 \\ -11 u_1 + (23 - \lambda_1) u_2 \end{cases}$$

This is eqⁿ to foll. 2 equations:

$$(14 - \lambda_1)u_1 - 11u_2 = 0$$

$$-11u_1 + (23 - \lambda_1)u_2 = 0$$

using theory of system of linear eqns, we note that these eqns are not independent & sol's are given by.

$$\begin{bmatrix} u_1 \\ u_2 \end{bmatrix} =$$

$$\begin{bmatrix} 8088.0 \\ 11u - 14\lambda_1 \end{bmatrix}$$

that is,

$$u_1 = 11t, u_2 = (14 - \lambda_1)t,$$

where t is any real number

Taking $t=1$, we get an eigenvector corr. to λ_1 as

$$v_1 = \begin{bmatrix} 11 \\ 14 - \lambda_1 \end{bmatrix} = \begin{bmatrix} 11 \\ 14 - 30 \end{bmatrix} = \begin{bmatrix} 11 \\ -16 \end{bmatrix}$$

To find a unit eigenvector, we compute the length of x_1 , which is given by

$$\|v_1\| = \sqrt{11^2 + (14 - 30)^2}$$

$$= \sqrt{11^2 + (-16)^2} = \sqrt{121 + 256} = \sqrt{377}$$

$$= 19.7348$$

$$= (19.7348) \times 10^{-2}$$

$$= (19.7348) \times 10^{-2}$$

\therefore , a unit eigenvector (not at Step 2), with corr. to λ_1 , is

$$\mathbf{e}_1 = \frac{\mathbf{v}_{11}}{\|\mathbf{v}_{11}\|} = \frac{\begin{bmatrix} 11 \\ 11 \\ 11 \end{bmatrix}}{\sqrt{(11-\lambda_1)^2 + (11-\lambda_1)^2 + (11-\lambda_1)^2}} = \frac{\begin{bmatrix} 11 \\ 11 \\ 11 \end{bmatrix}}{\sqrt{3(11-\lambda_1)^2}} = \frac{\begin{bmatrix} 11 \\ 11 \\ 11 \end{bmatrix}}{\sqrt{3(11-10.3849)^2}} = \frac{\begin{bmatrix} 11 \\ 11 \\ 11 \end{bmatrix}}{\sqrt{3(0.6151^2)}} = \frac{\begin{bmatrix} 11 \\ 11 \\ 11 \end{bmatrix}}{\sqrt{3(0.3781)}} = \frac{\begin{bmatrix} 11 \\ 11 \\ 11 \end{bmatrix}}{\sqrt{1.1343}} = \frac{\begin{bmatrix} 11 \\ 11 \\ 11 \end{bmatrix}}{1.0651} = \begin{bmatrix} 10.3849 \\ 10.3849 \\ 10.3849 \end{bmatrix}$$

$$\begin{aligned} \text{unit eigenvector corr. to } \lambda_1 &= \frac{\begin{bmatrix} 11 \\ 11 \\ 11 \end{bmatrix}}{\sqrt{1.1343}} = \begin{bmatrix} 11/\sqrt{1.1343} \\ 11/\sqrt{1.1343} \\ 11/\sqrt{1.1343} \end{bmatrix} = \begin{bmatrix} 11/19.7348 \\ 11/19.7348 \\ 11/19.7348 \end{bmatrix} \\ &= \begin{bmatrix} 0.5574 \\ 0.5574 \\ -0.8303 \end{bmatrix} \end{aligned}$$

By carrying out similar computations, the unit eigenvector \mathbf{e}_2 corr. to the eigenvalue $\lambda_2 = 1.0651$ can be shown to be $\mathbf{e}_2 = \begin{bmatrix} 0.8303 \\ 0.5574 \end{bmatrix}$

Step 5: Computation of 1st principal component

Let,

$$\begin{bmatrix} x_{1k} \\ x_{2k} \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \mathbf{u}$$

be the k^{th} sample in the above table (dataset).

The 1st principal component of this ex. is given

by (here " T " denotes the transpose of the

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

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

$$- 0.8303(x_{2k} - \bar{x}_2)$$

Code:

```
import pandas as pd  
  
from sklearn.decomposition import PCA  
  
from sklearn.preprocessing import StandardScaler  
import matplotlib.pyplot as plt  
  
  
# Load dataset  
  
data = pd.read_csv("your_data.csv") # Replace with your file  
X = data.select_dtypes(include=['float64', 'int64'])  
  
  
# Step 1: Standardize  
  
scaler = StandardScaler()  
  
X_scaled = scaler.fit_transform(X)  
  
  
# Step 2: Apply PCA  
  
pca = PCA(n_components=2)  
  
X_pca = pca.fit_transform(X_scaled)  
  
  
# Print explained variance ratio  
  
print("Explained variance ratio:", pca.explained_variance_ratio_)  
  
  
# Visualize  
  
plt.scatter(X_pca[:, 0], X_pca[:, 1], c='blue', alpha=0.5)  
plt.title("PCA - 2D Projection")  
plt.xlabel("Principal Component 1")
```

```
plt.ylabel("Principal Component 2")  
plt.show()
```



Choose Files Student.csv

- Student.csv(text/csv) - 41 bytes, last modified: 5/18/2025 - 100% done

Saving Student.csv to Student.csv

Uploaded file: Student.csv

	Hours_Studied	Passes	grid icon
0	1	0	grid icon
1	2	0	grid icon
2	3	1	grid icon
3	4	1	grid icon
4	5	1	grid icon

Numerical features found: ['Hours_Studied', 'Passes']

2D PCA Visualization

