

**VISVESVARAYA TECHNOLOGICAL
UNIVERSITY**

“JnanaSangama”, Belgaum -590014, Karnataka.



**LAB REPORT
on**

Machine Learning (23CS6PCMAL)

Submitted by

AKSHARA SINGA (1BM22CS029)

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 **Akshara Singa(1BM22CS029)**, who is bonafide student of **B.M.S. College of Engineering**. It is in partial fulfillment for the award of **Bachelor of Engineering in Computer Science and Engineering** of the Visvesvaraya Technological University, Belgaum. The Lab report has been approved as it satisfies the academic requirements in respect of an Machine Learning (23CS6PCMAL) work prescribed for the said degree.

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

Index

Sl. No.	Date	Experiment Title	Page No.
1	21-2-2025	Write a python program to import and export data using Pandas library functions	2-7
2	3-3-2025	Demonstrate various data pre-processing techniques for a given dataset	8-11
3	10-3-2025	Implement Linear and Multi-Linear Regression algorithm using appropriate dataset	12-16
4	17-3-2025	Build Logistic Regression Model for a given dataset	17-21
5	24-3-2025	Use an appropriate data set for building the decision tree (ID3) and apply this knowledge to classify a new sample.	22-24
6	7-4-2025	Build KNN Classification model for a given dataset.	25-28
7	21-4-2025	Build Support vector machine model for a given dataset	29-31
8	5-5-2025	Implement Random forest ensemble method on a given dataset.	32-33
9	5-5-2025	Implement Boosting ensemble method on a given dataset.	34-35
10	12-5-2025	Build k-Means algorithm to cluster a set of data stored in a .CSV file.	36-38
11	12-5-2025	Implement Dimensionality reduction using Principal Component Analysis (PCA) method.	39-42

Github Link:

<https://github.com/Singa-akshara/6A-ML-LAB-Batch2>

Program 1

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

Screenshot:

9/3/24 LAB-1

① Initializing values directly into DataFrame

```
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			
2			

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())
```

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

④ Downloading datasets from existing dataset repositories like kaggle

```
df = pd.read_csv("MobileDataset.csv")
print(df.head())
```

Output -

0		
1		

* BANK DATA ANALYSIS

- Using the code of "Stock Market Data Analysis", consider the following

- HDFC Bank Ltd, ICICI Bank Ltd, Kotak Mahindra Bank Ltd
 Tickers = ["HDFCBANK.NS", "ICICIBANK.NS", "KOTAKBANK.NS"]
- Start date : 2024-01-01, End-date : 2024-12-30
- Plot the closing price and 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', group_by='tickers')
 print(data.head())
 print(data.shape)

Output -

Ticker KOTAKBANK.NS

	Price	Open	High	Low	Close	Volume
2024-01-01	1906	1916	1891	1907.0	14259	
:	:	:	:	:	:	
:	:	:	:	:	:	

Code:

```
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("Sample data:")
print(df.head())
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())
file_path = 'data.csv'
df = pd.read_csv(file_path)
print("Sample data:")
print(df.head())
print("\n")
file_path = 'mobiles-dataset-2025.csv'
df = pd.read_csv(file_path, encoding='latin-1') # or 'cp1252' or other suitable
encoding
print("Sample data:")
print(df.head())
import pandas as pd

data = {
'USN': ['IS001','IS002','IS003','IS004','IS005'],
'Name': ['Alice', 'Bob', 'Charlie', 'David', 'Eve'],
'Marks': [25, 30, 35, 40,45]
}

df = pd.DataFrame(data)
print("Sample data:")
print(df.head())
from sklearn.datasets import load_diabetes
iris = load_diabetes()
df = pd.DataFrame(iris.data, columns=iris.feature_names)

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

df = pd.read_csv("/content/dataset-of-diabetes .csv",encoding='latin-1')
print("Sample data:")
print(df.head())
print("\n")

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

df.to_csv('output.csv',index=False)
print("Data saved to output.csv")

```

```

sales_df = pd.read_csv('sample_sales_data.csv')
print("Sample data:")
print(sales_df.head())
sales_by_region = sales_df.groupby('Region')['Sales'].sum()
print("\nTotal sales by region:")
print(sales_by_region)
best_selling_products
=sales_df.groupby('Product')['Quantity'].sum().sort_values(ascending=False)
print("\nBest-selling products by quantity:")
print(best_selling_products)
sales_by_region.to_csv('sales_by_region.csv')
best_selling_products.to_csv('best_selling_products.csv')
print("Data saved to sales_by_region.csv and best_selling_products.csv")

import yfinance as yf
import matplotlib.pyplot as plt
tickers = ["RELIANCE.NS", "TCS.NS", "INFY.NS"]
data = yf.download(tickers, start="2022-10-01", end="2023-10-01",
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)
print("\n")
reliance_data = data['RELIANCE.NS']
print("\nSummary statistics for Reliance Industries:")
print(reliance_data.describe())
reliance_data['Daily Return'] = reliance_data['Close'].pct_change()
print("\n")
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()
reliance_data.to_csv('reliance_stock_data.csv')

tickers = ["HDFCBANK.NS", "ICICI.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)
print("\n")
reliance_data = data['HDFCBANK.NS']
print("\nSummary statistics for Reliance Industries:")

```

```

print(reliance_data.describe())
reliance_data['Daily Return'] = reliance_data['Close'].pct_change()
print("\n")
plt.figure(figsize=(12, 6))
plt.subplot(2, 1, 1)
reliance_data['Close'].plot(title="HDFC Industries - Closing Price")
plt.subplot(2, 1, 2)
reliance_data['Daily Return'].plot(title="HDFCIndustries - Daily Returns",
color='red') plt.tight_layout()
plt.show()
reliance_data.to_csv('hdfc_stock_data.csv')
print("\nhdfc stock data saved to 'hdfc_stock_data.csv'.")

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)
print("\n")
reliance_data = data['ICICIBANK.NS']
print("\nSummary statistics for ICICI Industries:")
print(reliance_data.describe())
reliance_data['Daily Return'] = reliance_data['Close'].pct_change()
print("\n")
plt.figure(figsize=(12, 6))
plt.subplot(2, 1, 1)
reliance_data['Close'].plot(title="ICICI Industries - Closing Price")
plt.subplot(2, 1, 2)
reliance_data['Daily Return'].plot(title="ICICI Industries - Daily Returns",
color='BLACK') plt.tight_layout()

plt.show()
reliance_data.to_csv('icici_stock_data.csv')
print("\nicici stock data saved to 'icici_stock_data.csv'.")
```

tickers = ["HDFCBANK.NS", "ICICI.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)

```
print("\n")
reliance_data = data['KOTAKBANK.NS']
print("\nSummary statistics for Reliance Industries:")
print(reliance_data.describe())
reliance_data['Daily Return'] =
reliance_data['Close'].pct_change() print("\n")
plt.figure(figsize=(12, 6))
plt.subplot(2, 1, 1)
reliance_data['Close'].plot(title="KOTAK Industries - "
Closing Price") plt.subplot(2, 1, 2)
reliance_data['Daily Return'].plot(title="kotak Industries - Daily Returns",
color='red') plt.tight_layout()
plt.show()
reliance_data.to_csv('kotak_stock_data.csv')
print("\nkotak stock data saved to 'kotak_stock_data.csv'.")
```

Program 2

Demonstrate various data pre-processing techniques for a given dataset

Screenshot:

4/3/24 Lab-2

Write python code, consider filename "housing.csv"

Diagram: Data Page

i) To load .csv file into the data frame
ii) To display information of all columns
iii) To display statistical information of all numerical
iv) To display the count of unique labels for "Ocean Proximity" column
v) To display which attributes (columns) in a dataset have missing values count greater than zero

Housing Dataset

Code

i) import pandas as pd
df = pd.DataFrame()
df = pd.read_csv("housing.csv")
import sklearn.datasets import fetch_california_housing

ii) print(df.info())
iii) print(df.describe())
iv) print(df['Ocean-Proximity'].value_counts())
v) missing_values = df.isnull().sum()
print(missing_values[missing_values > 0])

i) housing_data = fetch_california_housing()
df = pd.DataFrame(housing_data.data, columns=housing_data.feature_names)

Output

# Column	Non-Null Count	Dtype
MedInc	20690	non-null float64
HouseAge	20690	int64
AveRooms	20690	float64
AveBedrms	20690	float64
Popol	20690	float64

count 20690 20690 20690 20690 20690

iv) 5

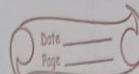
v) No missing values

* Diabetes dataset
* Adult dataset

i) Which columns in the dataset had missing values? How did you handle them?
ii) Which categorical columns did you identify in the dataset? How did you encode them?
iii) What is the difference between Min-Max Scaling & Standardisation? When would you use one over the other?

4/8/24

LAB - 2 (continuation)



* DIABETES DATASET

- ① Missing values are present in numerical columns if present, which are replaced by the mean of the respective column

- ② No categorical columns, hence no encoding

- ③ Min Max Scaling transforms 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 and need to be bound

Standardization transforms data to have zero mean and unit variance

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

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

* ADULT INCOME DATASET

- ① Missing values are represented by '?' which we replace with 'NaN'

for numerical columns - replace with mean
for categorical columns - replace with mode

- ② workclass, education, marital status, occupation, relationship, sex, native-country, income

Encoding method:

Label Encoding to convert categorical to numerical formats.

- ③ Same as prev.

Code:

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

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

```
df1=pd.read_csv("Dataset of Diabetes .csv")
```

```

df1.head()

df2=pd.read_csv("adult.csv")
df2.head()

df1.info()
df2.info()
df1.describe()
df2.describe()

missing_values1 = df1.isnull().sum()
print(missing_values1)
missing_values2 = df2.isnull().sum()
print(missing_values2)

df1['Gender'] = df1['Gender'].replace('f', 'F')
ordinal_encoder = OrdinalEncoder(categories=[["F", "M"]])
df1["Gender_Encoded"] =
ordinal_encoder.fit_transform(df1[["Gender"]]) onehot_encoder =
OneHotEncoder()
encoded_data =
onehot_encoder.fit_transform(df1[["CLASS"]]) encoded_array
= encoded_data.toarray()
encoded_df = pd.DataFrame(encoded_array,
columns=onehot_encoder.get_feature_names_out(["CLASS"])) df_encoded = pd.concat([df1,
encoded_df], axis=1)
df1 = pd.concat([df1, encoded_df], axis=1)
df1.drop("CLASS", axis=1, inplace=True)
df1.drop("Gender", axis=1, inplace=True)
print(df2.head())
from sklearn.preprocessing import OrdinalEncoder, OneHotEncoder
df_copy2 = df2
ordinal_encoder = OrdinalEncoder(categories=[["Male", "Female"]])
df_copy2["Gender_Encoded"] =
ordinal_encoder.fit_transform(df_copy2[["gender"]])
print(df_copy2[["gender", "Gender_Encoded"]])

onehot_encoder = OneHotEncoder()
encoded_data =
onehot_encoder.fit_transform(df2[["occupation", "workclass", "education
", "marital status", "relationship", "race", "native-country", "income"]])
encoded_array = encoded_data.toarray()
encoded_df =
pd.DataFrame(encoded_array,
columns=onehot_encoder.get_feature_names_out(["occupation", "workclass", "education
", "marital status", "relationship", "race", "native-country", "income"]))
df_encoded = pd.concat([df_copy2, encoded_df], axis=1)

df_encoded.drop("gender", axis=1, inplace=True)
df_encoded.drop("occupation", axis=1, inplace=True)
df_encoded.drop("workclass", axis=1, inplace=True)
df_encoded.drop("education", axis=1, inplace=True)

```

```

df_encoded.drop("marital-status", axis=1, inplace=True)
df_encoded.drop("relationship", axis=1, inplace=True)
df_encoded.drop("race", axis=1, inplace=True)
df_encoded.drop("native-country", axis=1, inplace=True)
df_encoded.drop("income", axis=1, inplace=True)
print(df_encoded.head())

normalizer = MinMaxScaler()
df_encoded[["fnlwgt","educational-num","capital-gain","capital-loss","hours-per-week"]] =
normalizer.fit_transform(df_encoded[["fnlwgt","educational-num","capital-gain","capital-
loss","hours-per week"]])
df_encoded.head()
normalizer = MinMaxScaler()
df1[["No_Pation","AGE","Urea","Cr" , "HbA1c" ,
"Chol","TG","HDL","LDL","VLDL","BMI"]] =
normalizer.fit_transform(df1[["No_Pation","AGE","Urea","Cr" , "HbA1c" ,
"Chol","TG","HDL","LDL","VLDL","BMI"]])
df1.head()

```

Program 3

Implement Linear and Multi-Linear Regression algorithm using appropriate dataset

Screenshot:

Date _____
Page _____

11/3/25	LAB - 4	
Q	Week	Sales
	1	2
	2	4
	3	5
	4	9
	5	9

$$x^T = \begin{bmatrix} 1 & 2 & 3 & 4 \end{bmatrix} \quad x = \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ 1 & 4 \end{bmatrix} \quad y = \begin{bmatrix} 2 \\ 4 \\ 5 \\ 9 \end{bmatrix}$$

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

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

compute $x^T x =$

$$\cdot \begin{bmatrix} 1 & 2 & 3 & 4 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \\ 2 \\ 1 \\ 3 \\ 1 \\ 4 \end{bmatrix} = \begin{bmatrix} 4 & 10 \\ 10 & 30 \end{bmatrix}$$

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

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

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

$$\beta = \begin{bmatrix} -0.5 \\ 2.2 \end{bmatrix}$$

Intercept
Slope

$$y = \beta_0 + \beta_1 x + e \quad \text{At } x=5 \quad y = -0.5 + 2.2(5) = 10.5$$

$$y = -0.5 + 2.2x$$

<p style="text-align: center;">Linear Regression</p> <pre> ① Canada_per_capita_income.csv import pandas as pd from sklearn.linear_model import LinearRegression import matplotlib.pyplot as plt data = pd.read_csv('per-capita-canada-per-capita-income.csv') X = data[['year']] y = data['per capita income (US\$)'] model = LinearRegression() model.fit(X, y) predicted_income = model.predict([[2020]]) print(f"Predicted per capita income in 2020: {predicted_income[0]} US\$") plt.scatter(data[['year']], data['per capita income (US\$)'], color='red', marker='+') plt.plot(data[['year']], model.predict(data[['year']]), color='blue') plt.xlabel('Year') plt.ylabel('Per capita Income (US\$)') plt.title('Canada Per Capita Income Prediction') plt.show() </pre>	<p style="text-align: right;">Date _____ Page _____</p> <p>Output</p> <p>Predicted per capita income in 2020: 41288.69 US\$ Coefficient (m): 828.47 Intercept (b): -1632210.76</p> <p>Salary.csv</p> <p>Predicted Salary for 12 years of experience: \$ 139049.67 .</p> <p>Coefficient (m): 9398.64 Intercept (b): 26265.99</p> <p style="text-align: center;">Multiple Regression</p> <p>① hiring.csv</p> <pre> import pandas as pd import numpy as np from sklearn import linear_model df = pd.readcsv('hiring.csv') word_to_number = {'one': 1, 'two': 2, 'three': 3, 'four': 4, 'five': 5, 'six': 6, 'seven': 7, 'eight': 8, 'nine': 9, 'ten': 10, 'eleven': 11} df['experience'] = df['experience'].apply(lambda n: word_to_number.get(n, n)) df['experience'] = pd.to_numeric(df['experience']) df['expenfce'] = df['experience'].fillna(df['experience'].median()) </pre>
---	--

```

df['test_score(out of 10)'] = df['test_score(out of
10)'].fillna(df['test_score(out of 10)'].median())
df['interview_score'] = df['interview_score'].fillna(
    df['interview_score'].median())
reg = linear_model.LinearRegression()
reg.fit(df[['experience', 'test_score(out of 10)', 'interview_score']], df['Salary'])
candidate1_Salary = reg.predict([[2, 9, 6]])
candidate2_Salary = reg.predict([[12, 10, 10]])
print(f"candidate1_Salary : ${candidate1_Salary[0]:.2f}")
print(f"candidate2_Salary : ${candidate2_Salary[0]:.2f}")
print(f"Coefficients : {reg.coef_}")
print(f"Intercept : {reg.intercept_}")

```

Output

candidate1_Salary : \$ 47056.91
 Candidate2_Salary : \$ 88227.64
 Coefficients : [2813.008 1333.333 2926.829]
 Intercept : 11869.918

Not
11869.918

1000_companies.csv

coefficients : [0.55389 1.02672 0.0805 46.6238]
 Intercept : -70214.4917
 Predicted profit : \$ 511209.20

Code:

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

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

```
import pandas as pd
```

```

import numpy as np
import matplotlib.pyplot as plt
from sklearn.model_selection import
train_test_split from sklearn.impute
import SimpleImputer
from sklearn.preprocessing import OrdinalEncoder,
OneHotEncoder from sklearn.preprocessing import
StandardScaler, MinMaxScaler from scipy import stats
from sklearn import linear_model

df1=pd.read_csv("canada_per_capita_inc
ome.csv") df1.head()

df2=pd.read_csv("salary.csv")
df2.head()
df2.YearsExperience.median()
df2.YearsExperience =
df2.YearsExperience.fillna(df2.YearsExperience.median()) df2

plt.xlabel("year")
plt.ylabel("per capita income (US$)")
plt.scatter(df1.year, df1['per capita income (US$)'])

plt.xlabel("YearsExperience")
plt.ylabel("Salary")
plt.scatter(df2.YearsExperience, df2.Salary)

reg1 = linear_model.LinearRegression()
reg1.intercept_
reg1.predict([[2020]])

reg2 = linear_model.LinearRegression()
reg2.fit(df2.drop('Salary', axis='columns'),
df2['Salary']) reg2.coef_
reg2.intercept_
reg2.predict([[12]])

from google.colab import files
hirings=files.upload()

from google.colab import files
companies=files.upload()

df3=pd.read_csv("hirings.csv")
df3.head()

df4=pd.read_csv("1000_Companies.csv")
df4.head()

df3.isnull().sum()
df4.isnull().sum()

df3_copy = df3.copy()

```

```

experience_mapping = {'two': 2, 'three': 3, 'five': 5, 'seven': 7, 'ten': 10,
'eleven': 11} df3_copy['experience'] =
df3_copy['experience'].map(experience_mapping)
median_experience = df3_copy['experience'].median()
df3_copy['experience'] = df3_copy['experience'].fillna(median_experience)
df3_copy
df3_copy['test_score(out of 10)'] = df3_copy['test_score(out of
10)'].fillna(df3_copy['test_score(out of 10)'].mean())
reg3 = linear_model.LinearRegression()
reg3.fit(df3_copy.drop('salary($)', axis='columns'),
df3_copy['salary($)']) reg3.coef_
reg3.intercept_
reg3.predict([[2,9,6]])
reg3.predict([[12,10,10]])
ohe = OneHotEncoder(sparse_output=False, handle_unknown='ignore')
state_encoded = ohe.fit_transform(df4[['State']])
state_encoded_df = pd.DataFrame(state_encoded, columns=ohe.get_feature_names_out(['State']))

df4 = pd.concat([df4, state_encoded_df], axis=1).drop(columns=['State'])
print(df4)
reg4 = linear_model.LinearRegression()
reg4.fit(df4.drop('Profit',axis='columns'),df4.Profit)
print(reg4.coef_)
print(reg4.intercept_)
reg4.predict([[91694.48, 515841.3, 11931.24,0,1,0]])

```

Program 4

Build Logistic Regression Model for a given dataset

Screenshot:

Date _____
Page _____

LAB - 3

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

a. Write the logistic regression equation

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

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

$$\rightarrow p(x) = \frac{1}{1 + e^{(-5 + 0.8x)}} \quad x = 7$$

$$p(x) = \frac{1}{1 + e^{(-5 + 0.8(7))}} = 0.6457$$

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

$$p(x) = 0.6457$$

$$= 0.6457 \geq 0.5$$

Thus $y = 1$ (pass)

18/3/18

2. Consider $z = [2, 1, 0]$ for 3 classes. Apply Softmax function to find probability values of 3 classes

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

$$\text{Softmax}(z_1) = \frac{e^2}{e^2 + e^1 + e^0} = 0.665$$

$$\text{Softmax}(z_2) = \frac{e^1}{e^2 + e^1 + e^0} = 0.249$$

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

The probabilities of 3 classes are approximately 66.5%, 24.9% and 9.1%.

* Binary Logistic Regression

```
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([[50]])

y_predicted

model.conf

model.intercept

import math

def sigmoid(z):

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

def prediction_function(age):

$z = 0.127 * age - 4.973$

$y = \text{sigmoid}(z)$

return y

age = 35

prediction_function(35)

→ 0.37 is less than 0.5 which means person with 35 will not buy insurance

* Logistic Regression Multiclass

```
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)
```

Output:

Accuracy of the Multinomial Logistic Regression model on the test set : 1.00

1. For dataset file "HR_Comma-Sep.csv"

i) Which variables did you identify as having a direct & clear impact on employee retention? → Key variables such as Satisfaction level, last evaluation and promotion last 5 years had a direct impact on retention

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 the model was 75.9%, which is reasonable but might be improved with more complex models or hyperparameter tuning

2. For 2nn dataset

i) Did you perform any data preprocessing steps? If yes, what were they, & why are they necessary? → Data preprocessing steps like merging datasets, handling missing data, and encoding categorical variables were necessary to prepare the data for logistic regression

- ii) Were there any missing or inconsistent values in 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 confusion matrix revealed how well the model predicted different class types, and misclassifications often occurred for similar animal classes
- iv) which class types were most frequently misclassified? Why do you think this happened?
 → The misclassifications are likely due to overlapping characteristics between certain classes which are challenging for logistic regression to separate effectively.

Code:

```

from google.colab import files
hr=files.upload()

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
  
```

```

from sklearn.impute import SimpleImputer
from sklearn.preprocessing import OrdinalEncoder,
OneHotEncoder from sklearn.preprocessing import
StandardScaler, MinMaxScaler
from scipy import stats
from sklearn import linear_model
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score

df1=pd.read_csv("HR_comma_sep.csv")
df1.head()
df1.isnull().sum()
plt.figure(figsize=(12, 6))
sns.barplot(x='Department', y='left', data=df1)
plt.title('Employee Retention Rate by Department')
plt.xlabel('Department')
plt.ylabel('Proportion of Employees Left')
plt.xticks(rotation=45, ha='right')
plt.show()

ohe = OneHotEncoder(handle_unknown='ignore',
sparse_output=False) department_encoded =
ohe.fit_transform(df1[['Department']])
department_encoded_df = pd.DataFrame(department_encoded,
columns=ohe.get_feature_names_out(['Department']))
df1 = pd.concat([df1, department_encoded_df], axis=1)
df1 = df1.drop('Department', axis=1)
ordinal_encoder = OrdinalEncoder(categories=[[ 'low', 'medium', 'high']],
dtype=np.int64) salary_encoded =
ordinal_encoder.fit_transform(df1[['salary']])
df1['salary_encoded'] = salary_encoded
df1 = df1.drop('salary', axis=1)
df1.head()

correlation_matrix = df1.corr()
plt.figure(figsize=(12, 10))
sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm',
fmt=".2f") plt.title('Correlation Matrix of Features')
plt.show()
plt.figure(figsize=(8, 6))
sns.barplot(x='salary_encoded', y='left', data=df1)
plt.title('Impact of Employee Salary on Retention')
plt.xlabel('Salary Level (Encoded)')
plt.ylabel('Proportion of Employees Left')
plt.show()

df_copy = df1[['number_project', 'average_monthly_hours', 'time_spend_company',
'left','salary_encoded', 'satisfaction_level','Work_accident']]
df_copy.head()

```

```

X = df_copy.drop('left', axis=1)
y = df_copy['left']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42) model = LogisticRegression(max_iter=1000)
model.fit(X_train, y_train)
y_pred = model.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy of the Logistic Regression model: {accuracy}")

from google.colab import files
zodata=files.upload()
zootype=files.upload()

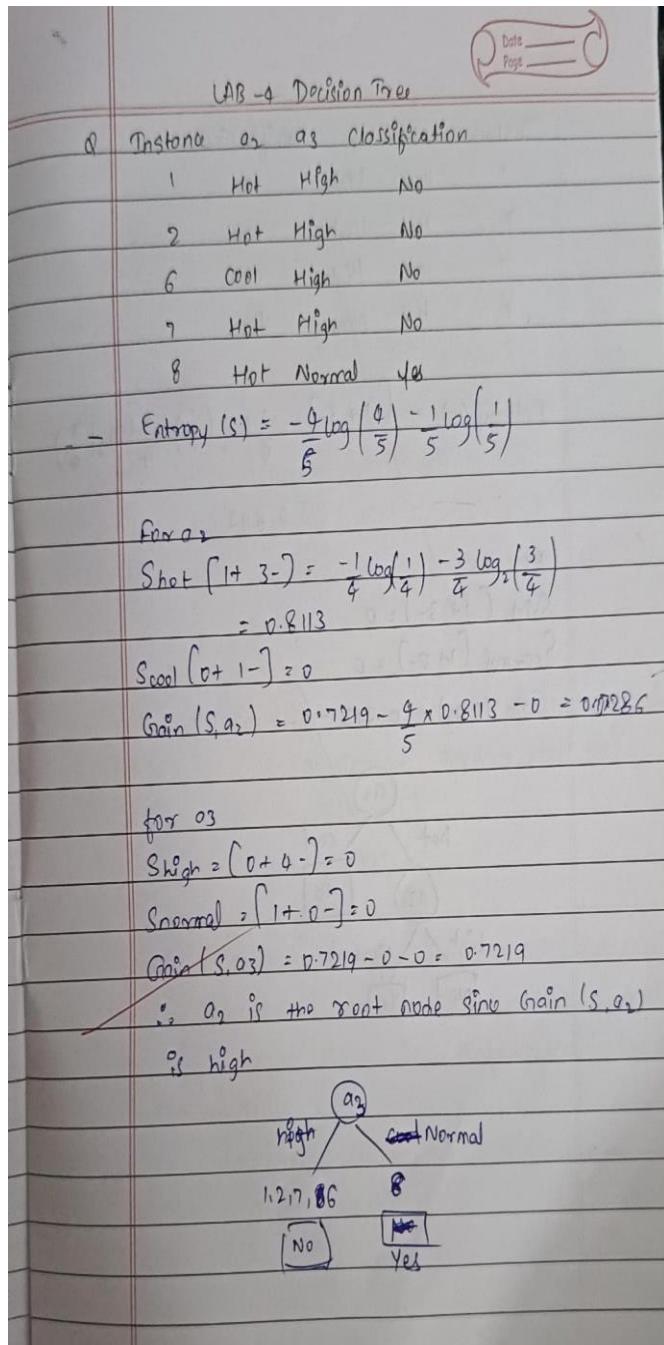
zoo_data = pd.read_csv('zoo-data.csv')
zoo_class = pd.read_csv('zoo-class-type.csv')
merged_data = pd.merge(zoo_data, zoo_class, left_on='class_type',
right_on='Class_Number') merged_data = merged_data.drop(['Animal_Names',
'Number_Of_Animal_Species_In_Class',
'Class_Number','class_type','animal_name'], axis=1)
X = merged_data.drop('Class_Type', axis=1)
y = merged_data['Class_Type']
print(merged_data.head())
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42) model = LogisticRegression(max_iter=1000)
model.fit(X_train, y_train)
y_pred = model.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy: {accuracy}")
cm = confusion_matrix(y_test, y_pred)
disp = ConfusionMatrixDisplay(confusion_matrix=cm,
display_labels=np.unique(y_test)) disp.plot(cmap="Blues", values_format="d")
plt.title("Confusion Matrix")
plt.show()

```

Program 5

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

Screenshot:



* Code for Building a Decision Tree classification to classify IRIS flower dataset.

- 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.datasets

```

iris = pd.read_csv("iris.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))

```

Output for IRIS dataset					
Accuracy: 1.0					
Confusion matrix					
$\begin{bmatrix} 10 & 0 & 0 \\ 0 & 9 & 0 \\ 0 & 0 & 11 \end{bmatrix}$					
Classification report:					
	Precision	Recall	F1-Score	Support	
Iris-setosa	1.00	1.00	1.00	10	
versicolor	1.00	1.00	1.00	9	
virginica	1.00	1.00	1.00	11	
accuracy			1.00	30	
micro avg.	1.00	1.00	1.00	30	
weighted avg.	1.00	1.00	1.00	30	

Code:

```

from google.colab import files
iris=files.upload()
df1=pd.read_csv("iris.csv")
df1.head()

df1.isnull().sum()

X = df1.drop('species', axis=1)
y = df1['species']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42) clf = DecisionTreeClassifier(criterion='entropy')
clf.fit(X_train, y_train)
y_pred = clf.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print(f'Accuracy: {accuracy:.2f}')
print(classification_report(y_test, y_pred))
plt.figure(figsize=(12, 8))
plot_tree(clf, filled=True, feature_names=X.columns,
class_names=y.unique()) plt.show()

cm = confusion_matrix(y_test, y_pred)
disp = ConfusionMatrixDisplay(confusion_matrix=cm,
display_labels=clf.classes_) cmap = plt.cm.get_cmap('PuBuGn')
disp.plot(cmap=cmap)
plt.show()

drug=files.upload()
df2=pd.read_csv("drug.csv")
df2.head()

```

```

df2.isnull().sum()

label_encoders = { }
for column in df2.columns:
    le = LabelEncoder()
    df2[column] = le.fit_transform(df2[column])
    label_encoders[column] = le
X = df2.drop('Drug', axis=1)
y = df2['Drug']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42) clf = DecisionTreeClassifier(criterion='entropy')
clf.fit(X_train, y_train)
y_pred = clf.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print(f'Accuracy: {accuracy:.2f}')
print(classification_report(y_test, y_pred))
plt.figure(figsize=(12, 8))
plot_tree(clf, filled=True, feature_names=X.columns, class_names=[str(c) for c in
y.unique()])
plt.show()

cm = confusion_matrix(y_test, y_pred)
disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=clf.classes_)
cmap = plt.cm.Blues
disp.plot(cmap=cmap)
plt.show()

pc=files.upload()
df3=pd.read_csv("petrol_consumption.csv")
df3.head()
df3.isnull().sum()
X = df3.drop('Petrol_Consumption', axis=1)
y = df3['Petrol_Consumption']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42) regressor = DecisionTreeRegressor(random_state=42)
regressor.fit(X_train, y_train)
y_pred =
regressor.predict(X_test)
mse = mean_squared_error(y_test, y_pred)
rmse = sqrt(mse)
mae = mean_absolute_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)
print(f'Mean Squared Error: {mse:.2f}')
print(f'Root Mean Squared Error:
{rmse:.2f}') print(f'Mean Absolute Error:
{mae:.2f}') print(f'R-squared: {r2:.2f}')
plt.figure(figsize=(30, 30))
plot_tree(regressor, filled=True, feature_names=X.columns,
fontsize=10) plt.show()

```

Program 6

Build KNN Classification model for a given dataset.

Screenshot :

KNN					
	Person	Age	Salary/k	Target	Distance Rank
Q	A	18	50	N	52.81
	B	23	55	N	46.57
	C	24	70	N	31.95 2
	D	41	60	Y	40.44 3
	E	43	70	Y	31.04 1
	F	38	40	Y	60.07
	X	35	100	?	.
→	The target is Y for data ($x, 35, 100$)				
Q	For Iris dataset how to choose the k value? Demonstrate using accuracy rate & error rate				
-	K=3 gives 100% accuracy here but generally take k=5, error rate = $1 - \text{accuracy} = 0.0$ (no misclassification), to find best k, test multiple values & plot the error rate				
Q	For diabetes data set what is the purpose of feature scaling? How to perform it? → It is needed bcz features have different ranges (eg glucose vs BMI), Standard Scaler ensure equal feature contribn by normalizing data Improves KNN performance (accuracy = 69.84% after Scaling)				

To do questions for decision tree

Q for Iris.csv dataset what was the accuracy score, what does the confusion matrix tell you about the model's performance? What were there any misclassifications? If so which classes were most confused?

→ Accuracy Score = 1.0 (100%)
 The confusion matrix showed perfect classification with no miss classifications
 Misclassification analysis: No classes were confused, all predictions were correct

Q for petrol.csv can you interpret the regression tree structure? What are the most important features for predicting the petrol consumption? How does the regression tree handle continuous target variables compared to the decision tree classifier?

→ The tree identifies the key features that impact petrol consumption & splits the dataset accordingly, the most influential features are likely variables such as road infrastructure, vehicle taxation & petrol price, unlike classification trees, regression trees predict continuous values by averaging observations in terminal nodes instead of assigning categorical labels

KNN code for iris dataset

```

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
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
iris = pd.read_csv("iris.csv")
diabetes = pd.read_csv("diabetes.csv")
heart = pd.read_csv("heart.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)
knn = KNeighborsClassifier(n_neighbors=3)
knn.fit(X_train, y_train)
y_pred_k = knn.predict(X_test)
print("KNN Classification for IRIS Dataset:")
print("Accuracy:", accuracy_score(y_test, y_pred_k))
print("Confusion matrix:\n", confusion_matrix)
print("Classification Report: classification_report")
  
```

- Some output

Code:

```

from google.colab import files
iris=files.upload()
df1=pd.read_csv("iris (2).csv")
df1.head()
df1.isnull().sum()
X = df1.drop('species', axis=1)
y = df1['species']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42) best_k = 1
best_accuracy = 0
for k in range(1,
11):
  knn = KNeighborsClassifier(n_neighbors=k)
  knn.fit(X_train, y_train)
  y_pred = knn.predict(X_test)
  accuracy = accuracy_score(y_test, y_pred)
  print(f"Accuracy for k={k}: {accuracy}, Error Rate for k={k}: {1-accuracy}") if accuracy >
  best_accuracy: best_accuracy = accuracy best_k = k
print(f"Best k value: {best_k}")
knn = KNeighborsClassifier(n_neighbors=3)
knn.fit(X_train, y_train)
y_pred = knn.predict(X_test)
  
```

```

print("Accuracy Score:", accuracy_score(y_test, y_pred))
print("\nConfusion Matrix:")
cm = confusion_matrix(y_test, y_pred)
print(cm)
print("\nClassification Report:")
print(classification_report(y_test, y_pred))
plt.figure(figsize=(8, 6))
sns.heatmap(cm, annot=True, fmt="d", cmap="Blues",
            xticklabels=knn.classes_, yticklabels=knn.classes_)
plt.xlabel('Predicted')
plt.ylabel('Actual')
plt.title('Confusion Matrix')
plt.show()

diabetes=files.upload()
df2=pd.read_csv("diabetes.csv")
df2.head()
df2.isnull().sum()
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
X_scaled = scaler.fit_transform(df2.drop('Outcome', axis=1))
X_train, X_test, y_train, y_test = train_test_split(X_scaled, df2['Outcome'], test_size=0.2,
random_state=42) best_k = 1
best_accuracy = 0
for k in range(1,
11):
    knn = KNeighborsClassifier(n_neighbors=k)
    knn.fit(X_train, y_train)
    y_pred = knn.predict(X_test)
    accuracy = accuracy_score(y_test, y_pred)
    print(f"Accuracy for k={k}: {accuracy}")
    if accuracy > best_accuracy:
        best_accuracy = accuracy
        best_k = k
print(f"Best k value: {best_k}")
knn = KNeighborsClassifier(n_neighbors=best_k)
knn.fit(X_train, y_train) y_pred = knn.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy:", accuracy)
cm = confusion_matrix(y_test, y_pred)
print("Confusion Matrix:")
print(cm)
plt.figure(figsize=(8, 6))
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues')
plt.xlabel("Predicted") plt.ylabel("Actual")
plt.title("Confusion Matrix")
plt.show()
print("\nClassification Report:")
print(classification_report(y_test, y_pred))

heart=files.upload()
df3=pd.read_csv("heart.csv")

```

```

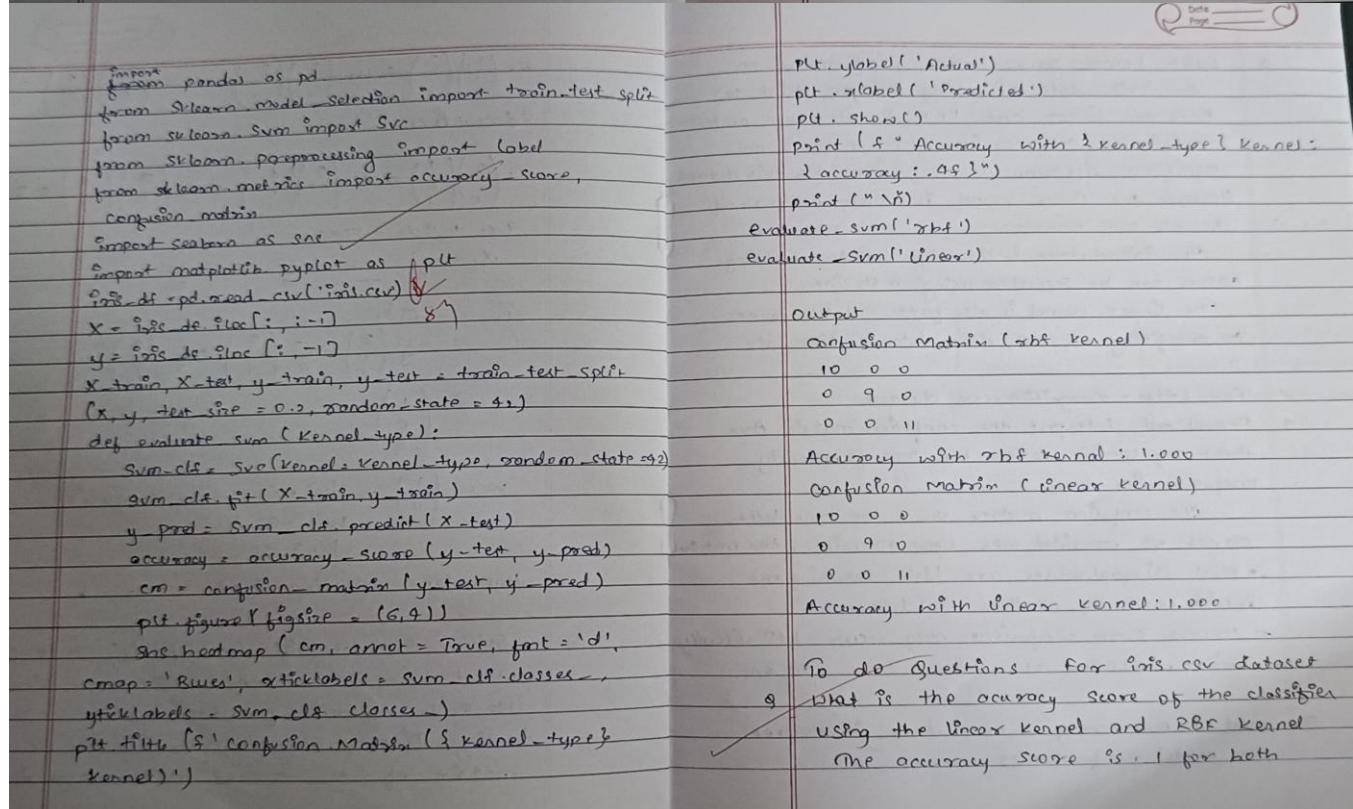
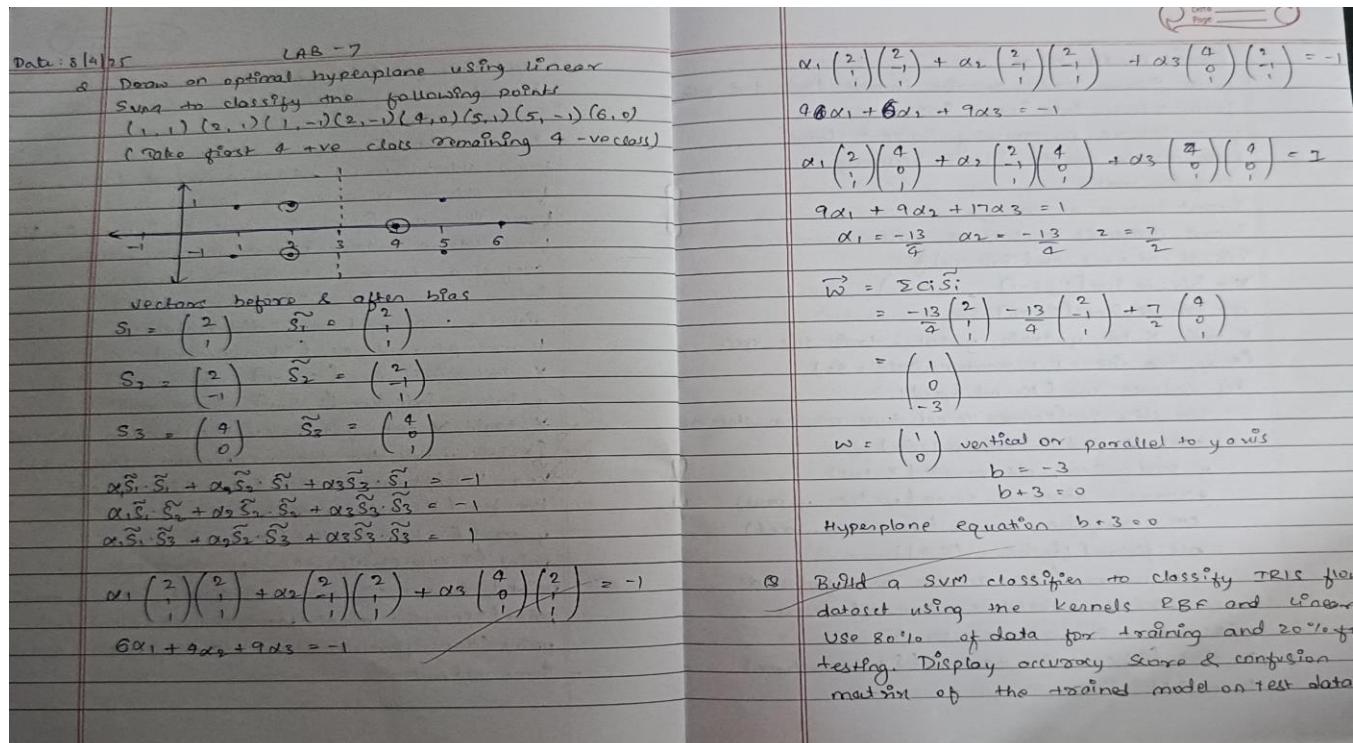
df3.head()
df3.isnull().sum()
X = df3.drop('target', axis=1)
y = df3['target']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42) best_k = 1
best_accuracy = 0
for k in range(1,
11):
    knn = KNeighborsClassifier(n_neighbors=k)
    knn.fit(X_train, y_train)
    y_pred = knn.predict(X_test)
    accuracy = accuracy_score(y_test, y_pred)
    print(f"Accuracy for k={k}: {accuracy}, Error Rate for k={k}: {1-
accuracy}") if accuracy > best_accuracy:
        best_accuracy = accuracy
        best_k = k
print(f"Best k value: {best_k}")
knn = KNeighborsClassifier(n_neighbors=optimal_k)
knn.fit(X_train, y_train)
y_pred = knn.predict(X_test)
print("Accuracy Score:", accuracy_score(y_test, y_pred))
print("\nConfusion Matrix:")
cm = confusion_matrix(y_test, y_pred)
print(cm)
print("\nClassification Report:")
print(classification_report(y_test, y_pred))
plt.figure(figsize=(8, 6))
sns.heatmap(cm, annot=True, fmt="d", cmap="Blues",
            xticklabels=knn.classes_, yticklabels=knn.classes_)
plt.xlabel('Predicted')
plt.ylabel('Actual')
plt.title('Confusion Matrix')
plt.show()

```

Program7

Build Support vector machine model for a given dataset

Screenshot:



Q Which kernel gave better performance on the Iris dataset? Why do you think so?
 Both the linear and RBF kernels gave perfect classification accuracy (100%) so there is no significant difference in performance for this dataset. However the Linear kernel is usually preferred in this case because:

- The Iris dataset is linearly separable, meaning classes can be separated using a straight line.
- The linear kernel is also faster and less prone to overfitting on small, clean datasets like Iris.

For letter-recognition.csv dataset

Q Present and interpret the confusion matrix. Are there any specific letters that are frequently confused with others?

- The confusion matrix shows the true letter classes vs predicted ones. Upon examining the matrix:
 - Most diagonal values are high, meaning the model correctly predicts most letters.
 - However some confusions may occur b/w visually similar letters such as D & O, T & I, M & N.

Q What is AUC score, & how does it reflect the model performance?

The AUC scores are as follows:

AUC	Model Performance
= 1	Perfect classifier
> 0.9	Excellent
between 0.7 - 0.9	Good
< 0.7	Model may struggle

How does the performance of the SVM model on this dataset compare to its performance on the Iris dataset?

Metric	Time	Letter-Recognition
Accuracy	$\sim 100\%$	lower, but still $\geq 90\%$
AUC score	N/A (3-class)	High AUC (~ 0.95)
Confusion Matrix	Clean, No overlap	Same confusion between letters

Date _____
Page _____

Code:

```

from google.colab import files
iris=files.upload()
df1=pd.read_csv("iris (1).csv")
df1.head()
X = df1.drop('species', axis=1)
y = df1['species']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)
rbf_svm = SVC(kernel='rbf')
rbf_svm.fit(X_train, y_train)
rbf_y_pred = rbf_svm.predict(X_test)
print("RBF Kernel SVM:")
print("Accuracy:", accuracy_score(y_test, rbf_y_pred))
cm = confusion_matrix(y_test, rbf_y_pred)
sns.heatmap(cm, annot=True, fmt='d', cmap="Blues")
plt.title('Confusion Matrix for RBF Kernel SVM')
plt.xlabel('Predicted')
plt.ylabel('True')
plt.show()
print(classification_report(y_test, rbf_y_pred))
linear_svm = SVC(kernel='linear')
linear_svm.fit(X_train, y_train)

```

```

linear_y_pred = linear_svm.predict(X_test)
print("\nLinear Kernel SVM:")
print("Accuracy:", accuracy_score(y_test, linear_y_pred))
cm = confusion_matrix(y_test, linear_y_pred)
sns.heatmap(cm, annot=True, fmt='d', cmap="Blues")
plt.title('Confusion Matrix for Linear Kernel SVM')
plt.xlabel('Predicted')
plt.ylabel('True')
plt.show()
print(classification_report(y_test, linear_y_pred))
letter=files.upload()
df2=pd.read_csv("letter-recognition.csv")
df2.head()
X = df2.drop('letter', axis=1)
y = df2['letter']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42) svm_classifier = SVC(kernel='linear', probability=True)
svm_classifier.fit(X_train, y_train)
y_pred =
svm_classifier.predict(X_test)
print("Accuracy:", accuracy_score(y_test, y_pred))
print(classification_report(y_test, y_pred))
cm = confusion_matrix(y_test, y_pred)
plt.figure(figsize=(10,10))
sns.heatmap(cm, annot=True, fmt='d', cmap="Blues")
plt.title('Confusion Matrix for SVM')
plt.xlabel('Predicted')
plt.ylabel('True')
plt.show()
lb = LabelBinarizer()
lb.fit(y_test)

```

```

y_test_lb = lb.transform(y_test)
y_pred_prob =
svm_classifier.predict_proba(X_test) fpr = { }
tpr = { }
thresh = { }
roc_auc = dict()
n_class = y_test_lb.shape[1]
for i in range(n_class):
    fpr[i], tpr[i], thresh[i] = roc_curve(y_test_lb[:,i], y_pred_prob[:,i])
    roc_auc[i] = auc(fpr[i], tpr[i])
plt.plot(fpr[0], tpr[0], linestyle='--', color='orange', label='SVM (AUC = %0.2f)' %
roc_auc[0]) plt.title('ROC Curve for Class 0')
plt.xlabel('False Positive Rate') plt.ylabel('True Positive rate')
plt.legend(loc='best')
plt.show()
print(f"AUC score for class 0: {roc_auc[0]}")

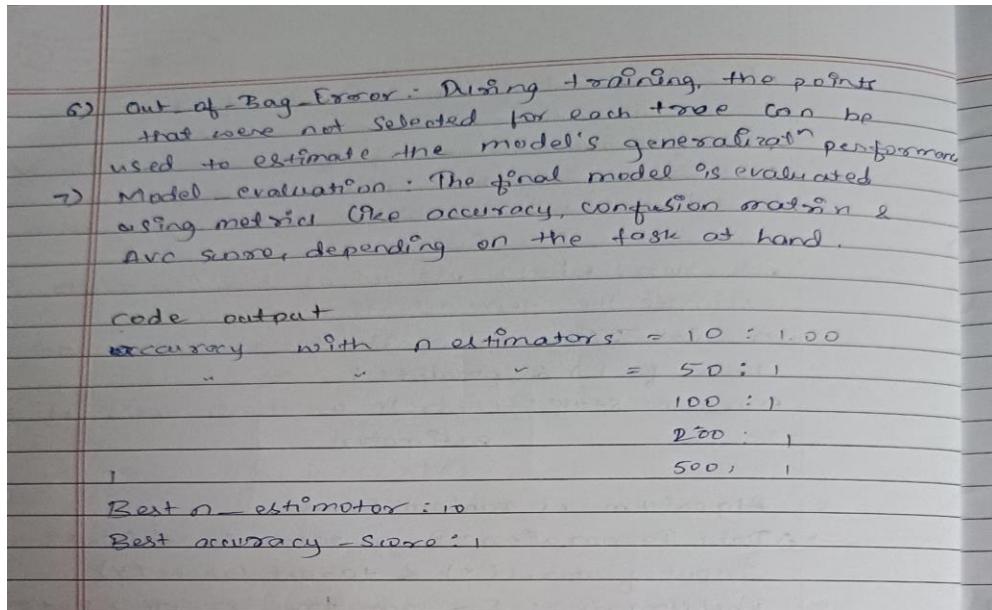
```

Program 8

Implement Random forest ensemble method on a given dataset

Screenshot:

<p>Q Get the difference in Decision tree & random forest</p> <p>→ Decision tree Random Forest</p> <ul style="list-style-type: none"> - A decision tree is a single tree structure where decision method that builds multiple trees by splitting the dataset at each node - These are highly prone to overfitting especially with complex datasets - A model has high variance and low bias - These are easy to interpret & visualize as they represent a sequence of decisions <p>Q Discuss all the parameters of RandomForestClassifier</p> <p>→ n_estimators: The no. of trees in the forest</p> <p>Criterion: The function to measure the quality of a split</p> <p>max_depth: The max depth of a tree</p> <p>min_samples_split: The min no. of samples required to split an internal node</p>	<p>min_samples_leaf: The min no. of samples required to be at a leaf node</p> <p>max_features: The no. of features to consider when looking for the best split</p> <p>Bootstrap: whether to use bootstrapping when creating trees.</p> <p>oob_score: whether use out of bag samples to estimate the generalization accuracy</p> <p>n_jobs: The no. of jobs to run in parallel of both fit() & predict()</p> <p>random_state: controls the randomness of the estimator</p> <p>1) Algorithm of random forest</p> <p>→ 1) Data Preparation: The datasets is divided into input features (X) & target labels (y)</p> <p>2) Bootstrapping: For each tree, create a bootstrap sample from the dataset</p> <p>3) Feature Selection: For each tree, randomly select a subset of features</p> <p>4) Building decision Trees: Build a decision tree on the bootstrap dataset using the selected subset of features</p> <p>5) Predicting: once all trees are built, make prediction by aggregating the predictions of all the trees</p>
--	--



Code:

```

from google.colab import files
iris=files.upload()
df1=pd.read_csv("iris (4).csv")
df1.head()
X = df1.drop('species', axis=1)
y = df1['species']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=0)
rf_classifier = RandomForestClassifier(random_state=0)
rf_classifier.fit(X_train, y_train)
y_pred =
rf_classifier.predict(X_test)
default_accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy with default n_estimators: {default_accuracy}")
best_accuracy = 0
best_n_estimators = 0
for n_estimators in range(1, 101):
    rf_classifier = RandomForestClassifier(n_estimators=n_estimators, random_state=0)
    rf_classifier.fit(X_train, y_train)
    y_pred = rf_classifier.predict(X_test)
    accuracy = accuracy_score(y_test, y_pred)
    if accuracy > best_accuracy:
        best_accuracy =
        best_n_estimators =
n_estimators
print(f"\nBest accuracy: {best_accuracy} achieved with n_estimators =
{best_n_estimators}")
cm = confusion_matrix(y_test, y_pred)
plt.figure(figsize=(8, 6))
sns.heatmap(cm, annot=True, fmt="d", cmap="Blues",
            xticklabels=np.unique(y_test), yticklabels=np.unique(y_test))
plt.xlabel("Predicted")
plt.ylabel("Actual")
plt.title("Confusion Matrix")
plt.show()

```

Program 9

Implement Boosting ensemble method on a given dataset

Screenshot:

LAB - 9

AdaBoost

Boosting combines multiple weak learners to create a strong learner. It works by training models sequentially where each model focuses on errors made by previous one.

Parameters

estimator	The base model
n_estimators	no. of weak learners
learning_rate	shinks contribut ⁿ of each learner
algorithm	SAMME or SAMME.K
random_state	for reproducibility

Algorithm

- Start with equal wts for all training sample
- Train a weak model
- later calculate error & update sample wts
- Add weak model to ensemble with a wt based on its accuracy
- Repeat for n-estimators
- Final prediction.

21/12/25

Code:

```
from google.colab import files  
income=files.upload()  
df1=pd.read_csv("income.csv")  
df1.head()  
X = df1.drop('income_level', axis=1)
```

```

y = df1['income_level'] X = pd.get_dummies(X)
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
random_state=42) abc = AdaBoostClassifier(n_estimators=10,
random_state=42)
abc.fit(X_train, y_train)
y_pred = abc.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print(f"Initial AdaBoost accuracy (10 trees): {accuracy}")
param_grid = {'n_estimators': [50, 100, 150, 200]}
grid_search = GridSearchCV(AdaBoostClassifier(random_state=42), param_grid, cv=5,
scoring='accuracy') grid_search.fit(X_train, y_train)
print(f"Best parameters: {grid_search.best_params_}")
print(f"Best cross-validation score: {grid_search.best_score_}")
best_abc = grid_search.best_estimator_
y_pred_best = best_abc.predict(X_test)
best_accuracy = accuracy_score(y_test,
y_pred_best)
print(f"Accuracy of the best model on the test set: {best_accuracy}")
cm = confusion_matrix(y_test, y_pred_best)
plt.figure(figsize=(8, 6))
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues',
xticklabels=['<=50K', '>50K'], yticklabels=['<=50K', '>50K'])
plt.xlabel('Predicted')
plt.ylabel('Actual')
plt.title('Confusion Matrix')
plt.show()

```

Program 10

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

Screenshot:

LAB - 10

Page

K means algorithm

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

How to determine no of clusters?

Common Methods are Elbow Method, Silhouette Score & Domain Knowledge

Formula for Sum of Squared errors

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

Elbow Technique - Plot SSE vs k
Choose the k where SSE starts to flatten - this is the "elbow"

~~Notes~~ Key parameters in Kmeans()

- n_clusters : NO of clusters to form
- init : Initialization method

n_init : Number of initialization
 max_iter : Max Iterations per run
 random_state : Controls randomness
 tol : convergence threshold
 algorithm : Algo to use ('loyd', 'elkan')

- cluster eight points (with (x, y) representing location) into 3 clusters. $A_1(2,10)$, $A_2(2,5)$, $A_3(8,4)$, $A_4(5,8)$, $A_5(7,5)$, $A_6(6,4)$, $A_7(1,2)$, $A_8(9,9)$. Initial clusters: $C_1(2,10)$, $C_2(5,8)$, $C_3(1,2)$

→ Calculation of distance b/w point $A_1(2,10)$ and $C_1(2,10)$:

$$D(A_1, C_1) = |x_2 - x_1| + |y_2 - y_1| = |2 - 2| + |10 - 10|$$

Calculating distance b/w $A_1(2,10)$ & $C_2(5,8)$:

$$D(A_1, C_2) = |5 - 2| + |8 - 10| = 3 + 2 = 5$$

Calculating distance b/w $A_1(2,10)$ & $C_3(1,2)$:

$$D(A_1, C_3) = |2 - 1| + |10 - 2| = 8 + 1 = 9$$

Now calculate the distance of other points from each of the center of 3 clusters.

For cluster - 01

- We have only one point $A_1(2,10)$ in cluster - 01
- So cluster center remains the same

Iteration-01	Given Pts	Distance from			Point Belongs to
		C ₁	C ₂	C ₃	
	$A_1(2,10)$	0	5	9	C ₁
	$A_2(2,5)$	5	6	4	C ₃
	$A_3(8,4)$	12	7	9	C ₂
	$A_4(5,8)$	5	0	10	C ₂
	$A_5(7,5)$	10	5	9	C ₂
	$A_6(6,4)$	10	5	7	C ₂
	$A_7(1,2)$	9	10	0	C ₃
	$A_8(9,9)$	2	2	10	C ₂

For cluster - 02

$$\text{Center of cluster - 02} = ((8+5+7+6)/4, (4+8+5+4)/4) = (6, 6)$$

For cluster - 03

$$\text{Center of cluster - 03} = ((2+1)/2, (5+2)/2) = (1.5, 3.5)$$

This is completion of Iteration - 01

Iteration-02	Given Pts	Distance from			Point Belongs to
		C ₁	C ₂	C ₃	
	$A_1(2,10)$	0	8	7	C ₁
	$A_2(2,5)$	5	5	2	C ₃
	$A_3(8,4)$	12	4	7	C ₂
	$A_4(5,8)$	5	3	8	C ₂
	$A_5(7,5)$	10	2	7	C ₂
	$A_6(6,4)$	10	2	5	C ₂
	$A_7(1,2)$	9	9	2	C ₃
	$A_8(9,9)$	3	5	8	C ₁

PAGE NO:	DATE:			
For c ₁	Center = $((2+0)/2, (10+9)/2) = (3, 9.5)$			
For c ₂	Center = $((8+5+7+6)/4, (4+8+5+4)/4) = (6, 6)$			
For c ₃	Center = $((2+1)/2, (5+2)/2) = (1.5, 3.5)$			

Iteration - 03

Given Pts	Distance from			Point belongs to cluster
	from C ₁	from C ₂	from C ₃	
$A_1(2,10)$	1.5	9.25	7	C ₁
$A_2(2,5)$	5	4.75	2	C ₃
$A_3(8,4)$	10.5	2.75	7	C ₂
$A_4(5,8)$	3.5	4.75	8	C ₁
$A_5(7,5)$	8.5	0.75	7	C ₂
$A_6(6,4)$	8.5	1.75	5	C ₂
$A_7(1,2)$	9.5	8.75	2	C ₃
$A_8(9,9)$	1.5	6.25	8	C ₁

$$\text{Centres } (2+5+4)/3 = 3.66, (10+8+9)/3 = 9$$

$$\text{For } C_2 \quad (8+7+6)/3 = 7, (4+8+5+4)/3 = 6.33$$

$$\text{For } C_3 \quad (2+1)/3 = 1.5, 3.5$$

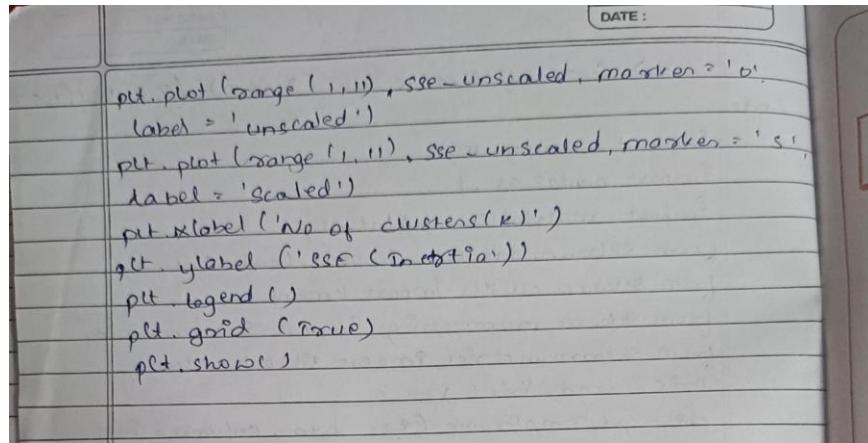
Given Pts	Distance from			Point
	D from C ₁	D from C ₂	D from C ₃	
$A_1(2,10)$	1.56	10.67	8	C ₁
$A_2(2,5)$	5.66	5.67	2	C ₃
$A_3(8,4)$	9.34	1.33	7	C ₂
$A_4(5,8)$	2.34	5.67	8	C ₁
$A_5(7,5)$	7.34	0.67	7	C ₂
$A_6(6,4)$	7.34	1.33	6	C ₂
$A_7(1,2)$	9.66	8.33	2	C ₃
$A_8(9,9)$	0.34	7.66	8	C ₁

```

PAGE NO:
DATE:
→ Center of clusters are (3.66, 9), (7, 4.33), (1.53, 3.5)

code:
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.datasets import load_iris
from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import silhouette_score
iris = load_iris()
df = pd.DataFrame(iris.data, columns=iris.feature_names)
ds['target'] = iris.target
X = df[['petal length (cm)', 'petal width (cm)']]
sse_unscaled = []
for k in range(1, 11):
    kmeans = KMeans(n_clusters=k, random_state=42)
    kmeans.fit(X)
    sse_unscaled.append(kmeans.inertia_)
sse_unscaled = np.array(sse_unscaled)
Scalen = StandardScaler()
X_Scaled = Scalen.fit_transform(X)
for k in range(1, 11):
    kmeans = KMeans(n_clusters=k, random_state=42)
    kmeans.fit(X_Scaled)
    sse_scaled.append(kmeans.inertia_)
sse_scaled = np.array(sse_scaled)
plt.figure(figsize=(10, 5))

```



Code:

```

from google.colab import files
iris=files.upload()
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from scipy import stats
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
from sklearn.metrics import classification_report, confusion_matrix,
accuracy_score from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler

df1=pd.read_csv("iris (4).csv")
df1.head()
df = df1.drop(['sepal_length','sepal_width','species'],axis=1)
scaler = StandardScaler()
scaled_df = scaler.fit_transform(df)
wcss = []
for i in range(1, 11):
    kmeans = KMeans(n_clusters=i, init='k-means++', max_iter=300, n_init=10,
                    random_state=0)
    kmeans.fit(scaled_df)
    wcss.append(kmeans.inertia_)
plt.plot(range(1, 11), wcss)
plt.title('Elbow Method')
plt.xlabel('Number of clusters')
plt.ylabel('WCSS')
plt.show()

kmeans = KMeans(n_clusters=3, init='k-means++', max_iter=300, n_init=10,
                random_state=0)
pred_y = kmeans.fit_predict(scaled_df)
df['cluster'] = pred_y
plt.scatter(df['petal_length'], df['petal_width'], c=df['cluster'])
plt.title('Clusters of Iris Flowers')
plt.xlabel('Petal Length')
plt.ylabel('Petal Width')
plt.show()

```

Program 11

Implement Dimensionality reduction using Principal Component Analysis (PCA)

method.

Screenshot:

Date _____
Page _____

LAB - 11
Build a PCA on a given dataset (Steps)

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

Feature	Ex 1	Ex 2	Ex 3	Ex 4	reduce 2D
x_{11}	4	8	13	7	to 1D
x_{12}	11	4	5	14	

1. $\bar{x}_1 = (4+8+13+7)/4 = 8$
2. $\bar{x}_2 = (11+4+5+14)/4 = 8.5$
3. $\text{cov}(x_1, x_2) = \frac{1}{N-1} \sum_{k=1}^N (x_{1k} - \bar{x}_1)(x_{2k} - \bar{x}_2)$
4. $= \frac{1}{3} ((4-8)^2 + (8-8)^2 + (13-8)^2 + (7-8)^2) = 14$
5. $\text{cov}(x_1, x_1) = \frac{1}{N-1} \sum_{k=1}^N (x_{1k} - \bar{x}_1)^2$
6. $= \frac{1}{3} ((4-8)^2 + (8-8)^2 + (13-8)^2 + (7-8)^2) = -11$
7. $\text{cov}(x_2, x_1) = \text{cov}(x_1, x_2) = -11$
8. $\text{cov}(x_2, x_2) = \frac{1}{N-1} \sum_{k=1}^N (x_{2k} - \bar{x}_2)^2$
9. $= \frac{1}{3} ((11-8.5)^2 + (4-8.5)^2 + (5-8.5)^2 + (14-8.5)^2) = 23$

$S = \begin{bmatrix} 14 & -11 \\ -11 & 23 \end{bmatrix}$

3. $0 = \det(S - \lambda I) \Rightarrow \begin{vmatrix} 14-\lambda & -11 \\ -11 & 23-\lambda \end{vmatrix}$

$\lambda^2 - 37\lambda + 201 \Rightarrow \lambda_1 = 30.3809$
 $\lambda_2 = 6.6151$

4. $\begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 14-\lambda_1 & -11 \\ -11 & 23-\lambda_1 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} = \begin{bmatrix} (14-\lambda_1)u_1 - 11u_2 \\ -11u_1 + (23-\lambda_1)u_2 \end{bmatrix}$

$(14-\lambda_1)u_1 - 11u_2 = 0$
 $-11u_1 + (23-\lambda_1)u_2 = 0$

$\frac{u_1}{11} = \frac{u_2}{14-\lambda_1} = t, \Rightarrow u_1 = 11t, u_2 = (14-\lambda_1)t$

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

Taking $t = 1$ we get eigen vector

for unit eigen vector $\|u_i\| = \sqrt{11^2 + (14-\lambda_1)^2}$

$\|u_1\| = 19.7348$

$\therefore e_1 = \begin{bmatrix} 11/19.7348 \\ (14-\lambda_1)/19.7348 \end{bmatrix} = \begin{bmatrix} 0.5574 \\ -0.8303 \end{bmatrix}$

$\|u_2\| = 0.8303$

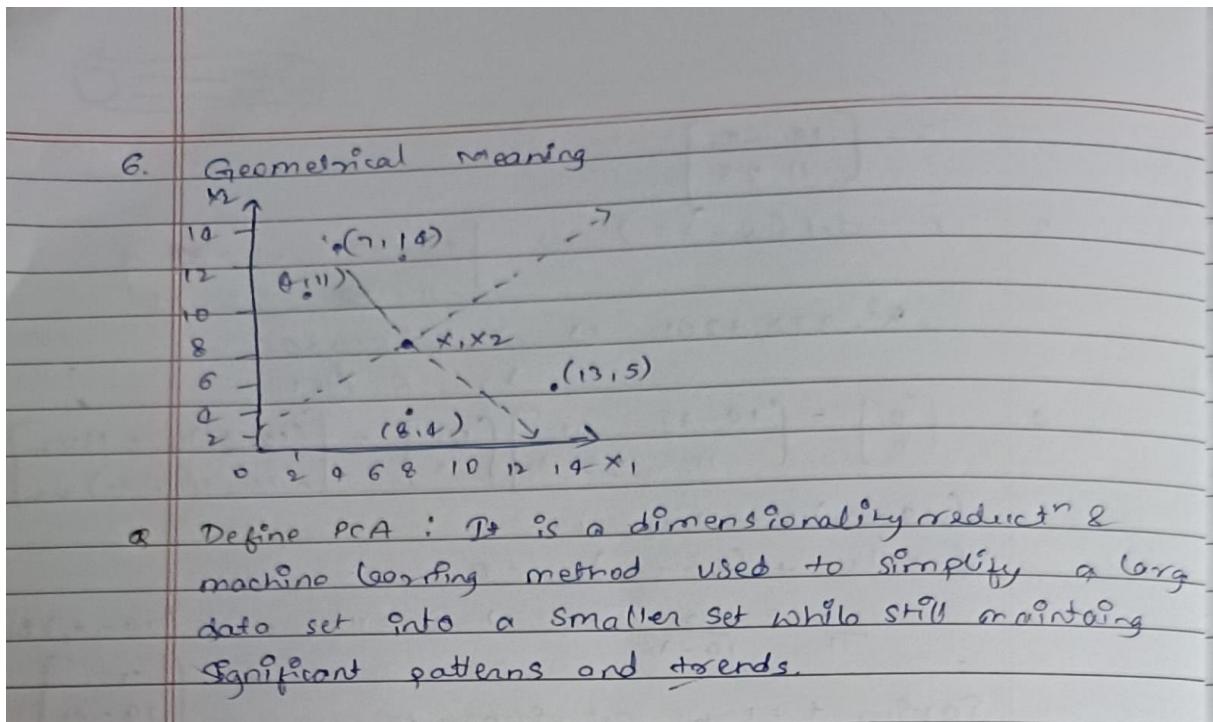
$e_2 = \begin{bmatrix} 0.8303 \\ 0.5574 \end{bmatrix}$

5. $PC \Rightarrow \begin{bmatrix} x_{1k} \\ x_{2k} \end{bmatrix}, e_1^T \begin{bmatrix} x_{1k} - \bar{x}_1 \\ x_{2k} - \bar{x}_2 \end{bmatrix} = 0.5574(x_{1k} - \bar{x}_1) - 0.8303(x_{2k} - \bar{x}_2)$

$\Rightarrow \begin{bmatrix} x_{11} \\ x_{21} \end{bmatrix} = \begin{bmatrix} 4 \\ 11 \end{bmatrix}$

After substituting $= -4.30535$

	4	8	13	7
x_1	4	8	13	7
x_2	11	4	5	14
First PC	-4.305	3.76	5.692	-5.1238



Code:

```

from google.colab import files
heart=files.upload()

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from scipy import stats
import seaborn as sns
from sklearn.preprocessing import LabelEncoder,
OneHotEncoder from sklearn.model_selection import
train_test_split
from sklearn.metrics import accuracy_score
from sklearn.metrics import classification_report, confusion_matrix,
accuracy_score from sklearn.preprocessing import StandardScaler
from sklearn.svm import SVC
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier
from sklearn.decomposition import PCA

df1=pd.read_csv("heart (1).csv")
df1.head()
text_cols = df1.select_dtypes(include=['object']).columns
label_encoder = LabelEncoder()
for col in text_cols:
    
```

```

df1[col] =
label_encoder.fit_transform(df1[col])
print(df1.head())
X = df1.drop('HeartDisease', axis=1)
y = df1['HeartDisease']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42) scaler = StandardScaler()
X_train =
scaler.fit_transform(X_train) X_test =
scaler.transform(X_test)
# Support Vector Machine
svm_model = SVC(kernel='linear', random_state=42)
svm_model.fit(X_train, y_train)
svm_predictions = svm_model.predict(X_test)
svm_accuracy = accuracy_score(y_test, svm_predictions)
print(f'SVM Accuracy: {svm_accuracy}')

# Logistic Regression
lr_model = LogisticRegression(random_state=42)
lr_model.fit(X_train, y_train) lr_predictions =
lr_model.predict(X_test) lr_accuracy =
accuracy_score(y_test, lr_predictions)
print(f'Logistic Regression Accuracy: {lr_accuracy}')

# Random Forest
rf_model = RandomForestClassifier(random_state=42)
rf_model.fit(X_train, y_train)
rf_predictions = rf_model.predict(X_test)
rf_accuracy = accuracy_score(y_test, rf_predictions)
print(f'Random Forest Accuracy: {rf_accuracy}')

models = {
    "SVM": svm_accuracy,
    "Logistic Regression": lr_accuracy,
    "Random Forest": rf_accuracy
}

best_model = max(models, key=models.get)
print(f"\nBest Model: {best_model} with accuracy {models[best_model]}")
pca = PCA(n_components=0.95)
X_train_pca = pca.fit_transform(X_train)
X_test_pca = pca.transform(X_test)

svm_model_pca = SVC(kernel='linear', random_state=42)
svm_model_pca.fit(X_train_pca, y_train)
svm_predictions_pca = svm_model_pca.predict(X_test_pca)
svm_accuracy_pca = accuracy_score(y_test, svm_predictions_pca)
print(f'SVM Accuracy (with PCA): {svm_accuracy_pca}')

lr_model_pca = LogisticRegression(random_state=42)

```

```

lr_model_pca.fit(X_train_pca, y_train)
lr_predictions_pca = lr_model_pca.predict(X_test_pca)
lr_accuracy_pca = accuracy_score(y_test, lr_predictions_pca)
print(f"Logistic Regression Accuracy (with PCA): {lr_accuracy_pca}")

rf_model_pca = RandomForestClassifier(random_state=42)
rf_model_pca.fit(X_train_pca, y_train)
rf_predictions_pca = rf_model_pca.predict(X_test_pca)
rf_accuracy_pca = accuracy_score(y_test, rf_predictions_pca)
print(f"Random Forest Accuracy (with PCA): {rf_accuracy_pca}")

models_pca = {
    "SVM": svm_accuracy_pca,
    "Logistic Regression": lr_accuracy_pca,
    "Random Forest": rf_accuracy_pca
}

best_model_pca = max(models_pca, key=models_pca.get)
print(f"\nBest Model (with PCA): {best_model_pca} with accuracy {models_pca[best_model_pca]}")

```