

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



LAB REPORT
on
Machine Learning (23CS6PCMAL)

Submitted by

Shreya C Y (1BM22CS265)

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 **Shreya CY (1BM22CS265)**, who is bonafide student of **B.M.S. College of Engineering**. It is in partial fulfillment for the award of **Bachelor of Engineering in Computer Science and Engineering** of the Visvesvaraya Technological University, Belgaum. The Lab report has been approved as it satisfies the academic requirements in respect of an Machine Learning (23CS6PCMAL) work prescribed for the said degree.

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

Index

Sl. No.	Date	Experiment Title	Page No.
1	21-2-2025	Write a python program to import and export data using Pandas library functions	1
2	3-3-2025	Demonstrate various data pre-processing techniques for a given dataset	4
3	10-3-2025	Implement Linear and Multi-Linear Regression algorithm using appropriate dataset	6
4	17-3-2025	Build Logistic Regression Model for a given dataset	9
5	24-3-2025	Use an appropriate data set for building the decision tree (ID3) and apply this knowledge to classify a new sample	13
6	7-4-2025	Build KNN Classification model for a given dataset	17
7	21-4-2025	Build Support vector machine model for a given dataset	21
8	5-5-2025	Implement Random forest ensemble method on a given dataset	22
9	5-5-2025	Implement Boosting ensemble method on a given dataset	23
10	12-5-2025	Build k-Means algorithm to cluster a set of data stored in a .CSV file	24
11	12-5-2025	Implement Dimensionality reduction using Principal Component Analysis (PCA) method	26

Github Link:

<https://github.com/Shreyacy/ML-LABB>

Program 1

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

Screenshot

The image shows handwritten notes on two pages of lined paper. The left page contains code for HDFC and ICICI stocks, while the right page contains code for KOTAK stock. The notes include comments explaining the steps and imports.

Left Page (HDFC and ICICI Stock Data):

- HDFC-data['Daily Returns'].plot(title="HDFC Industries - daily returns", color='orange')
plt.tight_layout()
plt.show()
print("In HDFC stock data saved")
- ICICI Bank:
1. ICICI-data = data[ICICIBANK.NS]
print("In ICICI stock data saved")
2. ICICI-data['Daily Returns'].plot()
print("In ICICI Industries - daily returns")
3. ICICI-data.to_csv("HDFC-stock-data.csv")
print("In HDFC stock data saved")

Right Page (KOTAK Stock Data):

- To do 2:
1. Import yfinance as yf
import pandas as pd
Import matplotlib.pyplot as plt
tickers = ["HDFCBANK.NS", "KOTAKBANK.NS",
ICICIBANK.NS] for i in range(3):
data = yf.download(tickers, start="2022-01-01",
end="2023-01-01", group_by='ticker')
print("First 5 rows of the dataset")
print(data.head())
2. print("In shape of the dataset")
print(data.shape)
3. print("In column names of dataset")
print(data.columns)
- HDFC-Bank:
1. HDFC-data = data['HDFCBANK.NS']
print("In summary statistics")
print(HDFC-data.describe())
HDFC-data['dly returns'] = HDFC-data['close'].pct_change()
2. plt.figure(figsize=(12,6))
plt.subplot(2,1,1)
HDFC-data['Close'].plot(title="HDFC Industries - closing Price")
plt.subplot(2,1,2)

PAGE NO.: 1
DATE: 5/3/25.

Lab-0

To do!

1. Import pandas as pd

```
data = pd.read_csv('USN.csv')
USN = [100, 101, 102, 103]
Names = ['Alice', 'Bob', 'Charlie', 'David'],
Marks = [90, 92, 95, 97]
```
2. from sklearn.datasets import load_diabetes

```
diabetes = load_diabetes()
df = pd.DataFrame(diabetes.data, columns=diabetes.feature_names)
df['target'] = diabetes.target
print("Sample data:")
print(df.head())
```
3. file_path = 'sample-sales-data.csv'

```
df = pd.read_csv(file_path)
print("Sample data:")
print(df.head())
print("n")
print(df.describe())
print("n")
print("n")
```
4. file_path = '10babies Dataset.zip'

```
df = pd.read_csv(file_path)
print("Sample data:")
print(df.head())
print("n")
print("n")
```

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())
```

```
import pandas as pd
data = {
    'USN': [100, 101, 102, 103],
    'Name': ['Alice', 'Bob', 'Charlie', 'David'],
    'Marks': [25, 30, 35, 40],
}
```

```

df=pd.DataFrame(data)
print(df.head())

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

file_path = '/content/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')
print("Sample data:")
print(df.head())


import yfinance as yf
import pandas as pd
import matplotlib.pyplot as plt
tickers = ["HDFCBANK.NS", "ICICIBANK.NS", "KOTAKBANK.NS"]
data = yf.download(tickers, start="2024-01-01", end="2024-12-30", group_by='ticker')

print("First 5 rows of the dataset:")
print(data.head())


hdfc_data=data['HDFCBANK.NS']
icici_data=data['ICICIBANK.NS']
kotak_data=data['KOTAKBANK.NS']

print("\nSummary statistics for Reliance Industries:")
print(hdfc_data.describe())
hdfc_data['Daily Return'] = hdfc_data['Close'].pct_change()
plt.figure(figsize=(12, 6))
plt.subplot(2, 1, 1)
hdfc_data['Close'].plot(title="HDFC - Closing Price")
plt.subplot(2, 1, 2)
hdfc_data['Daily Return'].plot(title="HDFC - Daily Returns", color='orange')
plt.tight_layout()
plt.show()
print("\nSummary statistics for Reliance Industries:")
print(icici_data.describe())

```

```

icici_data['Daily Return'] = icici_data['Close'].pct_change()
plt.figure(figsize=(12, 6))
plt.subplot(2, 1, 1)
icici_data['Close'].plot(title="ICICI - Closing Price")
plt.subplot(2, 1, 2)
icici_data['Daily Return'].plot(title="ICICI - Daily Returns", color='orange')
plt.tight_layout()
plt.show()

```

```

print("\nSummary statistics for Reliance Industries:")
print(kotak_data.describe())
kotak_data['Daily Return'] = kotak_data['Close'].pct_change()
plt.figure(figsize=(12, 6))
plt.subplot(2, 1, 1)
kotak_data['Close'].plot(title="kotak - Closing Price")
plt.subplot(2, 1, 2)
kotak_data['Daily Return'].plot(title="kotak - Daily Returns", color='orange')
plt.tight_layout()
plt.show()

```

Program 2

Demonstrate various data pre-processing techniques for

Lab-1	Lab-2	Lab-3
<pre> 1. import pandas as pd import numpy as np df = pd.read_csv('adult.csv') print(df) 2. df.describe() 3. df['Ocean_proximity'].value_counts() 4. df['Ocean_proximity'].nunique() 5. miv_count = df['Percuile'].sum() columns_miv = miv_count[miv_count > 0] print(columns_miv) 6. df = pd.read_csv('adult/adult.csv') print(df) 7. df = pd.read_csv('adult/adult.csv') df.head(5) 8. miv_by_col = adult['Age'].mean() print(j"Column with missing value :") missing_col = df['Age'].isnull().sum() print(missing_col) 9. categorical_col = adult['Age'].dtypes (Pandas is object type) </pre>	<pre> PAGE NO. 4 DATE 5/6/25 categorical - col is diabetes = df['diabetes'].dtypes Include 'object' in column. 1. If my column in dataset has missing value we can replace categorical value with mode and numerical value with median. 2. In diabetes column are gender & class value in adult categorical column are workers, off education, marital status, occupation, relationship, race, gender. We are ordinal encoding to encode categorical column. 3. Min-max is also called normalization, transform data to fit into a specific range (0 to 1) for example. Standardization reduce data by subtracting the mean and dividing by standard deviation. </pre>	<pre> PAGE NO. 5 DATE 5/6/25 8.2: categorical variables transformation Categorical variables are object type. 1. If my column in dataset has missing value we can replace categorical value with mode and numerical value with median. 2. In diabetes column are gender & class value in adult categorical column are workers, off education, marital status, occupation, relationship, race, gender. We are ordinal encoding to encode categorical column. 3. Min-max is also called normalization, transform data to fit into a specific range (0 to 1) for example. Standardization reduce data by subtracting the mean and dividing by standard deviation. 9.1: Numerical variables - continuous For numerical variable, standardization Normalizing Data - 0 to 1 Z-score = (x - mean)/std </pre>

given dataset

Screenshot

Code:

```
import pandas as pd
import numpy as np
df = pd.read_csv('/content/housing.csv')
print(df)

df.describe()

df['ocean_proximity'].value_counts()

df['ocean_proximity'].nunique()

mv=df.isnull().sum()
cmv=mv[mv>0]
print(cmv)

df4=pd.read_csv('/content/adult.csv')
df4.head()

df5=pd.read_csv('/content/Dataset_of_Diabetes.csv')
df5.head()

missing_cols_adult = df4.columns[df4.isnull().sum() > 0]
print(f"Columns with missing values in 'adult.csv': {missing_cols_adult.tolist()}")
missing_cols_diabetes = df5.columns[df5.isnull().sum() > 0]
print(f"Columns with missing values in 'Dataset of Diabetes .csv': {missing_cols_diabetes.tolist()")

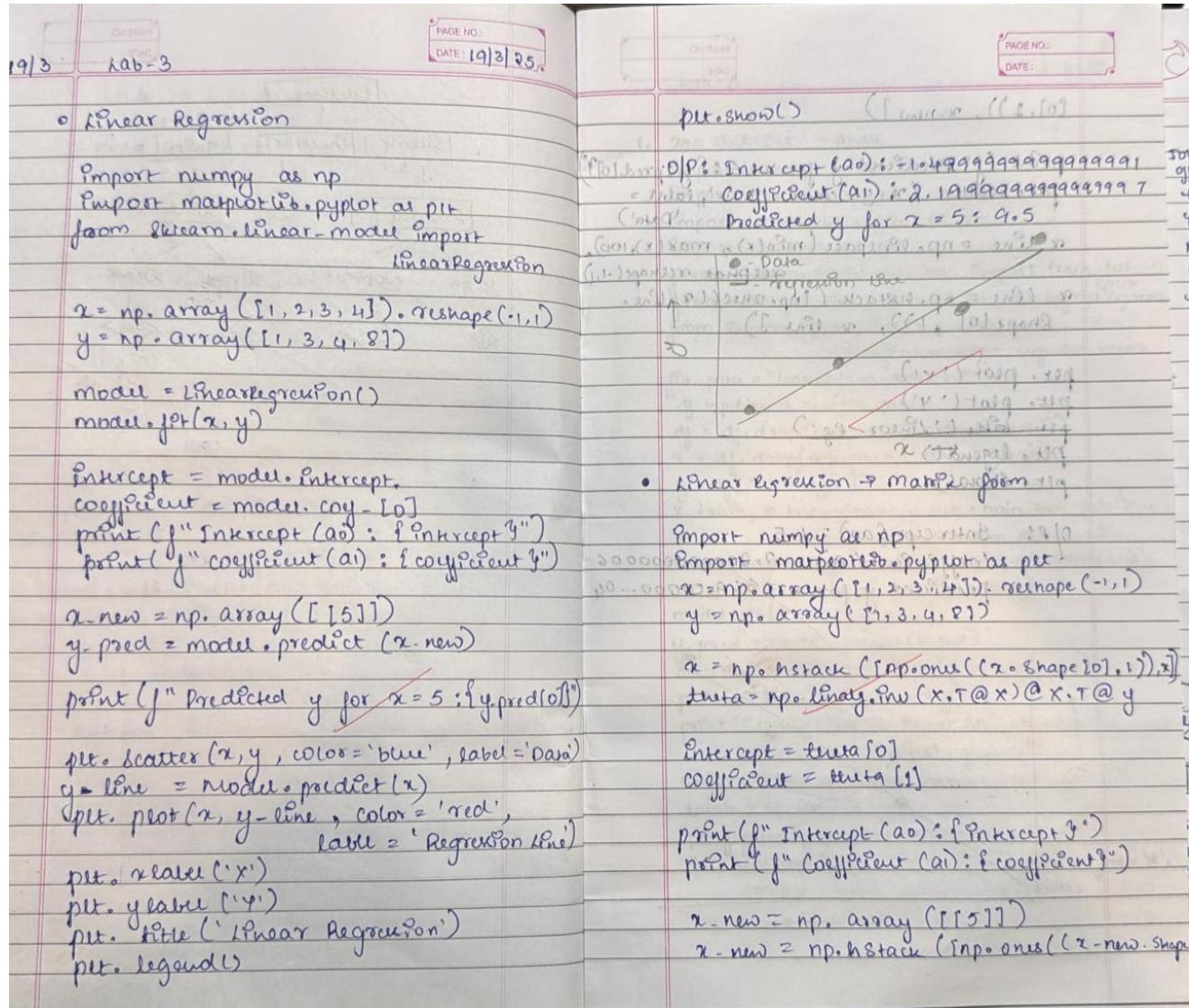
categorical_cols_adult = df4.select_dtypes(include=['object']).columns
print(f"Categorical columns in 'adult.csv': {categorical_cols_adult.tolist()")

categorical_cols_diabetes = df5.select_dtypes(include=['object']).columns
print(f"Categorical columns in 'Dataset of Diabetes .csv': {categorical_cols_diabetes.tolist()")
```

Program 3

Implement Linear and Multi-Linear Regression algorithm using appropriate dataset

Screenshot



[0], 1)), x-new])
 print(f"predicted y for x={x}: {y_pred[0]}")
 plt.scatter(x, y, color='blue', label='Data')
 x_line = np.linspace(min(x), max(x), 100)
 x_line = np.meshgrid(x_line, shape[0], 1), x_line])
 plt.plot(x_line)
 plt.plot(y)
 plt.title('Linear reg.')
 plt.legend()
 plt.show()

O/P:
 Intercept (a0) = 1.05
 coefficient (a1) = 2.0000000000000006
 predicted y for x=3.5 is 9.0000...04

~~$y = a_0 + a_1 x$
 $y = 1.05 + 2.0000000000000006 \cdot x$
 $y = 1.05 + 2.0000000000000006 \cdot 3.5$
 $y = 1.05 + 7.000000000000001$
 $y = 8.05$~~

~~$y = a_0 + a_1 x$
 $y = 1.05 + 2.0000000000000006 \cdot x$
 $y = 1.05 + 2.0000000000000006 \cdot 3.5$
 $y = 1.05 + 7.000000000000001$
 $y = 8.05$~~

~~$y = a_0 + a_1 x$
 $y = 1.05 + 2.0000000000000006 \cdot x$
 $y = 1.05 + 2.0000000000000006 \cdot 3.5$
 $y = 1.05 + 7.000000000000001$
 $y = 8.05$~~

~~$y = a_0 + a_1 x$
 $y = 1.05 + 2.0000000000000006 \cdot x$
 $y = 1.05 + 2.0000000000000006 \cdot 3.5$
 $y = 1.05 + 7.000000000000001$
 $y = 8.05$~~

Code:

```
import numpy as np  
import matplotlib.pyplot as plt  
from sklearn.linear_model import LinearRegression  
x = np.array([1, 2, 3, 4]).reshape(-1, 1)  
y = np.array([1, 3, 4, 8])
```

```

model = LinearRegression()
model.fit(x, y)
intercept = model.intercept_
coefficient = model.coef_[0]
print(f"Intercept (a0): {intercept}")
print(f"Coefficient (a1): {coefficient}")
x_new = np.array([[5]])
y_pred = model.predict(x_new)
print(f"Predicted y for x=5: {y_pred[0]}")
plt.scatter(x, y, color='blue', label='Data')
y_line = model.predict(x)
plt.plot(x, y_line, color='red', label='Regression Line')
plt.xlabel('X')
plt.ylabel('Y')
plt.title('Linear Regression')
plt.legend()
plt.show()

import numpy as np
import matplotlib.pyplot as plt
x = np.array([1, 2, 3, 4]).reshape(-1, 1)
y = np.array([1, 3, 4, 8])
X = np.hstack([np.ones((x.shape[0], 1)), x])
theta = np.linalg.inv(X.T @ X) @ X.T @ y
intercept = theta[0]
coefficient = theta[1]
print(f"Intercept (a0): {intercept}")
print(f"Coefficient (a1): {coefficient}")
x_new = np.array([[5]])
X_new = np.hstack([np.ones((x_new.shape[0], 1)), x_new])
y_pred = X_new @ theta

print(f"Predicted y for x=5: {y_pred[0]}")
plt.scatter(x, y, color='blue', label='Data')
x_line = np.linspace(min(x), max(x), 100).reshape(-1, 1)
X_line = np.hstack([np.ones((x_line.shape[0], 1)), x_line])
y_line = X_line @ theta

plt.plot(x_line, y_line, color='red', label='Regression Line')
plt.xlabel('X')
plt.ylabel('Y')
plt.title('Linear Regression')
plt.legend()
plt.show()

```

Program 4

Build Logistic Regression Model for a given dataset

Screenshot

PAGE NO.: 2 - 24 - 25

```

Lab - 4 Group 11
1. 200 dataset - multi.
Import numpy as np
Import pandas as pd
Import matplotlib.pyplot as plt
Import Seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score,
classification_report, confusion_matrix
file_path = "content/200-data.csv"
df = pd.read_csv(file_path)
df = df.drop(columns=['animal name'])
x = df.drop(columns=['class-type'])
y = df['class-type']
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',
solver='lbfgs', max_iter=5000)
model.fit(x_train, y_train)
y_pred = model.predict(x_test)
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy: ", accuracy)
print("Classification Report:\n", classification_report(y_test, y_pred))
conf_matrix = confusion_matrix(y_test, y_pred)
labels = np.unique(y)
plt.figure(figsize=(8, 6))
sns.heatmap(conf_matrix, annot=True, fmt='d',
cmap='Blues', xticklabels=labels,
yticklabels=labels)

```

Code:

```

import pandas as pd
from matplotlib import pyplot as plt
import math
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression

```

Load HR dataset

```

df = pd.read_csv("/content/HR_comma_sep (1).csv")
df.head()

```

```

plt.scatter(df.satisfaction_level, df.left, marker='+', color='red')

```

Splitting dataset into training and testing sets

```

X_train, X_test, y_train, y_test = train_test_split(df[['satisfaction_level']], df.left, train_size=0.9,
random_state=10)
X_train.shape

```

X_test

Training logistic regression model

```

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([[0.5]])

# model.coef_ indicates value of m in y=m*x + b equation
model.coef_

# model.intercept_ indicates value of b in y=m*x + b equation
model.intercept_

# Define sigmoid function and do the math manually
def sigmoid(x):
    return 1 / (1 + math.exp(-x))

def prediction_function(satisfaction_level):
    z = model.coef_[0][0] * satisfaction_level + model.intercept_[0]
    y = sigmoid(z)
    return y

satisfaction_level = 0.35
probability = prediction_function(satisfaction_level)
print(f"Probability of leaving: {probability:.2f}")

if probability >= 0.5:
    print("The employee will leave.")
else:
    print("The employee will stay.")

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

# Load dataset
file_path = "/content/zoo-data (1).csv"
df = pd.read_csv(file_path)

# Drop the animal_name column as it is not a feature
df = df.drop(columns=["animal_name"])

# Define features and target

```

```

X = df.drop(columns=["class_type"])
y = df["class_type"]

# Split data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Train multinomial logistic regression model
model = LogisticRegression(multi_class='multinomial', solver='lbfgs', max_iter=500)
model.fit(X_train, y_train)

# Make predictions
y_pred = model.predict(X_test)

# Evaluate model
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy:", accuracy)
print("Classification Report:\n", classification_report(y_test, y_pred))

# Confusion matrix
conf_matrix = confusion_matrix(y_test, y_pred)
labels = np.unique(y)
plt.figure(figsize=(8, 6))
sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues', xticklabels=labels, yticklabels=labels)
plt.xlabel('Predicted')
plt.ylabel('Actual')
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

PAGE NO. 6
DATE: 12/3/25

Lab-2

- Import pandas as pd
Import numpy as np
from collections import Counter
- file_path = "/content/WEATHER.csv"
df = pd.read_csv(file_path)
print("Dataset loaded successfully")
print(df.head(5))

O/P: Dataset loaded successfully
- Imports numpy as np
df_entropy = entropy(df)
label_count = Counter(df['label'])
total_instances = len(df)
entropy_value = -sum((label_count[i]/total_instances) * np.log2(label_count[i] / total_instances)) for i in label_count
values
return entropy_value
- dataset_entropy = entropy(df)
- def information_gain(data, attribute):
 total_entropy = entropy(data)
 attribute_values = data[attribute].unique()
 total_instances = len(data)

PAGE NO.:
DATE:

4: Decision: Yes

5: Decision: No

6: Decision: Yes

7: Decision: No

8: Decision: Yes

9: Decision: No

10: Decision: Yes

11: Decision: Yes

12: Decision: Yes

13: Decision: No

14: Decision: Yes

15: Decision: No

16: Decision: No

17: Decision: No

18: Decision: No

19: Decision: No

20: Decision: No

21: Decision: No

22: Decision: No

23: Decision: No

24: Decision: No

25: Decision: No

26: Decision: No

27: Decision: No

28: Decision: No

29: Decision: No

30: Decision: No

31: Decision: No

32: Decision: No

33: Decision: No

34: Decision: No

35: Decision: No

36: Decision: No

37: Decision: No

38: Decision: No

39: Decision: No

40: Decision: No

41: Decision: No

42: Decision: No

43: Decision: No

44: Decision: No

45: Decision: No

46: Decision: No

47: Decision: No

48: Decision: No

49: Decision: No

50: Decision: No

51: Decision: No

52: Decision: No

53: Decision: No

54: Decision: No

55: Decision: No

56: Decision: No

57: Decision: No

58: Decision: No

59: Decision: No

60: Decision: No

61: Decision: No

62: Decision: No

63: Decision: No

64: Decision: No

65: Decision: No

66: Decision: No

67: Decision: No

68: Decision: No

69: Decision: No

70: Decision: No

71: Decision: No

72: Decision: No

73: Decision: No

74: Decision: No

75: Decision: No

76: Decision: No

77: Decision: No

78: Decision: No

79: Decision: No

80: Decision: No

81: Decision: No

82: Decision: No

83: Decision: No

84: Decision: No

85: Decision: No

86: Decision: No

87: Decision: No

88: Decision: No

89: Decision: No

90: Decision: No

91: Decision: No

92: Decision: No

93: Decision: No

94: Decision: No

95: Decision: No

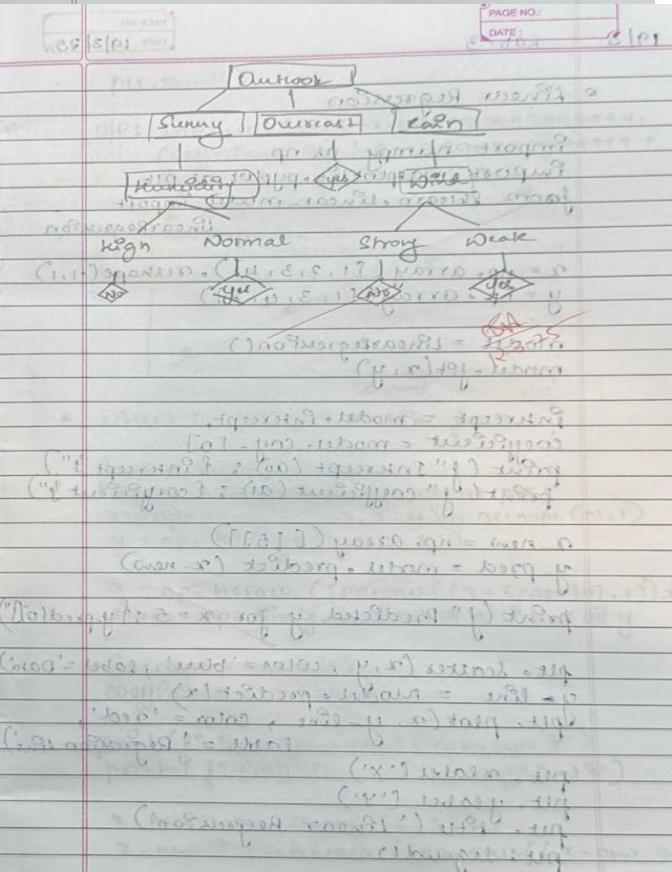
96: Decision: No

97: Decision: No

98: Decision: No

99: Decision: No

100: Decision: No



Code:

```
import numpy as np
import pandas as pd
from collections import Counter

class Node:
    def __init__(self, feature=None, value=None, label=None):
        self.feature = feature
        self.value = value
        self.label = label
        self.children = {}

def entropy(y):
    counts = np.bincount(y)
    probabilities = counts / len(y)
    return -np.sum([p * np.log2(p) for p in probabilities if p > 0])

def information_gain(X, y, feature):
    total_entropy = entropy(y)
    values, counts = np.unique(X[:, feature], return_counts=True)
    weighted_entropy = sum((counts[i] / sum(counts)) * entropy(y[X[:, feature] == v]) for i, v in enumerate(values))
    return total_entropy - weighted_entropy

def best_feature_to_split(X, y):
    gains = [information_gain(X, y, i) for i in range(X.shape[1])]
    return np.argmax(gains)

def id3(X, y, features):
    if len(set(y)) == 1:
        return Node(label=y[0])
    if len(features) == 0:
        return Node(label=Counter(y).most_common(1)[0][0])
    best_feature = best_feature_to_split(X, y)
    node = Node(feature=features[best_feature])
    feature_values = np.unique(X[:, best_feature])
    for value in feature_values:
        sub_X = X[X[:, best_feature] == value]
        sub_y = y[X[:, best_feature] == value]
        if len(sub_y) == 0:
            node.children[value] = Node(label=Counter(y).most_common(1)[0][0])
        else:
            node.children[value] = id3(np.delete(sub_X, best_feature, axis=1), sub_y, features[:best_feature] + features[best_feature+1:])
    return node

def print_tree(node, depth=0):
    if node.label is not None:
        print(f"{' ' * depth}Leaf: {node.label}")
        return
    print(f"{' ' * depth}Feature: {node.feature}")
    for value, child in node.children.items():
        print(f"{' ' * (depth+1)}Value: {value}")
        print_tree(child, depth+1)
```

```

print(f"{' ' * depth}Value: {value}")
print_tree(child, depth + 1)

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', 'Hot', 'Mild'],
    'Humidity': ['High', 'High', 'High', 'High', 'Normal', 'Normal', 'Normal', 'High', 'Normal', 'Normal', 'Normal',
    'High', 'Normal', 'High'],
    'Wind': ['Weak', 'Strong', 'Weak', 'Weak', 'Weak', 'Strong', 'Strong', 'Weak', 'Weak', 'Weak', 'Strong', 'Strong',
    'Weak', 'Strong'],
    'PlayTennis': ['No', 'No', 'Yes', 'Yes', 'Yes', 'No', 'Yes', 'No', 'Yes', 'Yes', 'Yes', 'Yes', 'Yes', 'No']
})

```

```

X = data.iloc[:, :-1].apply(lambda col: pd.factorize(col)[0]).to_numpy()
y = pd.factorize(data['PlayTennis'])[0]
features = list(data.columns[:-1])

decision_tree = id3(X, y, features)
print_tree(decision_tree)

```

Program 6

Build KNN Classification model for a given dataset

Screenshot

PAGE NO.: 2/4/25		PAGE NO.: 2/4/25	
KNN classification		1. For iris dataset	
Businessman's age and salary		o Various values of k are tested and for each, accuracy & error rate is calculated. This is done using grid search & optimal found k=3	
Person Age Salary K Target		2. Feature scaling ensures all features contribute equally to nearest neighbor. Scaling is done so that feature like gender, age, don't have their own with smaller ranges.	
A 18 50 N		(35-18)^2 + (100-50)^2 = 52.81	
B 23 55 N		d1 = $\sqrt{(35-23)^2 + (100-55)^2} = 66.57$	
C 24 70 N		d2 = $\sqrt{(35-24)^2 + (100-70)^2} = 31.75$	
D 41 60 4		d3 = $\sqrt{(35-41)^2 + (100-60)^2} = 60.45$	
E 43 70 4		d4 = $\sqrt{(35-43)^2 + (100-70)^2} = 31.65$	
F 38 40 4		d5 = $\sqrt{(35-38)^2 + (100-40)^2} = 60.07$	
X 35 100 ?		d6 = $\sqrt{(35-35)^2 + (100-100)^2} = 0$	
		d7 = Target = 4	

Code:

```

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

# Load the dataset
iris = pd.read_csv("/content/iris (2).csv")

# Split features and labels
X = iris.iloc[:, :-1] # Features: all columns except last
y = iris.iloc[:, -1] # Labels: last column (species)

# Split into training (80%) and testing (20%) sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

```

```

# Initialize and train the KNN model
k = 5 # Choosing k=5
knn = KNeighborsClassifier(n_neighbors=k)
knn.fit(X_train, y_train)

# Make predictions
y_pred = knn.predict(X_test)

# Evaluate the model
print("Accuracy Score:", accuracy_score(y_test, y_pred))
print("\nConfusion Matrix:\n", confusion_matrix(y_test, y_pred))
print("\nClassification Report:\n", classification_report(y_test, y_pred))

# Plot confusion matrix
plt.figure(figsize=(6, 5))
sns.heatmap(confusion_matrix(y_test, y_pred), annot=True, cmap="Blues", fmt="d",
            xticklabels=knn.classes_, yticklabels=knn.classes_)
plt.xlabel("Predicted")
plt.ylabel("Actual")
plt.title("Confusion Matrix - Iris Dataset")
plt.show()

```

```

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

# Load the dataset
diabetes = pd.read_csv("/content/diabetes (1).csv")

# Split features and labels
X = diabetes.iloc[:, :-1] # Features: all columns except last
y = diabetes.iloc[:, -1] # Labels: last column (diabetic or not)

# Split into training (80%) and testing (20%) sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Feature scaling (important for KNN)
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)

# Initialize and train the KNN model
k = 5 # Choosing k=5

```

```

knn = KNeighborsClassifier(n_neighbors=k)
knn.fit(X_train, y_train)

# Make predictions
y_pred = knn.predict(X_test)

# Evaluate the model
print("Accuracy Score:", accuracy_score(y_test, y_pred))
print("\nConfusion Matrix:\n", confusion_matrix(y_test, y_pred))

# Plot confusion matrix
plt.figure(figsize=(6, 5))
sns.heatmap(confusion_matrix(y_test, y_pred), annot=True, cmap="Blues", fmt="d",
            xticklabels=["Non-Diabetic", "Diabetic"], yticklabels=["Non-Diabetic", "Diabetic"])
plt.xlabel("Predicted")
plt.ylabel("Actual")
plt.title("Confusion Matrix - Diabetes Dataset")
plt.show()

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

# Load dataset
df = pd.read_csv('/content/heart (1).csv')

# Define features and target
X = df.drop(columns=['target']) # Assuming 'target' is the classification column
y = df['target']

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

# Feature scaling
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)

# Find the best K value
k_values = range(1, 21)
accuracy_scores = []
for k in k_values:
    model = KNeighborsClassifier(n_neighbors=k)
    model.fit(X_train, y_train)
    y_pred = model.predict(X_test)
    accuracy_scores.append(accuracy_score(y_test, y_pred))

```

```

best_k = k_values[np.argmax(accuracy_scores)]
print(f'Best K value: {best_k}')

# Train model with best K
best_model = KNeighborsClassifier(n_neighbors=best_k)
best_model.fit(X_train, y_train)
y_pred = best_model.predict(X_test)

# Evaluate model
accuracy = accuracy_score(y_test, y_pred)
print(f'Accuracy with best K ({best_k}): {accuracy:.4f}')
print("Classification Report:")
print(classification_report(y_test, y_pred))

# Confusion matrix
cm = confusion_matrix(y_test, y_pred)
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues')
plt.xlabel('Predicted')
plt.ylabel('Actual')
plt.title(f'Confusion Matrix - KNN (K={best_k})')
plt.show()

# Plot K values vs. Accuracy
plt.plot(k_values, accuracy_scores, marker='o')
plt.xlabel('K Value')
plt.ylabel('Accuracy')
plt.title('K Value vs Accuracy')
plt.show()

```

Program 7

Build Support vector machine model for a given dataset

Screenshot

PAGE NO.:
DATE:

Lab - 6

SVM

```

from sklearn.datasets import load_breast_cancer
import matplotlib.pyplot as plt
from sklearn.inspection import DecisionBoundaryDisplay

```

Cancer = load_breast_cancer()
X = Cancer.data[:, :2]
y = Cancer.target

Svm = SVC(kernel='rbf', gamma=0.5, C=1.0)
Svm.fit(X, y)

DecisionBoundaryDisplay.from_estimator(
 Svm,
 X,
 response_method='predict',
 alpha=0.8,
 xlabel=Cancer.feature_names[0],
 ylabel=Cancer.feature_names[1],

plt.scatter(X[:, 0], X[:, 1],
 c=y, s=20,
 edgecolors='k')
plt.show()

Code:

Load the important packages

```

from sklearn.datasets import load_breast_cancer
import matplotlib.pyplot as plt
from sklearn.inspection import DecisionBoundaryDisplay
from sklearn.svm import SVC

# Load the datasets
cancer = load_breast_cancer()
X = cancer.data[:, :2]
y = cancer.target

#Build the model
svm = SVC(kernel="rbf", gamma=0.5, C=1.0)
# Trained the model
svm.fit(X, y)

# Plot Decision Boundary
DecisionBoundaryDisplay.from_estimator(
    svm,
    X,
    response_method="predict",
    cmap=plt.cm.Spectral,
    alpha=0.8,
    xlabel=cancer.feature_names[0],
    ylabel=cancer.feature_names[1],
)

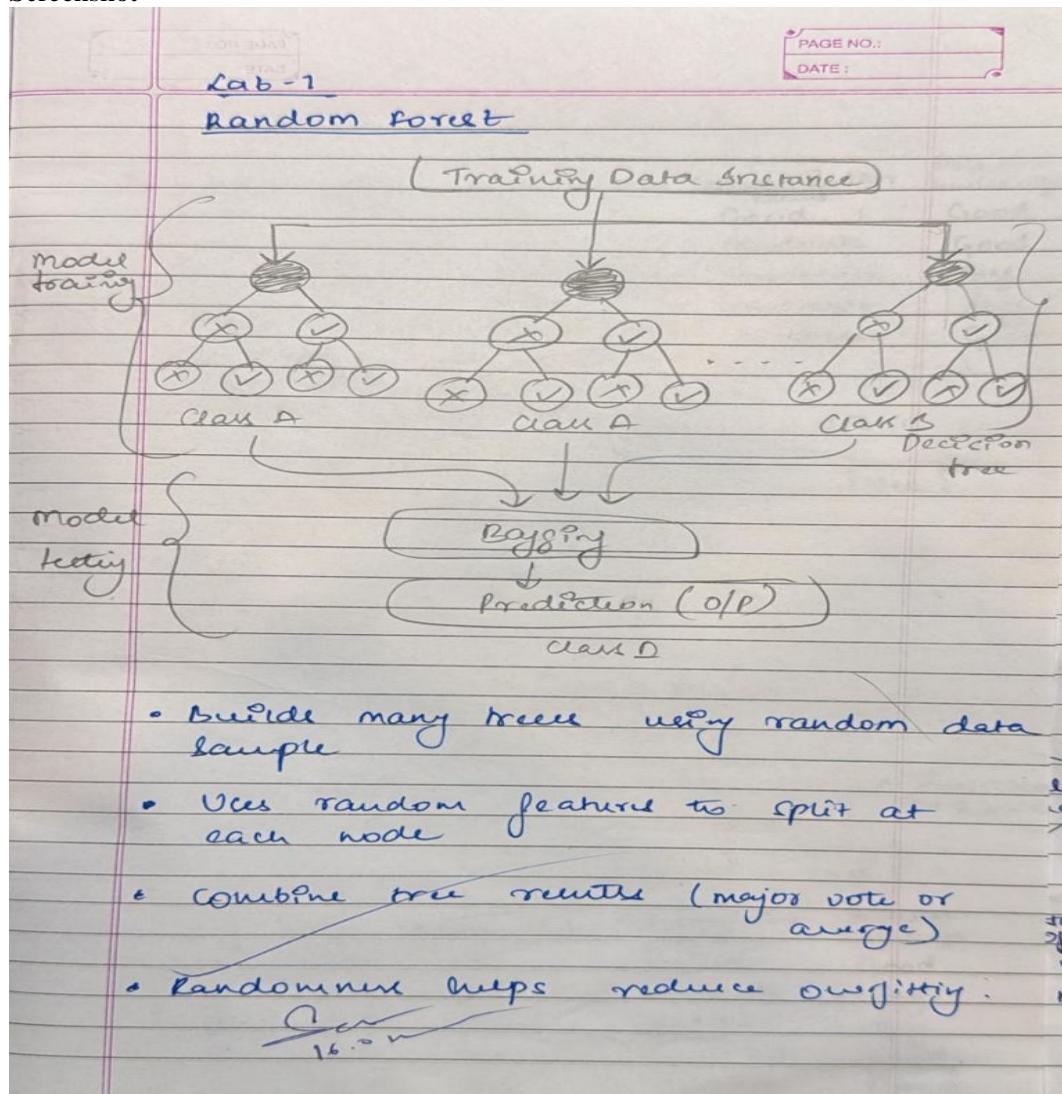
# Scatter plot
plt.scatter(X[:, 0], X[:, 1],
            c=y,
            s=20, edgecolors="k")
plt.show()

```

Program 8

Implement Random forest ensemble method on a given dataset

Screenshot



Code:

```
import pandas as pd
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, confusion_matrix
```

```
# Load the iris dataset
iris = load_iris()
```

```
# Convert to DataFrame
df = pd.DataFrame(data=iris.data, columns=iris.feature_names)
df['species'] = iris.target
```

```

# Split features and target
X = df.drop('species', axis=1)
y = df['species']

# Train-test split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

best_score = 0
best_n = 0
best_cm = None

# Try different numbers of trees
for n in range(1, 101):
    clf = RandomForestClassifier(n_estimators=n, random_state=42)
    clf.fit(X_train, y_train)
    y_pred = clf.predict(X_test)
    acc = accuracy_score(y_test, y_pred)
    if acc > best_score:
        best_score = acc
        best_n = n
        best_cm = confusion_matrix(y_test, y_pred)

print(f"Best Accuracy: {best_score}")
print(f"Number of Trees: {best_n}")
print("Confusion Matrix:")
print(best_cm)

```

Program 9

Implement Boosting ensemble method on a given dataset

Screenshot

Lab-8 - AdaBoost Algorithm

3. CGPA (Intraclumus knowledge)	Practical knowledge		Communication Skill	Job Profile
	Good	Bad		
≥ 9	Yes	Good	Good	Yes
< 9	No	Good	Moderate	Yes
≥ 9	No	Bad	Moderate	No
< 9	No	Bad	Good	No
≥ 9	Yes	Good	Moderate	Yes
≥ 9	Yes	Good	Moderate	Yes

CGPA	Predicted job offer	Actual Job Offer	Weight
≥ 9	Yes	Yes	1/6
< 9	No	Yes	1/6
≥ 9	Yes	No	1/6
< 9	No	No	1/6
≥ 9	Yes	Yes	1/6
≥ 9	Yes	Yes	1/6

$$\begin{aligned} E_{\text{CGPA}} &= \frac{1}{6}(1 - 0.333) = 0.6666666666666667 \cdot 0.667 \\ &\approx 0.4444444444444444 \\ \alpha_{\text{CGPA}} &= \frac{1}{2} \ln(1 - E_{\text{CGPA}}) \approx 0.5 \\ &\approx 0.5 \ln(1 - 0.6666666666666667) \approx 0.347. \end{aligned}$$

Date _____
Page _____

$$\begin{aligned} z_{\text{CGPA}} &= \frac{1}{6}(1 - e^{-0.347}) + \frac{1}{6}(1 + e^{-0.347}) \\ &\Rightarrow 0.9428. \end{aligned}$$

Weight of correct instance: $\frac{1}{6} \cdot e^{-0.347} = 0.9428$

$$\begin{aligned} &0.1 \quad 0.1 \quad 0.1 \\ &0.2 \quad 0.2 \quad 0.2 \\ &0.3 \quad 0.3 \quad 0.3 \\ &0.4 \quad 0.4 \quad 0.4 \\ &0.5 \quad 0.5 \quad 0.5 \\ &0.6 \quad 0.6 \quad 0.6 \\ &0.7 \quad 0.7 \quad 0.7 \\ &0.8 \quad 0.8 \quad 0.8 \\ &0.9 \quad 0.9 \quad 0.9 \\ &1.0 \quad 1.0 \quad 1.0 \end{aligned} = 0.1249.$$

CGPA	Predicted Job Offer	Actual Job Offer	Weight
≥ 9	Yes	Yes	0.1249
< 9	No	Yes	0.2501
≥ 9	Yes	No	0.1249
< 9	No	No	0.1249
≥ 9	Yes	No	0.1249
< 9	No	No	0.1249
≥ 9	Yes	Yes	0.1249

• Insurance 0.50% off.
Default Accuracy: 0.8182

Confusion Matrix

$$\begin{bmatrix} 16282 & 632 & 18 & 12 \\ 632 & 11448 & 12117 & 12 \\ 18 & 12117 & 12 & 12 \\ 12 & 12 & 12 & 12 \end{bmatrix}$$

$$\text{EPE} = \frac{1}{4}(\text{C}_{11} + \text{C}_{12} + \text{C}_{21} + \text{C}_{22}) = 0.0000000000000002$$

$$(C_{11}, C_{12}, C_{21}, C_{22}) = (16282, 632, 18, 12)$$

$$16282 + 632 + 18 + 12 = 17023$$

$$2 \cdot 81.13 = 162.26$$

$$2 \cdot 81.13 = 162.26$$

$$(81.13, 81.13) = (40.565, 40.565)$$

Code:

```
import pandas as pd
from sklearn.ensemble import AdaBoostClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, confusion_matrix
from sklearn.preprocessing import LabelEncoder
import matplotlib.pyplot as plt

# Load the dataset
df = pd.read_csv("/content/income (1).csv")

# Encode categorical variables (if any)
label_encoders = {}
for column in df.select_dtypes(include='object').columns:
    le = LabelEncoder()
    df[column] = le.fit_transform(df[column])
    label_encoders[column] = le

# Split dataset into features and target
X = df.drop("income_level", axis=1) # replace 'income' with the actual target column name if different
y = df["income_level"]

# Train-test split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Train AdaBoost with default n_estimators=10
ada_default = AdaBoostClassifier(n_estimators=10, random_state=42)
ada_default.fit(X_train, y_train)
y_pred_default = ada_default.predict(X_test)
default_accuracy = accuracy_score(y_test, y_pred_default)

print(f"Default Accuracy (10 estimators): {default_accuracy:.4f}")
print("Confusion Matrix:")
print(confusion_matrix(y_test, y_pred_default))

# Tune number of estimators
scores = []
n_estimators_range = range(1, 101)

for n in n_estimators_range:
```

```

ada = AdaBoostClassifier(n_estimators=n, random_state=42)
ada.fit(X_train, y_train)
y_pred = ada.predict(X_test)
acc = accuracy_score(y_test, y_pred)
scores.append(acc)

# Best accuracy and corresponding n_estimators
best_score = max(scores)
best_n_estimators = n_estimators_range[scores.index(best_score)]

# Final model with best n_estimators
ada_best = AdaBoostClassifier(n_estimators=best_n_estimators, random_state=42)
ada_best.fit(X_train, y_train)
y_pred_best = ada_best.predict(X_test)
conf_matrix_best = confusion_matrix(y_test, y_pred_best)

print(f"\nBest Accuracy: {best_score:.4f} using {best_n_estimators} estimators")
print("Best Confusion Matrix:")
print(conf_matrix_best)

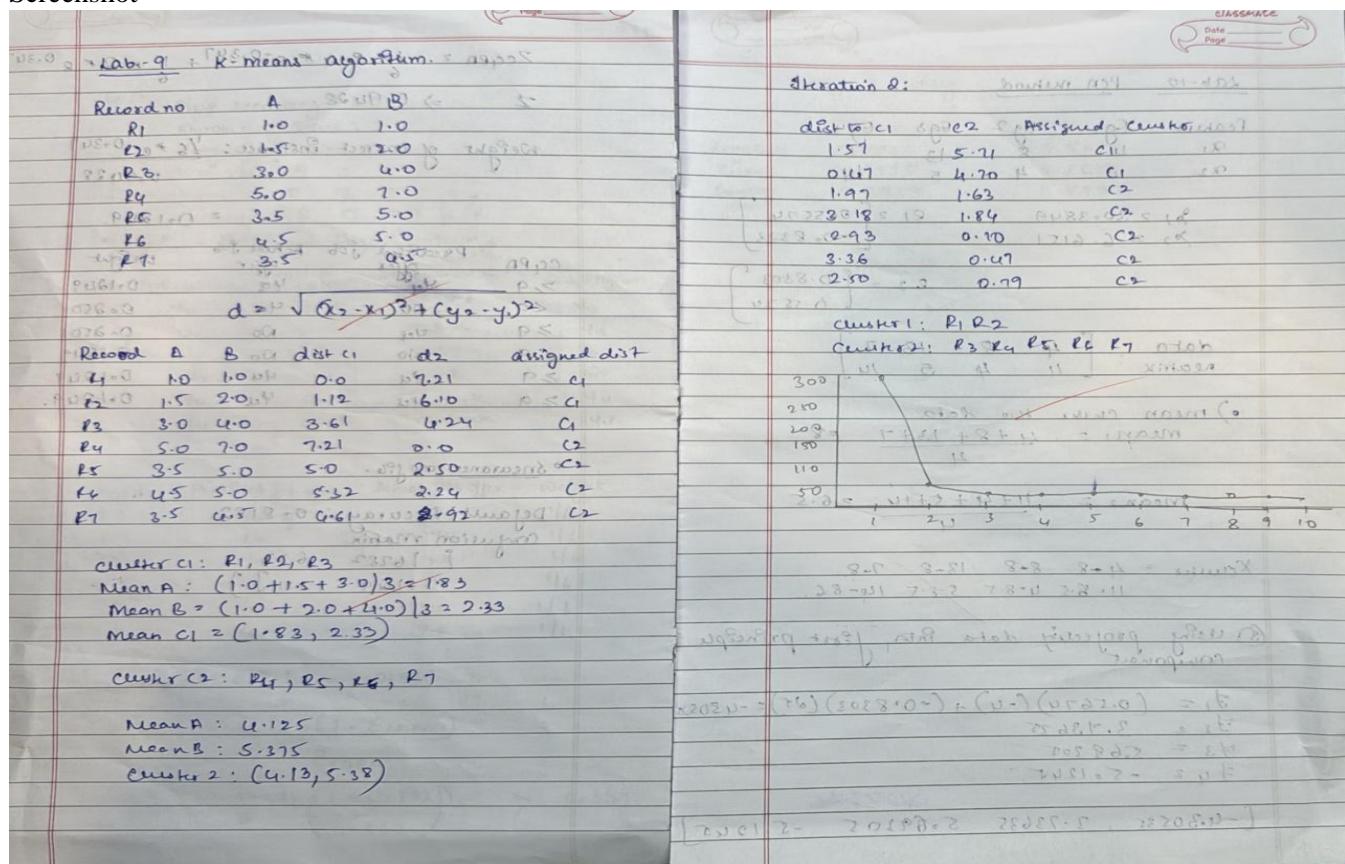
# Optional: Plot accuracy vs number of estimators
plt.figure(figsize=(10, 6))
plt.plot(n_estimators_range, scores, marker='o')
plt.title("AdaBoost Accuracy vs Number of Estimators")
plt.xlabel("Number of Estimators")
plt.ylabel("Accuracy")
plt.grid(True)
plt.show()

```

Program 10

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

Screenshot



Code:

```
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
```

```

# Load the Iris dataset
iris = load_iris()
df = pd.DataFrame(iris.data, columns=iris.feature_names)
df['target'] = iris.target
# Select only petal length and width
X = df[['petal length (cm)', 'petal width (cm)']]
# Optional: Scale the features (important for distance-based algorithms like K-Means)
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)
# Elbow method to determine optimal number of clusters
inertia = []
k_range = range(1, 11)

for k in k_range:
    kmeans = KMeans(n_clusters=k, random_state=42)
    kmeans.fit(X_scaled)
    inertia.append(kmeans.inertia_)

# Plot the elbow curve
plt.figure(figsize=(8, 5))
plt.plot(k_range, inertia, marker='o')
plt.title('Elbow Method for Optimal k')
plt.xlabel('Number of clusters (k)')
plt.ylabel('Inertia')
plt.grid(True)
plt.show()
# From the elbow plot, we choose k = 3 (typical for Iris dataset)
optimal_k = 3
kmeans_final = KMeans(n_clusters=optimal_k, random_state=42)
clusters = kmeans_final.fit_predict(X_scaled)

# Add cluster labels to original data
df['Cluster'] = clusters

# Plotting the clusters
plt.figure(figsize=(8, 5))
plt.scatter(X_scaled[:, 0], X_scaled[:, 1], c=clusters, cmap='viridis', s=50)
plt.scatter(kmeans_final.cluster_centers_[:, 0], kmeans_final.cluster_centers_[:, 1],
            c='red', s=200, alpha=0.7, marker='X', label='Centroids')
plt.title('K-Means Clustering (k=3) on Iris Petal Features')
plt.xlabel('Petal Length (scaled)')
plt.ylabel('Petal Width (scaled)')
plt.legend()
plt.grid(True)
plt.show()

```

Program 11

Implement Dimensionality reduction using Principal Component Analysis (PCA) method

Screenshot

Lab-10 PCA method		Node Accuracy	
Features		Before	
x_1	4	0.8833	0.8667
x_2	11	0.8718	0.8556
x_3	13	0.8889	0.8722
x_4	7		
$\mu_1 = \begin{bmatrix} 30.3849 \\ 6.6151 \end{bmatrix}$		54 75.205	
$\mu_2 = \begin{bmatrix} 0.5574 \\ -0.8303 \end{bmatrix}$			
$e_1 = \begin{bmatrix} 0.5574 \\ 0.8303 \end{bmatrix}$			
$e_2 = \begin{bmatrix} 0.8303 \\ 0.5574 \end{bmatrix}$			
data : $\begin{bmatrix} 4 & 8 & 13 & 7 \\ 11 & 4 & 5 & 14 \end{bmatrix}$			
mean center the data			
$\text{mean}_1 = \frac{4+8+13+7}{4} = 8$			
$\text{mean}_2 = \frac{11+4+5+14}{4} = 6.5$			
$X_{center} = \begin{bmatrix} 4-8 & 8-8 & 13-8 & 7-8 \\ 11-8.5 & 4-8.5 & 5-8.5 & 14-8.5 \end{bmatrix}$			
③ using projecting data onto first principle component			
$f_1 = (0.5574)(-4) + (-0.8303)(0.5) = -4.8058$			
$f_2 = 3.73625$			
$f_3 = 5.69305$			
$f_4 = -5.1245$			
$\begin{bmatrix} -4.8058 & 3.73625 & 5.69305 & -5.1245 \end{bmatrix}$			

Code:

```

import pandas as pd
from sklearn.preprocessing import LabelEncoder, StandardScaler
from sklearn.model_selection import train_test_split
from sklearn.svm import SVC
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier
from sklearn.decomposition import PCA
from sklearn.metrics import accuracy_score

# Load dataset
df = pd.read_csv("/content/heart (2).csv") # Update to match your file path if needed

# Define features and target
X = df.drop('HeartDisease', axis=1)
y = df['HeartDisease']

# Identify categorical columns
categorical_cols = X.select_dtypes(include=['object']).columns.tolist()

# Encode categorical columns
for col in categorical_cols:

```

```

if X[col].nunique() == 2:
    X[col] = LabelEncoder().fit_transform(X[col])
else:
    X = pd.get_dummies(X, columns=[col])

# Scale features
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

# Train-test split
X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=0.2, random_state=42)

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

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

# Apply PCA (reduce to 5 components)
pca = PCA(n_components=5)
X_pca = pca.fit_transform(X_scaled)
X_train_pca, X_test_pca, y_train_pca, y_test_pca = train_test_split(X_pca, y, test_size=0.2, random_state=42)

# Train and evaluate models (with PCA)
print("\n📝 Accuracy with PCA:")
for name, model in models.items():
    model.fit(X_train_pca, y_train_pca)
    y_pred_pca = model.predict(X_test_pca)
    print(f"{name}: {accuracy_score(y_test_pca, y_pred_pca):.4f}")

```