

# VISVESVARAYA TECHNOLOGICAL UNIVERSITY

“JnanaSangama”, Belgaum -590014, Karnataka.



**LAB REPORT**  
**on**

## **Machine Learning (23CS6PCMAL)**

*Submitted by*

**NAVNEETH K S (1BM22CS174)**

*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 **Navneeth K S (1BM22CS174)**, 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 a Machine Learning (23CS6PCMAL) work prescribed for the said degree.

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

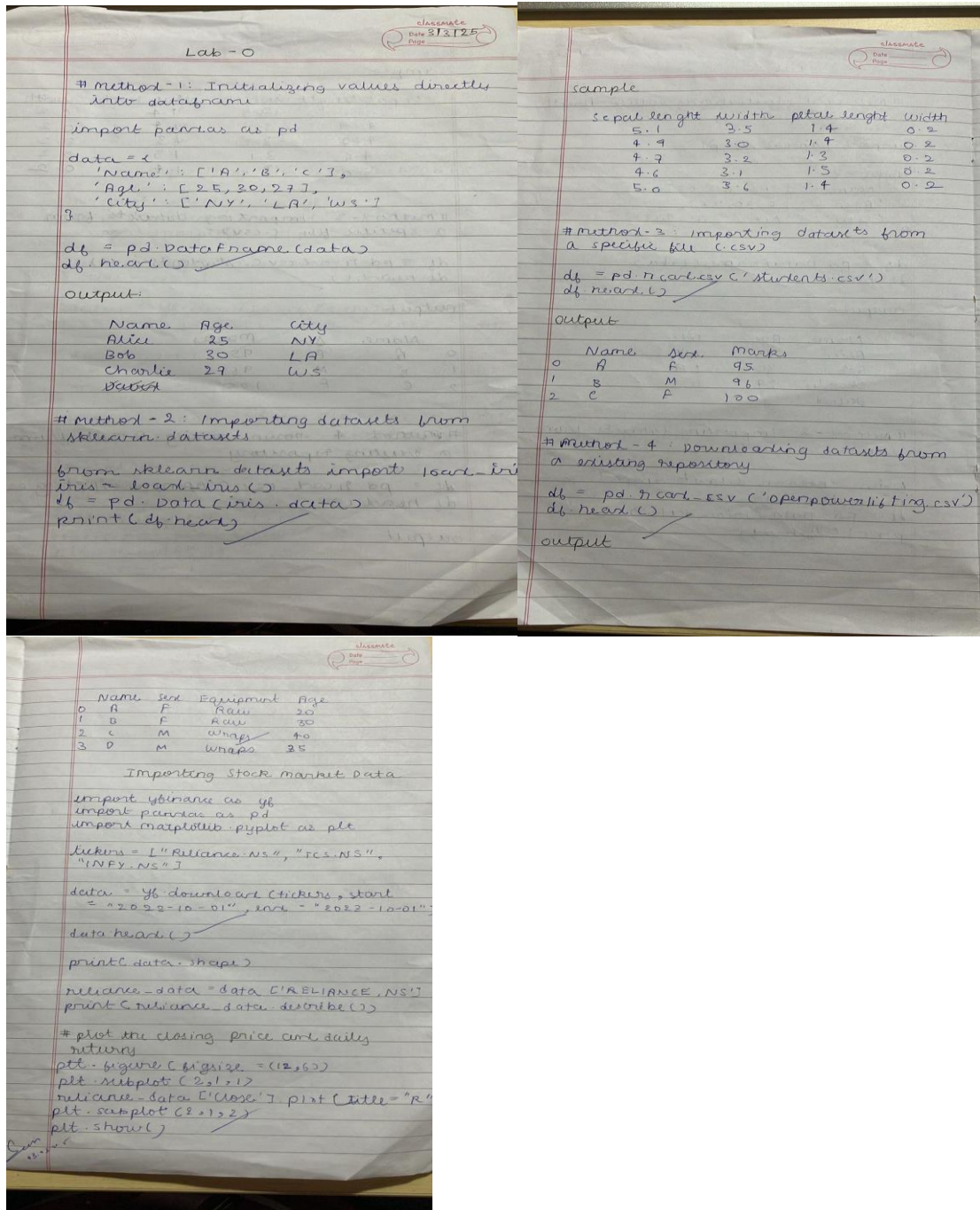
## Index

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

## Program 1

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

Screenshot



Code:

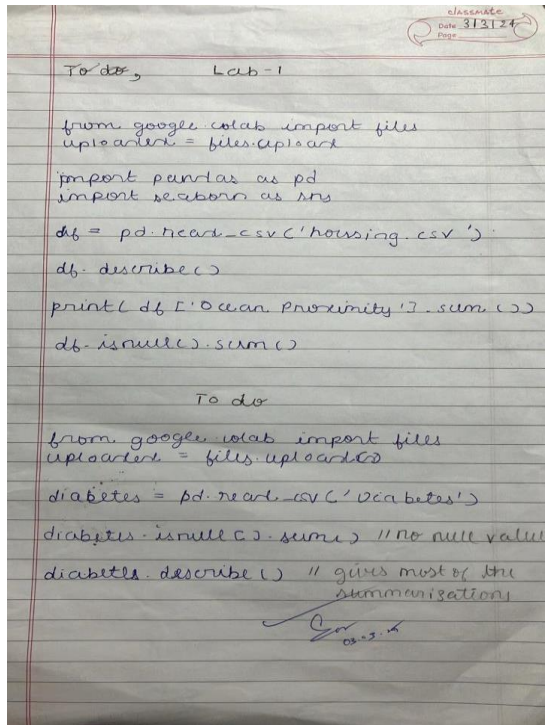
```
import pandas as pd

try:
    df = pd.read_csv('input.csv')
    print("Data imported successfully!\n")
    print(df)
except FileNotFoundError:
    print("The file 'input.csv' was not found.")
df["Processed"] = True
df.to_csv('output.csv', index=False)
print("\nData exported successfully to 'output.csv'.")
```

## Program 2

Demonstrate various data pre-processing techniques for a given dataset

### Screenshots



### Code

```
import pandas as pd
```

```
import numpy as np
```

```
from sklearn.preprocessing import LabelEncoder, StandardScaler, MinMaxScaler
```

```
data = {
    'Name': ['Alice', 'Bob', 'Charlie', 'David', 'Eve', None],
    'Age': [25, 30, np.nan, 35, 29, 40],
```

```

'Department': ['HR', 'IT', 'Finance', 'IT', 'HR', 'Finance'],
'Salary': [50000, 60000, 58000, 62000, np.nan, 52000]
}

df = pd.DataFrame(data)

print("Original DataFrame:\n", df)

df['Age'].fillna(df['Age'].mean(), inplace=True)
df['Salary'].fillna(df['Salary'].median(), inplace=True)
df['Name'].fillna('Unknown', inplace=True)

le = LabelEncoder()
df['Department_Encoded'] = le.fit_transform(df['Department'])

df.drop_duplicates(inplace=True)

df.rename(columns={'Salary': 'Monthly_Salary'}, inplace=True)

df['Age'] = df['Age'].astype(int)

scaler = MinMaxScaler()
df['Salary_Normalized'] = scaler.fit_transform(df[['Monthly_Salary']])

standard_scaler = StandardScaler()

```

```
df['Age_Standardized'] = standard_scaler.fit_transform(df[['Age']])
```

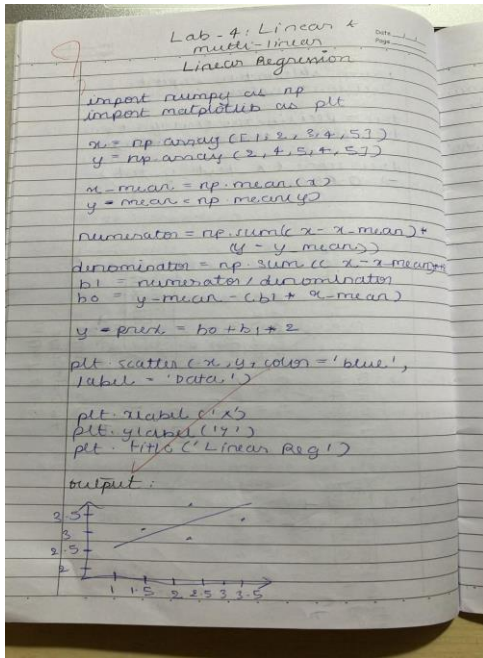
```
print("\nPreprocessed DataFrame:\n", df)
```

### **Program 3**



Implement Linear and Multi-Linear Regression algorithm using appropriate dataset

## Screenshots



## Code

### # Linear Regression

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.datasets import load_boston
from sklearn.linear_model import LinearRegression
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error, r2_score

# Load dataset
boston = load_boston()
df = pd.DataFrame(boston.data, columns=boston.feature_names)
df['PRICE'] = boston.target

# Use only one feature for simple linear regression (e.g., RM = average number of rooms)
X = df[['RM']]
y = df['PRICE']

# Split dataset
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=1)

# Train model
lr = LinearRegression()
lr.fit(X_train, y_train)
```

```

# Predict
y_pred = lr.predict(X_test)

# Output
print("Linear Regression Results")
print("Coefficients:", lr.coef_)
print("Intercept:", lr.intercept_)
print("MSE:", mean_squared_error(y_test, y_pred))
print("R2 Score:", r2_score(y_test, y_pred))

# Plot
plt.scatter(X_test, y_test, color='blue')
plt.plot(X_test, y_pred, color='red')
plt.xlabel('Average Number of Rooms (RM)')
plt.ylabel('House Price')
plt.title('Simple Linear Regression')
plt.show()

# Multiple Linear Regression

# Use all features
X = df.drop('PRICE', axis=1)
y = df['PRICE']

# Split dataset
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=1)

# Train model
mlr = LinearRegression()
mlr.fit(X_train, y_train)

# Predict
y_pred = mlr.predict(X_test)

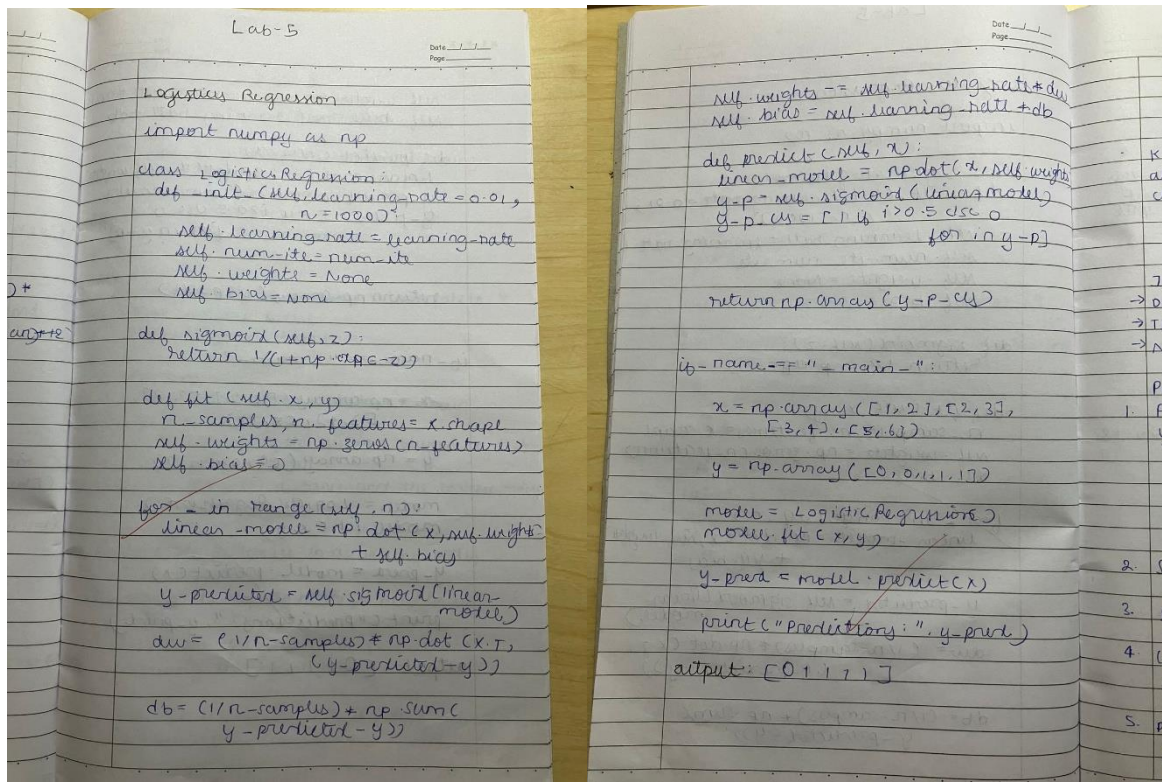
# Output
print("\nMultiple Linear Regression Results")
print("Coefficients:", mlr.coef_)
print("Intercept:", mlr.intercept_)
print("MSE:", mean_squared_error(y_test, y_pred))
print("R2 Score:", r2_score(y_test, y_pred))

```

### **Program 4**

Build Logistic Regression Model for a given dataset

## Screenshot's



## Code

```

import pandas as pd

from sklearn.datasets import load_iris

from sklearn.linear_model import LogisticRegression

from sklearn.model_selection import train_test_split

from sklearn.metrics import confusion_matrix, accuracy_score, classification_report

iris = load_iris()

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

df['species'] = iris.target
    
```

```

df_binary = df[df['species'] != 2] # Remove class 2 (Virginica)

X = df_binary.iloc[:, :-1] # Features
y = df_binary['species'] # Target (0 or 1)

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

model = LogisticRegression()
model.fit(X_train, y_train)

# Step 5: Predict and evaluate
y_pred = model.predict(X_test)

print("Confusion Matrix:\n", confusion_matrix(y_test, y_pred))
print("\nAccuracy:", accuracy_score(y_test, y_pred))
print("\nClassification Report:\n", classification_report(y_test, y_pred))

```

### **Program 5**

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

## Screenshots

Lab 3: ID3 algorithm

What is ID3 algorithm?

- The ID3 is a popular decision tree algorithm used in machine learning.
- It aims to build a decision tree by iteratively selecting the best attribute to split the data, based on the info gain.
- It is a greedy algorithm and best for features that categorical compared to continuous.

ID3 Metrics

The metrics are entropy and information gain.

Entropy

- also referred to as gini impurity, it measures randomness in the datasets.
- high entropy data sets are evenly distributed across all categories, where low entropy ones have data concentrated at few particular points.

$$H(S) = - \sum (P_i \cdot \log_2(P_i))$$

$S \leftarrow$  current dataset  
 $i \rightarrow$  set of classes in  $S$

Information gain

- it assesses how much value info an attribute can provide.
- we select attribute with highest info gain.

$$I(A;B,D) = H(S) - \sum_{i=1}^n \frac{|S_i|}{|S|} \cdot H(S_i)$$

Pseudocode

def ID3(D, A):

- if  $D$  is pure or  $A$  is empty: return a leaf node with the majority class in  $D$
- else:
  - $A_{best} = \text{argmax}(\text{info}(D, A))$
  - $root = \text{Node}(A_{best})$
  - for  $v$  in values( $A_{best}$ ):
    - $D_v = \text{subset}(D, A_{best}, v)$
    - $child = \text{ID3}(D_v, A - A_{best})$
    - $root.add\_child(v, child)$
- return root

Code

```
import pandas as pd
import numpy as np
import sklearn as sk
from sklearn.preprocessing import LabelEncoder

df = pd.read_csv('play-tennis.csv')

# label encoding
from sklearn.preprocessing import LabelEncoder
le = LabelEncoder()

for col in df.columns:
    if col != 'day':
        df[col] = le.fit_transform(df[col])

X = df.drop('day', axis=1)
y = df['play']

from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
```

from sklearn import DecisionTreeClassifier

tree = DecisionTreeClassifier()

tree.fit(X\_train, y\_train)

tree.score(X\_test, y\_test)

$\rightarrow 0.666$

## Code

```
import pandas as pd
```

```

from sklearn.preprocessing import LabelEncoder

from sklearn.tree import DecisionTreeClassifier

from sklearn import tree

data = {

    'Outlook': ['Sunny', 'Sunny', 'Overcast', 'Rain', 'Rain', 'Rain', 'Overcast',
               'Sunny', 'Sunny', 'Rain', 'Sunny', 'Overcast', 'Overcast', 'Rain'],
    'Temperature': ['Hot', 'Hot', 'Hot', 'Mild', 'Cool', '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']
}

df = pd.DataFrame(data)

le = LabelEncoder()

for column in df.columns:

    df[column] = le.fit_transform(df[column])

# Step 3: Separate features and label
X = df.drop('PlayTennis', axis=1)

y = df['PlayTennis']

clf = DecisionTreeClassifier(criterion='entropy') # ID3 uses 'entropy'

```

```

clf = clf.fit(X, y)

print("\nDecision Tree Rules:")

tree_text = tree.export_text(clf, feature_names=X.columns.tolist())

print(tree_text)

# Example: Outlook=Rain, Temperature=Mild, Humidity=High, Wind=Weak

# Encode input sample with same label encoding order used earlier

sample = pd.DataFrame({

    'Outlook': [le.transform(['Rain'])[0]],

    'Temperature': [le.transform(['Mild'])[0]],

    'Humidity': [le.transform(['High'])[0]],

    'Wind': [le.transform(['Weak'])[0]]

})

# Predict

prediction = clf.predict(sample)

result = 'Yes' if prediction[0] == 1 else 'No'

print(f"\nPrediction for new sample (Rain, Mild, High, Weak): {result}")

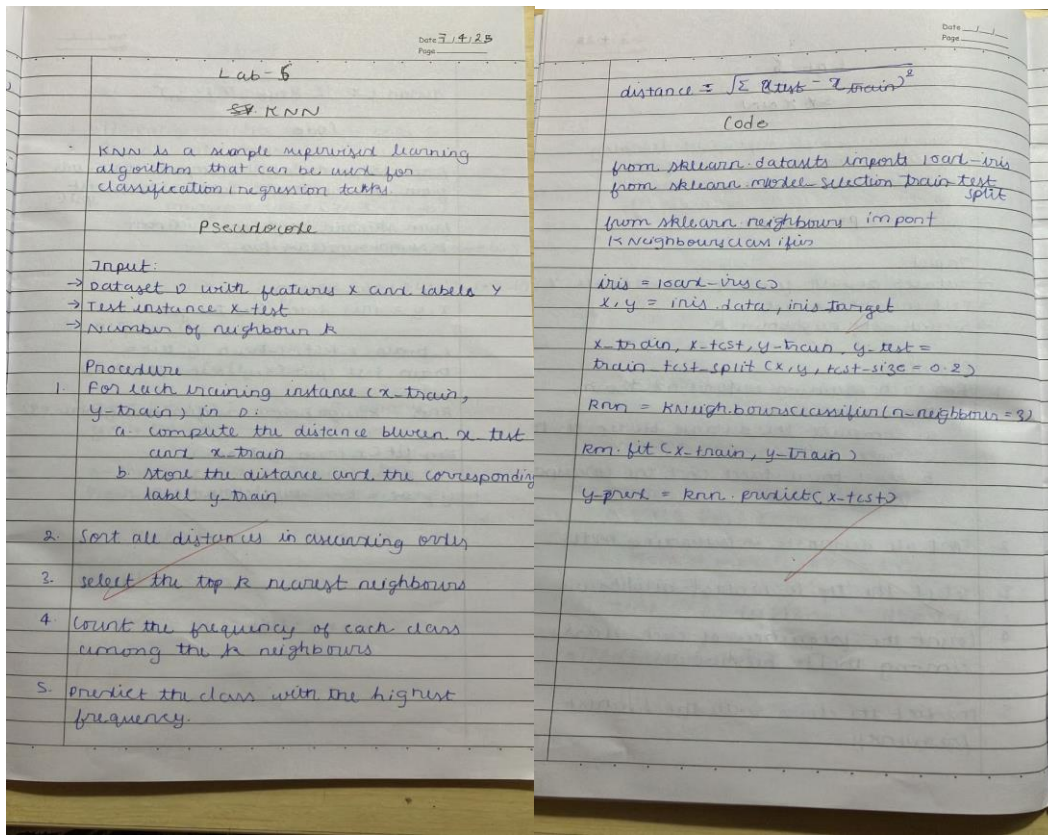
```

### **Program 6**

Build KNN Classification model for a given dataset



## Screenshots



## Code

```
import pandas as pd

from sklearn.datasets import load_iris

from sklearn.model_selection import train_test_split

from sklearn.neighbors import KNeighborsClassifier

from sklearn.metrics import accuracy_score, classification_report, confusion_matrix

# Step 1: Load the Iris dataset

iris = load_iris()

X = pd.DataFrame(iris.data, columns=iris.feature_names)
```



```

y = pd.Series(iris.target)

# Step 2: Split the data (80% training, 20% testing)
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Step 3: Build KNN model (k = 3)
knn = KNeighborsClassifier(n_neighbors=3)
knn.fit(X_train, y_train)

# Step 4: Predict and evaluate
y_pred = knn.predict(X_test)

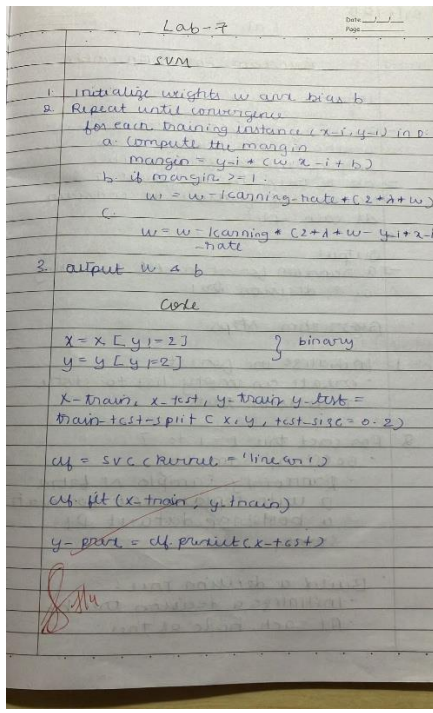
# Results
print("Confusion Matrix:\n", confusion_matrix(y_test, y_pred))
print("\nClassification Report:\n", classification_report(y_test, y_pred))
print("Accuracy Score:", accuracy_score(y_test, y_pred))

```

### **Program 7**

Build Support vector machine model for a given dataset

Screenshots



## Code

```

import pandas as pd

from sklearn.datasets import load_iris

from sklearn.model_selection import train_test_split

from sklearn.svm import SVC

from sklearn.metrics import confusion_matrix, classification_report, accuracy_score

# Step 1: Load the Iris dataset

iris = load_iris()

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

y = pd.Series(iris.target)

# Step 2: Split the data into train and test sets (80% train, 20% test)

```

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

# Step 3: Build and train the SVM model (linear kernel)

svm_model = SVC(kernel='linear')

svm_model.fit(X_train, y_train)

# Step 4: Predict and evaluate

y_pred = svm_model.predict(X_test)

# Output results

print("Confusion Matrix:\n", confusion_matrix(y_test, y_pred))

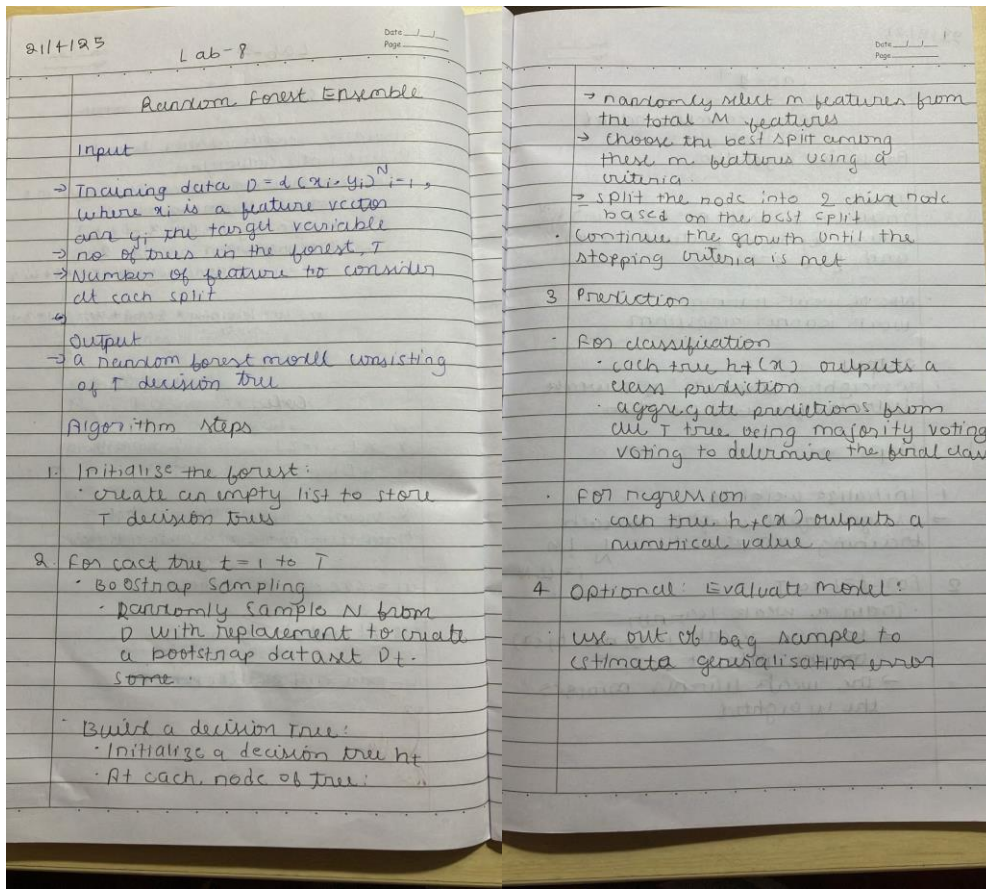
print("\nClassification Report:\n", classification_report(y_test, y_pred))

print("Accuracy Score:", accuracy_score(y_test, y_pred))
```

## **Program 8**

Implement Random Forest ensemble method on a given dataset

Screenshots



## 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 classification_report, confusion_matrix, accuracy_score

# Step 1: Load Iris dataset

iris = load_iris()

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

y = pd.Series(iris.target)
```

```

# Step 2: Split into train and test sets (80% train, 20% test)

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


# Step 3: Train Random Forest model

rf = RandomForestClassifier(n_estimators=100, random_state=42) # 100 trees

rf.fit(X_train, y_train)


# Step 4: Predict and evaluate

y_pred = rf.predict(X_test)


print("Confusion Matrix:\n", confusion_matrix(y_test, y_pred))

print("\nClassification Report:\n", classification_report(y_test, y_pred))

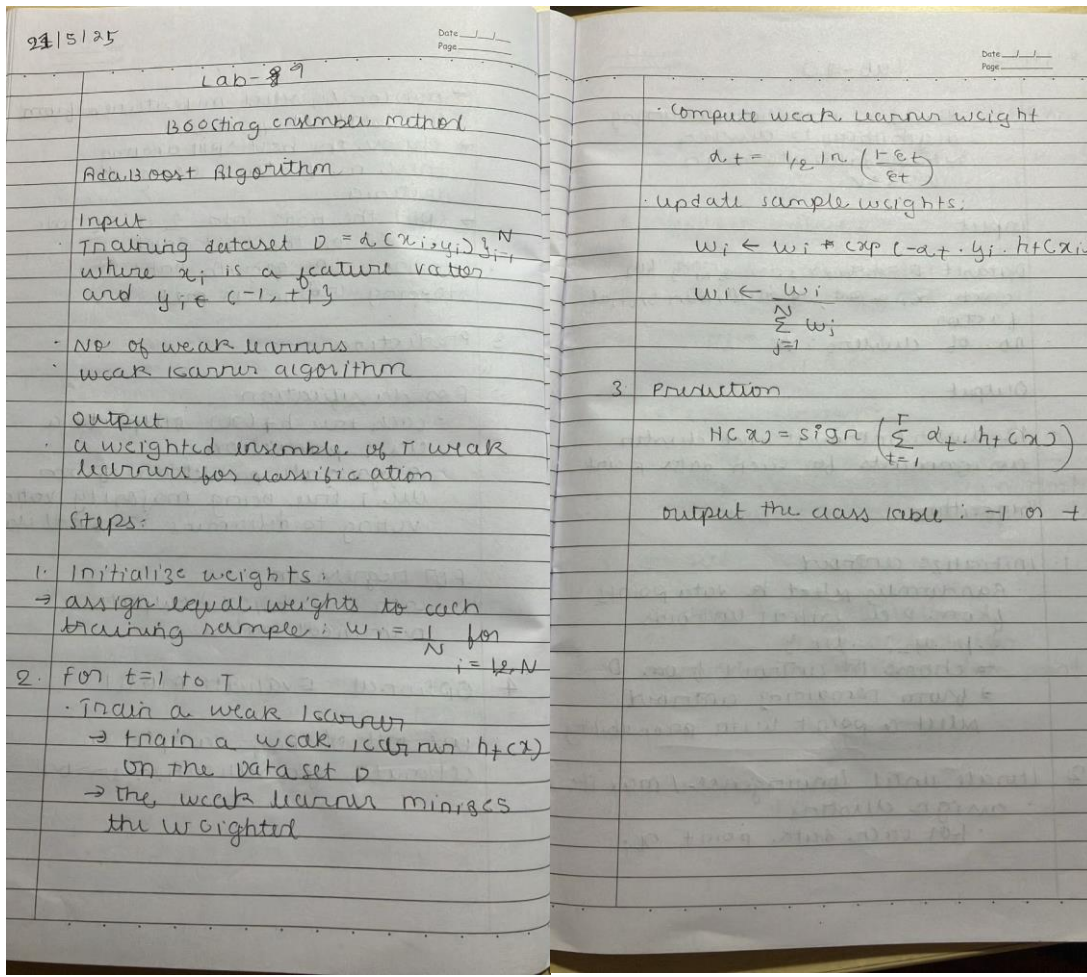
print("Accuracy Score:", accuracy_score(y_test, y_pred))

```

### **Program 9**

Implement Boosting ensemble method on a given dataset

Screenshots



## Code

```
import pandas as pd

from sklearn.datasets import load_iris

from sklearn.model_selection import train_test_split

from sklearn.ensemble import AdaBoostClassifier

from sklearn.metrics import classification_report, confusion_matrix, accuracy_score

from sklearn.tree import DecisionTreeClassifier

# Step 1: Load the Iris dataset

iris = load_iris()
```

```

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

y = pd.Series(iris.target)


# Step 2: Split into train and test sets (80% train, 20% test)

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


# Step 3: Build AdaBoost model with DecisionTreeClassifier as base estimator

base_estimator = DecisionTreeClassifier(max_depth=1)

model = AdaBoostClassifier(base_estimator=base_estimator, n_estimators=50, learning_rate=1.0,
random_state=42)

model.fit(X_train, y_train)


# Step 5: Predict and evaluate

y_pred = model.predict(X_test)

print("Confusion Matrix:\n", confusion_matrix(y_test, y_pred))

print("\nClassification Report:\n", classification_report(y_test, y_pred))

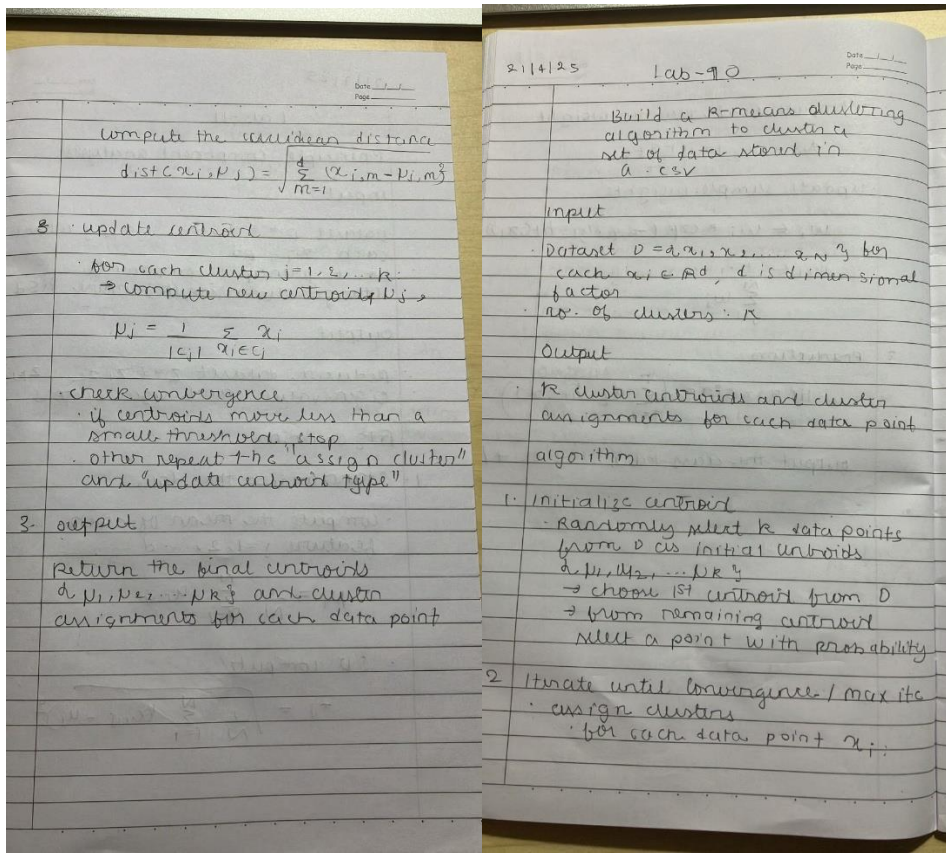
print("Accuracy Score:", accuracy_score(y_test, y_pred))

```

### **Program 10**

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

Screenshots



## Code

```
import pandas as pd

from sklearn.cluster import KMeans

import matplotlib.pyplot as plt

from sklearn.preprocessing import StandardScaler

# Step 1: Load dataset from CSV

df = pd.read_csv('your_dataset.csv') # Replace with your file path

# Optional: View first few rows

print("Data Preview:\n", df.head())
```



```

# Step 2: Select relevant numeric columns for clustering

# You can specify specific columns like: df[['column1', 'column2']]

X = df.select_dtypes(include='number')


# Step 3: Scale the data (important for K-Means)

scaler = StandardScaler()

X_scaled = scaler.fit_transform(X)


# Step 4: Apply K-Means clustering (e.g., 3 clusters)

kmeans = KMeans(n_clusters=3, random_state=42)

df['Cluster'] = kmeans.fit_predict(X_scaled)


# Step 5: Print cluster centers

print("Cluster Centers:\n", kmeans.cluster_centers_)


# Optional Step 6: Visualize (works well for 2D or PCA-reduced data)

if X.shape[1] >= 2:

    plt.scatter(X_scaled[:, 0], X_scaled[:, 1], c=df['Cluster'], cmap='viridis')

    plt.title("K-Means Clustering")

    plt.xlabel("Feature 1")

    plt.ylabel("Feature 2")

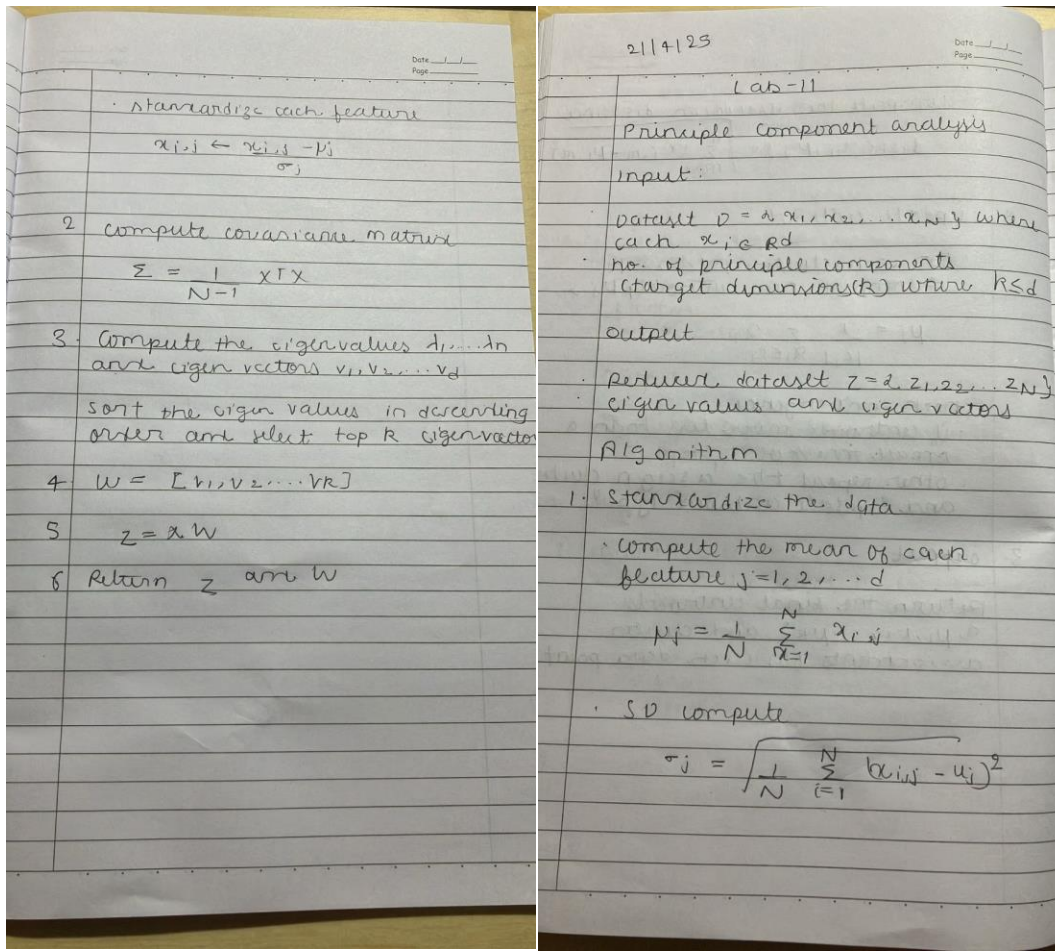
    plt.show()

```

### **Program 11**

Implement Dimensionality reduction using Principal Component Analysis (PCA) method

Screenshots



Code

```
import pandas as pd

from sklearn.datasets import load_iris

from sklearn.preprocessing import StandardScaler

from sklearn.decomposition import PCA

import matplotlib.pyplot as plt

# Step 1: Load the Iris dataset

iris = load_iris()

X = iris.data
```

```

y = iris.target

feature_names = iris.feature_names


# Step 2: Standardize the data

scaler = StandardScaler()

X_scaled = scaler.fit_transform(X)


# Step 3: Apply PCA (reduce to 2 components for visualization)

pca = PCA(n_components=2)

X_pca = pca.fit_transform(X_scaled)


# Step 4: Plot the 2D PCA result

plt.figure(figsize=(8, 6))

scatter = plt.scatter(X_pca[:, 0], X_pca[:, 1], c=y, cmap='viridis', edgecolor='k', s=60)

plt.xlabel("Principal Component 1")

plt.ylabel("Principal Component 2")

plt.title("PCA - Iris Dataset")

plt.legend(handles=scatter.legend_elements()[0], labels=iris.target_names)

plt.grid(True)

plt.show()


# Explained variance

print("Explained variance ratio:", pca.explained_variance_ratio_)

```