# VISVESVARAYA TECHNOLOGICAL UNIVERSITY

**"JnanaSangama", Belgaum- 590014, Karnataka.**



## LAB REPORT on

# Machine Learning (23CS6PCMAL)

### *Submitted by*

**Pranav R Hegde (1BM22CS202)**

*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)**

# INDEX

| 11 | 12-05-2025 | Implement Dimensionality reduction using Principle Component Analysis (PCA) method. | 44 |
| --- | --- | --- | --- |

Github Link:  https://github.com/pranav0hegde/ML

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

**Screenshot:**

① file-path = 'data.csv'
df = pd.read_csv(file-path)
print("sample_data:")
print(df.head(5))
print("\n")

Output:
Sample Data

| | Sepal length | Sepal Width | Petal length |
|---|---|---|---|
| 0 | 5.1 | 3.5 | 1.4 |
| 1 | 4.9 | 3 | 1.4 |
| 2 | 4.7 | 3.2 | 1.3 |
| 3 | 4.6 | 3.1 | 1.5 |
| 4 | 5 | 3.6 | 1.4 |

② df = pd.read_csv(mobile, dataset=2025.
   csv')
   print(df.head())

To do 2

② 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("\n Shape of dataset:")
print(data.shape)
print("\n column names:")
print(data.columns)

summary_data = data["RELIANCE.NS"]

print("\n Summary Statistics for
   Reliance Industries:")
print(summary_data.describe())

Lab 0

To Do 1

ⓐ import pandas as pd
data = {
'Name' : ['Alice', 'Bob', 'Charlie', 'David'],
'USN' : [25, 30, 35, 40],
'Marks' : ['New York', '88', '97', '62']
}

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

Output.

Sample data.

| | Name | USN | Marks |
|---|---|---|---|
| 0 | Alice | 25 | New York |
| 1 | Bob | 30 | 88 |
| 2 | Charlie | 35 | 97 |
| 3 | David | 40 | 62 |

ⓑ 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 ())
Sample data.

| | Sepal Length (cm) | Sepal width (cm) | petal length |
|---|---|---|---|
| 0 | | 3.0 | 1.4 |
| 1 | 4.9 | 3.2 | 1.3 |
| 2 | 4.7 | 3.1 | 1.5 |
| 3 | 4.6 | 3.6 | 1.4 |
| 4 | 5 | | |

**Code:**

```
from sklearn.datasets import load_iris import

pandas as pd

iris = load_iris()

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

df.head()


df['target'] = iris.target

df

import kagglehub


# Download latest version

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


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


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

df.head()

df['Company Name']


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

df = pd.DataFrame(data)

df
```

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

df.columns

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

df.head() import yfinance as yf import pandas as pd

import matplotlib.pyplot as plt


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

# Fetch historical data for the last 1 year

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


# Display the first 5 rows of the dataset

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


print(data.head())

print("\nShape of the dataset:")

print(data.shape)
# Summary statistics for a specific stock (e.g., Reliance)

reliance_data = data['RELIANCE.NS']

print("\nSummary statistics for Reliance Industries:")

print(reliance_data.describe())
```

```python
# Calculate daily returns

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


# Plot the closing price and daily returns

plt.figure(figsize=(12, 6))

plt.subplot(2, 1, 1)

reliance_data['Close'].plot(title="Reliance Industries - Closing Price")

plt.subplot(2, 1, 2)

reliance_data['Daily Return'].plot(title="Reliance Industries - Daily Returns",
color='orange') plt.tight_layout() plt.show()
```

# Program 2

Demonstrate various data pre-processing techniques for a given dataset

**Screenshot:**

**Code:**

```python
import pandas as pd import
numpy as np


# Load dataset
df       =       pd.read_csv("data.csv")
print(df.head())
#       Check       missing       values
print(df.isnull().sum())


# Drop rows with missing values df_cleaned
= df.dropna()


# Or fill missing values with mean/median
df['Age'].fillna(df['Age'].mean(),       inplace=True)
df['Salary'].fillna(df['Salary'].median(),
inplace=True)


# For nominal categories
df = pd.get_dummies(df, columns=['Gender', 'Country'], drop_first=True)


# For ordinal categories
```

```python
from sklearn.preprocessing import OrdinalEncoder encoder

= OrdinalEncoder()

df[['Education_Level']] = encoder.fit_transform(df[['Education_Level']])



from sklearn.preprocessing import StandardScaler, MinMaxScaler



# Standardization (Z-score) scaler

=    StandardScaler()    df[['Age',

'Salary']]                        =

scaler.fit_transform(df[['Age',

'Salary']])



# Min-Max Normalization minmax

= MinMaxScaler()

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



# Using IQR method

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

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

IQR = Q3 - Q1
```

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


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



# Drop irrelevant columns

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


# Correlation-based filtering correlation_matrix

= df.corr() print(correlation_matrix)

from sklearn.model_selection import train_test_split


X = df.drop('Purchased', axis=1) y

= df['Purchased']


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

Program 3

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

**Screenshot:**

**Code:**

```python
import pandas as pd import

numpy as np

from graphviz import Digraph


# Calculate Entropy def
entropy(data):

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

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


# Calculate Information Gain

def    information_gain(data,    feature):

    total_entropy    =    entropy(data)
```

```python
        feature_values = data[feature].unique()

        weighted_entropy = 0

        for value in feature_values:

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

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

        total_entropy - weighted_entropy


# Find the best feature to split the data def

best_feature(data):

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

    gains = {feature: information_gain(data, feature) for feature in features} return max(gains,
    key=gains.get)


# Create the decision tree

def id3(data, features=None):

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

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


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

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


    best = best_feature(data) tree

    = {best: {}}
```

```python
    new_features = features.copy()
    new_features.remove(best)

    for value in data[best].unique(): subset
        = data[data[best] == value]
        tree[best][value] = id3(subset, new_features)

    return tree


# Function to classify new examples based on the decision tree def
classify(tree, example):
    if not isinstance(tree, dict):
    return    tree    feature    =
    list(tree.keys())[0]   value   =
    example[feature]
    return classify(tree[feature][value], example)


# Function to visualize the decision tree using Graphviz
def create_tree_diagram(tree, dot=None, parent_name="Root", parent_value=""):
    if dot is None:
        dot = Digraph(format="png", engine="dot")
```

```python
    if isinstance(tree, dict): # Tree node for

        feature, branches in tree.items():

            feature_name = f"{parent_name}_{feature}" dot.node(feature_name,

            feature)

            dot.edge(parent_name, feature_name, label=parent_value)


            for value, subtree in branches.items():

                value_name        =        f"{feature_name}_{value}"

                dot.node(value_name,      f"{feature}:      {value}")

                dot.edge(feature_name,                    value_name,

                label=str(value))

                # Recurse for each subtree create_tree_diagram(subtree,

    dot, value_name, str(value)) else: # Leaf node

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

        parent_name + "_class", label="Leaf")

    return dot
# Example usage

data = pd.DataFrame({

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

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

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

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

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

})

# Train the decision tree

tree = id3(data, features=list(data.columns[:-1])) print("Decision

Tree:", tree)

# Classify a new example

example = {'Outlook': 'Sunny', 'Temperature': 'Cool', 'Humidity': 'Low', 'Wind': 'Strong'} prediction

= classify(tree, example)

print("Prediction for the example:", prediction)


# Visualize the decision tree dot = create_tree_diagram(tree) dot.render("decision_tree",

view=True) # This will generate and open the tree diagram **Program 4**

Implement Linear and Multi-Linear Regression algorithm for appropriate dataset


**Screenshot:**

# LAB-4

## Linear Regression

- A relation between variable & the outcome was used to predict and were so most suitable variable & the predictions.

A assumption of a linear relationship between 2 independent variable and dependent variable.

To predict and only on independent variable "To find we if $mx + c$ is minimum and the error between the prediction value ($y$) and actual values.

Let data points be $(x, y)$ ... $(x_n, y_n)$ Expressed the data points in matrix form

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} \quad x = \begin{bmatrix} x_1 & \varepsilon_1 \\ \vdots & \\ x_n & \varepsilon_n \end{bmatrix}$$

$$y = mx + c \quad \text{can be represented as}$$

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} + \varepsilon$$

where $\varepsilon$ is error.

Then $\beta_0$ and $\beta_1$ can be determined by
$$\beta = ((x^T x)^{-1} . x^T) y$$

the best fit line and can be used to predict future values

## Multiple Linear Regression

In multiple linear regression best fit line $y = b_0 + b_1 x_1 + b_2 x_2 \ldots b_n x_n + \varepsilon$

Initializing to $b_0$ to 0

Let data points be $(x_1, x_2, \ldots x_n)$
$x_i \in [0 \ldots m]$ where $x_1^i, x_2^i \in [0 \ldots m]$ represents independent variable and $y$ values are dependent variables.

In the form of matrix

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & x_{21} & \ldots & x_{n1} \\ 1 & x_{12} & x_{22} & \ldots & x_{n2} \\ \vdots & & & & \vdots \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \end{bmatrix} + \varepsilon$$

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

Used to predict future values

**Code:**

**Linear Regression**

```python
import pandas as pd

df = pd.read_csv("/content/tvmarketing.csv") df
```

```python
# Visualise the relationship between the features and the response using scatterplots
df.plot(x='TV',y='Sales',kind='scatter')
```

```python
from sklearn.model_selection import train_test_split

x_train, x_test, y_train, y_test = train_test_split(df['TV'], df['Sales'], test_size=0.2, random_state=42)

from sklearn.linear_model import LinearRegression model = LinearRegression()
model.fit(x_train.values.reshape(-1,  1),  y_train)  y_train  model.coef_
model.intercept_
```

**MultiLinearRegression**

```python
import pandas as pd #
```

```python
Step 2 : import data
```

```python
house = pd.read_csv('https://github.com/YBIFoundation/Dataset/raw/main/Boston.csv')
```

```python
# display first 5 rows
house.head()
```

y = house['MEDV']

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

# Step 4 : train test split

from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(X,y, train_size=0.7, random_state=2529)

# Step 5 : select model from sklearn.linear_model

import    LinearRegression    model    =

LinearRegression() # Step 6 : train or fit model

model.fit(X_train,y_train) model.intercept_

model.coef_

## Program 5

Build Logistic Regression Model for a given dataset

**Screenshot:**

**Code:**

```
from sklearn.linear_model import LogisticRegression

from sklearn.datasets import load_iris from

sklearn.model_selection import train_test_split from

sklearn.metrics import accuracy_score


# Load sample dataset (binary classification - Iris with only 2 classes)

iris = load_iris() X = iris.data[iris.target != 2] y = iris.target[iris.target

!= 2]


# Train/Test split

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

```
# Logistic Regression model

model = LogisticRegression()

model.fit(X_train, y_train)


# Predict and evaluate

y_pred = model.predict(X_test)

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

**Program 6**

Build KNN Classification model for a given dataset

**Screenshot:**

Lab 5    KNN

→ Input
   Dataset   $D = \{(x_1, y_1), (x_2, y_2) \ldots \ldots (x_n, y_n)\}$  Test data X.
   Number of neighbours K

→ Output
   Predicted labels. Y used for test data
   Classification accuracy.

Algorithm:
① Load the Dataset.
   Input feature X
   Target labels y
② Split the Dataset
   Divide the dataset into:
   Training set : $X_{train}$, $y_{train}$
   Test set: $X_{test}$, $Y_{test}$

   Use a fixed random seed for reproducibility

③ Initialize the KNN classifier
   set the number of neighbours K.

④ Train the classifier
   Store the trainer data $X_{train}$, $y_{train}$

⑤ Predict for each Test Instance.

For each test sample $x \in x_{test}$:

ⓐ Compute Euclidian Distance to all training samples.

$$d(x, x_i) = \sqrt{\sum_{j=1}^{n} (x_j - x_{ij})^2}$$

ⓑ Identify the $k$ nearest neighbour Select the $k$ smallest distance.

ⓒ Extract labels of the $k$ nearest neigh

ⓓ Determine the majority class among these labels.

ⓔ Assign the majority class as the predicted label for $x$.

⑥ Evaluate Accuracy
Compare actual and predicted label
Compute Accuracy

$$Accuracy = \frac{No \ of \ correct \ predictions}{Total \ test \ samples}$$

7) Display Results.

Print :
Predicted labels.
Actual labels
Classification accuracy

**Code:**

**KNN**

```python
import numpy as np

from collections import Counter


class KNN:

def _init_(self, k=3): self.k = k

    def fit(self, X, y):

        self.X_train = np.array(X) self.y_train

        = np.array(y)


    def euclidean_distance(self, x1, x2):

        return np.sqrt(np.sum((x1 - x2) ** 2))


    def predict(self, X):

        predictions = [self._predict(x) for x in X] return

        np.array(predictions)


    def _predict(self, x):

        # Compute distances to all training points

        distances = [self.euclidean_distance(x, x_train) for x_train in self.X_train]
```

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


        # Get the labels of those neighbors

        k_nearest_labels = [self.y_train[i] for i in k_indices]

        #  Return  the  most  common  label  most_common  =

        Counter(k_nearest_labels).most_common(1)            return

        most_common[0][0]


# Sample dataset (like a mini version of Iris) X_train

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

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


# Test data

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


# Using the KNN modelh knn

= KNN(k=3)

knn.fit(X_train, y_train)

predictions = knn.predict(X_test)


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

**Program 7**

Build Support vector machine model for a given dataset

**Screenshot:**

# Support Vector Machine

→ Input.

Dataset $D = \{(x_1, y_1), (x_2, y_2) \ldots (x_n, y_n)\}$

Test dataset $S_{test}$. Regularization
parameters (Max no of iteration)

→ Output

Predicted class labels
Classification accuracy.

① Data loading and preprocessing
- Load the tries dataset.
- Apply Z-score normalization to
  standardize feature. $x' = \dfrac{x - \mu}{\sigma}$

② Split the Dataset.
- Split the data into
  Training set (70)%. , Test set (70%.)

③ Initalize the SVM classifier
- Set the following parameter
  C : Regularization constant.
  $K(x, x') = x \cdot x'$

④ One vs Rest Training strategy
For each class C in the set of uniques
classes :
   Convert labels into
   Binary format.

2/5/25

$$y_{binary} = \begin{cases} 1 & \text{if } y = c \\ -1 & \text{otherwise} \end{cases}$$

Train a binary SVM classifier using the complified SMO algorithm.

Prediction Phase.

For each test sample $x$:
For each attained binary classifier
Compute decision tr score:

$$f(x) = \sum_{j} a_i y_i \, K(x_i, x_j) + b$$

Store the score, predict the class with the maximum decision tree.

Evaluation.

Compare predicted labels $y$ pred with true labels $y$ test.

Calculate accuracy.

$$Accuracy = \frac{No \text{ of correct predictions}}{Total \text{ test samples}}$$

**Code:**

```
from sklearn import datasets

from          sklearn.model_selection          import

train_test_split   from   sklearn.svm   import   SVC

import matplotlib.pyplot as plt

from sklearn.decomposition import PCA


# Load dataset

iris = datasets.load_iris()

X = iris.data y

= iris.target


# For binary classification (class 0 vs 1) X

= X[y != 2]

y = y[y != 2]


# Train-test split

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


# Train SVM

clf = SVC(kernel='linear') # Try 'rbf', 'poly', etc.

clf.fit(X_train, y_train)
```

# Accuracy

print("Test Accuracy:", clf.score(X_test, y_test))

**<u>Program 8</u>**

Implement Random forest ensemble method on a given dataset

**Screenshot:**

# Random Forest

1) Input :

Dataset $D = \{(x_1, y_1), (x_2, y_2) \dots (x_n, y_n)\}$
where : $x_i$ are vectors
$y_i$ are corresponding labels.
Number of trees $T$ to be created in forest

2) Create Random subsets
For each tree $t = 1, 2 \dots T$ :
Randomly select a subset $D_t$ from dataset
$D$ with Replacement (bootstrapping).

3) Build Decision Tree.
For each tree $t$, build a decision tree :
. Start with the whole subset $D_t$.
. At each node in the tree,
ⓘ Randomly select a subset of features
$F_{subset} \subseteq F$ where $F$ is the set of all
features.
ⓘⓘ Find the best split among the featu
-res in $F_{subset}$.
ⓘⓘⓘ Split the node based on best feature.
Repeat this process recursively until one of
stopping criteria is met

4) Aggregate the Results.

5) Output :

Predict class label or the predicted value

based on the majority vote of all the trees

**Boot strap Sampling:**

Random feature Selection:

**Aggregation:**

Random Forest Algorithm

import numpy as np
from collections import import counter
from random import sample
del train

import pandas as pd
from sklearn.model_selection import
train-test-split.
from sklearn.ensemble import Random
forest classification
from graph.colab import files

uploaded = files.upload()

for filename in uploaded.key():
  df = pd.ml = csv (filename)
  print (pd Data loaded : filename(?)
  display (df.head())

x = df.iloc[:, :-1]

---

y = df.iloc[:,-1]

x_train, x_test, y_train, y_test = train
test = train-test split (x, y, file-size = train
0.2, random-state = 42)

rf model = Random forest classification (n-
estimator = 100, random-state = 3)
rf model. fit (x-train, y-train)

y-pred = rf model. predict (x_test)

accuracy: accuracy score (y_test, y_pred)
print ("Accuracy of Random Forest:
0.97234.")
print ("classification Report:")

**Output:**
Studied CSV
Data loaded from student.csv

Accuracy of Random Forest: 100.00%
classification Report:

| | Precision | Recall | f1-score | support |
|---|---|---|---|---|
| 0 | 1.00 | 1.00 | 1.00 | 1 |
| 1 | 1.00 | 1.00 | 1.00 | 1 |
| macro avg | 1.00 | 1.00 | 1.00 | 2 |
| weighted | 1.00 | 1.00 | 1.00 | 2 |

avg

**Code:**

```
from sklearn.datasets import load_iris

from sklearn.model_selection import train_test_split

from sklearn.ensemble import RandomForestClassifier

from sklearn.metrics import accuracy_score

# Load sample dataset iris

= load_iris()

X, y = iris.data, iris.target


# Train/test split

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


# Initialize Random Forest

rf = RandomForestClassifier(n_estimators=100, random_state=42) rf.fit(X_train,

y_train)


# Predict and evaluate y_pred = rf.predict(X_test)

print("Accuracy:", accuracy_score(y_test, y_pred))
```
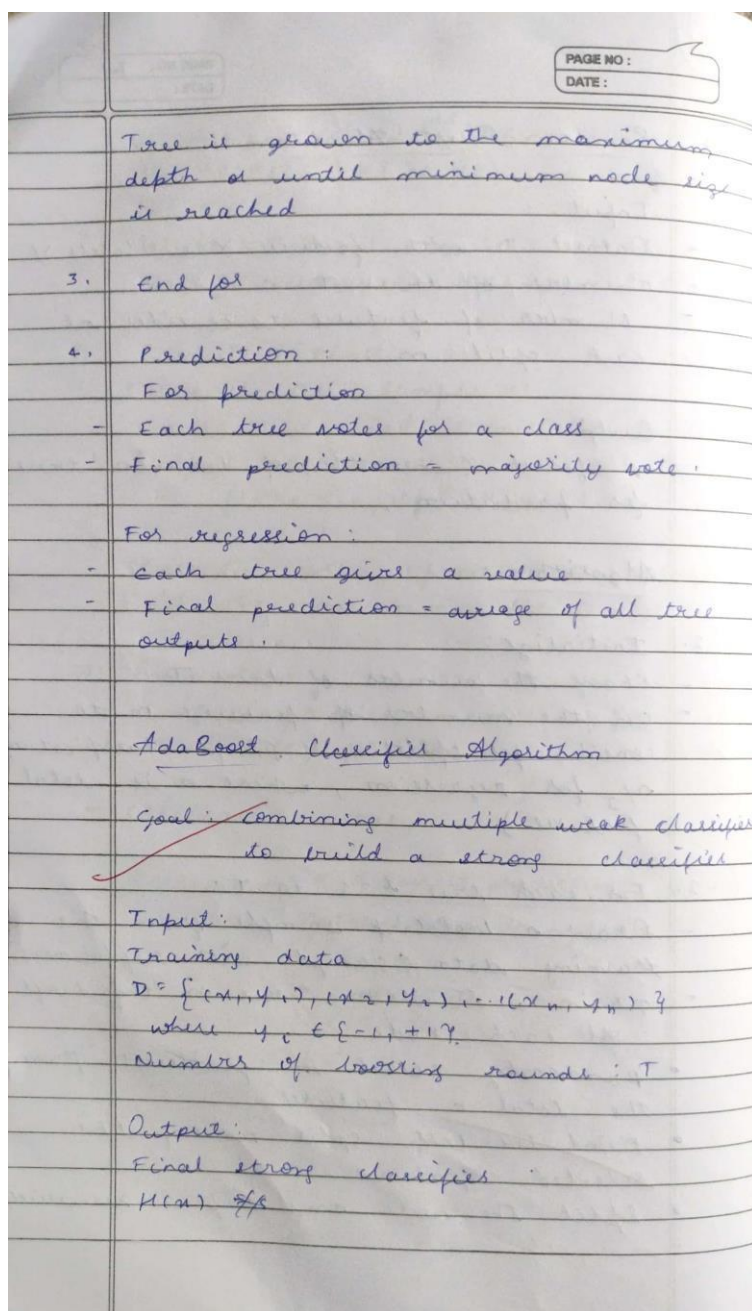
**Program 9**

Implement Boosting ensemble method on a given dataset

**Screenshot:**

Tree is grown to the maximum depth or until minimum node size is reached

3. End for

4. Prediction :
For prediction
- Each tree votes for a class
- Final prediction = majority vote.

For regression :
- Each tree gives a value
- Final prediction = average of all tree outputs.

Adaboost Classifier Algorithm

Goal : combining multiple weak classifier to build a strong classifier

Input :
Training data
$D = \{(x_1, y_1), (x_2, y_2), \ldots (x_n, y_n)\}$
where $y_i \in \{-1, +1\}$
Number of boosting rounds : T

Output :
Final strong classifier
$H(x)$

**Code:**

from



sklearn.ensemble import AdaBoostClassifier from

sklearn.datasets import load_iris from

sklearn.model_selection import train_test_split

from sklearn.metrics import accuracy_score


\# Load Iris dataset iris

= load_iris()

X, y = iris.data, iris.target
\# For AdaBoost, we'll use binary classification \#

Convert to binary (setosa vs. not-setosa)

y = (y == 0).astype(int)

# Split data

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


# Train AdaBoost

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

model.fit(X_train, y_train)


# Predict and evaluate

y_pred = model.predict(X_test)

print("AdaBoost Accuracy (sklearn):", accuracy_score(y_test, y_pred))
```

**Program 10**

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

**Screenshot:**

# K-Mean Clustering

**Input:**

A dataset with n data points
$$X = \{ x_1, x_2, x_3 \ldots x_n \}$$
numbers of clusters = k.

**Algorithm**

1) Load the Dataset
   Read the csv file to extract data point
   $X = \ldots$
   Store the data in suitable structure

2) Preprocess the Data
   Handle missing values
   Normalize or standardize the features
   to bring them to similar scale.

3) Initialize Centroids
   Randomly select k clusters data points
   from the dataset as initial centroids

4) Repeat until convergence or a set number
   of iterations

   Assign Cluster
   For each data point, calculate the
   distance to each centroids.
   Assign the point to cluster of the
   nearest centroid

   Update Centroids
   For each cluster, compute the mean
   of the point assigned to it
   Update the centroid with this new

5) Convergence check
   If cluster assignment do not change
   or centroids remain the same step

6) Output:
   Final centroids of all clusters
   Cluster labels for each data point

**Code:**

```
import pandas as pd

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

from sklearn.datasets import load_iris # Import load_iris


# Step 1: Load the Iris dataset directly iris

= load_iris()

# Create a DataFrame from the data and target

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

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

= iris.target




# Step 2: Extract only numeric columns (or select required features)

# All features in the Iris dataset are numeric

X = data[iris.feature_names].values # Use the feature names to select columns



# Step 3: Apply KMeans

# Adjust n_clusters based on the expected number of clusters in your data (3 for Iris)

kmeans = KMeans(n_clusters=3, random_state=42, n_init=10) # Added n_init to suppress future
warnings

data['Cluster'] = kmeans.fit_predict(X)
```

# Step 4: Plot clusters (for 2D data)

# Iris data has 4 features. We will plot the first two features for visualization. if

X.shape[1] >= 2:

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

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

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

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

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

plt.show()
```

else:

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

# Program 11

Implement Dimensionality reduction using Principal Component Analysis (PCA) method

**Screenshot:**

# Principal Component Analysis

1) Standardize the Data.
   - Given a dataset $X$ of $n \times d$ ord.
     $n \rightarrow$ no of samples
     $d \rightarrow$ number of features

   $$X_{centered} = X - \alpha.$$

2) Compute the covariance matrix.
   - Calculate the covariance matrix of the centred data. This matrix captures the relationship between different features.

   $$C = \frac{1}{n-1} \times X_{centered}^T \cdot X_{centered}$$

3) Compute the Eigen Values & Eigen Vectors.
   - Find the eigenvalue and eigen vector.
   (Directions of the new feature space)
   (variance)

4) Sort Eigen values and Eigenvectors.
   - Sort the eigen values in descending order
     $$\lambda_1 > \lambda_2 > \lambda_3 > \lambda_4 \ldots \lambda_d$$

5) Select the Top $k$ Eigen Vectors
   - Choose the top $k$ Eigen Vectors
   corresponding to the $k$ largest eigenvalue

6) Construction of Projection Matrix.
   $$W = [V_1, V_2, V_3 \ldots V_k]$$

7) Project the Data.
   $$X_{reduced} = X_{centered} \cdot W$$

   Multiply the centered data $X_{centered}$ by the projection matrix $W$ to obtain the reduced dataset in the new $k$ dimension $k$ space.

**Code:**

```
import pandas as pd

from sklearn.decomposition import PCA

from sklearn.preprocessing import StandardScaler
import matplotlib.pyplot as plt



# Load dataset

data = pd.read_csv("your_data.csv") # Replace with your file

X = data.select_dtypes(include=['float64', 'int64'])



# Step 1: Standardize scaler

= StandardScaler()

X_scaled = scaler.fit_transform(X)



# Step 2: Apply PCA

pca = PCA(n_components=2)

X_pca = pca.fit_transform(X_scaled)



# Print explained variance ratio

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



# Visualize
```

```
plt.scatter(X_pca[:, 0], X_pca[:, 1], c='blue', alpha=0.5) plt.title("PCA

- 2D Projection")

plt.xlabel("Principal Component 1")

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

```
⮞  📋  Accuracy Before PCA:
    Logistic Regression: 0.9016
    SVM: 0.8525
    Random Forest: 0.8361

    📉  Accuracy After PCA (n_components=5):
    Logistic Regression: 0.8689
    SVM: 0.8689
    Random Forest: 0.8852
```