

# Final Lab Report

## Machine Learning

Subject Code -SE204a



**Submitted to Ms Anjali Bansal**

Submitted by Gaurav Meena

Roll Number - 2K22/SE/65

Department of software engineering

DEPARTMENT OF SOFTWARE ENGINEERING

DELHI TECHNOLOGICAL UNIVERSITY

SHAHBAD, DAULATPUR

MAIN BAWANA ROAD,

Delhi 110042

# INDEX

Serial Number	Experiment	Date	Page Number	Signature
1.	To implement basics of python			
2.	To learn and implement data representation in Python			
3.	Implementation of ID3 Decision Tree Classifier			
4.	Implementation of C 4.5 Decision Tree Classifier			
5.	Implementation of CART Decision Tree Classifier			
6.	Build an Artificial Neural Network by implementing the Backpropagation algorithm and test the same using appropriate data sets			
7.	Write a program to implement the Naïve Bayesian classifier for the appropriate dataset and compute the performance measures of the model			
8.	Write a program to implement the k-nearest Neighbour algorithm to classify any dataset and compute the performance measures of the model			
9.	Apply the k-Means Clustering algorithm on suitable datasets and comment on the quality of clustering			
10.	Write a program to implement ensemble algorithms - AdaBoost and Bagging using the appropriate dataset and evaluate their performance on that dataset			

# INDEX

Serial Number	Experiment	Page Number	Signature
1.	To implement basics of python		
2.	To learn and implement data representation in Python		
3.	Implementation of ID3 Decision Tree Classifier		
4.	Implementation of C 4.5 Decision Tree Classifier		
5.	Implementation of CART Decision Tree Classifier		
6.	Build an Artificial Neural Network by implementing the Backpropagation algorithm and test the same using appropriate data sets		
7.	Write a program to implement the Naïve Bayesian classifier for the appropriate dataset and compute the performance measures of the model		
8.	Write a program to implement the k-nearest Neighbour algorithm to classify any dataset and compute the performance measures of the model		
9.	Apply the k-Means Clustering algorithm on suitable datasets and comment on the quality of clustering		
10.	Write a program to implement ensemble algorithms - AdaBoost and Bagging using the appropriate dataset and evaluate their performance on that dataset		

# Experiment – 1

Aim: To implement basics of python

Theory:

PYTHON IS A HIGH-LEVEL, GENERAL-PURPOSE PROGRAMMING LANGUAGE. ITS DESIGN PHILOSOPHY EMPHASIZES CODE READABILITY WITH THE USE OF SIGNIFICANT INDENTATION.

FEATURES OF PYTHON LANGUAGE:

- HIGH-LEVEL LANGUAGE
- DYNAMICALLY-TYPED
- GARBAGE-COLLECTED
- IT SUPPORTS MULTIPLE PROGRAMMING PARADIGMS, SUCH AS STRUCTURED (PARTICULARLY PROCEDURAL PROGRAMMING), OBJECT-ORIENTED PROGRAMMING, AND FUNCTIONAL PROGRAMMING.
- CONTAINS A COMPREHENSIVE STANDARD LIBRARY.
- MOST SUITABLE PROGRAMMING LANGUAGE FOR DATA SCIENCE
- ONE OF THE MOST POPULAR PROGRAMMING LANGUAGES

NUMPY: NUMPY IS A LIBRARY FOR THE PYTHON PROGRAMMING LANGUAGE, ADDING SUPPORT FOR LARGE, MULTIDIMENSIONAL ARRAYS AND MATRICES, ALONG WITH A LARGE COLLECTION OF HIGH-LEVEL MATHEMATICAL FUNCTIONS TO OPERATE ON THESE MULTI-DIMENSIONAL ARRAYS.

MATPLOTLIB: MATPLOTLIB IS A PLOTTING LIBRARY FOR THE PYTHON PROGRAMMING LANGUAGE. IT CAN BE USED TO DRAW A WIDE VARIETY OF PLOTS INCLUDING LINE PLOTS, SCATTER PLOTS, BAR CHARTS, AND HISTOGRAMS.

PANDAS: PANDAS (DERIVED FROM THE TERM "PANEL DATA") IS A SOFTWARE LIBRARY FOR THE PYTHON PROGRAMMING LANGUAGE FOR DATA MANIPULATION AND ANALYSIS. IN PARTICULAR, IT OFFERS DATA STRUCTURES AND OPERATIONS FOR MANIPULATING NUMERICAL TABLES AND TIME SERIES.

SCIPY: SCIPY IS A PYTHON LIBRARY USED FOR SCIENTIFIC COMPUTING AND TECHNICAL COMPUTING. SCIPY CONTAINS MODULES FOR OPTIMIZATION, LINEAR ALGEBRA, INTEGRATION, INTERPOLATION, SPECIAL FUNCTIONS, FFT, SIGNAL AND IMAGE PROCESSING,

SCIKIT-LEARN: SCIKIT-LEARN (ALSO CALLED SKLEARN) IS A FREE SOFTWARE MACHINE-LEARNING LIBRARY FOR THE PYTHON PROGRAMMING LANGUAGE. IT FEATURES VARIOUS CLASSIFICATION, REGRESSION, AND CLUSTERING ALGORITHMS INCLUDING SUPPORT--VECTOR MACHINES, RANDOM FORESTS, GRADIENT BOOSTING, K-MEANS, AND DBSCAN.

Learning:

IN THIS EXPERIMENT, I WAS ABLE TO LEARN:

- THE CONCEPTS OF PYTHON LANGUAGE
- THE USE OF VARIOUS PYTHON LIBRARIES
- THE USE OF SEVERAL ML LIBRARIES FOR USING BUILT-IN ML MODELS.

## Experiment – 2

Aim: To learn and implement data representation in Python

Theory:

THE LIBRARY KNOWN AS NUMPY, OR "NUMERICAL PYTHON," CONTAINS MULTIDIMENSIONAL ARRAY OBJECTS AND A SELECTION OF PROCEDURES FOR HANDLING SUCH ARRAYS. ARRAYS MAY BE SUBJECTED TO LOGICAL AND MATHEMATICAL OPERATIONS USING NUMPY. IT IS A LIBRARY MADE UP OF FUNCTIONS FOR HANDLING ARRAYS AND MULTIDIMENSIONAL ARRAY OBJECTS. LARGE, MULTI-DIMENSIONAL ARRAYS AND MATRICES ARE NOW SUPPORTED BY NUMPY, ALONG WITH A VAST VARIETY OF SOPHISTICATED MATHEMATICAL OPERATIONS THAT MAY BE PERFORMED ON THESE ARRAYS.

Code:

```
#importing libraries
import numpy as np
import pandas as pd
import matplotlib as plt
import seaborn as sns

#Loading the dataset
dataset = pd.read_csv("Bank Customer Churn Prediction.csv")

#printing first 5 instances of the dataset
dataset.head(5)

#number of dimensions in the dataset
print("Total number of dimension in the given dataset:")
print(dataset.ndim)

#shape of the dataset
print("Shape of the dataset:")
print(dataset.shape)

#removing null or missing values
dataset = dataset.dropna()

#Null vs Not null values graph
null_counts = dataset.isnull().sum() # Count null values for each column
nonnull_counts = dataset.notnull().sum() # Count non-null values for each column
# Create a plot to visualize the counts
plt.figure(figsize=(10, 6))
# Bar plot for non-null values
sns.barplot(x=nonnull_counts.index, y=nonnull_counts.values, color='blue',
alpha=0.6, label='Non-Null Values')
# Bar plot for null values
sns.barplot(x=null_counts.index, y=null_counts.values, color='red', alpha=0.6,
label='Null Values')
# Adding labels and title
plt.xlabel('Columns')
plt.ylabel('Count')
plt.title('Null vs Non-Null Values Count for Each Column')
plt.xticks(rotation=90) # Rotate x-labels for better readability
plt.legend()
# Show plot
plt.show()
```

Output:

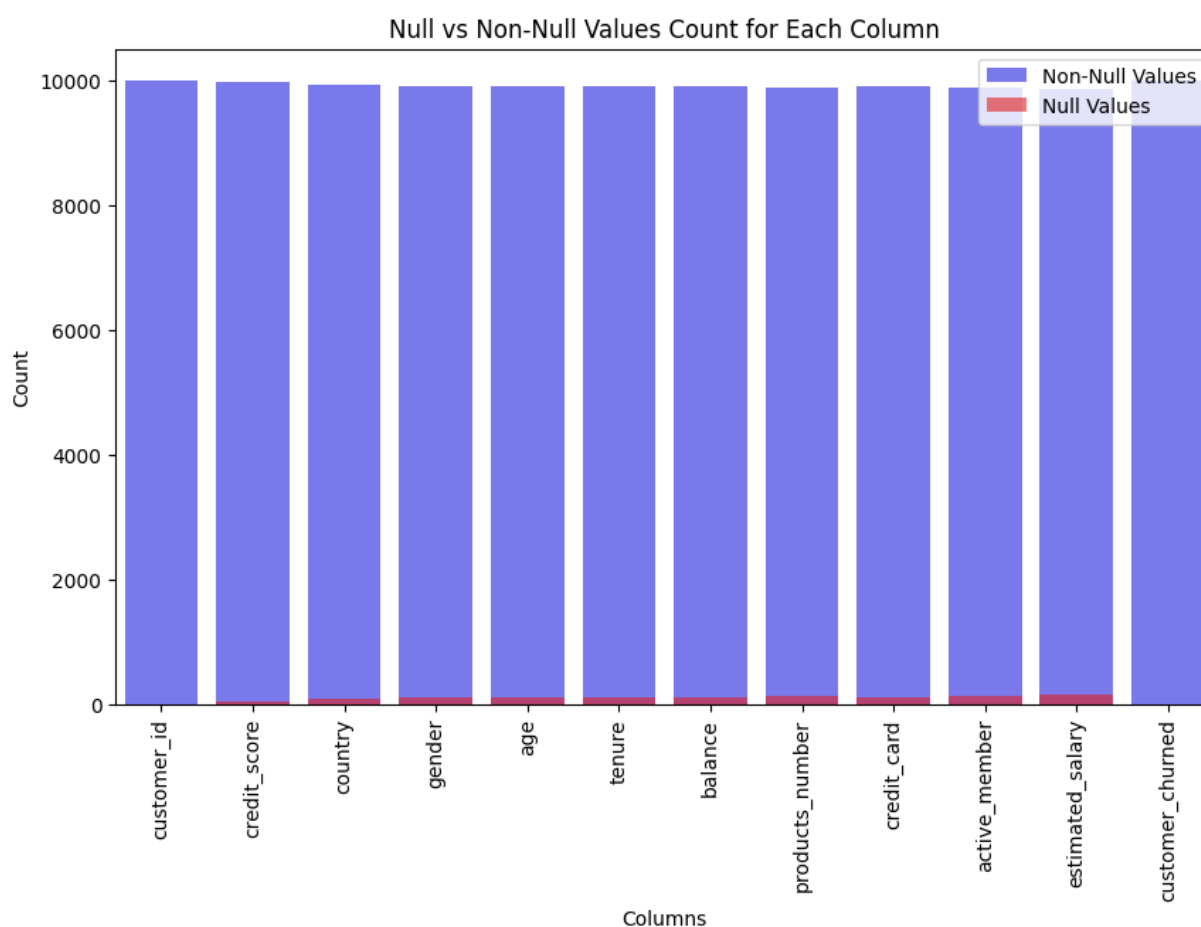
	customer_id	credit_score	country	gender	age	tenure	balance	products_number	credit_card	active_member	estimated_salary	customer_churned
0	15634602.0	619.0	France	Female	42.0	2.0	0.00	1.0	1.0	1.0	101348.88	1
1	15647311.0	608.0	Spain	Female	41.0	1.0	83807.86	1.0	0.0	1.0	112542.58	0
2	15619304.0	502.0	France	Female	42.0	8.0	159660.80	3.0	1.0	0.0	113931.57	1
3	15701354.0	NaN	France	Female	39.0	1.0	0.00	2.0	0.0	0.0	93826.63	0
4	15737888.0	850.0	Spain	Female	43.0	2.0	125510.82	1.0	1.0	1.0	79084.10	0

Total number of dimension in the given dataset:

2

Shape of the dataset:

(9995, 12)



Learning:

NUMPY IS OFTEN USED ALONG WITH PACKAGES LIKE SCIPY (SCIENTIFIC PYTHON) AND MATPLOTLIB (PLOTING LIBRARY). USING NUMPY, A DEVELOPER CAN PERFORM THE FOLLOWING OPERATIONS

- MATHEMATICAL AND LOGICAL OPERATIONS ON ARRAYS.
- FOURIER TRANSFORMS AND ROUTINES FOR SHAPE MANIPULATION.
- OPERATIONS RELATED TO LINEAR ALGEBRA. NUMPY HAS IN-BUILT FUNCTIONS FOR LINEAR ALGEBRA AND RANDOM NUMBER GENERATION.

## Experiment – 3

### Aim: Implementation of ID3 Decision Tree Classifier

#### Theory:

A DECISION TREE IS A TREE-BASED DATA-STRUCTURE THAT CONTAINS NODES (RECTANGULAR BOXES) AND EDGES (ARROWS) AND IS BUILT FROM A DATASET (TABLE OF COLUMNS REPRESENTING FEATURES/ATTRIBUTES AND ROWS CORRESPONDS TO RECORDS). EACH NODE IS EITHER USED TO MAKE A DECISION (KNOWN AS DECISION NODE) OR REPRESENT AN OUTCOME (KNOWN AS LEAF NODE).

ID3 STANDS FOR ITERATIVE DICHOTOMISE 3 AND IS NAMED SUCH BECAUSE THE ALGORITHM ITERATIVELY (REPEATEDLY) DICHOTOMIZES(DIVIDES) FEATURES INTO TWO OR MORE GROUPS AT EACH STEP. INVENTED BY ROSS QUINLAN, ID3 USES A TOP-DOWN GREEDY APPROACH TO BUILD A DECISION TREE. IN SIMPLE WORDS, THE TOP-DOWN APPROACH MEANS THAT WE START BUILDING THE TREE FROM THE TOP AND THE GREEDY APPROACH MEANS THAT AT EACH ITERATION WE SELECT THE BEST FEATURE AT THE PRESENT MOMENT TO CREATE A NODE. MOST GENERALLY ID3 IS ONLY USED FOR CLASSIFICATION PROBLEMS WITH NOMINAL FEATURES ONLY.

#### METRICS IN ID3:

ENTROPY: ENTROPY IS THE MEASURE OF DISORDER AND THE ENTROPY OF A DATASET IS THE MEASURE OF DISORDER IN THE TARGET FEATURE OF THE DATASET. IN THE CASE OF BINARY CLASSIFICATION (WHERE THE TARGET COLUMN HAS ONLY TWO TYPES OF CLASSES) ENTROPY IS 0 IF ALL VALUES IN THE TARGET COLUMN ARE HOMOGENOUS(SIMILAR) AND WILL BE 1 IF THE TARGET COLUMN HAS EQUAL NUMBER VALUES FOR BOTH THE CLASSES.

$$\text{ENTROPY}(S) = - \sum p_i * \log_2(p_i); i = 1 \text{ TO } N$$

WHERE, N IS THE TOTAL NUMBER OF CLASSES IN THE TARGET COLUMN;  $p_i$  IS THE PROBABILITY OF CLASS 'i' OR THE RATIO OF "NUMBER OF ROWS WITH CLASS i IN THE TARGET COLUMN" TO THE "TOTAL NUMBER OF ROWS" IN THE DATASET.

INFORMATION GAIN: INFORMATION GAIN CALCULATES THE REDUCTION IN THE ENTROPY AND MEASURES HOW WELL A GIVEN FEATURE SEPARATES OR CLASSIFIES THE TARGET CLASSES. THE FEATURE WITH THE HIGHEST INFORMATION GAIN IS SELECTED AS THE BEST ONE.

$$\text{IG}(S, A) = \text{ENTROPY}(S) - \sum (|S_v| / |S|) * \text{ENTROPY}(S_v)$$

WHERE,  $S_v$  IS THE SET OF ROWS IN S FOR WHICH THE FEATURE COLUMN A HAS VALUE v,  $|S_v|$  IS THE NUMBER OF ROWS IN  $S_v$  AND LIKEWISE  $|S|$  IS THE NUMBER OF ROWS IN S.

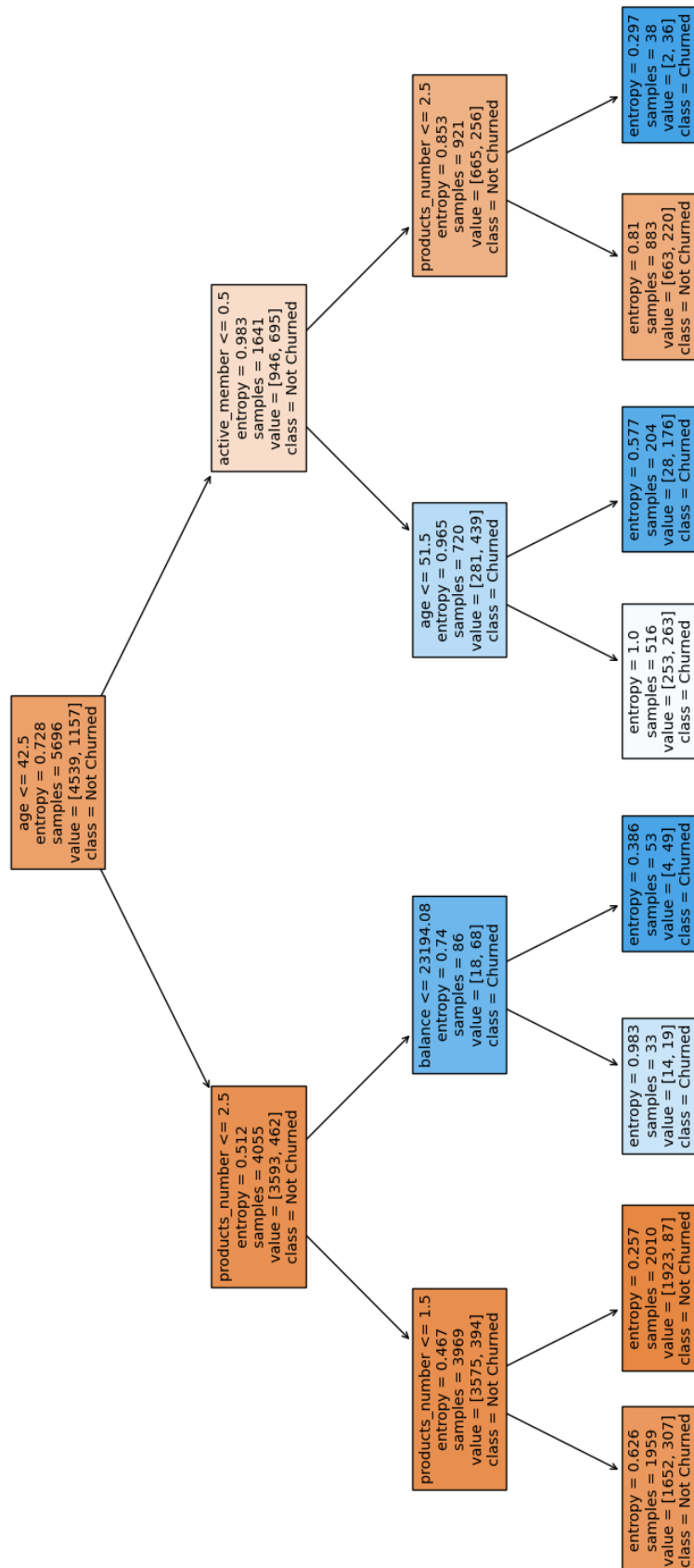
#### ID3 STEPS:

1. CALCULATE THE INFORMATION GAIN OF EACH FEATURE.
2. CONSIDERING THAT ALL ROWS DON'T BELONG TO THE SAME CLASS, SPLIT THE DATASET S INTO SUBSETS USING THE FEATURE FOR WHICH THE INFORMATION GAIN IS MAXIMUM.
3. MAKE A DECISION TREE NODE USING THE FEATURE WITH THE MAXIMUM INFORMATION GAIN.
4. IF ALL ROWS BELONG TO THE SAME CLASS, MAKE THE CURRENT NODE AS A LEAF NODE WITH THE CLASS AS ITS LABEL.
5. REPEAT FOR THE REMAINING FEATURES UNTIL WE RUN OUT OF ALL FEATURES, OR THE DECISION TREE HAS ALL LEAF NODES.

#### Code:

```
dt_id3 = DecisionTreeClassifier(criterion='entropy', max_depth=3)
dt_id3.fit(X_train, Y_train)
plt.figure(figsize=(20,10))
plot_tree(dt_id3, feature_names=list(X.columns), class_names=['Not Churned',
'Churned'], filled=True)
plt.show()
```

Output:



Learning:

THE EXPERIMENT SHOULD HELP THE LEARNERS UNDERSTAND THE SIGNIFICANCE OF THE IMPLEMENTATION OF THE ID3 DECISION TREE CLASSIFIER.



## Experiment – 4

### Aim: Implementation of C 4.5 Decision Tree Classifier

#### Theory:

C4.5 IS A DECISION TREE ALGORITHM DEVELOPED BY ROSS QUINLAN THAT IS COMMONLY USED FOR MACHINE LEARNING AND DATA MINING TASKS. IT IS AN IMPROVEMENT OVER THE EARLIER ID3 ALGORITHM AND CAN HANDLE BOTH CONTINUOUS AND DISCRETE ATTRIBUTES.

THE C4.5 ALGORITHM BUILDS A DECISION TREE BY RECURSIVELY SPLITTING THE DATASET BASED ON THE ATTRIBUTE THAT PROVIDES THE MOST INFORMATION GAIN. INFORMATION GAIN MEASURES THE REDUCTION IN ENTROPY (OR IMPURITY) OF THE DATA WHEN A GIVEN ATTRIBUTE IS USED FOR SPLITTING. THE ALGORITHM STOPS WHEN ALL INSTANCES IN A NODE BELONG TO THE SAME CLASS OR WHEN NO MORE USEFUL ATTRIBUTES CAN BE FOUND.

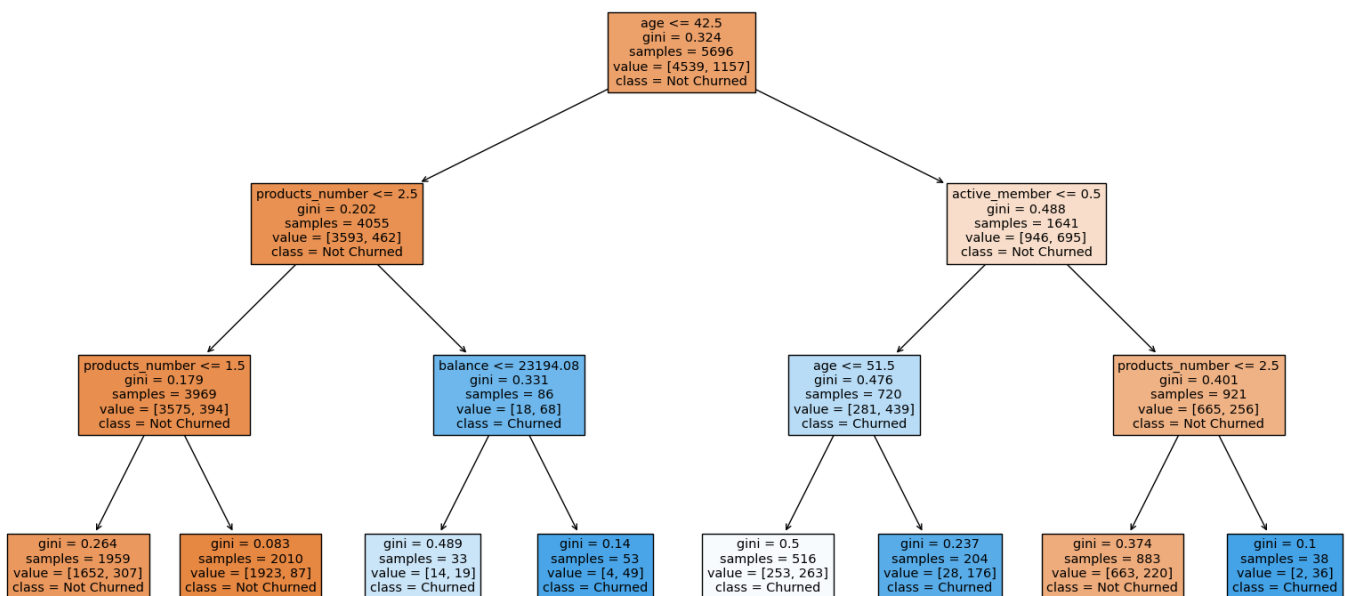
C4.5 ALSO INCLUDES A MECHANISM FOR HANDLING MISSING ATTRIBUTE VALUES AND FOR PRUNING THE DECISION TREE TO AVOID OVERFITTING TO THE TRAINING DATA.

ONCE THE DECISION TREE IS BUILT, IT CAN BE USED TO CLASSIFY NEW INSTANCES BY TRAVERSING THE TREE FROM THE ROOT TO A LEAF NODE, WHERE THE PREDICTED CLASS IS ASSIGNED BASED ON THE MAJORITY CLASS OF THE TRAINING INSTANCES IN THAT LEAF NODE.

#### Code:

```
dt_c45 = DecisionTreeClassifier(criterion='gini', max_depth=3)
dt_c45.fit(X_train, Y_train)
plt.figure(figsize=(20,10))
plot_tree(dt_c45, feature_names=list(X.columns), class_names=['Not Churned',
'Churned'], filled=True)
plt.show()
```

#### Output:



#### Learning:

THE EXPERIMENT SHOULD HELP THE LEARNERS UNDERSTAND THE SIGNIFICANCE OF THE IMPLEMENTATION OF THE CART DECISION TREE CLASSIFIER.

## Experiment – 5

### Aim: Implementation of CART Decision Tree Classifier

#### Theory:

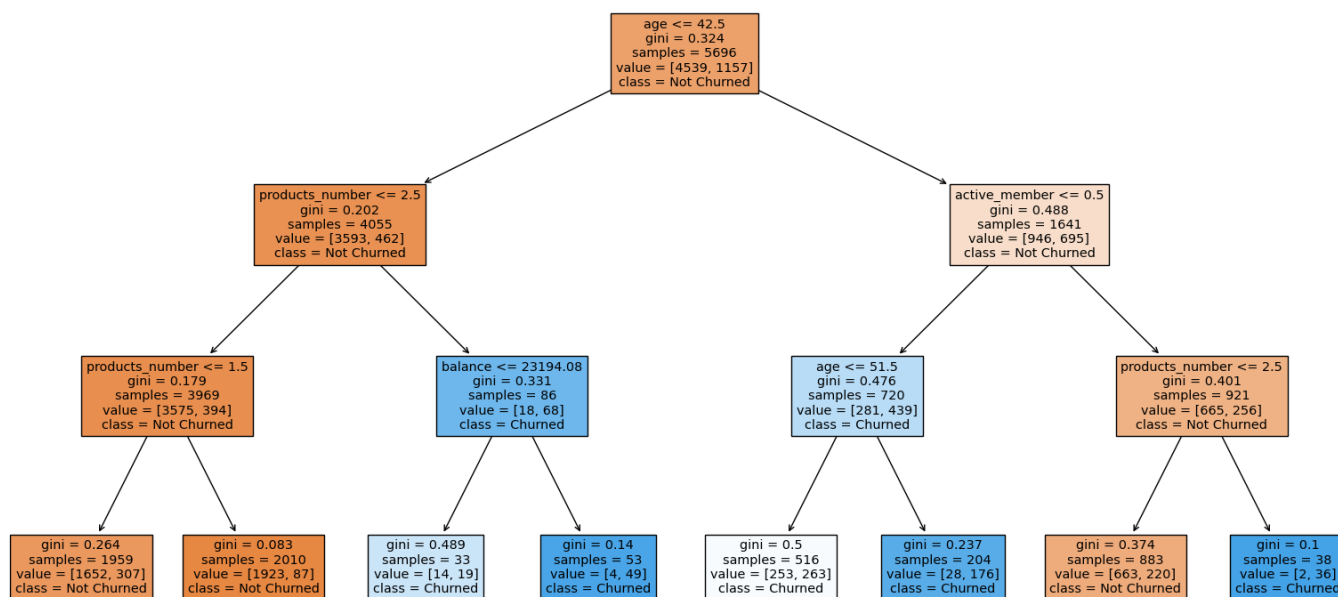
DECISION TREE LEARNING USES A DECISION TREE (AS A PREDICTIVE MODEL) TO GO FROM OBSERVATIONS ABOUT AN ITEM (REPRESENTED IN THE BRANCHES) TO CONCLUSIONS ABOUT THE ITEM'S TARGET VALUE (REPRESENTED IN THE LEAVES). IT IS ONE OF THE PREDICTIVE MODELLING APPROACHES USED IN STATISTICS, DATA MINING AND MACHINE LEARNING. TREE MODELS WHERE THE TARGET VARIABLE CAN TAKE A DISCRETE SET OF VALUES ARE CALLED CLASSIFICATION TREES; IN THESE TREE STRUCTURES, LEAVES REPRESENT CLASS LABELS AND BRANCHES REPRESENT CONJUNCTIONS OF FEATURES THAT LEAD TO THOSE CLASS LABELS. DECISION TREES WHERE THE TARGET VARIABLE CAN TAKE CONTINUOUS VALUES (TYPICALLY REAL NUMBERS) ARE CALLED REGRESSION TREES. GINI IMPURITY IS A MEASURE OF HOW OFTEN A RANDOMLY CHOSEN ELEMENT FROM THE SET WOULD BE INCORRECTLY LABELLED IF IT WAS RANDOMLY LABELLED ACCORDING TO THE DISTRIBUTION OF LABELS IN THE SUBSET.

THE GINI IMPURITY CAN BE COMPUTED BY SUMMING THE PROBABILITY  $P(i)$  OF AN ITEM WITH LABEL  $i$  BEING CHOSEN TIMES THE PROBABILITY OF A MISTAKE IN CATEGORIZING THAT ITEM. IT REACHES ITS MINIMUM (ZERO) WHEN ALL CASES IN THE NODE FALL INTO A SINGLE TARGET CATEGORY.

#### Code:

```
dt_clf = DecisionTreeClassifier(criterion='gini', max_depth=3)
plt.figure(figsize=(20, 10))
plot_tree(dt_clf, feature_names=X.columns, class_names=['Not Churned', 'Churned'],
          filled=True)
plt.show()
```

#### Output:



## Experiment – 6

Aim: Build an Artificial Neural Network by implementing the Backpropagation algorithm and test the same using appropriate data sets

Theory:

AN ARTIFICIAL NEURAL NETWORK (ANN) IS A TYPE OF MACHINE-LEARNING MODEL THAT IS INSPIRED BY THE HUMAN BRAIN. IT IS COMPOSED OF LAYERS OF INTERCONNECTED NEURONS THAT PROCESS AND TRANSMIT INFORMATION. THE BACKPROPAGATION ALGORITHM IS ONE OF THE MOST WIDELY USED ALGORITHMS FOR TRAINING ANN. IT IS A SUPERVISED LEARNING ALGORITHM THAT IS USED TO ADJUST THE WEIGHTS OF THE NETWORK SO THAT THE NETWORK CAN LEARN TO MAP INPUTS TO OUTPUTS.

THE BACKPROPAGATION ALGORITHM WORKS BY CALCULATING THE ERROR BETWEEN THE PREDICTED OUTPUT OF THE NETWORK AND THE ACTUAL OUTPUT. IT THEN ADJUSTS THE WEIGHTS OF THE NETWORK IN SUCH A WAY THAT THE ERROR IS REDUCED. THIS PROCESS IS REPEATED OVER MULTIPLE ITERATIONS UNTIL THE ERROR IS MINIMIZED.

Code:

```
import numpy as np

class NeuralNetwork:
    def __init__(self, input_size, hidden_size, output_size, learning_rate):
        self.input_size = input_size
        self.hidden_size = hidden_size
        self.output_size = output_size
        self.learning_rate = learning_rate

        # Initialize weights and biases randomly
        self.weights_input_hidden = np.random.randn(input_size, hidden_size)
        self.biases_input_hidden = np.random.randn(hidden_size)
        self.weights_hidden_output = np.random.randn(hidden_size, output_size)
        self.biases_hidden_output = np.random.randn(output_size)

    def sigmoid(self, x):
        return 1 / (1 + np.exp(-x))

    def sigmoid_derivative(self, x):
        return x * (1 - x)

    def forward_pass(self, inputs):
        # Calculate hidden layer activations
        hidden_inputs = np.dot(inputs, self.weights_input_hidden) +
self.biases_input_hidden
        hidden_outputs = self.sigmoid(hidden_inputs)

        # Calculate output layer activations
        final_inputs = np.dot(hidden_outputs, self.weights_hidden_output) +
self.biases_hidden_output
        final_outputs = self.sigmoid(final_inputs)

        return final_outputs, hidden_outputs

    def backward_pass(self, inputs, targets, outputs, hidden_outputs):
        # Compute output layer error
        output_errors = targets - outputs
        output_gradients = output_errors * self.sigmoid_derivative(outputs)
```

```

        # Update weights and biases for hidden to output layer
        self.weights_hidden_output += np.dot(hidden_outputs.T, output_gradients) *
self.learning_rate
        self.biases_hidden_output += np.sum(output_gradients, axis=0) *
self.learning_rate

        # Compute hidden layer error
        hidden_errors = np.dot(output_gradients, self.weights_hidden_output.T)
        hidden_gradients = hidden_errors * self.sigmoid_derivative(hidden_outputs)

        # Update weights and biases for input to hidden layer
        self.weights_input_hidden += np.dot(inputs.T, hidden_gradients) *
self.learning_rate
        self.biases_input_hidden += np.sum(hidden_gradients, axis=0) *
self.learning_rate

    def train(self, inputs, targets, epochs):
        for epoch in range(epochs):
            outputs, hidden_outputs = self.forward_pass(inputs)
            self.backward_pass(inputs, targets, outputs, hidden_outputs)

            if epoch % 100 == 0:
                loss = np.mean(np.square(targets - outputs))
                print(f"Epoch {epoch}, Loss: {loss}")

# Example usage
# Define input, hidden, and output layer sizes
input_size = 10
hidden_size = 5
output_size = 1

# Create a neural network
nn = NeuralNetwork(input_size, hidden_size, output_size, learning_rate=0.1)

# Generate synthetic dataset
X_train = np.random.rand(100, input_size)
y_train = np.random.randint(0, 2, size=(100, output_size))

# Train the neural network
nn.train(X_train, y_train, epochs=1000)

```

Output:

```

Epoch 0, Loss: 0.37320200391211117
Epoch 100, Loss: 0.23476138652802273
Epoch 200, Loss: 0.21115199101453416
Epoch 300, Loss: 0.19607072200874393
Epoch 400, Loss: 0.18606742759657494
Epoch 500, Loss: 0.17158831366066388
Epoch 600, Loss: 0.1595023603992833
Epoch 700, Loss: 0.14707143430528563
Epoch 800, Loss: 0.13667280895709216
Epoch 900, Loss: 0.12822250773232358

```

## Learning:

IMPLEMENTING THE BACKPROPAGATION ALGORITHM TO TRAIN AN ANN CAN HELP IN ACHIEVING HIGH ACCURACY IN PREDICTING OUTCOMES. BY TESTING THE ANN WITH APPROPRIATE DATASETS, ONE CAN ANALYSE THE EFFECTIVENESS OF THE ALGORITHM IN DIFFERENT SCENARIOS. THE EXPERIMENT CAN ALSO HELP IN UNDERSTANDING THE IMPACT OF DIFFERENT HYPERPARAMETERS LIKE THE NUMBER OF LAYERS, NUMBER OF NEURONS, AND LEARNING RATE ON THE PERFORMANCE OF THE ANN. ADDITIONALLY, THE EXPERIMENT CAN ALSO HELP IN IDENTIFYING THE LIMITATIONS OF THE ALGORITHM AND THE DATASET USED.

## Experiment – 7

Aim: Write a program to implement the Naïve Bayesian classifier for the appropriate dataset and compute the performance measures of the model

Theory:

THE NAÏVE BAYESIAN CLASSIFIER IS A PROBABILISTIC ALGORITHM USED IN MACHINE LEARNING FOR CLASSIFICATION TASKS. IT IS BASED ON BAYES' THEOREM AND ASSUMES THAT THE PRESENCE OF A PARTICULAR FEATURE IN A CLASS IS INDEPENDENT OF THE PRESENCE OF ANY OTHER FEATURE IN THE SAME CLASS. THE NAÏVE BAYESIAN CLASSIFIER CALCULATES THE PROBABILITY OF A GIVEN DATA POINT BELONGING TO A CERTAIN CLASS BY COMPUTING THE CONDITIONAL PROBABILITY OF EACH FEATURE GIVEN THE CLASS AND MULTIPLYING THEM TOGETHER.

TO IMPLEMENT THE NAÏVE BAYESIAN CLASSIFIER, WE FIRST NEED TO TRAIN THE MODEL USING A DATASET WITH LABELLED DATA. THE LABELLED DATA IS USED TO ESTIMATE THE CONDITIONAL PROBABILITIES OF EACH FEATURE GIVEN IN EACH CLASS. THEN, THE MODEL CAN BE USED TO PREDICT THE CLASS OF NEW DATA POINTS BASED ON THEIR FEATURES.

PERFORMANCE MEASURES OF THE MODEL CAN BE COMPUTED USING METRICS SUCH AS ACCURACY, PRECISION, RECALL, AND F1 SCORE. ACCURACY IS THE RATIO OF CORRECTLY PREDICTED DATA POINTS TO THE TOTAL NUMBER OF DATA POINTS. PRECISION IS THE RATIO OF TRUE POSITIVE PREDICTIONS TO THE TOTAL NUMBER OF POSITIVE PREDICTIONS, WHILE RECALL IS THE RATIO OF TRUE POSITIVE PREDICTIONS TO THE TOTAL NUMBER OF ACTUAL POSITIVE INSTANCES. THE F1 SCORE IS A WEIGHTED AVERAGE OF PRECISION AND RECALL.

Code:

```
import numpy as np

class NaiveBayesClassifier:
    def __init__(self):
        self.class_probabilities = None
        self.feature_probabilities = None

    def fit(self, X_train, y_train):
        num_samples, num_features = X_train.shape
        self.classes = np.unique(y_train)
        num_classes = len(self.classes)

        # Calculate class probabilities
        self.class_probabilities = np.zeros(num_classes)
        for i, c in enumerate(self.classes):
            self.class_probabilities[i] = np.sum(y_train == c) / num_samples

        # Calculate feature probabilities
        self.feature_probabilities = {}
        for c in self.classes:
            self.feature_probabilities[c] = []
            for feature_index in range(num_features):
                feature_values = X_train[y_train == c, feature_index]
                probabilities = {}
                for value in np.unique(X_train[:, feature_index]):
                    probabilities[value] = np.sum(feature_values == value) /
len(feature_values)
                self.feature_probabilities[c].append(probabilities)

    def predict(self, X_test):
        num_samples = X_test.shape[0]
        predictions = np.zeros(num_samples, dtype=int)
```

```

        for i in range(num_samples):
            probabilities = []
            for j, c in enumerate(self.classes):
                class_probability = np.log(self.class_probabilities[j])
                for feature_index, value in enumerate(X_test[i]):
                    if value in self.feature_probabilities[c][feature_index]:
                        class_probability +=
np.log(self.feature_probabilities[c][feature_index][value])
                probabilities.append(class_probability)
            predictions[i] = self.classes[np.argmax(probabilities)]
        return predictions

# Example usage
X_train = np.array([[1, 1], [1, 0], [0, 1], [0, 0]])
y_train = np.array([1, 1, 0, 0])

X_test = np.array([[1, 1], [0, 0]])

nb_classifier = NaiveBayesClassifier()
nb_classifier.fit(X_train, y_train)
predictions = nb_classifier.predict(X_test)

print("Predictions:", predictions)

```

Output:

```
Predictions: [1 0]
```

Learning:

IMPLEMENTING THE NAÏVE BAYESIAN CLASSIFIER FOR A DATASET CAN HELP IN ACHIEVING ACCURATE CLASSIFICATION OF DATA POINTS. THE PERFORMANCE MEASURES COMPUTED FOR THE MODEL CAN HELP IN ANALYSING THE EFFECTIVENESS OF THE ALGORITHM IN DIFFERENT SCENARIOS. THE EXPERIMENT CAN ALSO HELP IN UNDERSTANDING THE IMPACT OF DIFFERENT FEATURES ON THE PERFORMANCE OF THE CLASSIFIER.

ADDITIONALLY, THE EXPERIMENT CAN ALSO PROVIDE INSIGHTS INTO THE IMPORTANCE OF FEATURE SELECTION AND DATA PREPROCESSING FOR THE NAÏVE BAYESIAN CLASSIFIER. BY TESTING THE CLASSIFIER WITH DIFFERENT DATASETS, ONE CAN IDENTIFY THE STRENGTHS AND LIMITATIONS OF THE ALGORITHM AND IMPROVE ITS PERFORMANCE.

## Experiment – 8

**Aim:** Write a program to implement the k-nearest Neighbour algorithm to classify any dataset and compute the performance measures of the model

**Theory:**

THE K-NEAREST NEIGHBOUR (K-NN) ALGORITHM IS A SIMPLE YET POWERFUL ALGORITHM USED FOR CLASSIFICATION TASKS. IT IS A NON-PARAMETRIC ALGORITHM THAT STORES ALL THE AVAILABLE DATA POINTS IN MEMORY AND CLASSIFIES NEW DATA POINTS BASED ON THEIR PROXIMITY TO THE EXISTING DATA POINTS. THE K-NN ALGORITHM WORKS BY FIRST CALCULATING THE DISTANCE BETWEEN THE NEW DATA POINT AND ALL THE EXISTING DATA POINTS IN THE DATASET. IT THEN SELECTS THE K NEAREST DATA POINTS AND ASSIGNS THE CLASS TO THE NEW DATA POINT BASED ON THE MAJORITY CLASS OF THESE K NEAREST DATA POINTS.

TO IMPLEMENT THE K-NN ALGORITHM, WE FIRST NEED TO TRAIN THE MODEL USING A DATASET WITH LABELLED DATA. THE LABELLED DATA IS USED TO STORE THE DATA POINTS AND THEIR CORRESPONDING CLASS LABELS IN MEMORY. THEN, THE MODEL CAN BE USED TO PREDICT THE CLASS OF NEW DATA POINTS BASED ON THEIR PROXIMITY TO THE EXISTING DATA POINTS.

PERFORMANCE MEASURES OF THE MODEL CAN BE COMPUTED USING METRICS SUCH AS ACCURACY, PRECISION, RECALL, AND F1 SCORE. ACCURACY IS THE RATIO OF CORRECTLY PREDICTED DATA POINTS TO THE TOTAL NUMBER OF DATA POINTS. PRECISION IS THE RATIO OF TRUE POSITIVE PREDICTIONS TO THE TOTAL NUMBER OF POSITIVE PREDICTIONS, WHILE RECALL IS THE RATIO OF TRUE POSITIVE PREDICTIONS TO THE TOTAL NUMBER OF ACTUAL POSITIVE INSTANCES. THE F1 SCORE IS A WEIGHTED AVERAGE OF PRECISION AND RECALL.

**Code:**

```
import numpy as np
from collections import Counter

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

    def fit(self, X_train, y_train):
        self.X_train = X_train
        self.y_train = y_train

    def predict(self, X_test):
        predictions = []
        for x in X_test:
            distances = np.linalg.norm(self.X_train - x, axis=1)
            nearest_indices = distances.argsort()[:self.k]
            nearest_labels = self.y_train[nearest_indices]
            most_common_label = Counter(nearest_labels).most_common(1)[0][0]
            predictions.append(most_common_label)
        return np.array(predictions)

def accuracy(y_true, y_pred):
    return np.mean(y_true == y_pred)

# Example usage
X_train = np.array([[1, 2], [2, 3], [3, 4], [4, 5]])
y_train = np.array([0, 0, 1, 1])
X_test = np.array([[1, 2.5], [3.5, 4]])

knn_classifier = KNNClassifier(k=2)
knn_classifier.fit(X_train, y_train)
predictions = knn_classifier.predict(X_test)
```



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

# Compute accuracy
y_test = np.array([0, 1])
print("Accuracy:", accuracy(y_test, predictions))
```

Output:

```
Predictions: [0 1]
Accuracy: 1.0
```

Learning:

IMPLEMENTING THE K-NN ALGORITHM FOR A DATASET CAN HELP IN ACHIEVING ACCURATE CLASSIFICATION OF DATA POINTS. THE PERFORMANCE MEASURES COMPUTED FOR THE MODEL CAN HELP IN ANALYSING THE EFFECTIVENESS OF THE ALGORITHM IN DIFFERENT SCENARIOS. THE EXPERIMENT CAN ALSO HELP IN UNDERSTANDING THE IMPACT OF DIFFERENT VALUES OF K ON THE PERFORMANCE OF THE ALGORITHM.

ADDITIONALLY, THE EXPERIMENT CAN ALSO PROVIDE INSIGHTS INTO THE IMPORTANCE OF FEATURE SELECTION AND DATA PREPROCESSING FOR THE K-NN ALGORITHM. BY TESTING THE ALGORITHM WITH DIFFERENT DATASETS, ONE CAN IDENTIFY THE STRENGTHS AND LIMITATIONS OF THE ALGORITHM AND IMPROVE ITS PERFORMANCE.

## Experiment – 9

Aim: Apply the k-Means Clustering algorithm on suitable datasets and comment on the quality of clustering

Theory:

THE K-MEANS CLUSTERING ALGORITHM IS A POPULAR UNSUPERVISED MACHINE LEARNING ALGORITHM USED FOR CLUSTERING TASKS. IT IS A SIMPLE YET EFFECTIVE ALGORITHM THAT PARTITIONS THE GIVEN DATASET INTO K CLUSTERS BASED ON THE SIMILARITY OF THE DATA POINTS. THE ALGORITHM WORKS BY FIRST SELECTING K INITIAL CENTROIDS RANDOMLY FROM THE DATASET. IT THEN ASSIGNS EACH DATA POINT TO THE NEAREST CENTROID AND RECALCULATES THE CENTROID FOR EACH CLUSTER. THIS PROCESS IS REPEATED UNTIL CONVERGENCE, I.E., WHEN THE ASSIGNMENT OF DATA POINTS TO CLUSTERS NO LONGER CHANGES.

TO APPLY THE K-MEANS CLUSTERING ALGORITHM, WE FIRST NEED TO SELECT A SUITABLE DATASET WITH UNLABELLED DATA. THE ALGORITHM WILL AUTOMATICALLY IDENTIFY PATTERNS IN THE DATA AND CLUSTER SIMILAR DATA POINTS TOGETHER. WE CAN THEN EVALUATE THE QUALITY OF CLUSTERING BY USING METRICS SUCH AS THE SILHOUETTE SCORE, CALINSKI-HARABASZ INDEX, OR DAVIES-BOULDIN INDEX.

THE SILHOUETTE SCORE MEASURES THE SIMILARITY OF DATA POINTS WITHIN A CLUSTER COMPARED TO OTHER CLUSTERS. THE CALINSKI-HARABASZ INDEX MEASURES THE RATIO OF THE BETWEEN-CLUSTER DISPERSION TO THE WITHIN-CLUSTER DISPERSION. THE DAVIES-BOULDIN INDEX MEASURES THE AVERAGE SIMILARITY BETWEEN EACH CLUSTER AND ITS MOST SIMILAR CLUSTER.

Code:

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn.datasets import load_iris
from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA

# Load the Iris dataset
iris = load_iris()
X = iris.data
y = iris.target

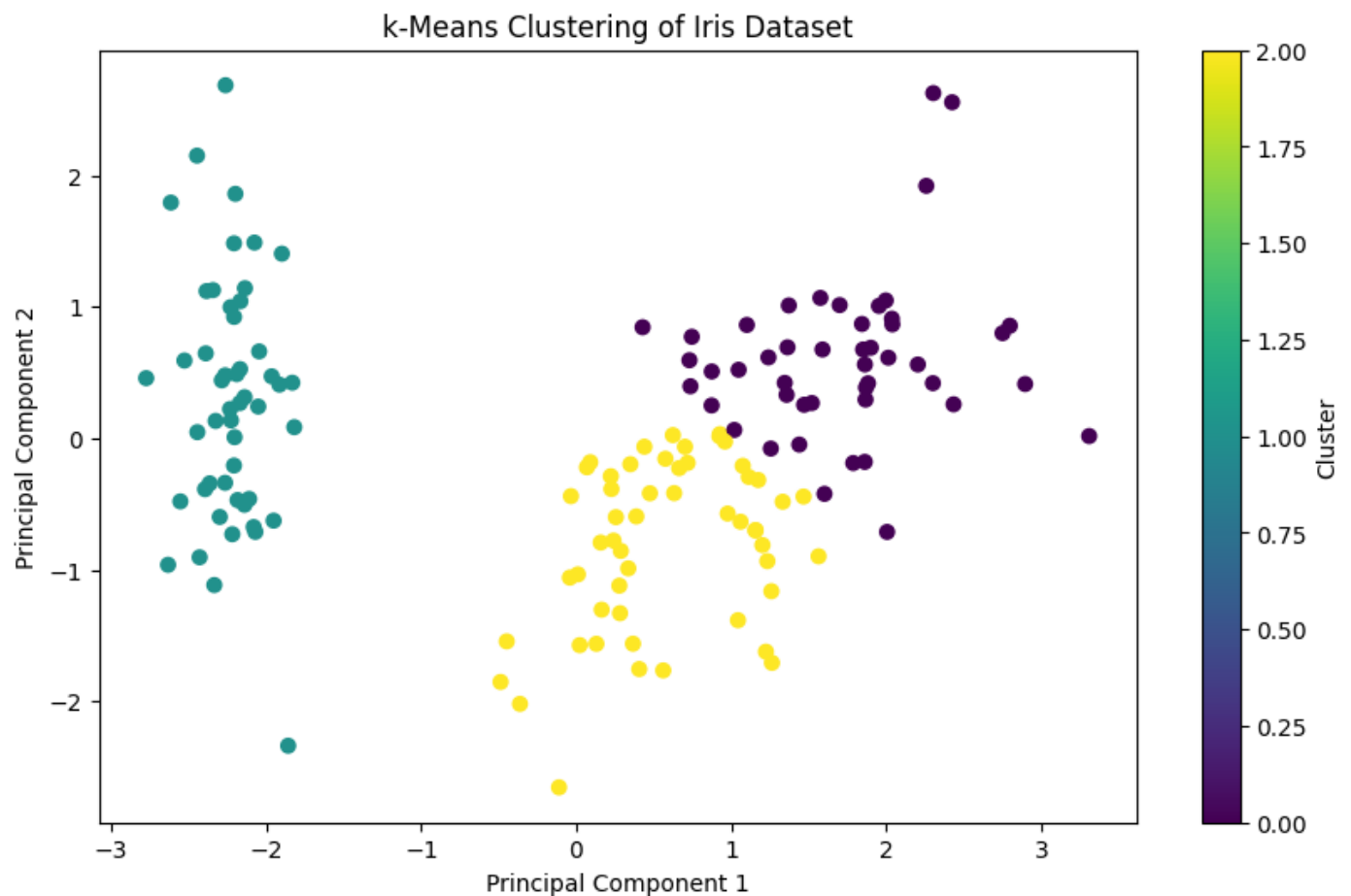
# Standardize the features
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

# Apply PCA for visualization
pca = PCA(n_components=2)
X_pca = pca.fit_transform(X_scaled)

# Apply k-Means clustering
kmeans = KMeans(n_clusters=3, random_state=42)
kmeans.fit(X_scaled)
y_pred = kmeans.predict(X_scaled)

# Visualize the clusters
plt.figure(figsize=(10, 6))
plt.scatter(X_pca[:, 0], X_pca[:, 1], c=y_pred, cmap='viridis')
plt.title('k-Means Clustering of Iris Dataset')
plt.xlabel('Principal Component 1')
plt.ylabel('Principal Component 2')
plt.colorbar(label='Cluster')
plt.show()
```

Output:



Learning:

APPLYING THE K-MEANS CLUSTERING ALGORITHM TO A SUITABLE DATASET CAN PROVIDE INSIGHTS INTO THE STRUCTURE OF THE DATA AND THE SIMILARITY OF THE DATA POINTS. THE QUALITY OF CLUSTERING CAN BE EVALUATED USING VARIOUS METRICS, WHICH CAN HELP IN SELECTING THE OPTIMAL VALUE OF K AND IMPROVING THE PERFORMANCE OF THE ALGORITHM.

ADDITIONALLY, THE EXPERIMENT CAN ALSO PROVIDE INSIGHTS INTO THE IMPORTANCE OF FEATURE SELECTION AND DATA PREPROCESSING FOR THE K-MEANS CLUSTERING ALGORITHM. BY TESTING THE ALGORITHM WITH DIFFERENT DATASETS, ONE CAN IDENTIFY THE STRENGTHS AND LIMITATIONS OF THE ALGORITHM AND IMPROVE ITS PERFORMANCE.

## Experiment – 10

Aim: Write a program to implement ensemble algorithms- AdaBoost and Bagging using the appropriate dataset and evaluate their performance on that dataset

Theory:

ENSEMBLE LEARNING IS A POPULAR MACHINE LEARNING TECHNIQUE USED TO IMPROVE THE ACCURACY OF PREDICTIVE MODELS. IT INVOLVES COMBINING MULTIPLE MODELS TO MAKE MORE ACCURATE PREDICTIONS THAN A SINGLE MODEL. TWO COMMON ENSEMBLE LEARNING ALGORITHMS ARE ADABOOST AND BAGGING.

ADABOOST (ADAPTIVE BOOSTING) IS A BOOSTING ALGORITHM THAT TRAINS MULTIPLE WEAK CLASSIFIERS SEQUENTIALLY. EACH WEAK CLASSIFIER IS TRAINED ON A SUBSET OF THE DATA AND ASSIGNED A WEIGHT. IN EACH SUBSEQUENT ROUND, THE ALGORITHM FOCUSES ON THE MISCLASSIFIED SAMPLES FROM THE PREVIOUS ROUND AND ADJUSTS THE WEIGHTS OF THE SAMPLES. THE FINAL PREDICTION IS MADE BY TAKING A WEIGHTED AVERAGE OF THE PREDICTIONS OF ALL THE WEAK CLASSIFIERS.

BAGGING (BOOTSTRAP AGGREGATING) IS A PARALLEL ENSEMBLE ALGORITHM THAT TRAINS MULTIPLE INSTANCES OF A BASE CLASSIFIER ON RANDOM SUBSETS OF THE DATA. THE FINAL PREDICTION IS MADE BY TAKING A MAJORITY VOTE ON THE PREDICTIONS OF ALL THE BASE CLASSIFIERS.

TO IMPLEMENT ADABOOST AND BAGGING, WE FIRST NEED TO SELECT A SUITABLE DATASET WITH LABELLED DATA. WE CAN THEN TRAIN MULTIPLE INSTANCES OF THE BASE CLASSIFIER USING ADABOOST OR BAGGING AND EVALUATE THEIR PERFORMANCE USING METRICS SUCH AS ACCURACY, PRECISION, RECALL, AND F1 SCORE.

Code:

```
from sklearn.datasets import load_iris
from sklearn.model_selection import train_test_split
from sklearn.ensemble import AdaBoostClassifier, BaggingClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy_score

# Load the Iris dataset
iris = load_iris()
X = iris.data
y = iris.target

# Split the dataset 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)

# Define base classifier
base_classifier = DecisionTreeClassifier(max_depth=1)

# AdaBoost
ada_boost = AdaBoostClassifier(base_estimator=base_classifier, n_estimators=50,
random_state=42)
ada_boost.fit(X_train, y_train)
y_pred_ada = ada_boost.predict(X_test)
accuracy_ada = accuracy_score(y_test, y_pred_ada)
print("Accuracy of AdaBoost:", accuracy_ada)

# Bagging
bagging = BaggingClassifier(base_estimator=base_classifier, n_estimators=50,
random_state=42)
bagging.fit(X_train, y_train)
y_pred_bagging = bagging.predict(X_test)
accuracy_bagging = accuracy_score(y_test, y_pred_bagging)
```

```
print("Accuracy of Bagging:", accuracy_bagging)
```

Output:

```
Accuracy of AdaBoost: 1.0
```

```
Accuracy of Bagging: 1.0
```

Learning:

IMPLEMENTING ENSEMBLE ALGORITHMS SUCH AS ADABOOST AND BAGGING CAN HELP IN IMPROVING THE ACCURACY AND ROBUSTNESS OF PREDICTIVE MODELS. THE PERFORMANCE MEASURES COMPUTED FOR THE MODELS CAN HELP IN ANALYSING THE EFFECTIVENESS OF THE ALGORITHMS IN DIFFERENT SCENARIOS. THE EXPERIMENT CAN ALSO HELP IN UNDERSTANDING THE IMPACT OF DIFFERENT HYPERPARAMETERS ON THE PERFORMANCE OF THE ALGORITHMS.

ADDITIONALLY, THE EXPERIMENT CAN PROVIDE INSIGHTS INTO THE IMPORTANCE OF FEATURE SELECTION AND DATA PREPROCESSING FOR ENSEMBLE LEARNING ALGORITHMS. BY TESTING THE ALGORITHMS WITH DIFFERENT DATASETS, ONE CAN IDENTIFY THE STRENGTHS AND LIMITATIONS OF THE ALGORITHMS AND IMPROVE THEIR PERFORMANCE