# ASSIGNMENT 6

1. Compute the split point for each attribute in the dataset using the following strategies:
   a. Information Gain
   b. Gini Indexes
   c. Gain Ratio

```python
import pandas as pd
import numpy as np
from sklearn.tree import DecisionTreeClassifier
# Load the dataset
df = pd.read_csv("vehicle.csv")



# Preprocessing (if required)
# Check for missing values
print(df.isnull().sum())



# Handle missing values (if any)
# For simplicity, let's drop rows with missing values
df.dropna(inplace=True)



# Separate features (X) and target variable (y)
X = df.drop(columns=['class'])  # Assuming 'class' is the target variable
y = df['class']



# Initialize the decision tree classifier
clf = DecisionTreeClassifier()



# Fit the classifier to the data
clf.fit(X, y)
```

```python
# Function to compute Information Gain
def information_gain(y, y_splits):
    entropy_parent = entropy(y)
    total_instances = len(y)
    weighted_entropy_children = sum((len(y_split) / total_instances) *
entropy(y_split) for y_split in y_splits)
    return entropy_parent - weighted_entropy_children
# Function to compute Gini Index
def gini_index(y, y_splits):
    gini_parent = gini_impurity(y)
    total_instances = len(y)
    weighted_gini_children = sum((len(y_split) / total_instances) *
gini_impurity(y_split) for y_split in y_splits)
    return gini_parent - weighted_gini_children


# Function to compute entropy
def entropy(y):
    classes, counts = np.unique(y, return_counts=True)
    probabilities = counts / len(y)
    return -sum(p * np.log2(p) for p in probabilities)


# Function to compute Gini impurity
def gini_impurity(y):
    classes, counts = np.unique(y, return_counts=True)
    probabilities = counts / len(y)
    return 1 - sum(p**2 for p in probabilities)
# Function to compute split points using Information Gain, Gini Index, and
Gain Ratio
def compute_split_points(X, y):
    split_points = {}
    for col in X.columns:
        # Assuming each attribute is numeric
        values = X[col].unique()
        for value in values:
            # Split the dataset based on the attribute value
            left_indices = X[col] < value
```

```python
            right_indices = ~left_indices
            y_splits = [y[left_indices], y[right_indices]]
            # Compute metrics for split points
            info_gain = information_gain(y, y_splits)
            gini_idx = gini_index(y, y_splits)
            gain_ratio = info_gain / (entropy(X[col]) + 1e-10)  # Add small
value to avoid division by zero

            split_points[(col, value)] = {'Information Gain': info_gain,
'Gini Index': gini_idx, 'Gain Ratio': gain_ratio}

    return split_points


# Compute split points for each attribute
split_points = compute_split_points(X, y)


# Print split points
for key, value in split_points.items():
    print(f"Split Point for {key}: {value}")
```
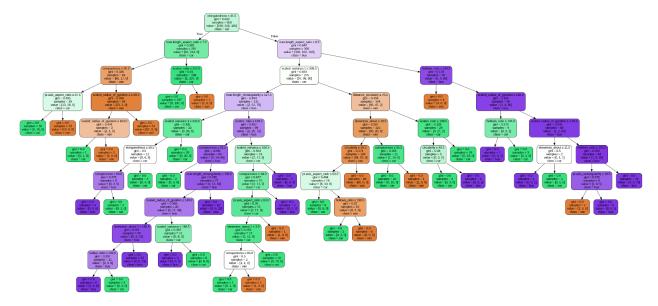
2. Design module for creating the decision tree and its representation in graphical format for the following cases:
      a. Binary Tree (each node split into exactly two branches).
      b. General Tree (each node may split into more than two branches depending on
         count nominal labels corresponding attributes).

```python
import pandas as pd
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, classification_report
from sklearn.tree import export_graphviz
import graphviz


# Load the dataset
df = pd.read_csv("vehicle.csv")
```

```python
# Check for missing values before preprocessing
print("Missing values before preprocessing:")
print(df.isnull().sum())


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


# Check for missing values after preprocessing
print("\nMissing values after preprocessing:")
print(df.isnull().sum())


# Split the dataset into features and target variable
X = df.drop(columns=['class'])
y = df['class']


# 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)


# Initialize the decision tree classifier
clf = DecisionTreeClassifier()


# Fit the classifier to the training data
clf.fit(X_train, y_train)


# Predict the labels of the test set
y_pred = clf.predict(X_test)


# Evaluate the model
accuracy = accuracy_score(y_test, y_pred)
print("\nAccuracy:", accuracy)
```

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



# Visualize the decision tree
dot_data = export_graphviz(clf, out_file=None,
                           feature_names=X.columns,
                           class_names=y.unique(),
                           filled=True, rounded=True,
                           special_characters=True)
graph = graphviz.Source(dot_data)
graph.render("vehicle_decision_tree", format='png', cleanup=True)
```



3. Design module which predicts the class label of unknown and unseen data using tree traversal or any other techniques.

```
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy_score
```

```python
# Load the dataset
df = pd.read_csv('vehicle.csv')
# Drop rows with missing values
df.dropna(inplace=True)
# Split dataset into features and target variable
X = df.drop(columns=['class'])
y = df['class']



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



# Create a decision tree classifier
clf = DecisionTreeClassifier()



# Train the classifier on the training data
clf.fit(X_train, y_train)



# Predict the classes of testing data
y_pred = clf.predict(X_test)



# Evaluate accuracy
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy:", accuracy)
```