| |
|---|
| Experiment No. 4 |
| Apply Random Forest Algorithm on Adult Census Income Dataset and analyze the performance of the model |
| Date of Performance: |
| Date of Submission: |

**Aim:** Apply Random Forest Algorithm on Adult Census Income Dataset and analyze the performance of the model.

**Objective:** Able to perform various feature engineering tasks, apply Random Forest Algorithm on the given dataset and maximize the accuracy, Precision, Recall, F1 score.
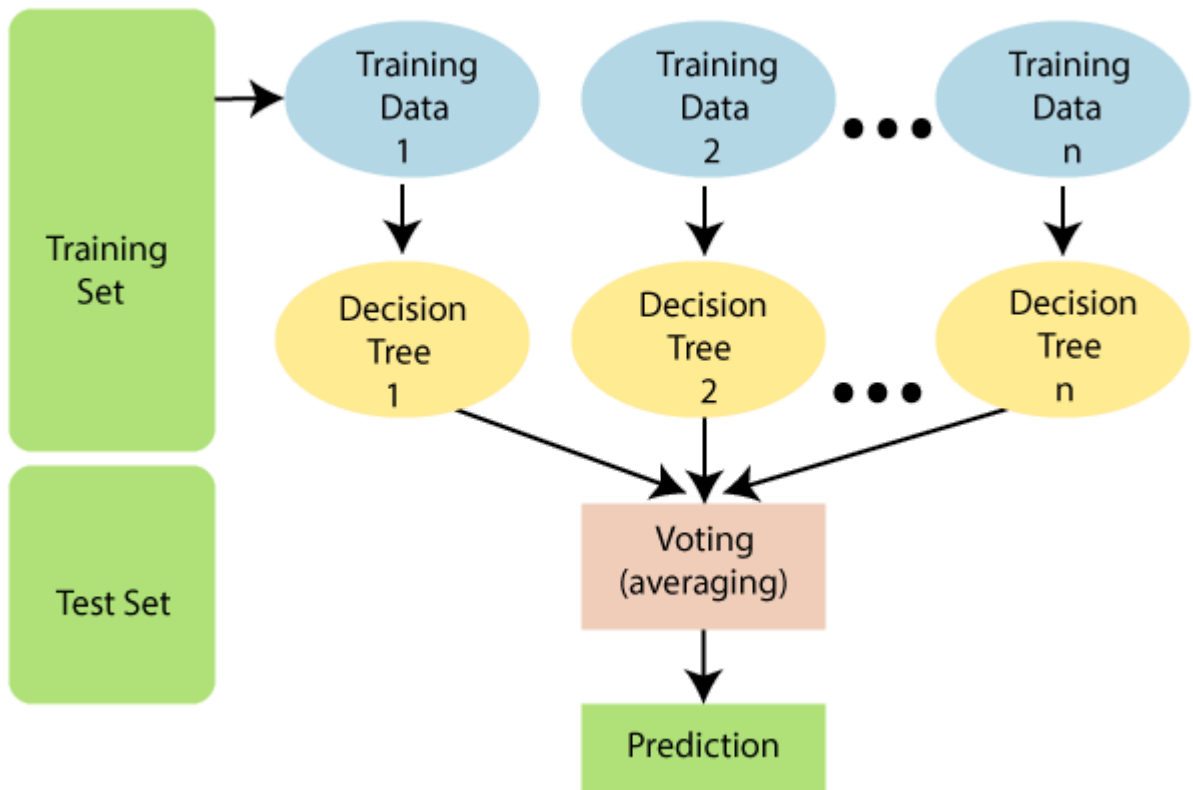
**Theory:**

Random Forest is a popular machine learning algorithm that belongs to the supervised learning technique. It can be used for both Classification and Regression problems in ML. It is based on the concept of ensemble learning, which is a process of combining multiple classifiers to solve a complex problem and to improve the performance of the model.

As the name suggests, "Random Forest is a classifier that contains a number of decision trees on various subsets of the given dataset and takes the average to improve the predictive accuracy of that dataset." Instead of relying on one decision tree, the random forest takes the prediction from each tree and based on the majority votes of predictions, and it predicts the final output.

The greater number of trees in the forest leads to higher accuracy and prevents the problem of overfitting.

The below diagram explains the working of the Random Forest algorithm:

**Dataset:**

Predict whether income exceeds $50K/yr based on census data. Also known as "Adult" dataset.

Attribute Information:

Listing of attributes:

>50K, <=50K.

age: continuous.

workclass: Private, Self-emp-not-inc, Self-emp-inc, Federal-gov, Local-gov, State-gov, Without-pay, Never-worked.

fnlwgt: continuous.

education: Bachelors, Some-college, 11th, HS-grad, Prof-school, Assoc-acdm, Assoc-voc, 9th, 7th-8th, 12th, Masters, 1st-4th, 10th, Doctorate, 5th-6th, Preschool.

education-num: continuous.

marital-status: Married-civ-spouse, Divorced, Never-married, Separated, Widowed, Married-spouse-absent, Married-AF-spouse.

occupation: Tech-support, Craft-repair, Other-service, Sales, Exec-managerial, Prof-specialty, Handlers-cleaners, Machine-op-inspct, Adm-clerical, Farming-fishing, Transport-moving, Priv-house-serv, Protective-serv, Armed-Forces.

relationship: Wife, Own-child, Husband, Not-in-family, Other-relative, Unmarried.

race: White, Asian-Pac-Islander, Amer-Indian-Eskimo, Other, Black.

sex: Female, Male.

capital-gain: continuous.

capital-loss: continuous.

hours-per-week: continuous.

native-country: United-States, Cambodia, England, Puerto-Rico, Canada, Germany, Outlying-US(Guam-USVI-etc), India, Japan, Greece, South, China, Cuba, Iran, Honduras, Philippines, Italy, Poland, Jamaica, Vietnam, Mexico, Portugal, Ireland, France, Dominican-Republic, Laos, Ecuador, Taiwan, Haiti, Columbia, Hungary, Guatemala, Nicaragua, Scotland, Thailand, Yugoslavia, El-Salvador, Trinadad&Tobago, Peru, Hong, Holand-Netherlands.

**Code:**

```
import pandas as pd
import numpy as np
from sklearn.ensemble import RandomForestClassifier
```

CSL701: Machine Learning Lab

```
df = pd.read_csv("/content/adult.csv")
print(df.isnull().sum())

print(df.info())

from sklearn.preprocessing import LabelEncoder
le = LabelEncoder()
categorical_columns = df.select_dtypes(include=['object']).columns
for column in categorical_columns:
    df[column] = le.fit_transform(df[column])

from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report, accuracy_score, confusion_matrix

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

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

rf_classifier = RandomForestClassifier(n_estimators=100, random_state=42)

rf_classifier.fit(X_train, y_train)

y_pred = rf_classifier.predict(X_test)

accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy: {accuracy * 100:.2f}%")

print("\nClassification Report:\n", classification_report(y_test, y_pred))
print("\nConfusion Matrix:\n", confusion_matrix(y_test, y_pred))
```

**Conclusion:**

The evaluation of the Random Forest Classifier on the dataset yielded the following insights based on accuracy, confusion matrix, precision, recall, and F1 score:

**Accuracy (86.54%):**
The model has a high overall accuracy, indicating that it correctly classifies a substantial portion of the instances in the test set. However, accuracy alone can be misleading, especially in the context of class imbalance.

**Confusion Matrix:**
The confusion matrix reveals the distribution of correct and incorrect predictions:

**True Negatives (10,493):** The model successfully identifies a large number of class 0 instances.

CSL701: Machine Learning Lab

**True Positives (2,187)**: While the model identifies a considerable number of class 1 instances, the number of False Negatives (1,233) suggests that many actual class 1 instances are misclassified as class 0.

**Precision**:

For class 0, the precision of **0.89** indicates a high level of confidence in the model's positive predictions. However, for class 1, the precision of **0.75** indicates that while the model predicts class 1 correctly 75% of the time, there are still significant false positives, suggesting that the model is less reliable in identifying this class.

**Recall**:

The recall for class 0 is high at **0.93**, demonstrating the model's effectiveness in capturing actual instances of this class. Conversely, the recall for class 1 is lower at **0.64**, indicating that the model fails to identify a substantial portion of actual class 1 instances, pointing to a potential weakness in detecting this minority class.

**F1 Score**:

The F1 score for class 0 (**0.91**) reflects a strong balance between precision and recall, signifying robust performance. However, the F1 score for class 1 (**0.69**) indicates challenges in achieving a similar balance, revealing that improvements are necessary for better classification of this class.