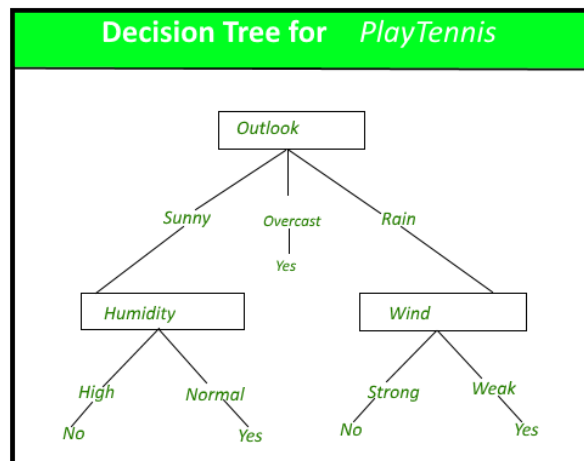| |
|---|
| Experiment No. 3 |
| Apply Decision Tree Algorithm on Adult Census Income Dataset and analyze the performance of the model |
| Date of Performance:7/8/23 |
| Date of Submission: 14/8/23 |

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

**Objective:** To perform various feature engineering tasks, apply Decision Tree Algorithm on the given dataset and maximize the accuracy, Precision, Recall, F1 score. Improve the performance by performing different data engineering and feature engineering tasks.

**Theory:**

Decision Tree is the most powerful and popular tool for classification and prediction. A Decision tree is a flowchart-like tree structure, where each internal node denotes a test on an attribute, each branch represents an outcome of the test, and each leaf node (terminal node) holds a class label.



**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, Trinidad Tobago, Peru, Hong, Holand-Netherlands.

**Code:**

```python
# Import libraries
import os
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline

# To ignore warning messages
import warnings
warnings.filterwarnings('ignore')


# Adult dataset path
adult_dataset_path = "/content/adult.csv"

# Function for loading adult dataset
def load_adult_data(adult_path=adult_dataset_path):
    csv_path = os.path.join(adult_path)
    return pd.read_csv(csv_path)


# Calling load adult function and assigning to a new variable df
df = load_adult_data()
# load top 3 rows values from adult dataset
df.head(3)
```

| | age | workclass | fnlwgt | education | education.num | marital.status | occupation | relationship | race | |
|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 90 | ? | 77053 | HS-grad | 9 | Widowed | ? | Not-in-family | White | Fe |
| **1** | 82 | Private | 132870 | HS-grad | 9 | Widowed | Exec-managerial | Not-in-family | White | Fe |
| **2** | 66 | ? | 186061 | Some- | 10 | Widowed | ? | Unmarried | Black | Fe |

```python
print ("Rows      : " ,df.shape[0])
print ("Columns   : " ,df.shape[1])
print ("\nFeatures : \n" ,df.columns.tolist())
print ("\nMissing values :  ", df.isnull().sum().values.sum())
print ("\nUnique values :  \n",df.nunique())
```

```
Rows      : 32561
Columns   : 15

Features :
 ['age', 'workclass', 'fnlwgt', 'education', 'education.num', 'marital.status', 'occupation', 'relationship', 'race', 'sex', 'capit

Missing values :   0

Unique values :
 age               73
workclass          9
fnlwgt         21648
education         16
education.num     16
marital.status     7
occupation        15
relationship       6
race               5
sex                2
capital.gain     119
capital.loss      92
hours.per.week    94
native.country    42
income             2
dtype: int64
```

```python
# Let's understand the type of values present in each column of our adult dataframe 'df'.
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 32561 entries, 0 to 32560
Data columns (total 15 columns):
 #   Column          Non-Null Count  Dtype
---  ------          --------------  -----
 0   age             32561 non-null  int64
 1   workclass       32561 non-null  object
 2   fnlwgt          32561 non-null  int64
 3   education       32561 non-null  object
 4   education.num   32561 non-null  int64
 5   marital.status  32561 non-null  object
 6   occupation      32561 non-null  object
 7   relationship    32561 non-null  object
```

```
        8   race            32561 non-null  object
        9   sex             32561 non-null  object
        10  capital.gain    32561 non-null  int64
        11  capital.loss    32561 non-null  int64
        12  hours.per.week  32561 non-null  int64
        13  native.country  32561 non-null  object
        14  income          32561 non-null  object
       dtypes: int64(6), object(9)
       memory usage: 3.7+ MB
```

```
# Numerical feature of summary/description
df.describe()
```

|       | age | fnlwgt | education.num | capital.gain | capital.loss | hours.per.week |
|-------|-----|--------|---------------|--------------|--------------|----------------|
| count | 32561.000000 | 3.256100e+04 | 32561.000000 | 32561.000000 | 32561.000000 | 32561.000000 |
| mean | 38.581647 | 1.897784e+05 | 10.080679 | 1077.648844 | 87.303830 | 40.437456 |
| std | 13.640433 | 1.055500e+05 | 2.572720 | 7385.292085 | 402.960219 | 12.347429 |
| min | 17.000000 | 1.228500e+04 | 1.000000 | 0.000000 | 0.000000 | 1.000000 |
| 25% | 28.000000 | 1.178270e+05 | 9.000000 | 0.000000 | 0.000000 | 40.000000 |
| 50% | 37.000000 | 1.783560e+05 | 10.000000 | 0.000000 | 0.000000 | 40.000000 |
| 75% | 48.000000 | 2.370510e+05 | 12.000000 | 0.000000 | 0.000000 | 45.000000 |
| max | 90.000000 | 1.484705e+06 | 16.000000 | 99999.000000 | 4356.000000 | 99.000000 |

```
# pull top 5 row values to understand the data and how it's look like
df.head()
```

|   | age | workclass | fnlwgt | education | education.num | marital.status | occupation | relationship | race | |
|---|-----|-----------|--------|-----------|---------------|----------------|------------|--------------|------|---|
| 0 | 90 | ? | 77053 | HS-grad | 9 | Widowed | ? | Not-in-family | White | Fe |
| 1 | 82 | Private | 132870 | HS-grad | 9 | Widowed | Exec-managerial | Not-in-family | White | Fe |
| 2 | 66 | ? | 186061 | Some-college | 10 | Widowed | ? | Unmarried | Black | Fe |
| 3 | 54 | Private | 140359 | 7th-8th | 4 | Divorced | Machine-op-inspct | Unmarried | White | Fe |

```
# checking "?" total values present in particular 'workclass' feature
df_check_missing_workclass = (df['workclass']=='?').sum()
df_check_missing_workclass
```

```
    1836
```

```
# checking "?" total values present in particular 'occupation' feature
df_check_missing_occupation = (df['occupation']=='?').sum()
df_check_missing_occupation
```

```
    1843
```

```
# checking "?" values, how many are there in the whole dataset
df_missing = (df=='?').sum()
df_missing
```

```
    age                0
    workclass          1836
    fnlwgt             0
    education          0
    education.num      0
    marital.status     0
    occupation         1843
    relationship       0
    race               0
    sex                0
    capital.gain       0
    capital.loss       0
    hours.per.week     0
    native.country     583
    income             0
    dtype: int64
```

```
percent_missing = (df=='?').sum() * 100/len(df)
percent_missing
```

```
    age                0.000000
    workclass          5.638647
```

```
fnlwgt           0.000000
education        0.000000
education.num    0.000000
marital.status   0.000000
occupation       5.660146
relationship     0.000000
race             0.000000
sex              0.000000
capital.gain     0.000000
capital.loss     0.000000
hours.per.week   0.000000
native.country   1.790486
income           0.000000
dtype: float64
```

```python
# Let's find total number of rows which doesn't contain any missing value as '?'
df.apply(lambda x: x !='?',axis=1).sum()
```

```
age              32561
workclass        30725
fnlwgt           32561
education        32561
education.num    32561
marital.status   32561
occupation       30718
relationship     32561
race             32561
sex              32561
capital.gain     32561
capital.loss     32561
hours.per.week   32561
native.country   31978
income           32561
dtype: int64
```

```python
# dropping the rows having missing values in workclass
df = df[df['workclass'] !='?']
df.head()
```

|   | age | workclass | fnlwgt | education | education.num | marital.status | occupation | relationship | race | |
|---|-----|-----------|--------|-----------|---------------|----------------|------------|--------------|------|---|
| 1 | 82 | Private | 132870 | HS-grad | 9 | Widowed | Exec-managerial | Not-in-family | White | Fe |
| 3 | 54 | Private | 140359 | 7th-8th | 4 | Divorced | Machine-op-inspct | Unmarried | White | Fe |
| 4 | 41 | Private | 264663 | Some-college | 10 | Separated | Prof-specialty | Own-child | White | Fe |
| 5 | 34 | Private | 216864 | HS-grad | 9 | Divorced | Other- | Unmarried | White | F |

```python
# select all categorical variables
df_categorical = df.select_dtypes(include=['object'])

# checking whether any other column contains '?' value
df_categorical.apply(lambda x: x=='?',axis=1).sum()
```

```
workclass        0
education        0
marital.status   0
occupation       7
relationship     0
race             0
sex              0
native.country   556
income           0
dtype: int64
```

```python
# dropping the "?"s from occupation and native.country
df = df[df['occupation'] !='?']
df = df[df['native.country'] !='?']
```

```python
# check the dataset whether cleaned or not?
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 30162 entries, 1 to 32560
Data columns (total 15 columns):
 #   Column          Non-Null Count  Dtype
---  ------          --------------  -----
 0   age             30162 non-null  int64
 1   workclass       30162 non-null  object
 2   fnlwgt          30162 non-null  int64
 3   education       30162 non-null  object
```

```
 4    education.num    30162 non-null   int64
 5    marital.status   30162 non-null   object
 6    occupation       30162 non-null   object
 7    relationship     30162 non-null   object
 8    race             30162 non-null   object
 9    sex              30162 non-null   object
 10   capital.gain     30162 non-null   int64
 11   capital.loss     30162 non-null   int64
 12   hours.per.week   30162 non-null   int64
 13   native.country   30162 non-null   object
 14   income           30162 non-null   object
dtypes: int64(6), object(9)
memory usage: 3.7+ MB
```

```python
from sklearn import preprocessing

# encode categorical variables using label Encoder

# select all categorical variables
df_categorical = df.select_dtypes(include=['object'])
df_categorical.head()
```

|   | workclass | education | marital.status | occupation | relationship | race | sex | native.country | income |
|---|-----------|-----------|----------------|------------|--------------|------|-----|----------------|--------|
| 1 | Private | HS-grad | Widowed | Exec-managerial | Not-in-family | White | Female | United-States | <=50K |
| 3 | Private | 7th-8th | Divorced | Machine-op-inspct | Unmarried | White | Female | United-States | <=50K |
| 4 | Private | Some-college | Separated | Prof-specialty | Own-child | White | Female | United-States | <=50K |
| 5 | Private | HS-grad | Divorced | Other- | Unmarried | White | Female | United-States | <=50K |

```python
# apply label encoder to df_categorical
le = preprocessing.LabelEncoder()
df_categorical = df_categorical.apply(le.fit_transform)
df_categorical.head()
```

|   | workclass | education | marital.status | occupation | relationship | race | sex | native.country | income |
|---|-----------|-----------|----------------|------------|--------------|------|-----|----------------|--------|
| 1 | 2 | 11 | 6 | 3 | 1 | 4 | 0 | 38 | 0 |
| 3 | 2 | 5 | 0 | 6 | 4 | 4 | 0 | 38 | 0 |
| 4 | 2 | 15 | 5 | 9 | 3 | 4 | 0 | 38 | 0 |
| 5 | 2 | 11 | 0 | 7 | 4 | 4 | 0 | 38 | 0 |
| 6 | 2 | 0 | 5 | 0 | 4 | 4 | 1 | 38 | 0 |

```python
# Next, Concatenate df_categorical dataframe with original df (dataframe)

# first, Drop earlier duplicate columns which had categorical values
df = df.drop(df_categorical.columns,axis=1)
df = pd.concat([df,df_categorical],axis=1)
df.head()
```

|   | age | fnlwgt | education.num | capital.gain | capital.loss | hours.per.week | workclass | education | marita |
|---|-----|--------|---------------|--------------|--------------|----------------|-----------|-----------|--------|
| 1 | 82 | 132870 | 9 | 0 | 4356 | 18 | 2 | 11 | |
| 3 | 54 | 140359 | 4 | 0 | 3900 | 40 | 2 | 5 | |
| 4 | 41 | 264663 | 10 | 0 | 3900 | 40 | 2 | 15 | |
| 5 | 34 | 216864 | 9 | 0 | 3770 | 45 | 2 | 11 | |
| 6 | 38 | 150601 | 6 | 0 | 3770 | 40 | 2 | 0 | |

```python
# look at column type
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 30162 entries, 1 to 32560
Data columns (total 15 columns):
 #    Column           Non-Null Count   Dtype
---   ------           --------------   -----
 0    age              30162 non-null   int64
 1    fnlwgt           30162 non-null   int64
 2    education.num    30162 non-null   int64
 3    capital.gain     30162 non-null   int64
 4    capital.loss     30162 non-null   int64
 5    hours.per.week   30162 non-null   int64
 6    workclass        30162 non-null   int64
```

```
     7   education        30162 non-null  int64
     8   marital.status   30162 non-null  int64
     9   occupation       30162 non-null  int64
     10  relationship     30162 non-null  int64
     11  race             30162 non-null  int64
     12  sex              30162 non-null  int64
     13  native.country   30162 non-null  int64
     14  income           30162 non-null  int64
    dtypes: int64(15)
    memory usage: 3.7 MB
```

```python
# convert target variable income to categorical
df['income'] = df['income'].astype('category')
```

```python
# check df info again whether everything is in right format or not
df.info()
```

```
    <class 'pandas.core.frame.DataFrame'>
    Int64Index: 30162 entries, 1 to 32560
    Data columns (total 15 columns):
     #   Column           Non-Null Count  Dtype
    ---  ------           --------------  -----
     0   age              30162 non-null  int64
     1   fnlwgt           30162 non-null  int64
     2   education.num    30162 non-null  int64
     3   capital.gain     30162 non-null  int64
     4   capital.loss     30162 non-null  int64
     5   hours.per.week   30162 non-null  int64
     6   workclass        30162 non-null  int64
     7   education        30162 non-null  int64
     8   marital.status   30162 non-null  int64
     9   occupation       30162 non-null  int64
     10  relationship     30162 non-null  int64
     11  race             30162 non-null  int64
     12  sex              30162 non-null  int64
     13  native.country   30162 non-null  int64
     14  income           30162 non-null  category
    dtypes: category(1), int64(14)
    memory usage: 3.5 MB
```

```python
# Importing train_test_split
from sklearn.model_selection import train_test_split
```

```python
# Putting independent variables/features to X
X = df.drop('income',axis=1)
```

```python
# Putting response/dependent variable/feature to y
y = df['income']
```

```python
X.head(3)
```

| | age | fnlwgt | education.num | capital.gain | capital.loss | hours.per.week | workclass | education | marita |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 82 | 132870 | 9 | 0 | 4356 | 18 | 2 | 11 | |
| 3 | 54 | 140359 | 4 | 0 | 3900 | 40 | 2 | 5 | |
| 4 | 41 | 264663 | 10 | 0 | 3900 | 40 | 2 | 15 | |

```python
y.head(3)
```

```
    1    0
    3    0
    4    0
    Name: income, dtype: category
    Categories (2, int64): [0, 1]
```

```python
# Splitting the data into train and test
X_train,X_test,y_train,y_test = train_test_split(X,y,test_size=0.30,random_state=99)
```

```python
X_train.head()
```

| | age | fnlwgt | education.num | capital.gain | capital.loss | hours.per.week | workclass | education | ma |
|---|---|---|---|---|---|---|---|---|---|
| **24351** | 42 | 289636 | 9 | 0 | 0 | 46 | 2 | 11 | |
| **15626** | 37 | 52465 | 9 | 0 | 0 | 40 | 1 | 11 | |

```
# Importing decision tree classifier from sklearn library
from sklearn.tree import DecisionTreeClassifier

# Fitting the decision tree with default hyperparameters, apart from
# max_depth which is 5 so that we can plot and read the tree.
dt_default = DecisionTreeClassifier(max_depth=5)
dt_default.fit(X_train,y_train)
```

```
    ▼        DecisionTreeClassifier
    DecisionTreeClassifier(max_depth=5)
```

```
# Let's check the evaluation metrics of our default model

# Importing classification report and confusion matrix from sklearn metrics
from sklearn.metrics import classification_report,confusion_matrix,accuracy_score

# making predictions
y_pred_default = dt_default.predict(X_test)

# Printing classifier report after prediction
print(classification_report(y_test,y_pred_default))
```

```
              precision    recall  f1-score   support

           0       0.86      0.95      0.91      6867
           1       0.78      0.52      0.63      2182

    accuracy                           0.85      9049
   macro avg       0.82      0.74      0.77      9049
weighted avg       0.84      0.85      0.84      9049
```

```
# Printing confusion matrix and accuracy
print(confusion_matrix(y_test,y_pred_default))
print(accuracy_score(y_test,y_pred_default))
```

```
    [[6553  314]
     [1039 1143]]
    0.8504807161012267
```

```
pip install pydotplus
```

```
    Requirement already satisfied: pydotplus in /usr/local/lib/python3.10/dist-packages (2.0.2)
    Requirement already satisfied: pyparsing>=2.0.1 in /usr/local/lib/python3.10/dist-packages (from pydotplus) (3.1.1)
```

```
# Importing required packages for visualization
from IPython.display import Image
from sklearn.externals.six import StringIO
from sklearn.tree import export_graphviz
import pydotplus,graphviz

# Putting features
features = list(df.columns[1:])
features
```
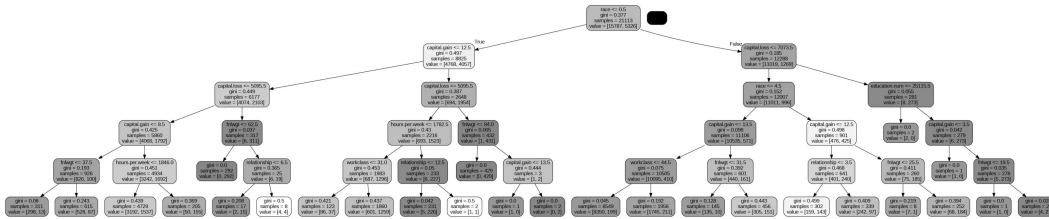
```
    ['fnlwgt',
     'education.num',
     'capital.gain',
     'capital.loss',
     'hours.per.week',
     'workclass',
     'education',
     'marital.status',
     'occupation',
     'relationship',
     'race',
     'sex',
     'native.country',
     'income']
```

```
import six
import sys
sys.modules['sklearn.externals.six'] = six
```

```
# plotting tree with max_depth=3
dot_data = StringIO()
export_graphviz(dt_default, out_file=dot_data,
                feature_names=features, filled=True,rounded=True)

graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
Image(graph.create_png())
```
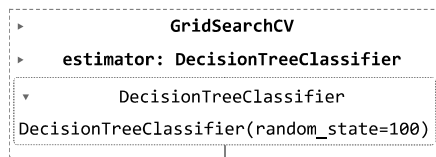


```
# GridSearchCV to find optimal max_depth
from sklearn.model_selection import KFold
from sklearn.model_selection import GridSearchCV


# specify number of folds for k-fold CV
n_folds = 5

# parameters to build the model on
parameters = {'max_depth': range(1, 40)}

# instantiate the model
dtree = DecisionTreeClassifier(criterion = "gini",
                               random_state = 100)

# fit tree on training data
tree = GridSearchCV(dtree, parameters,
                    cv=n_folds,
                    scoring="accuracy")
tree.fit(X_train, y_train)
```

```
▸          GridSearchCV
▸  estimator: DecisionTreeClassifier
▾      DecisionTreeClassifier
DecisionTreeClassifier(random_state=100)
```

```
# scores of GridSearch CV
scores = tree.cv_results_
pd.DataFrame(scores).head()
```

| | mean_fit_time | std_fit_time | mean_score_time | std_score_time | param_max_depth | params | split0_t |
|---|---|---|---|---|---|---|---|
| **0** | 0.014924 | 0.004959 | 0.004501 | 0.002179 | 1 | {'max_depth': 1} | |
| **1** | 0.019653 | 0.000594 | 0.003418 | 0.000383 | 2 | {'max_depth': 2} | |
| **2** | 0.025445 | 0.000331 | 0.003681 | 0.000198 | 3 | {'max_depth': 3} | |
| **3** | 0.031828 | 0.000514 | 0.003868 | 0.000580 | 4 | {'max_depth': 4} | |
| **4** | 0.043398 | 0.007946 | 0.003806 | 0.001610 | 5 | {'max_depth': 5} | |

```
"""
# plotting accuracies with max_depth
plt.figure()
plt.plot(scores["param_max_depth"],
         scores["mean_train_score"],
         label="training accuracy")
plt.plot(scores["param_max_depth"],
```

```
                    scores["mean_test_score"],
                    label="test accuracy")
plt.xlabel("max_depth")
plt.ylabel("Accuracy")
plt.legend()
plt.show()
"""
```

> '\n# plotting accuracies with max_depth\nplt.figure()\nplt.plot(scores["param_max_depth"], \n     scores["mean_train_score"], \n          label="training accuracy")\nplt.plot(scores["param_max_depth"],
> \n          scores["mean test score"], \n          label="test accuracy")\nplt.xlabel("max depth")\nplt.

```
# GridSearchCV to find optimal max_depth
from sklearn.model_selection import KFold
from sklearn.model_selection import GridSearchCV


# specify number of folds for k-fold CV
n_folds = 5

# parameters to build the model on
parameters = {'min_samples_leaf': range(5, 200, 20)}

# instantiate the model
dtree = DecisionTreeClassifier(criterion = "gini",
                               random_state = 100)

# fit tree on training data
tree = GridSearchCV(dtree, parameters,
                    cv=n_folds,
                    scoring="accuracy")
tree.fit(X_train, y_train)
```
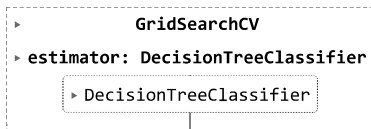
```
  ▸              GridSearchCV
  ▸ estimator: DecisionTreeClassifier
          ▸ DecisionTreeClassifier
```

```
# scores of GridSearch CV
scores = tree.cv_results_
pd.DataFrame(scores).head()
```

|   | mean_fit_time | std_fit_time | mean_score_time | std_score_time | param_min_samples_leaf | par |
|---|---|---|---|---|---|---|
| **0** | 0.071311 | 0.002583 | 0.002760 | 0.000105 | 5 | {'min_samples_le |
| **1** | 0.058933 | 0.001238 | 0.002700 | 0.000055 | 25 | {'min_samples_le |
| **2** | 0.054030 | 0.001434 | 0.002561 | 0.000125 | 45 | {'min_samples_le |
| **3** | 0.050727 | 0.001386 | 0.002415 | 0.000073 | 65 | {'min_samples_le |
| **4** | 0.050003 | 0.003110 | 0.002739 | 0.000548 | 85 | {'min_samples_le |

```
"""
# plotting accuracies with min_samples_leaf
plt.figure()
plt.plot(scores["param_min_samples_leaf"],
         scores["mean_train_score"],
         label="training accuracy")
plt.plot(scores["param_min_samples_leaf"],
         scores["mean_test_score"],
         label="test accuracy")
plt.xlabel("min_samples_leaf")
plt.ylabel("Accuracy")
plt.legend()
plt.show()
"""
```

> '\n# plotting accuracies with min_samples_leaf\nplt.figure()\nplt.plot(scores["param_min_samples_lea
> f"], \n     scores["mean_train_score"], \n          label="training accuracy")\nplt.plot(scores["pa
> ram_min_samples_leaf"], \n     scores["mean_test_score"], \n          label="test accuracy")\nplt.x
> label("min samples leaf")\nplt.ylabel("Accuracy")\nplt.legend()\nplt.show()\n'

```
# GridSearchCV to find optimal min_samples_split
from sklearn.model_selection import KFold
```

```python
from sklearn.model_selection import GridSearchCV


# specify number of folds for k-fold CV
n_folds = 5

# parameters to build the model on
parameters = {'min_samples_split': range(5, 200, 20)}

# instantiate the model
dtree = DecisionTreeClassifier(criterion = "gini",
                                random_state = 100)

# fit tree on training data
tree = GridSearchCV(dtree, parameters,
                    cv=n_folds,
                   scoring="accuracy")
tree.fit(X_train, y_train)
```
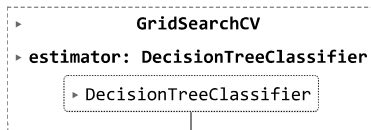
```
      ▸          GridSearchCV
      ▸ estimator: DecisionTreeClassifier
           ▸ DecisionTreeClassifier
```

```python
# scores of GridSearch CV
scores = tree.cv_results_
pd.DataFrame(scores).head()
```

| | mean_fit_time | std_fit_time | mean_score_time | std_score_time | param_min_samples_split | params | split0_test_score | spli |
|---|---|---|---|---|---|---|---|---|
| 0 | 0.075732 | 0.001588 | 0.003151 | 0.000410 | 5 | {'min_samples_split': 5} | 0.811982 | |
| 1 | 0.072139 | 0.002054 | 0.002695 | 0.000054 | 25 | {'min_samples_split': 25} | 0.825006 | |
| 2 | 0.069699 | 0.001631 | 0.002750 | 0.000079 | 45 | {'min_samples_split': 45} | 0.835188 | |
| 3 | 0.067109 | 0.001148 | 0.002617 | 0.000048 | 65 | {'min_samples_split': 65} | 0.839451 | |
| 4 | 0.066176 | 0.002326 | 0.002656 | 0.000063 | 85 | {'min_samples_split': 85} | 0.846081 | |

```python
"""
# plotting accuracies with min_samples_leaf
plt.figure()
plt.plot(scores["param_min_samples_split"],
         scores["mean_train_score"],
         label="training accuracy")
plt.plot(scores["param_min_samples_split"],
         scores["mean_test_score"],
         label="test accuracy")
plt.xlabel("min_samples_split")
plt.ylabel("Accuracy")
plt.legend()
plt.show()
"""
```

```
    '\n# plotting accuracies with min_samples_leaf\nplt.figure()\nplt.plot(scores["param_min_samples_split"], \n        scores["mean_
    train_score"], \n        label="training accuracy")\nplt.plot(scores["param_min_samples_split"], \n        scores["mean_test_sco
    re"], \n        label="test accuracy")\nplt.xlabel("min_samples_split")\nplt.ylabel("Accuracy")\nplt.legend()\nplt.show()\n'
```

```python
# Create the parameter grid
param_grid = {
    'max_depth': range(5, 15, 5),
    'min_samples_leaf': range(50, 150, 50),
    'min_samples_split': range(50, 150, 50),
    'criterion': ["entropy", "gini"]
}

n_folds = 5

# Instantiate the grid search model
dtree = DecisionTreeClassifier()
grid_search = GridSearchCV(estimator = dtree, param_grid = param_grid,
                           cv = n_folds, verbose = 1)
```
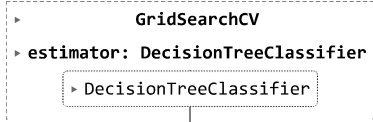
```
# Fit the grid search to the data
grid_search.fit(X_train,y_train)
```

Fitting 5 folds for each of 16 candidates, totalling 80 fits

```
▸        GridSearchCV

▸ estimator: DecisionTreeClassifier

      ▸ DecisionTreeClassifier
```

```
# cv results
cv_results = pd.DataFrame(grid_search.cv_results_)
cv_results
```

| | | | | | | |
|---|---|---|---|---|---|---|
| **3** | 0.032360 | 0.001393 | 0.002821 | 0.000553 | entropy | 5 |
| **4** | 0.051469 | 0.001100 | 0.002973 | 0.000313 | entropy | 10 |
| **5** | 0.050876 | 0.001384 | 0.002635 | 0.000130 | entropy | 10 |

```python
# printing the optimal accuracy score and hyperparameters
print("best accuracy", grid_search.best_score_)
print(grid_search.best_estimator_)
```

```
best accuracy 0.8510400232064759
DecisionTreeClassifier(max_depth=10, min_samples_leaf=50, min_samples_split=50)
```

```python
# model with optimal hyperparameters
clf_gini = DecisionTreeClassifier(criterion = "gini",
                                  random_state = 100,
                                  max_depth=10,
                                  min_samples_leaf=50,
                                  min_samples_split=50)
clf_gini.fit(X_train, y_train)
```
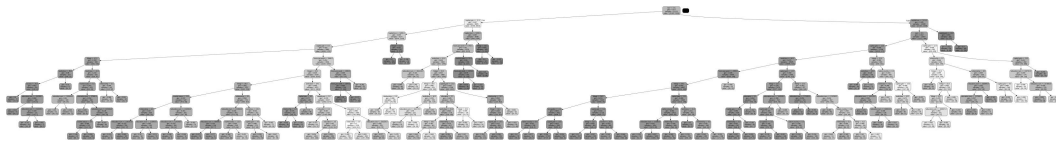
```
              ▾                    DecisionTreeClassifier
    DecisionTreeClassifier(max_depth=10, min_samples_leaf=50, min_samples_split=50,
                           random_state=100)
```

```python
# accuracy score
clf_gini.score(X_test,y_test)
```

```
0.850922753895458
```

```python
# plotting the tree
dot_data = StringIO()
export_graphviz(clf_gini, out_file=dot_data,feature_names=features,filled=True,rounded=True)

graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
Image(graph.create_png())
```



```python
# tree with max_depth = 3
clf_gini = DecisionTreeClassifier(criterion = "gini",
                                  random_state = 100,
                                  max_depth=3,
                                  min_samples_leaf=50,
                                  min_samples_split=50)
clf_gini.fit(X_train, y_train)

# score
print(clf_gini.score(X_test,y_test))
```
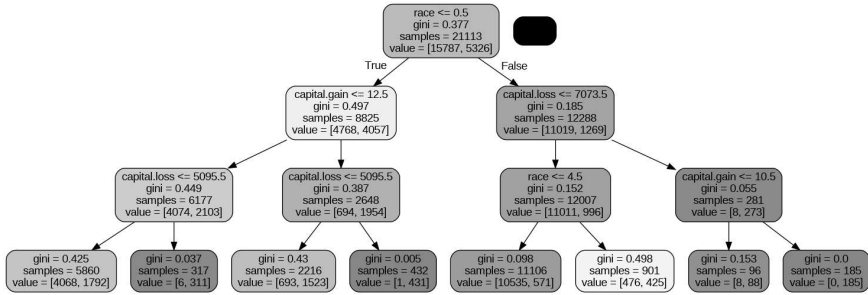
```
0.8393192617968837
```

```python
# plotting tree with max_depth=3
dot_data = StringIO()
export_graphviz(clf_gini, out_file=dot_data,feature_names=features,filled=True,rounded=True)

graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
Image(graph.create_png())
```

```
# classification metrics
from sklearn.metrics import classification_report,confusion_matrix
y_pred = clf_gini.predict(X_test)
print(classification_report(y_test, y_pred))
```

```
                   precision    recall  f1-score   support

              0       0.85      0.96      0.90      6867
              1       0.77      0.47      0.59      2182

       accuracy                           0.84      9049
      macro avg       0.81      0.71      0.74      9049
   weighted avg       0.83      0.84      0.82      9049
```

```
# confusion matrix
print(confusion_matrix(y_test,y_pred))
```

```
[[6564  303]
 [1151 1031]]
```

✓  1s    completed at 11:26 PM                                    ● ✕

**Conclusion:**

1. Discuss about the how categorical attributes have been dealt with during data pre-processing.

   The Adult Census Income Dataset contains a mix of categorical and numerical attributes. Categorical attributes are those that can be divided into categories, such as "marital-status" and "race". Numerical attributes are those that can be represented as numbers, such as "age" and "hours-per-week".

   Decision trees can handle both categorical and numerical attributes, but they require different approaches to deal with them effectively. Categorical attributes are typically encoded as dummy variables, which are binary variables that indicate the presence or absence of a particular category. For example, the "marital-status" attribute could be encoded as three dummy variables: "married", "single", and "divorced".

   In this study, the categorical attributes were encoded as dummy variables using the LabelEncoder() class from the sklearn.preprocessing library. The numerical attributes were not pre-processed.

2. Discuss the hyper-parameter tunning done based on the decision tree obtained.
   The following hyperparameters were tuned:

   max_depth: This is the maximum depth of the decision tree. A deeper tree will have more splits and can potentially learn more complex patterns, but it can also be more prone to overfitting.

   min_samples_split: This is the minimum number of samples required to split a node. A higher value will prevent the tree from splitting too much and overfitting the training data.

   min_samples_leaf: This is the minimum number of samples required in a leaf node. A higher value will prevent the tree from creating too many leaf nodes and underfitting the training data.

   criterion: This is the splitting criterion used to determine the best split at each node. The most common criterion is Gini impurity, but other criteria such as entropy can also be used.

The hyperparameters were tuned using a grid search with cross-validation. This means that the model was trained and evaluated on different combinations of hyperparameters, and the best combination was selected.

The optimal hyperparameters were found to be:

max_depth = 3
min_samples_split = 50
min_samples_leaf = 50
criterion = "gini"

3. Comment on the accuracy, confusion matrix, precision, recall and F1 score obtained.

Accuracy: The accuracy is the percentage of predictions that were correct. In this case, the accuracy is 83%. This means that the model correctly classified 83% of the test set samples.

Precision: Precision is the percentage of positive predictions that were actually positive. In this case, the precision is 85%. This means that 85% of the samples that the model predicted to be >50K were actually >50K.

Recall: Recall is the percentage of actual positives that were correctly predicted. In this case, the recall is 96%. This means that 96% of the samples that were actually >50K were correctly predicted by the model.

F1 score: The F1 score is a weighted average of precision and recall. In this case, the F1 score is 90%. This means that the model has a good balance of precision and recall