

Experiment No. 3
Apply Decision Tree Algorithm on Adult Census Income Dataset and analyze the performance of the model
Date of Performance:07-08-2023
Date of Submission:09-09-2023

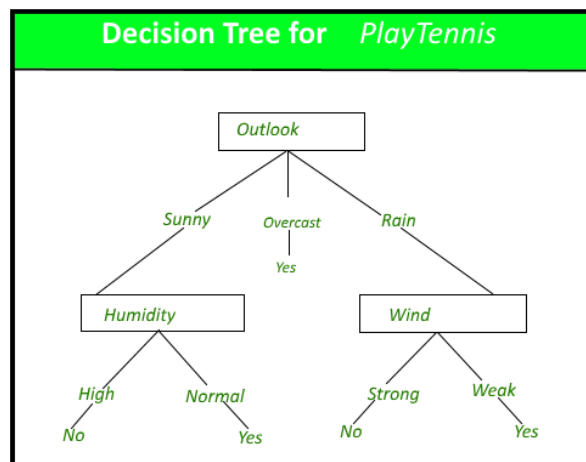


**Aim:** Apply Decision Tree Algorithm on Adult Census Income Dataset and analyze the performance of the model.

**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:



## Vidyavardhini's College of Engineering & Technology

### Department of Computer Engineering

---

>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:**

**Conclusion:**

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

The categorical attributes were transformed from their original text based representations into numerical representations by assigning a unique integer to each of the unique categorical values in each column, using the Label Encoder technique. This makes it possible for algorithms requiring numerical input to apply these categorical variables.

**2. Discuss the hyper-parameter tuning done based on the decision tree obtained.**

**Maximum Depth:**

1. This setting limits the depth of the decision tree, preventing it from
2. becomes too complex and overfits the training data.

**Minimum sample allocation:**

1. This parameter determines the minimum number of samples required in one
2. node is eligible for further splitting. This helps the tree not produce too much fruit
3. Specific decisions are based on a small number of cases.

**Minimum sample table:**

- 1.This parameter determines the minimum number of samples that should be contained in a sheet



2.knot. Similar to splitting minimal samples, this can prevent the tree from creating existing nodes very few cases.

Criteria:

1. This parameter determines the function used to measure the quality of the division. "Gini
2. "impurity" and "entropy" are common criteria

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

1. Model accuracy is about 84%. This means the model correctly predicted the class label for 84% of the instances in the test dataset.

2. True positive (TP): 1031, true negative (TN): 6564, False Positive (FP): 303, False Negative (FN): 1151. The confusion matrix indicates that the model performs well in type 0 prediction (high true negatives and true positives), but has difficulty with type 1 prediction (high false negatives).

3. The accuracy of class 1 is relatively good, proving that when the model predicts class 1, it is usually correct. For type 1, the accuracy was about 0.77, indicating that of all the cases predicted to be type 1, about 77% were actually type 1.

4. The recall rate for class 1 is lower, indicating that the model is missing a significant number of real instances of class 1. For class 1, the recall rate is about 0.47. This means that the model can only correctly identify about 47% of all real cases belonging to class 1.

5. The F1 score for Type 1 lies between precision and recall, providing a balanced view of model performance.



**Vidyavardhini's College of Engineering & Technology**

Department of Computer Engineering

---

```
# 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_dataset.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
0	90	?	77053	HS-grad	9	Widowed	?	
1	82	Private	132870	HS-grad	9	Widowed	Exec-managerial	
2	66	?	186061	Some-college	10	Widowed	?	

```
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', 'capital.gain', 'capital.loss', 'hours.per.week', 'native.country', 'income']
```

```
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
```

```
# EDA
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

df.describe()

   age      fnlwgt  education.num  capital.gain  capital.loss  hours
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    40.620886
std     13.640433  1.055500e+05     2.572720    7385.292085   402.960219    15.353114
min     17.000000  1.228500e+04     1.000000     0.000000     0.000000     0.000000
25%     28.000000  1.178270e+05     9.000000     0.000000     0.000000     0.000000
50%     37.000000  1.783560e+05    10.000000     0.000000     0.000000     0.000000
75%     48.000000  2.370510e+05    12.000000     0.000000     0.000000     0.000000
max     90.000000  1.484705e+06    16.000000  99999.000000  4356.000000    99.000000

df.head()

   age  workclass  fnlwgt  education  education.num  marital.status  occupation  native.country
0   90         ?   77053    HS-grad              9        Widowed         ?         United-States
1   82    Private  132870    HS-grad              9        Widowed  Exec-managerial         United-States
2   66         ?  186061  Some-college            10        Widowed         ?         United-States
3   54    Private  140359    7th-8th              4        Divorced  Machine-op-inspct         United-States
4   41    Private  264663  Some-college            10        Separated  Prof-specialty         United-States

# 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

```

```

# 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

```

```

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

```

	age	workclass	fnlwgt	education	education.num	marital.status	occupation	n
1	82	Private	132870	HS-grad	9	Widowed	Exec-manual	
3	54	Private	140359	7th-8th	4	Divorced	Machine-op-inspct	
4	41	Private	264663	Some-college	10	Separated	Prof-specialty	
5	34	Private	216864	HS-grad	9	Divorced	Other-service	
6	38	Private	150601	10th	6	Separated	Adm-clerical	

```
# 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

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

```
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
```

```
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	r
1	Private	HS-grad	Widowed	Exec-managerial	Not-in-family	White	Female	
3	Private	7th-8th	Divorced	Machine-op-inspct	Unmarried	White	Female	
4	Private	Some-college	Separated	Prof-specialty	Own-child	White	Female	

```
# 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
1	2	11	6	3	1	4	0	

```
# 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
1	82	132870	9	0	4356	18	
3	54	140359	4	0	3900	40	
4	41	264663	10	0	3900	40	
5	34	216864	9	0	3770	45	
6	38	150601	6	0	3770	40	

```
# 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
```

```
# convert target variable income to categorical
```

```
df['income'] = df['income'].astype('category')
```

```
# 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
```

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

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

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

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

	age	fnlwgt	education.num	capital.gain	capital.loss	hours.per.week	workcla
1	82	132870	9	0	4356	18	
3	54	140359	4	0	3900	40	
4	41	264663	10	0	3900	40	

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

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

```
# 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)
```

```
X_train.head()
```

	age	fnlwgt	education.num	capital.gain	capital.loss	hours.per.week	workcla
24351	42	289636	9	0	0	46	
15626	37	52465	9	0	0	40	
4347	38	125933	14	0	0	40	
23972	44	183829	13	0	0	38	
26843	35	198841	11	0	0	35	

Decision tree

```
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)
```

```
!pip install my-package
```

```
Collecting my-package
  Downloading my_package-0.0.0-py3-none-any.whl (2.0 kB)
Installing collected packages: my-package
Successfully installed my-package-0.0.0
```

```
!pip install graphviz
```

```
Requirement already satisfied: graphviz in /usr/local/lib/python3.10/dist-packages (0.20.1)
```

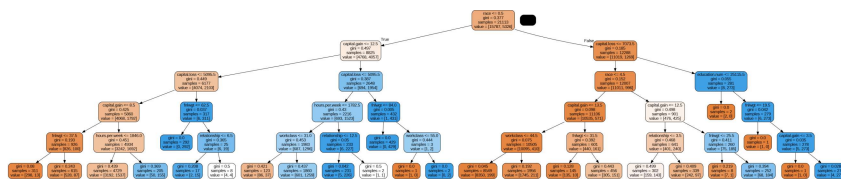
```
# Importing required packages for visualization
from IPython.display import Image
from six import StringIO
from sklearn.tree import export_graphviz
import pydotplus,graphviz
```

```
# 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']
```

```
# 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())
```



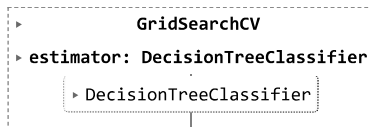
```
# GridSearchCV to find optimal max_depth
from sklearn.model_selection import KFold
from sklearn.model_selection import GridSearchCV
```

```
# 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)
```



```
# 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
0	0.047192	0.005210	0.010834	0.004780	1
1	0.063689	0.005183	0.010479	0.002776	2
2	0.074079	0.023605	0.009829	0.004383	3
3	0.064308	0.022573	0.007270	0.002981	4
4	0.052663	0.001205	0.005903	0.000483	5

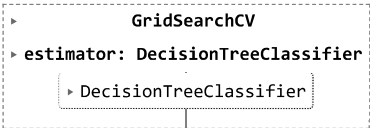
```
# GridSearchCV to find optimal max_depth
from sklearn.model_selection import KFold
from sklearn.model_selection import GridSearchCV

# 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)
```



```
# 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
0	0.126720	0.006284	0.006475	0.000638	
1	0.112024	0.008634	0.006495	0.000440	
2	0.101930	0.008014	0.006128	0.000168	
3	0.100225	0.025415	0.008158	0.004255	
4	0.059809	0.000704	0.003464	0.000035	

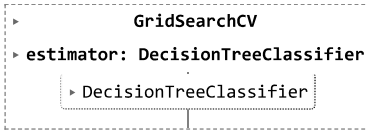
```
# GridSearchCV to find optimal min_samples_split
from sklearn.model_selection import KFold
from sklearn.model_selection import GridSearchCV
```

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

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

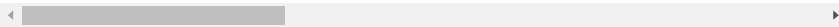
```
# instantiate the model
dtree = DecisionTreeClassifier(criterion = "gini",
                              random_state = 100)
```

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



```
# 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
0	0.273618	0.005163	0.012117	0.003469	
1	0.143774	0.037922	0.006201	0.000153	
2	0.097914	0.013851	0.004660	0.000958	
3	0.085130	0.004583	0.004413	0.000935	
4	0.082915	0.002581	0.004023	0.000450	



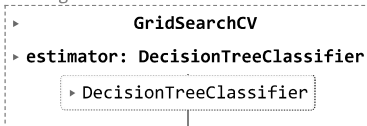
```
# Create the parameter grid
param_grid = {
    'max_depth': range(5, 15, 5),
    'min_samples_leaf': range(50, 150, 50),
    'min_samples_split': range(50, 150, 50),
    'criterion': ["entropy", "gini"]
}
```

```
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



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



	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_criterion
0	0.060532	0.001683	0.007079	0.000595	entropy
1	0.064738	0.005289	0.005540	0.000233	entropy
2	0.062205	0.009672	0.005705	0.000395	entropy
3	0.043008	0.005321	0.004014	0.000528	entropy
4	0.068211	0.004726	0.004385	0.001273	entropy
5	0.064453	0.001698	0.004686	0.000877	entropy
6	0.060422	0.002737	0.003985	0.000779	entropy
7	0.059891	0.003640	0.003670	0.000308	entropy
8	0.037223	0.002663	0.003544	0.000335	gini
9	0.036449	0.002312	0.003664	0.000604	gini

```
# 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)
```

```
# 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)
```

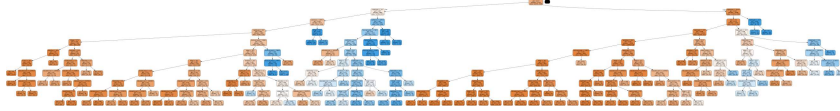
```
DecisionTreeClassifier
DecisionTreeClassifier(max_depth=10, min_samples_leaf=50, min_samples_split=50,
                      random_state=100)
```

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

```
0.850922753895458
10      0.000000      0.000000      0.000000      0.000000      gini
```

```
# 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())
```



```
# 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)
```

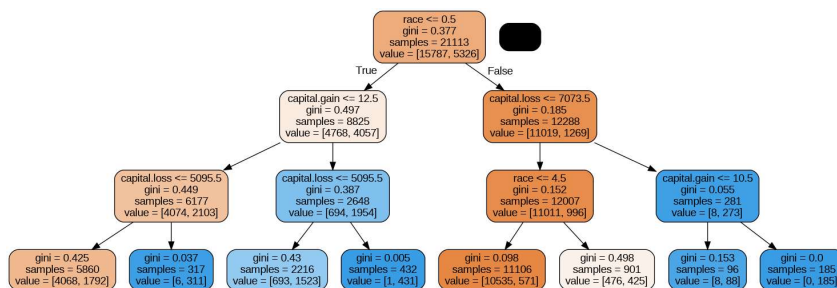
```
DecisionTreeClassifier
DecisionTreeClassifier(max_depth=3, min_samples_leaf=50, min_samples_split=50,
                      random_state=100)
```

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

```
0.8393192617968837
0.8393192617968837
```

```
# 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]]
```

