



Vidyavardhini's College of Engineering & Technology

Department of Computer Engineering

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: 20-08-2023



Vidyavardhini's College of Engineering & Technology

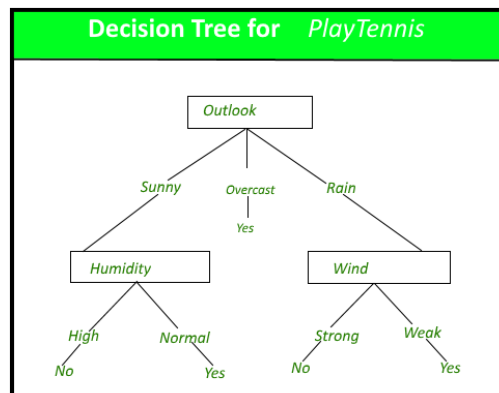
Department of Computer Engineering

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.



Vidyavardhini's College of Engineering & Technology

Department of Computer Engineering

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.

```
# 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 = "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	race	Fe
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

```
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 :

age73

workclass9

fnlwgt21648

education16

education.num16

marital.status7

occupation15

relationship6

race5

sex2

capital.gain119

capital.loss92

hours.per.week94

native.country42

income2

dtype: int64

```
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.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

df.head()

	age	workclass	fnlwgt	education	education.num	marital.status	occupation	relationship	race	sex	capital.gain	capital.loss
0	90	?	77053	HS-grad	9	Widowed	?	Not-in-family	White	Female	0	?
1	82	Private	132870	HS-grad	9	Widowed	Exec-managerial	Not-in-family	White	Female	0	?
2	66	?	186061	Some-college	10	Widowed	?	Unmarried	Black	Female	0	?
3	54	Private	140359	7th-8th	4	Divorced	Machine-op-inspct	Unmarried	White	Female	0	?
4	41	Private	264663	Some-college	10	Separated	Prof-specialty	Own-child	White	Female	0	?

```
# 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

# 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  relationship  race  sex  capital.gain  capital.loss
1   82    Private  132870    HS-grad             9         Widowed    Exec-managerial  Not-in-family  White  Female             0             0
3   54    Private  140359    7th-8th             4         Divorced    Machine-op-inspct  Unmarried    White  Female             0             0
4   41    Private  264663    Some-college        10         Separated    Prof-specialty  Own-child    White  Female             0             0
5   34    Private  216864    HS-grad             9         Divorced    Other-service  Unmarried    White  Female             0             0

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

```
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	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-service	Unmarried	White	Female	United-States	<=50K
6	Private	10th	Separated	Adm-clerical	Unmarried	White	Male	United-States	<=50K

```
# 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

```
# 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	marital.status	occupation	relation
1	82	132870	9	0	4356	18	2	11	6	3	
3	54	140359	4	0	3900	40	2	5	0	6	
4	41	264663	10	0	3900	40	2	15	5	9	
5	34	216864	9	0	3770	45	2	11	0	7	
6	38	150601	6	0	3770	40	2	0	5	0	

```
# 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

```

```

# 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	workclass	education	marital.status	occupation	relatio
1	82	132870	9	0	4356	18	2	11	6	3	
3	54	140359	4	0	3900	40	2	5	0	6	
4	41	264663	10	0	3900	40	2	15	5	9	

```
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	fnlwtg	education.num	capital.gain	capital.loss	hours.per.week	workclass	education	marital.status	occupation	rel
24351	42	289636	9	0	0	46	2	11	2	13	
15626	37	52465	9	0	0	40	1	11	4	7	

```
# 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)
```

```
# 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      0.85      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 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 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 six import StringIO
from sklearn.tree import export_graphviz
import pydotplus,graphviz
```

```
# Putting features
features = list(df.columns[1:])
features
```

```
['fnlwtg',
 'education.num',
 'capital.gain',
 'capital.loss',
 'hours.per.week',
 'workclass',
 'education',
 'marital.status',
 'occupation',
 'relationship',
 'race',
 'sex',
```

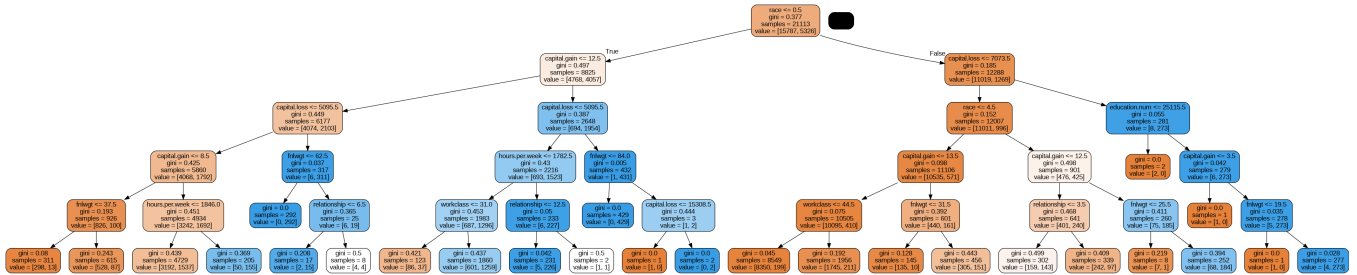
```
'native.country',
'income']
```

```
!pip install graphviz
```

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

```
# 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())
```



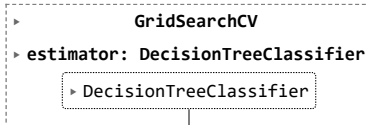
```
# GridSearchCV to find optimal max_depth
from sklearn.model_selection import KFold
from sklearn.model_selection import GridSearchCV
```

```
# 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	params	split0_test_score	split1_test_score
0	0.017606	0.003309	0.004545	0.000778	1	{'max_depth': 1}	0.747810	0.747810
1	0.019553	0.002510	0.003388	0.000431	2	{'max_depth': 2}	0.812219	0.818612
2	0.024394	0.001255	0.003224	0.000196	3	{'max_depth': 3}	0.828558	0.834241
3	0.031005	0.003047	0.003377	0.000224	4	{'max_depth': 4}	0.832583	0.840871
4	0.038414	0.004430	0.003606	0.000638	5	{'max_depth': 5}	0.834241	0.844897

```
"""
# 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\n         -"test accuracy"\n\nplt.xlabel("max_depth")\nplt.ylabel("Accuracy")\nplt.legend()\nplt.show()\n'

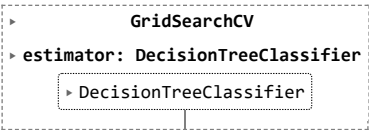
# GridSearchCV to find optimal max_depth
from sklearn.model_selection import KFold
from sklearn.model_selection import GridSearchCV

# 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_leaf	params	split0_test_score	split
0	0.206307	0.049675	0.006364	0.001320	5	{'min_samples_leaf': 5}	0.825716	
1	0.130202	0.019363	0.006910	0.002911	25	{'min_samples_leaf': 25}	0.841819	
2	0.109869	0.021408	0.005161	0.000183	45	{'min_samples_leaf': 45}	0.843003	
3	0.106612	0.017281	0.008429	0.006759	65	{'min_samples_leaf': 65}	0.841108	
4	0.116716	0.018254	0.009106	0.005065	85	{'min_samples_leaf': 85}	0.838030	

```
"""
# 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_leaf"], \n      scores["mean_train_score"], \n      label="training accuracy")\nplt.plot(scores["param_min_samples_leaf"], \n      scores["mean_test_score"], \n      label="test accuracy")\nplt.xlabel("min_samples_leaf")\nplt.ylabel("Accuracy")\nplt.legend()\nplt.show()\n'

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

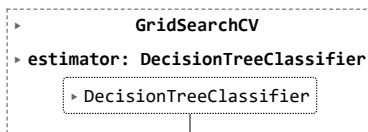
# specify number of folds for k-fold CV
n_folds = 5

# 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_split	params	split0_test_score	spli
0	0.135789	0.007802	0.005833	0.000620	5	{'min_samples_split': 5}	0.811982	
1	0.129377	0.003017	0.005809	0.000127	25	{'min_samples_split': 25}	0.825006	
2	0.127084	0.003051	0.005815	0.000104	45	{'min_samples_split': 45}	0.835188	
3	0.124147	0.005110	0.006897	0.001300	65	{'min_samples_split': 65}	0.839451	
4	0.090817	0.015012	0.004403	0.000369	85	{'min_samples_split': 85}	0.846081	

```

"""
# 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_score"], \n      label="test accuracy")\nplt.xlabel("min_samples_split")\nplt.ylabel("Accuracy")\nplt.legend()\nplt.show()\n'

# 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
```

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_criterion	param_max_depth	param_min_samples_leaf	param_m...
0	0.040022	0.005139	0.003463	0.000684	entropy	5	50	
1	0.038898	0.002205	0.003150	0.000035	entropy	5	50	
2	0.039383	0.001463	0.003649	0.000590	entropy	5	100	
3	0.038811	0.001748	0.003635	0.000666	entropy	5	100	

```
# 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

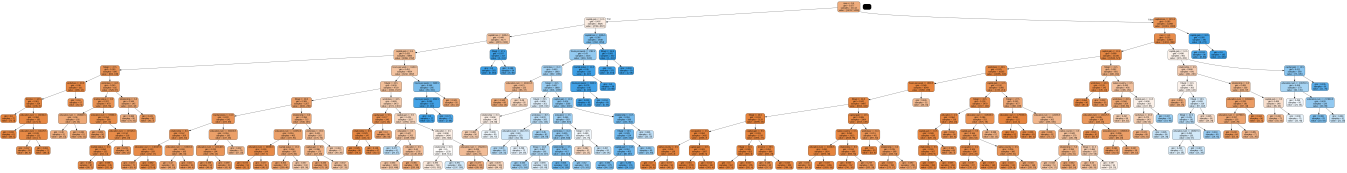
DecisionTreeClassifier(max_depth=10, min_samples_leaf=50, min_samples_split=50, random_state=100)

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

0.850922753895458
```

```
#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)

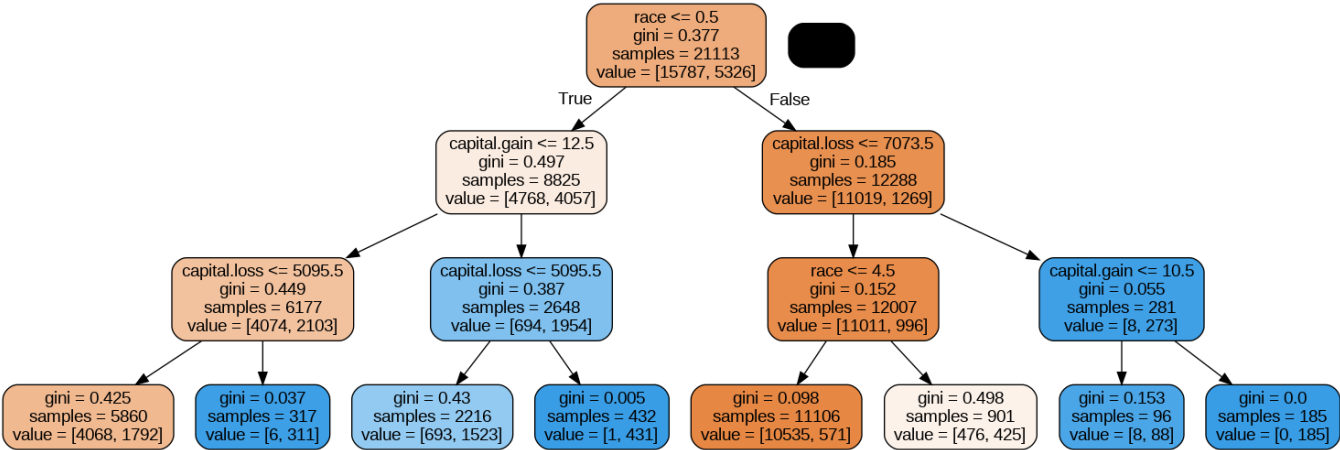
clf_gini.fit(X_train, y_train)

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

0.8393192617968837	0.000000	0.000000	0.000000	0.000000	gini	10	100	
--------------------	----------	----------	----------	----------	------	----	-----	--

```
# 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]]
```



Conclusion:

1. By using the Label Encoder technique, categorical attributes have been transformed from their original text-based representations into numerical representations by assigning a unique integer to each unique categorical value within each column. This makes it possible to use these categorical variables in machine learning algorithms that require numerical input.
2. Hyper-parameter tuning done based on the decision tree obtained:
 - Max Depth: This parameter restricts the depth of the decision tree, preventing it from becoming too complex and overfitting the training data.
 - Min Samples Split: This parameter sets the minimum number of samples required in a node to be eligible for further splitting. It helps prevent the tree from making overly specific decisions based on a small number of instances.
 - Min Samples Leaf: This parameter sets the minimum number of samples to be in a leaf node. Similar to min samples split, this can prevent the tree from creating nodes with very few instances.
 - Criterion: This parameter defines the function used to measure the quality of a split. "Gini impurity" and "entropy" are common criteria.
3. The accuracy of the model is approximately 84%. This means that the model correctly predicted the class labels for 84% of the instances in the test dataset.
4. True Positive (TP): 1031, True Negative (TN): 6564, False Positive (FP): 303, False Negative (FN): 1151. The confusion matrix indicates that the model is performing well in predicting class 0 (high true negatives and true positives), but it struggles with class 1 prediction (high false negatives).
5. The precision for class 1 is relatively good, indicating that when the model predicts class 1, it's often correct. For class 1, the precision is approximately 0.77, indicating that out of all instances predicted as class 1, around 77% are actually class 1.
6. The recall for class 1 is lower, suggesting that the model misses a significant number of actual class 1 instances. For class 1, the recall is approximately 0.47. This means that the model was able to correctly identify only about 47% of all actual instances that belong to class 1.
7. The F1 score for class 1 is in between precision and recall, providing a balanced view of the model's performance.