

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

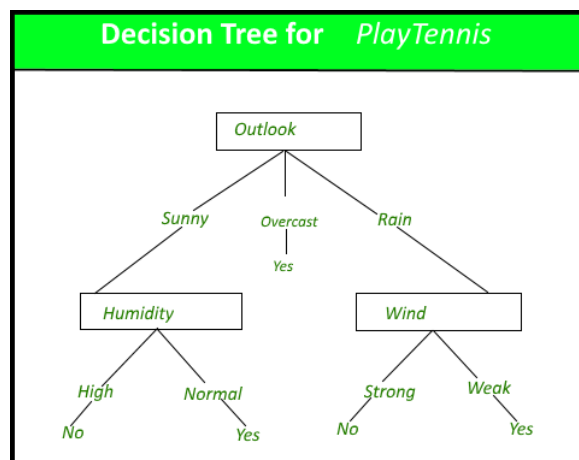


**Aim:** Apply Decision Tree Algorithm on Adult Census Income Dataset and analyze the performance of the model.

**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.



## Vidyavardhini's College of Engineering & Technology

### Department of Computer Engineering

---

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.

Handling categorical & numerical features during data preparation is crucial for creating precise models. There are various preprocessing methods, including:

- a. Label Encoding, which converts qualities like education, occupation, marital status, relationship, and sex into numerical values. Decision tree algorithm can use these attributes since label encoding gives each category a distinct number.
- b. Missing values: For categorical values with missing values, the median imputation method was used, or the mode method for ordinal values.
- c. One Hot Encoder: One Hot Encoder is used to convert categorical values into binary vectors and each category is represented as a binary column for characteristics with nominal categories like native country and occupation.

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

The Decision Tree model's performance is optimized by the hyper-parameter tweaking, which makes sure that the model is not under- or overfitting the data.

To investigate various combinations of factors, such as the maximum tree depth and splitting criteria, a systematic grid search method was used. Finding the appropriate hyper-parameter to maximize the model's performance is the goal of this method.



3. Comment on the accuracy, confusion matrix, precision, recall and F1 score obtained. The ratio of results that were accurately predicted to all instances was used to determine the model's accuracy. It gives an overview of the model's performance. The accuracy of this model is about 85%. Confusion Matrix: The confusion matrix offers a thorough analysis of the performance of the model by displaying true positives, true negatives, false positives, and false negatives. For this model, the confusion matrix is as follows:

```
[[6353 334]
 [1059 1243]]
```

The harmonic mean of recall and precision is known as the F1 score. It offers a fair evaluation of a model's performance. This model's f1 score falls between recall & precision.

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

```
pima=pd.read_csv("adult.csv")
```

```
pima.head(3)
```

	age	workclass	fnlwgt	education	education.num	marital.status	occupation	relationship	race	sex	capital.gain	capital.loss	hours.per.week	native.co
0	90	?	77053	HS-grad	9	Widowed	?	Not-in-family	White	Female	0	4356	40	United-
1	82	Private	132870	HS-grad	9	Widowed	Exec-managerial	Not-in-family	White	Female	0	4356	18	United-
2	66	?	186061	Some-college	10	Widowed	?	Unmarried	Black	Female	0	4356	40	United-

```
print("shape")
pima.shape
```

shape
(32561, 15)

```
pima.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
```

```
pima.describe()
```

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

```
print("? in workclass")
df_check_missing_workclass = (pima['workclass']=='?').sum()
print(df_check_missing_workclass)
print("? in occupation")
df_check_missing_occupation = (pima['occupation']=='?').sum()
print(df_check_missing_occupation)
```

```
? in workclass
1836
? in occupation
1843
```

```
df_missing = (pima=='?').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
```

```
pima = pima[pima['workclass'] != '?']
pima = pima[pima['occupation'] != '?']
pima = pima[pima['native.country'] != '?']
pima.head()
```

```

age workclass fnlwgt education education.num marital.status occupation relationship race sex capital.gain capital.loss hours.per.week native.co
1 82 Private 132870 HS-grad 9 Widowed Exec-managerial Not-in-family White Female 0 4356 18 United-
3 54 Private 140359 7th-8th 4 Divorced Machine-op-inspct Unmarried White Female 0 3900 40 United-
4 41 Private 264663 Some-college 10 Separated Prof-specialty Own-child White Female 0 3900 40 United-
5 34 Private 216864 HS-grad 9 Divorced Other-service Unmarried White Female 0 3770 45 United-
6 50 Private 158129 10th 5 Separated Adm-clerical Unmarried White Male 0 3450 40 United-

```

```
from sklearn import preprocessing
```

```
# encode categorical variables using label Encoder
```

```
# select all categorical variables
```

```
df_categorical = pima.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

```
pima = pima.drop(df_categorical.columns,axis=1)
```

```
pima = pd.concat([pima,df_categorical],axis=1)
```

```
pima.head()
```



	age	fnlwgt	education.num	capital.gain	capital.loss	hours.per.week	workclass	education	marital.status	occupation	relationship	race	sex	native.countr
1	82	132870	9	0	4356	18	2	11	6	3	1	4	0	3
3	54	140359	4	0	3900	40	2	5	0	6	4	4	0	3
4	41	264663	10	0	3900	40	2	15	5	9	3	4	0	3
5	34	216864	9	0	3770	45	2	11	0	7	4	4	0	3

```
pima.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
```

```
pima['income'] = pima['income'].astype('category')
pima.head(10)
```

	age	fnlwgt	education.num	capital.gain	capital.loss	hours.per.week	workclass	education	marital.status	occupation	relationship	race	sex	native.country	income
1	82	132870	9	0	4356	18	2	11	6	3	1	4	0	38	0
3	54	140359	4	0	3900	40	2	5	0	6	4	4	0	38	0
4	41	264663	10	0	3900	40	2	15	5	9	3	4	0	38	0
5	34	216864	9	0	3770	45	2	11	0	7	4	4	0	38	0
6	38	150601	6	0	3770	40	2	0	5	0	4	4	1	38	0
7	74	88638	16	0	3683	20	5	10	4	9	2	4	0	38	1
8	68	422013	9	0	3683	40	0	11	0	9	1	4	0	38	0
10	45	172274	16	0	3004	35	2	10	0	9	4	2	0	38	1
11	38	164526	15	0	2824	45	4	14	4	9	1	4	1	38	1
12	52	129177	13	0	2824	20	2	9	6	7	1	4	0	38	1

```
pima.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

```

```

from sklearn.model_selection import train_test_split
# Putting independent variables/features to X
X = pima.drop('income',axis=1)

```

```

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

```

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

	age	fnlwgt	education.num	capital.gain	capital.loss	hours.per.week	workclass	education	marital.status	occupation	relationship	race	sex	native.countr
1	82	132870	9	0	4356	18	2	11	6	3	1	4	0	3
3	54	140359	4	0	3900	40	2	5	0	6	4	4	0	3
4	41	264663	10	0	3900	40	2	15	5	9	3	4	0	3

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

```

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

```

```

X_train,X_test,y_train,y_test = train_test_split(X,y,test_size=0.30,random_state=99)
# 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)
```

```
from sklearn.metrics import classification_report,confusion_matrix,accuracy_score
y_pred_default = dt_default.predict(X_test)
```

```
print(confusion_matrix(y_test,y_pred_default))
print(accuracy_score(y_test,y_pred_default))
```

```
[[6553 314]
 [1039 1143]]
0.849
```