

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
```

```
df = pd.read_csv('Decision_Tree_Income_Prediction.csv')
```

```
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	4356
1	82	Private	132870	HS-grad	9	Widowed	Exec-managerial	Not-in-family	White	Female	0	4356
2	66	?	186061	Some-college	10	Widowed	?	Unmarried	Black	Female	0	4356
3	54	Private	140359	7th-8th	4	Divorced	Machine-op-inspct	Unmarried	White	Female	0	3900
4	41	Private	264663	Some-college	10	Separated	Prof-specialty	Own-child	White	Female	0	3900



```
df.tail()
```

	age	workclass	fnlwgt	education	education.num	marital.status	occupation	relationship	race	sex	capital.gain	capital.loss
32556	22	Private	310152	Some-college	10	Never-married	Protective-serv	Not-in-family	White	Male	0	
32557	27	Private	257302	Assoc-acdm	12	Married-civ-spouse	Tech-support	Wife	White	Female	0	
32558	40	Private	154374	HS-grad	9	Married-civ-spouse	Machine-op-inspct	Husband	White	Male	0	
32559	58	Private	151910	HS-grad	9	Widowed	Adm-clerical	Unmarried	White	Female	0	
32560	22	Private	201490	HS-grad	9	Never-married	Adm-clerical	Own-child	White	Male	0	

```
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	

```
df.replace('?',pd.NA, inplace=True)
```

```
df.isnull().sum()
```

```
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
```

```
df.dropna(inplace=True)
```

```
df.isnull().sum()
```

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

```
from sklearn.preprocessing import LabelEncoder
LEOBJ = LabelEncoder()
columns = ['workclass','education','marital.status','occupation','relationship','race','sex','native.country','income']
for col in columns:
    df[col] = LEOBJ.fit_transform(df[col])
```

```
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  int64
2   fnlwgt                30162 non-null  int64
3   education             30162 non-null  int64
4   education.num         30162 non-null  int64
5   marital.status        30162 non-null  int64
6   occupation            30162 non-null  int64
7   relationship          30162 non-null  int64
8   race                  30162 non-null  int64
9   sex                   30162 non-null  int64
```

```

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 int64
14 income          30162 non-null int64
dtypes: int64(15)
memory usage: 3.7 MB

```

```

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

```

```

from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X,y, test_size=0.2, random_state=42)

```

```

from sklearn.tree import DecisionTreeClassifier
from sklearn.tree import plot_tree

```

```
dtclf = DecisionTreeClassifier(max_leaf_nodes = 8)
```

```
dtclf.fit(X_train, y_train)
```

```

▼      DecisionTreeClassifier
DecisionTreeClassifier(max_leaf_nodes=8)

```

```
y_pred = dtclf.predict(X_test)
```

```
from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
```

```

print("Confusion Matrix : \n", confusion_matrix(y_test, y_pred))
print()
print("Accuracy score: \n ",accuracy_score(y_test,y_pred))
print()
print("Classification report: \n ", classification_report(y_test,y_pred))

```

```

Confusion Matrix :
[[4303  230]
 [ 750  750]]

```

```

Accuracy score:
0.8375600861926074

```

```

Classification report:

```

	precision	recall	f1-score	support
0	0.85	0.95	0.90	4533
1	0.77	0.50	0.60	1500
accuracy			0.84	6033
macro avg	0.81	0.72	0.75	6033
weighted avg	0.83	0.84	0.82	6033

```

plt.figure(figsize=(12,10))
plot_tree(dtclf,filled=True, feature_names=X.columns, class_names=list(map(str,dtclf.classes_)))
plt.show

```

```
<function matplotlib.pyplot.show(close=None, block=None)>
```

