

Assignment-14-Decision Trees [Fraud Check]

In [289]:

```
1 #Importing necessary Libraries
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from sklearn import datasets
5 import numpy as np
6 from sklearn.model_selection import train_test_split
7 from sklearn.tree import DecisionTreeClassifier
8 from sklearn import tree
9 from sklearn.metrics import classification_report
10 from sklearn import preprocessing
11
12 df = pd.read_csv("Fraud_check.csv")
```

In [290]:

```
1 # Viewing top 5 rows of dataset
2 df.head()
```

Out[290]:

|   | Undergrad | Marital.Status | Taxable.Income | City.Population | Work.Experience | Urban |
|---|-----------|----------------|----------------|-----------------|-----------------|-------|
| 0 | NO        | Single         | 68833          | 50047           | 10              | YES   |
| 1 | YES       | Divorced       | 33700          | 134075          | 18              | YES   |
| 2 | NO        | Married        | 36925          | 160205          | 30              | YES   |
| 3 | YES       | Single         | 50190          | 193264          | 15              | YES   |
| 4 | NO        | Married        | 81002          | 27533           | 28              | NO    |

In [291]:

```
1 df.tail()
```

Out[291]:

|     | Undergrad | Marital.Status | Taxable.Income | City.Population | Work.Experience | Urban |
|-----|-----------|----------------|----------------|-----------------|-----------------|-------|
| 595 | YES       | Divorced       | 76340          | 39492           | 7               | YES   |
| 596 | YES       | Divorced       | 69967          | 55369           | 2               | YES   |
| 597 | NO        | Divorced       | 47334          | 154058          | 0               | YES   |
| 598 | YES       | Married        | 98592          | 180083          | 17              | NO    |
| 599 | NO        | Divorced       | 96519          | 158137          | 16              | NO    |

In [292]:

```
1 # Creating dummy variables for ['Undergrad','marital.status','Urban'] dropping first dummy variable
2 df=pd.get_dummies(df,columns=['Undergrad','Marital.Status','Urban'], drop_first=True)
```

In [293]:

```
1 # Creating new cols TaxInc and dividing 'Taxable.Income' cols on the basic of [10002,30000,99620] for Risky and Good
2 df["TaxInc"] = pd.cut(df["Taxable.Income"], bins = [10002,30000,99620], labels = ["Risky", "Good"])
```

In [294]:

```
1 print(df)
```

|     | Taxable.Income         | City.Population       | Work.Experience | Undergrad_YES | \ |
|-----|------------------------|-----------------------|-----------------|---------------|---|
| 0   | 68833                  | 50047                 | 10              | 0             |   |
| 1   | 33700                  | 134075                | 18              | 1             |   |
| 2   | 36925                  | 160205                | 30              | 0             |   |
| 3   | 50190                  | 193264                | 15              | 1             |   |
| 4   | 81002                  | 27533                 | 28              | 0             |   |
| ..  | ...                    | ...                   | ...             | ...           |   |
| 595 | 76340                  | 39492                 | 7               | 1             |   |
| 596 | 69967                  | 55369                 | 2               | 1             |   |
| 597 | 47334                  | 154058                | 0               | 0             |   |
| 598 | 98592                  | 180083                | 17              | 1             |   |
| 599 | 96519                  | 158137                | 16              | 0             |   |
|     |                        |                       |                 |               |   |
|     | Marital.Status_Married | Marital.Status_Single | Urban_YES       | TaxInc        |   |
| 0   | 0                      | 1                     | 1               | Good          |   |
| 1   | 0                      | 0                     | 1               | Good          |   |
| 2   | 1                      | 0                     | 1               | Good          |   |
| 3   | 0                      | 1                     | 1               | Good          |   |
| 4   | 1                      | 0                     | 0               | Good          |   |
| ..  | ...                    | ...                   | ...             | ...           |   |
| 595 | 0                      | 0                     | 1               | Good          |   |
| 596 | 0                      | 0                     | 1               | Good          |   |
| 597 | 0                      | 0                     | 1               | Good          |   |
| 598 | 1                      | 0                     | 0               | Good          |   |
| 599 | 0                      | 0                     | 0               | Good          |   |

[600 rows x 8 columns]

Lets assume: taxable\_income <= 30000 as “Risky=0” and  
others are “Good=1”

In [295]:

```
1 # After creation of new col. TaxInc also made its dummies var concating right side of df.
2 df = pd.get_dummies(df, columns = ["TaxInc"], drop_first=True)
```

In [296]:

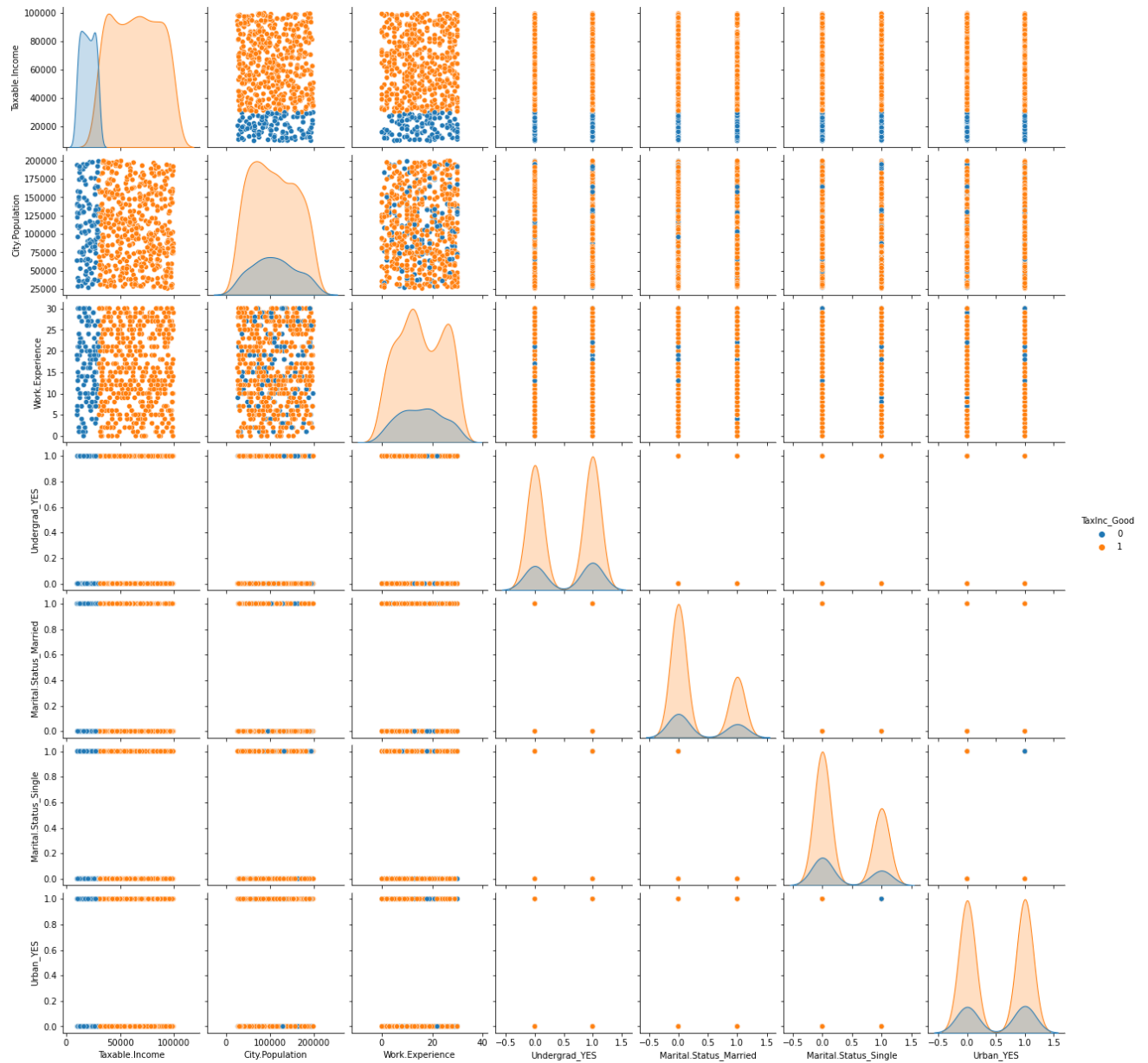
```
1 # Viewing buttom 10 observations
2 df.tail(10)
```

Out[296]:

|     | Taxable.Income | City.Population | Work.Experience | Undergrad_YES | Marital.Status_Married | Marital.Status_Single | Urban_YES | TaxInc_Good |
|-----|----------------|-----------------|-----------------|---------------|------------------------|-----------------------|-----------|-------------|
| 590 | 43018          | 85195           | 14              | 0             | 1                      | 0                     | 1         | 1           |
| 591 | 27394          | 132859          | 18              | 1             | 0                      | 1                     | 1         | 0           |
| 592 | 68152          | 75143           | 16              | 1             | 0                      | 1                     | 0         | 1           |
| 593 | 84775          | 131963          | 10              | 0             | 0                      | 0                     | 1         | 1           |
| 594 | 47364          | 97526           | 9               | 0             | 1                      | 0                     | 1         | 1           |
| 595 | 76340          | 39492           | 7               | 1             | 0                      | 0                     | 1         | 1           |
| 596 | 69967          | 55369           | 2               | 1             | 0                      | 0                     | 1         | 1           |
| 597 | 47334          | 154058          | 0               | 0             | 0                      | 0                     | 1         | 1           |
| 598 | 98592          | 180083          | 17              | 1             | 1                      | 0                     | 0         | 1           |
| 599 | 96519          | 158137          | 16              | 0             | 0                      | 0                     | 0         | 1           |

```
In [297]: 1 # Lets plot pair plot to visualise the attributes all at once
2 import seaborn as sns
3 sns.pairplot(data=df, hue = 'TaxInc_Good')
```

Out[297]: <seaborn.axisgrid.PairGrid at 0x260780d6af0>



```
In [298]: 1 # Normalization function
2 def norm_func(i):
3     x = (i-i.min())/(i.max()-i.min())
4     return (x)
```

```
In [299]: 1 # Normalized dataframe (considering the numerical part of data)
          2 df_norm = norm_func(df.iloc[:,1:])
          3 df_norm.tail(10)
```

```
Out[299]:
```

|     | City.Population | Work.Experience | Undergrad_YES | Marital.Status_Married | Marital.Status_Single | Urban_YES | TaxInc_Good |
|-----|-----------------|-----------------|---------------|------------------------|-----------------------|-----------|-------------|
| 590 | 0.341473        | 0.466667        | 0.0           | 1.0                    | 0.0                   | 1.0       | 1.0         |
| 591 | 0.615406        | 0.600000        | 1.0           | 0.0                    | 1.0                   | 1.0       | 0.0         |
| 592 | 0.283703        | 0.533333        | 1.0           | 0.0                    | 1.0                   | 0.0       | 1.0         |
| 593 | 0.610256        | 0.333333        | 0.0           | 0.0                    | 0.0                   | 1.0       | 1.0         |
| 594 | 0.412341        | 0.300000        | 0.0           | 1.0                    | 0.0                   | 1.0       | 1.0         |
| 595 | 0.078811        | 0.233333        | 1.0           | 0.0                    | 0.0                   | 1.0       | 1.0         |
| 596 | 0.170058        | 0.066667        | 1.0           | 0.0                    | 0.0                   | 1.0       | 1.0         |
| 597 | 0.737240        | 0.000000        | 0.0           | 0.0                    | 0.0                   | 1.0       | 1.0         |
| 598 | 0.886810        | 0.566667        | 1.0           | 1.0                    | 0.0                   | 0.0       | 1.0         |
| 599 | 0.760683        | 0.533333        | 0.0           | 0.0                    | 0.0                   | 0.0       | 1.0         |

```
In [300]: 1 # Declaring features & target
          2 X = df_norm.drop(['TaxInc_Good'], axis=1)
          3 y = df_norm['TaxInc_Good']
```

```
In [301]: 1 from sklearn.model_selection import train_test_split
```

```
In [302]: 1 # Splitting data into train & test
          2 Xtrain, Xtest, ytrain, ytest = train_test_split(X, y, test_size=0.2, random_state=0)
```

```
In [303]: 1 # converting the Taxable income variable to bucketing.
          2 df_norm["income"] = "<=30000"
          3 df_norm.loc[df["Taxable.Income"] >= 30000, "income"] = "Good"
          4 df_norm.loc[df["Taxable.Income"] <= 30000, "income"] = "Risky"
```

```
In [304]: 1 #Dropping the Taxable income variables
          2 df.drop(["Taxable.Income"], axis=1, inplace=True)
```

```
In [305]: 1 df.rename(columns={"Undergrad": "undergrad", "Marital.Status": "marital", "City.Population": "population", "Work.Experience": "expe
          2 # As we are getting error as "Value Error: could not convert string to float: 'Yes'"
          3 # Model.fit doesn't not consider string. So, we encoder
```

```
Out[305]:
```

|     | population | experience | Undergrad_YES | Marital.Status_Married | Marital.Status_Single | Urban_YES | TaxInc_Good |
|-----|------------|------------|---------------|------------------------|-----------------------|-----------|-------------|
| 0   | 50047      | 10         | 0             | 0                      | 1                     | 1         | 1           |
| 1   | 134075     | 18         | 1             | 0                      | 0                     | 1         | 1           |
| 2   | 160205     | 30         | 0             | 1                      | 0                     | 1         | 1           |
| 3   | 193264     | 15         | 1             | 0                      | 1                     | 1         | 1           |
| 4   | 27533      | 28         | 0             | 1                      | 0                     | 0         | 1           |
| ... | ...        | ...        | ...           | ...                    | ...                   | ...       | ...         |
| 595 | 39492      | 7          | 1             | 0                      | 0                     | 1         | 1           |
| 596 | 55369      | 2          | 1             | 0                      | 0                     | 1         | 1           |
| 597 | 154058     | 0          | 0             | 0                      | 0                     | 1         | 1           |
| 598 | 180083     | 17         | 1             | 1                      | 0                     | 0         | 1           |
| 599 | 158137     | 16         | 0             | 0                      | 0                     | 0         | 1           |

600 rows × 7 columns

```
In [306]: 1 from sklearn import preprocessing
          2 le=preprocessing.LabelEncoder()
          3 for column_name in df.columns:
          4     if df[column_name].dtype == object:
          5         df[column_name] = le.fit_transform(df[column_name])
          6     else:
          7         pass
```

```
In [307]: 1 # Splitting the data into feature and labels
          2 features = df.iloc[:,0:5]
          3 labels = df.iloc[:,5]
```

```
In [308]: 1 # Collecting the columns names
          2 colnames = list(df.columns)
          3 predictors = colnames[0:5]
          4 target = colnames[5]
          5 ## Splitting the data into train and test
```

```
In [309]: 1 from sklearn.model_selection import train_test_split
          2 x_train,x_test,y_train,y_test = train_test_split(features,labels,test_size = 0.2,stratify = labels)
```

```
In [310]: 1 # Model Building
          2 from sklearn.ensemble import RandomForestClassifier as RF
          3 model = RF(n_jobs = 3,n_estimators = 15, oob_score = True, criterion = "entropy")
          4 model.fit(x_train,y_train)
```

```
Out[310]: RandomForestClassifier(criterion='entropy', n_estimators=15, n_jobs=3,
                                oob_score=True)
```

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.  
On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

```
In [311]: 1 import warnings
          2 warnings.filterwarnings('ignore')
```

```
In [312]: 1 model.estimators_
          2 model.classes_
          3 model.n_features_
          4 model.n_classes_
```

```
Out[312]: 2
```

```
In [313]: 1 model.n_outputs_
```

```
Out[313]: 1
```

```
In [314]: 1 model.oob_score_
          2 ##74.7833%
```

```
Out[314]: 0.55
```

```
In [315]: 1 ## Predictions on train data
          2 prediction = model.predict(x_train)
```

```
In [316]: 1 ## Accuracy
          2 ## For Accuracy
          3 from sklearn.metrics import accuracy_score
          4 accuracy = accuracy_score(y_train,prediction)
          5 ## 98.33%
```

```
In [317]: 1 np.mean(prediction == y_train)
          2 ## 98.33%
```

```
Out[317]: 0.9958333333333333
```

```
In [318]: 1 # Confusion Matrix
          2 from sklearn.metrics import confusion_matrix
          3 confusion = confusion_matrix(y_train,prediction)
```

```
In [319]: 1 # predictions on test data
          2 pred_test = model.predict(x_test)
```

```
In [320]: 1 # Accuracy
          2 acc_test = accuracy_score(y_test,pred_test)
          3 ## 78.33%
```

```
In [321]: 1 !pip install pydotplus
```

Requirement already satisfied: pydotplus in c:\users\admin\anaconda3\lib\site-packages (2.0.2)  
Requirement already satisfied: pyparsing>=2.0.1 in c:\users\admin\anaconda3\lib\site-packages (from pydotplus) (3.0.4)

```
In [322]: 1 # in random forest we can plot a decision tree present in Random Forest
2 from sklearn.tree import export_graphviz
3 import pydotplus
4 from six import StringIO
```

```
In [323]: 1 tree = model.estimators_[5]
```

```
In [324]: 1 dot_data = StringIO()
2 export_graphviz(tree,out_file = dot_data, filled = True, rounded = True, feature_names = predictors, class_names = target)
```

```
In [325]: 1 graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
```

## Building Decision Tree Classifier Using Entropy Criteria

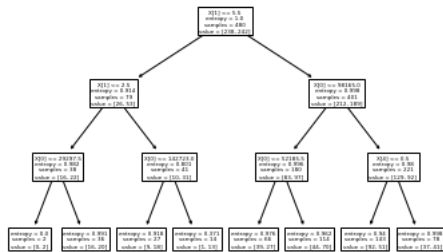
```
In [326]: 1 model = DecisionTreeClassifier(criterion = 'entropy',max_depth=3)
2 model.fit(x_train,y_train)
```

Out[326]: DecisionTreeClassifier(criterion='entropy', max\_depth=3)

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.  
On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

```
In [327]: 1 from sklearn import tree
```

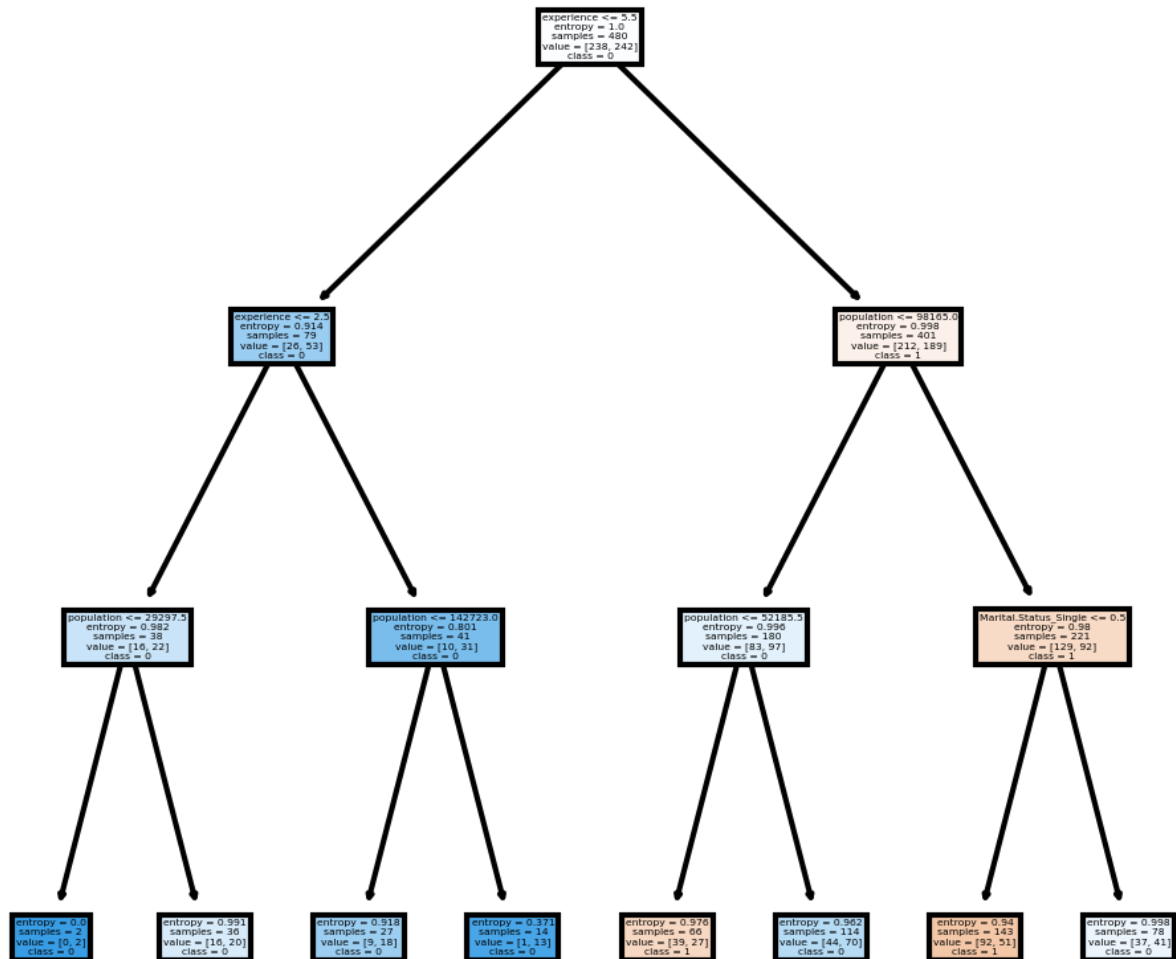
```
In [328]: 1 # Plot the import tree
2 tree.plot_tree(model);
```



```
In [329]: 1 colnames = list(df.columns)
2 colnames
```

Out[329]: ['City.Population',  
'Work.Experience',  
'Undergrad\_YES',  
'Marital.Status\_Married',  
'Marital.Status\_Single',  
'Urban\_YES',  
'TaxInc\_Good']

```
In [330]: 1 fn=['population','experience','Undergrad_YES','Marital.Status_Married','Marital.Status_Single','Urban_YES']
2 cn=['1', '0']
3 fig, axes = plt.subplots(nrows = 1,ncols = 1,figsize = (4,4), dpi=300)
4 tree.plot_tree(model,
5                 feature_names = fn,
6                 class_names=cn,
7                 filled = True);
```



```
In [331]: 1 # Predicting on test data
2 preds = model.predict(x_test) # Predicting on test dataset
3 pd.Series(preds).value_counts() # getting the count of each category
```

```
Out[331]: 1 73
0 47
dtype: int64
```

```
In [332]: 1 preds
```

```
Out[332]: array([1, 1, 1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0,
0, 1, 1, 0, 1, 1, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 0, 1, 1,
1, 1, 1, 1, 0, 0, 0, 0, 1, 1, 1, 0, 1, 0, 0, 0, 1, 1, 1, 1, 1, 1,
1, 1, 0, 1, 0, 1, 1, 0, 0, 1, 0, 1, 0, 1, 0, 0, 0, 1, 1, 1, 0, 1, 1,
1, 1, 0, 1, 1, 0, 0, 1, 1, 1, 1, 0, 1, 0, 0, 1, 1, 0, 0, 0, 0,
1, 1, 1, 0, 0, 1, 0, 1, 1, 1], dtype=uint8)
```

```
In [333]: 1 pd.crosstab(y_test,preds) # Getting the 2 way table to understand the correct and wrong predictions
```

```
Out[333]:      col_0  0  1
      Urban_YES
      0  23  37
      1  24  36
```

```
In [334]: 1 # Accuracy
          2 np.mean(preds==y_test)
```

```
Out[334]: 0.49166666666666664
```

## Building Decision Tree Classifier (CART) using Gini Criteria

```
In [335]: 1 from sklearn.tree import DecisionTreeClassifier
          2 model_gini = DecisionTreeClassifier(criterion='gini', max_depth=3)
```

```
In [336]: 1 model_gini.fit(x_train, y_train)
```

```
Out[336]: DecisionTreeClassifier(max_depth=3)
```

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.  
On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

```
In [337]: 1 # Prediction and computing the accuracy
          2 pred=model.predict(x_test)
          3 np.mean(preds==y_test)
```

```
Out[337]: 0.49166666666666664
```

## Decision Tree Regression Example

```
In [338]: 1 # Decision Tree Regression
          2 from sklearn.tree import DecisionTreeRegressor
```

```
In [339]: 1 array = df.values
          2 X = array[:,0:3]
          3 y = array[:,3]
```

```
In [340]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, random_state=1)
```

```
In [341]: 1 model = DecisionTreeRegressor()
          2 model.fit(X_train, y_train)
```

```
Out[341]: DecisionTreeRegressor()
```

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.  
On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

```
In [342]: 1 # Find the accuracy
          2 model.score(X_test,y_test)
```

```
Out[342]: -0.8931902985074629
```

```
In [ ]: 1
```