In [1]:	Decision Trees Assignment  Data Set - Fraud_check  1. Import Necessary libraries  import pandas as pd import numpy as np
In [2]:	<pre>from matplotlib import pyplot as plt import seaborn as sns  from sklearn import datasets from sklearn.metrics import classification_report import warnings warnings.filterwarnings('ignore')  2. Import Data  fraud_details = pd.read_csv('Fraud_check.csv') fraud_details</pre>
Out[2]:	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                   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           600 rows × 6 columns    3. Data Understanding
<pre>In [3]: Out[3]:</pre>	3.1 Initial Analysis:  fraud_details.head()  Undergrad Marital.Status Taxable.Income City.Population Work.Experience Urban  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
<pre>In [4]: Out[4]: In [5]:</pre>	<pre>fraud_details.info()  <class 'pandas.core.frame.dataframe'=""> RangeIndex: 600 entries, 0 to 599 Data columns (total 6 columns): # Column Non-Null Count Dtype 0 Undergrad 600 non-null object</class></pre>
<pre>In [6]: Out[6]:</pre>	1 Marital.Status 600 non-null object 2 Taxable.Income 600 non-null int64 3 City.Population 600 non-null int64 4 Work.Experience 600 non-null int64 5 Urban 600 non-null object dtypes: int64(3), object(3) memory usage: 28.2+ KB  fraud_details.isna().sum()  Undergrad 0 Marital.Status 0 Taxable.Income 0 City.Population 0 Work.Experience 0
<pre>In [7]: Out[7]:</pre>	Urban dtype: int64  fraud_details.describe()  Taxable.Income City.Population Work.Experience  count 600.000000 600.000000 600.000000  mean 55208.375000 108747.368333 15.558333  std 26204.827597 49850.075134 8.842147  min 10003.000000 25779.000000 0.0000000  25% 32871.500000 66966.750000 8.000000
In [8]: Out[8]:	50%       55074.500000       106493.500000       15.000000         75%       78611.750000       150114.250000       24.000000         max       99619.00000       199778.00000       30.000000                 Fraud_details.dtypes               Undergrad
<pre>In [9]: Out[9]: In [10]: Out[10]:</pre>	<pre>fraud_details.columns  Index(['Undergrad', 'Marital.Status', 'Taxable.Income', 'City.Population',</pre>
In [11]:	3.2 Correlation Matrix:  plt.figure(figsize = (12,8)) sns.heatmap(fraud_details.corr(), annot = True) plt.show()  1 -0.064 -0.0018 -0.8
	- 0.064 1 0.013 - 0.4  - 0.0018 0.0013 1 - 0.2
In [12]: In [13]:	3.3 Label Encoder:  from sklearn import preprocessing  label_encoder = preprocessing.LabelEncoder() label_encoder
out[13]: in [14]: in [15]: out[15]:	LabelEncoder()  fraud_details['Undergrad'] = label_encoder.fit_transform(fraud_details['Undergrad']) fraud_details['Marital.Status'] = label_encoder.fit_transform(fraud_details['Marital.Status']) fraud_details['Urban'] = label_encoder.fit_transform(fraud_details['Urban'])  fraud_details  Undergrad Marital.Status Taxable.Income City.Population Work.Experience Urban  0 0 2 68833 50047 10 1  1 1 0 33700 134075 18 1
	2         0         1         36925         160205         30         1           3         1         2         50190         193264         15         1           4         0         1         81002         27533         28         0                   595         1         0         76340         39492         7         1           596         1         0         69967         55369         2         1           597         0         0         47334         154058         0         1           598         1         1         98592         180083         17         0           599         0         0         96519         158137         16         0
	Sns.pairplot(fraud_details) plt.show()  10
	0.2
	80000 60000 20000 175000 8150000
	100000 250000 2500000 25000 25000 25000 25000 25000 25000 25000 25000 25000 25000 25
	0.0
n [17]: ut[17]:	3.4 Adding New Column:  fraud_details['Status'] = fraud_details['Taxable.Income'].apply(lambda Income: 'Risky' if Income <= 30000 else 'Good') fraud_details    Undergrad   Marital.Status   Taxable.Income   City.Population   Work.Experience   Urban   Status
	4         0         1         81002         27533         28         0         Good                     595         1         0         76340         39492         7         1         Good           596         1         0         69967         55369         2         1         Good           597         0         0         47334         154058         0         1         Good           598         1         1         98592         180083         17         0         Good           599         0         0         96519         158137         16         0         Good           600 rows × 7 columns
[18]: t[18]: [19]: t[19]:	<pre>fraud_details['Status'] = label_encoder.fit_transform(fraud_details['Status']) fraud_details</pre>
	3 1 2 50190 193264 15 1 0 4 0 1 81002 27533 28 0 0
[20]: [21]: [22]: t[22]:	4. Train Test Split  from sklearn.model_selection import train_test_split  x = fraud_details.iloc[:,0:4] y = fraud_details['Status']  x  Undergrad Marital.Status Taxable.Income City.Population
	0         0         2         68833         50047           1         1         0         33700         134075           2         0         1         36925         160205           3         1         2         50190         193264           4         0         1         81002         27533                  595         1         0         76340         39492           596         1         0         69967         55369           597         0         0         47334         154058
n [23]: ut[23]:	
[24]: t[24]: [25]: t[25]:	<pre>fraud_details.Status.value_counts()</pre>
n [26]: nt[26]: n [27]:	<pre>Name: Status, dtype: int64  list(fraud_details.columns)  ['Undergrad',     'Marital.Status',     'Taxable.Income',     'City.Population',     'Work.Experience',     'Urban',     'Status']  x_train, x_test,y_train,y_test = train_test_split(x,y, test_size = 0.2 , random_state = 40)</pre>
n [28]: n [29]: nt[29]: n [30]:	5. Building Model for Decision Tree Classifier using Entropy Criteria  from sklearn.tree import pecisionTreeClassifier from sklearn import tree  model = DecisionTreeClassifier(criterion = 'entropy', max_depth = 3) model.fit(x_train,y_train)  DecisionTreeClassifier(criterion='entropy', max_depth=3)  tree.plot_tree(model);
	X[2] <= 29949.5 entropy = 0.722 samples = 480 value = [384, 96]  entropy = 0.0 samples = 96 value = [0, 96]  entropy = 0.0 samples = 384 value = [384, 0]
	<pre>f_n = ['Undergrad', 'Marital.Status', 'Taxable.Income', 'City.Population'] c_n = ['Good', 'Risky']  fig, axes = plt.subplots(nrows = 1,ncols = 1,figsize = (2.5,2.5), dpi = 300) tree.plot_tree(model, feature_names = f_n,class_names = c_n, filled = True);  Taxable.Income &lt;= 29949.5</pre>
	entropy = 0.722 samples = 480 value = [384, 96] class = Good
	entropy = 0.0 samples = 96 value = [0, 96] class = Risky  entropy = 0.0 samples = 384 value = [384, 0] class = Good
[33]: t[33]: [34]:	
	array([0, 0, 0, 0, 1, 1, 0, 1, 1, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
[36]: t[36]: [37]: t[37]:	5.2 Feature Importance  model.feature_importances_ array([0., 0., 1., 0.])
[38]: : :[38]:	<pre>feature_imp = pd.Series(model.feature_importances_,index = f_n).sort_values(ascending = False) feature_imp  Taxable.Income    1.0 Undergrad    0.0 Marital.Status    0.0 City.Population    0.0 dtype: float64  sns.barplot(x = feature_imp, y = feature_imp.index) plt.xlabel('Feature Importance Score') plt.ylabel('Features') plt.title("Visualizing Important Features") plt.show()</pre>
	Visualizing Important Features  Taxable.Income -  Undergrad -  City.Population -
[40]: t[40]:	6. Building Model For Decision Tree Classifier (CART) using Gini Criteria  model_1 = DecisionTreeClassifier(criterion = 'gini', max_depth = 3) model_1.fit(x_train, y_train)  DecisionTreeClassifier(max_depth=3)
[41]:	<pre>tree.plot_tree(model_1);  X[2] &lt;= 29949.5     gini = 0.32     samples = 480     value = [384, 96]  gini = 0.0     samples = 96     value = [0, 96]</pre> gini = 0.0 samples = 384 value = [384, 0]
	Taxable.Income <= 29949.5 gini = 0.32 samples = 480 value = [384, 96] class = Good
	gini = 0.0 samples = 96 value = [0, 96] class = Risky  gini = 0.0 samples = 384 value = [384, 0] class = Good
[44]: t[44]:	<pre>pred_1 = model_1.predict(x_test) pd.Series(pred_1).value_counts()  0</pre>
	<pre>pred_1 array([0, 0, 0, 0, 1, 1, 0, 1, 1, 0, 0, 1, 1, 0, 1, 1, 0, 0, 0, 0, 0, 1,</pre>
[47]: t[47]: [48]:	<pre>0 92 0 1 0 28  np.mean(pred_1 == y_test)</pre>
[48]: t[48]: [49]: t[49]:	<pre>array([0., 0., 1., 0.])  feature_imp_1 = pd.Series(model_1.feature_importances_,index = f_n).sort_values(ascending = False) feature_imp_1  Taxable.Income</pre>
	plt.xlabel('Features') plt.title("Visualizing Important Features") plt.show()  Visualizing Important Features  Undergrad  Marital.Status
[51]: [52]:	7. Decision Tree Regression  from sklearn.tree import DecisionTreeRegressor
[52]: n [53]: n [54]: n [54]:	<pre>model_3.score(X_test,y_test)</pre>
out[55]:	Conclusion:  Model Accuracy for Entropy Criteria: 1.0 Model Accuracy for Gini Criteria: 1.0 Model Accuracy for Decision Tree Regression: 0.999