

In [1]:

```
import pandas as pd
import seaborn as sns
```

executed in 12.7s, finished 19:09:21 2022-01-05

In [2]:

```
data_set=pd.read_csv('Fraud_check.csv')
data_set
```

executed in 192ms, finished 19:09:28 2022-01-05

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

Initial investigation

In [3]:

```
data_set.shape
```

Out[3]:

(600, 6)

In [4]:

```
data_set.dtypes
```

Out[4]:

```
Undergrad      object
Marital.Status  object
Taxable.Income  int64
City.Population int64
Work.Experience int64
Urban          object
dtype: object
```

In [5]:

```
data_set.isnull().sum()
```

Out[5]:

```
Undergrad      0
Marital.Status  0
Taxable.Income  0
City.Population 0
Work.Experience 0
Urban          0
dtype: int64
```

Number of features and records in the given data set is 6 and 600 respectively

There is no null values in the data set

The categorical data can be converted into numeric data type by using encoder so that the model can learn the things more easily

Data preprocessing

In [6]:

```
from sklearn.preprocessing import LabelEncoder
le=LabelEncoder()
```

In [7]:

```
data_set['Undergrad']=le.fit_transform(data_set['Undergrad'])
data_set['Marital.Status']=le.fit_transform(data_set['Marital.Status'])
data_set['Urban']=le.fit_transform(data_set['Urban'])
data_set.dtypes
```

Out[7]:

```
Undergrad      int32
Marital.Status  int32
Taxable.Income  int64
City.Population int64
Work.Experience int64
Urban          int32
dtype: object
```

In [8]:

```
data_set.insert(6, 'tax_category', '')
data_set
```

Out[8]:

	Undergrad	Marital.Status	Taxable.Income	City.Population	Work.Experience	Urban	tax_cat
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	

600 rows × 7 columns

In [9]:

```
import warnings
warnings.filterwarnings('ignore')
```

Converting taxable income to category of 0 and 1

In [10]:

```
for i in range(0, len(data_set['tax_category']), 1):
    if data_set['Taxable.Income'][i] <= 30000:
        data_set['tax_category'][i] = '0'
    else:
        data_set['tax_category'][i] = '1'
```

In [11]:

```
data_set['tax_category'].unique()
```

Out[11]:

```
array(['1', '0'], dtype=object)
```

Model building

In [12]:

```
x=data_set.loc[:,('Undergrad','Marital.Status','City.Population','Work.Experience','Urban')]
y=data_set['tax_category']
```

In [13]:

```
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)
```

Model training

In [14]:

```
from sklearn.ensemble import RandomForestClassifier
rf_model=RandomForestClassifier().fit(x_train,y_train)
```

Model testing

In [15]:

```
y_pred_train=rf_model.predict(x_train)
y_pred_test=rf_model.predict(x_test)
```

Model evaluation

In [16]:

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

In [17]:

```
print(classification_report(y_test,y_pred_test))
```

	precision	recall	f1-score	support
0	0.10	0.04	0.06	25
1	0.78	0.91	0.84	95
accuracy			0.73	120
macro avg	0.44	0.47	0.45	120
weighted avg	0.64	0.72	0.68	120

In [18]:

```
print(accuracy_score(y_test,y_pred_test))
```

0.725

In [19]:

```
print(confusion_matrix(y_test,y_pred_test))
```

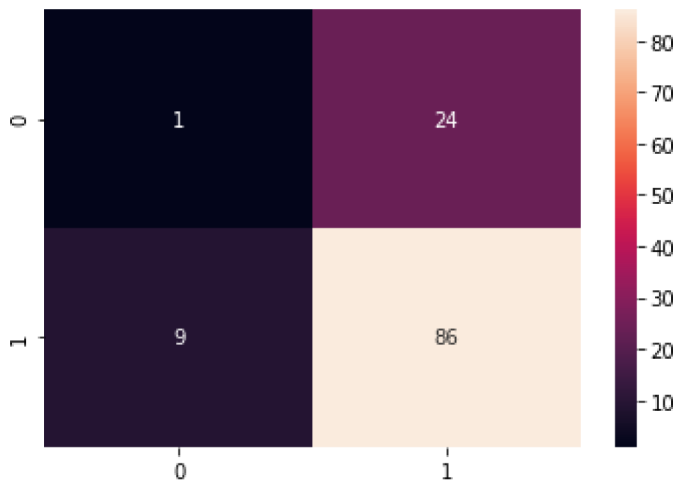
```
[[ 1 24]
 [ 9 86]]
```

In [23]:

```
confusion_matrix_test=confusion_matrix(y_test,y_pred_test)
sns.heatmap(confusion_matrix_test,annot=True)
```

Out[23]:

<AxesSubplot:>



In [24]:

```
auc_test= roc_auc_score(y_test, y_pred_test)
print('auc value for test data',auc_test)
```

```
auc value for test data 0.47263157894736846
```

In [25]:

```
from sklearn.model_selection import GridSearchCV
grid_model=GridSearchCV(estimator = rf_model,param_grid={'criterion':['entropy','gini'],
                                                           'max_depth':[2,4,8,10],
                                                           'min_samples_split':[2,4,6,8],
                                                           'min_samples_leaf':[1,2,3,4],
                                                           'n_estimators' : [20,50,70,100,150]}

grid_model.fit(x_train,y_train)
print(grid_model.best_params_)
print(grid_model.best_score_)
```

```
{'criterion': 'gini', 'max_depth': 10, 'min_samples_leaf': 1, 'min_samples_s
plit': 6, 'n_estimators': 150}
0.8
```

In [26]:

```
rf_model_tweak=RandomForestClassifier( n_estimators=150,min_samples_split=6,max_depth=10,mi
```

In [27]:

```
y_pred_test_tweak=rf_model_tweak.predict(x_test)
```

In [28]:

```
print(classification_report(y_test,y_pred_test_tweak))
```

	precision	recall	f1-score	support
0	0.33	0.04	0.07	25
1	0.79	0.98	0.88	95
accuracy			0.78	120
macro avg	0.56	0.51	0.47	120
weighted avg	0.70	0.78	0.71	120

In [29]:

```
print(accuracy_score(y_test,y_pred_test_tweak))
```

```
0.7833333333333333
```

In [30]:

```
auc_test= roc_auc_score(y_test, y_pred_test_tweak)
```

```
print('auc value for test data',auc_test)
```

```
auc value for test data 0.5094736842105263
```

In [31]:

```
confusion_matrix_test=confusion_matrix(y_test,y_pred_test_tweak)
sns.heatmap(confusion_matrix_test,annot=True)
```

Out[31]:

<AxesSubplot:>

