

```
In [69]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt

In [70]: df.head()

Out[70]:
   State-gov  7736  Bachelors  12  Never-married  Adm-clerical  Non-in-kennly  White  Male  2174  0  40  United-States  <DOK>
0  50  Self-emp-inc  10231  Bachelors  15  Married-cv-spouse  Exec-mgmt  Husband  White  Male  0  0  13  United-States  <DOK>
1  38  Private  22666  HS-grad  9  Divorced  Handwr-clerical  Husband  White  Male  0  0  40  United-States  <DOK>
2  53  Private  22471  11th  7  Married-cv-spouse  Handwr-clerical  Husband  Black  Male  0  0  40  United-States  <DOK>
3  28  Private  20840  Bachelors  12  Married-cv-spouse  Prof-specalty  Wife  Black  Female  0  0  40  Cuba  <DOK>
4  37  Private  20602  Masters  14  Married-cv-spouse  Exec-mgmt  Wife  White  Female  0  0  40  United-States  <DOK>

In [71]: df.shape
(2569, 15)

Out[71]:
Renaming the Columns.

In [72]: df.columns = ['Age', 'Workclass', 'Frtnght', 'Education', 'education_num', 'marital_status', 'occupation', 'relationship', 'race', 'sex', 'capital_gain', 'capital_loss', 'hours_per_week', 'native_country', 'income']

Out[72]:
df.columns = ['Age', 'Workclass', 'Frtnght', 'Education', 'education_num', 'marital_status', 'occupation', 'relationship', 'race', 'sex', 'capital_gain', 'capital_loss', 'hours_per_week', 'native_country', 'income']
dtype: object

In [73]: df.columns

Out[73]:
Index(['Age', 'Workclass', 'Frtnght', 'Education', 'education_num', 'marital_status', 'occupation', 'relationship', 'race', 'sex', 'capital_gain', 'capital_loss', 'hours_per_week', 'native_country', 'income'],
      dtype='object')

In [74]: len(df.columns)

Out[74]:
15

In [75]: df.dtypes

Out[75]:
Age                int64
Workclass          object
Frtnght            int64
Education          object
education_num      int64
marital_status     object
occupation         object
relationship       object
race              object
sex               object
capital_gain       int64
capital_loss       int64
hours_per_week     int64
native_country     object
income            object
dtype: object

Checking whether the dataset has Null Values.

In [76]: df.isnull().sum()

Out[76]:
Age                0
Workclass          0
Frtnght            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

Note:
Since there is no null values so that we can move to other steps.

Converting all 'object' parameters into 'integer' values.

In [124]: from sklearn.preprocessing import LabelEncoder

In [78]: lb = LabelEncoder()
df['Workclass'] = lb.fit_transform(df['Workclass'])
df['Education'] = lb.fit_transform(df['Education'])
df['marital_status'] = lb.fit_transform(df['marital_status'])
df['relationship'] = lb.fit_transform(df['relationship'])
df['race'] = lb.fit_transform(df['race'])
df['sex'] = lb.fit_transform(df['sex'])
df['native_country'] = lb.fit_transform(df['native_country'])
df['income'] = lb.fit_transform(df['income'])

In [79]: print(df['Workclass'].value_counts())
print(df['Education'].value_counts())
print(df['marital_status'].value_counts())
print(df['relationship'].value_counts())
print(df['race'].value_counts())
print(df['sex'].value_counts())
print(df['native_country'].value_counts())
print(df['occupation'].value_counts())
print(df['income'].value_counts())

4    22696
2    2993
3    1597
7    1759
5    969
1     34
Name: Workclass, dtype: int64
11    10591
15    7291
9     5354
12    1723
6     1392
8     1151
9     933
6     646
14    576
2     433
19    413
4     333
1     168
13     51
Name: Education, dtype: int64
2    14970
4    58652
6     442
5     1825
10    1119
3     428
Name: marital_status, dtype: int64
4     8364
5     5869
4     3446
3     3589
2     681
Name: relationship, dtype: int64
4    27915
2    12154
9     311
1     181
Name: race, dtype: int64
0    18775
1     584
Name: sex, dtype: int64
39    29169
16    613
6     593
19    198
21    127
13    111
39    114
9     106
5     95
6     90
23    81
16    60
2     75
12    71
6     70
40    47
13    64
14    62
21    60
6     59
16    41
10    41
39    34
32    37
27    34
29    31
18    29
12    29
7     29
21    24
27    20
19    19
38    19
13    19
37    18
41    18
28    14
18    13
18    13
16    12
15     1
Name: native_country, dtype: int64
19    4169
3     4090
4     4065
12    3959
9     3295
6     2897
0     1843
4     1379
5     954
13    928
11    949
9     149
2     6
Name: occupation, dtype: int64
0    28718
2     7841
Name: income, dtype: int64

In [80]: df.dtypes

Out[80]:
Age                int64
Workclass          int32
Frtnght            int64
Education          int32
education_num      int64
marital_status     int32
occupation         int32
relationship       int32
race              int32
sex               int32
capital_gain       int64
capital_loss       int64
hours_per_week     int64
native_country     int32
income            int32
dtype: object

Splitting the dataset into Training Data and Test Data.

In [81]: x = df.iloc[:,1:13]
y = df.iloc[:,14]
print(y.shape)
print(y.dtype)
print(y.shape)

<class 'pandas.core.frame.DataFrame'>
<class 'pandas.core.series.Series'>
(2569, 14)
(2569, 1)

In [82]: from sklearn.model_selection import train_test_split

In [83]: x_train,x_test,y_train,y_test = train_test_split(x,y,test_size= 0.25)
print(x_train.shape)
print(x_test.shape)
print(y_train.shape)
print(y_test.shape)

(2429, 14)
(140, 14)
(2429, 1)
(140, 1)

In [84]: print(x_train.shape)
print(y_train.shape)

(2429, 14)
(2569, 1)

In [85]: x_train.head()

Out[85]:
   Age  Workclass  Frtnght  Education  education_num  marital_status  occupation  relationship  race  sex  capital_gain  capital_loss  hours_per_week  native_country
0  44  50  10231  Bachelors  15  Married-cv-spouse  Exec-mgmt  Husband  White  Male  0  0  13  United-States
1  38  38  22666  HS-grad  9  Divorced  Handwr-clerical  Husband  White  Male  0  0  40  United-States
2  53  38  22471  11th  7  Married-cv-spouse  Handwr-clerical  Husband  Black  Male  0  0  40  United-States
3  28  38  20840  Bachelors  12  Married-cv-spouse  Prof-specalty  Wife  Black  Female  0  0  40  Cuba
4  37  38  20602  Masters  14  Married-cv-spouse  Exec-mgmt  Wife  White  Female  0  0  40  United-States

In [86]: y_train.head()

Out[86]:
0  28869  1
1  8778  9
2  2989  0
3  10008  0
4  1963  0
Name: income, dtype: int32

Different Training Models.
a. Decision Tree Classifier.

In [87]: from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import confusion_matrix,classification_report

In [88]: s1 = DecisionTreeClassifier(criterion = 'gini', max_depth = 15)
s1.fit(x_train,y_train)

Out[88]:
DecisionTreeClassifier(max_depth=15)

In [89]: print('Training Score', s1.score(x_train,y_train))
print('Testing Score', s1.score(x_test,y_test))

Training Score 0.91502865030905
Testing Score 0.832066339966329

Income Prediction for Test Data.

In [90]: ypredict_s1 = s1.predict(x_test)
print(ypredict_s1)

[0 0 0 ... 0 0 0]

Confusion Matrix and Classification Report.

In [127]: conf_mat = confusion_matrix(y_test,ypredict_s1)
print('Confusion Matrix :\n')
print(conf_mat)
print('\n')
print('Classification Report :\n')
print(classification_report(y_test,ypredict_s1))

Confusion Matrix :
[[593 129]
 [ 761 1275]]

Classification Report :

      precision    recall  f1-score   support

    0      0.88      0.90      0.89      6154
    1      1.00      0.97      0.98      1395

 accuracy      0.83      0.88      0.84
 macro avg      0.78      0.76      0.77      8149
weighted avg      0.83      0.89      0.83      8149

Validating the Result for Precision, Recall, F1-Score and Accuracy.
Precision = TP/(TP+FP), TN/(TN+FN)
Recall = TP/(TP+FN), TN/(TN+FP)
F1-Score = (2*Precision Recall)/(Precision + Recall)
Accuracy = (TP+FN)/(TP+FN+FP+TN)

In [118]: prec = 593/(593+129)
pre = 1255/(1255+129)
rec = 593/(593+129)
print('Precision 0 : ', prec)
print('\n')
print('Precision 1 : ', pre)
print('\n')
print('Recall 0 : ', rec)
print('\n')
print('Recall 1 : ', pre)
print('\n')
F0 = 2*prec*rec/(prec + rec)
F1 = 2*pre*pre/(pre + rec)
print('F1-score 0 : ', F0)
print('\n')
print('F1-score 1 : ', F1)
print('\n')
acc = (593+1255)/(593+129+1255+761)
print('Accuracy : ', acc)

Precision 0 : 0.87864879234058
Precision 1 : 0.6768483206303911

Recall 0 : 0.90394989871478
Recall 1 : 0.6225198412086112

F1-score 0 : 0.8916492827728802
F1-score 1 : 0.6486225208414294

Accuracy : 0.8305155238878091

Percentage of Misclassification.
Percentage of Misclassification = (FP+TN)/(TP+TN+FP+FN)

In [119]: per_0 = (761+593)/(593+129+1255+761)
print('Percentage of Misclassification : ', per_mis*100)

Percentage of Misclassification : 18.548347613219993

b. Random Forest Classifier.

In [94]: from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import confusion_matrix,classification_report

Testing Score 0.797665878658477

In [95]: s2 = RandomForestClassifier(n_estimators = 80, criterion = 'gini', max_depth = 15)
s2.fit(x_train,y_train)

Out[95]:
RandomForestClassifier(max_depth=15, n_estimators=80)

In [96]: print('Training Score', s2.score(x_train,y_train))
print('Testing Score', s2.score(x_test,y_test))

Training Score 0.916878262478224
Testing Score 0.860313530330320

Income Prediction for Test Data.

In [97]: ypredict_s2 = s2.predict(x_test)
print(ypredict_s2)

[0 0 0 ... 0 0 0]

Confusion Matrix and Classification Report.

In [98]: conf_mat = confusion_matrix(y_test,ypredict_s2)
print('Confusion Matrix :\n')
print(conf_mat)
print('\n')
print('Classification Report :\n')
print(classification_report(y_test,ypredict_s2))

Confusion Matrix :
[[596 148]
 [ 762 1224]]

Classification Report :

      precision    recall  f1-score   support

    0      0.88      0.94      0.91      6154
    1      1.00      0.78      0.89      1395

 accuracy      0.83      0.86      0.84
 macro avg      0.88      0.86      0.84
weighted avg      0.88      0.89      0.88      8149

Validating the Result for Precision, Recall, F1-Score and Accuracy.

In [117]: prec = 596/(596+148)
pre = 1224/(1224+148)
rec = 596/(596+148)
print('Precision 0 : ', prec)
print('\n')
print('Precision 1 : ', pre)
print('\n')
print('Recall 0 : ', rec)
print('\n')
print('Recall 1 : ', pre)
print('\n')
F0 = 2*prec*rec/(prec + rec)
F1 = 2*pre*pre/(pre + rec)
print('F1-score 0 : ', F0)
print('\n')
print('F1-score 1 : ', F1)
print('\n')
acc = (596+1224)/(596+148+1224+762)
print('Accuracy : ', acc)

Precision 0 : 0.805863895620505
Precision 1 : 0.7781389496941417

Recall 0 : 0.9423889177778056
Recall 1 : 0.6183141939397784

F1-score 0 : 0.8126449887343762
F1-score 1 : 0.687838611487698

Accuracy : 0.8035136121335358

Percentage of Misclassification.
per_mis = (762+549)/(596+148+1224+762)
print('Percentage of Misclassification : ', per_mis*100)

Percentage of Misclassification : 18.648648648648647

c. KNN Classifier.

In [101]: from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import confusion_matrix,classification_report

In [102]: s3 = KNeighborsClassifier(n_neighbors = 18)
s3.fit(x_train,y_train)

Out[102]:
KNeighborsClassifier(n_neighbors=18)

In [103]: print('Training Score', s3.score(x_train,y_train))
print('Testing Score', s3.score(x_test,y_test))

Training Score 0.89548354483055
Testing Score 0.798314863148664

Income Prediction for Test Data.

In [104]: ypredict_s3 = s3.predict(x_test)
print(ypredict_s3)

[0 0 0 ... 0 0 0]

Confusion Matrix and Classification Report.

In [105]: conf_mat = confusion_matrix(y_test,ypredict_s3)
print('Confusion Matrix :\n')
print(conf_mat)
print('\n')
print('Classification Report :\n')
print(classification_report(y_test,ypredict_s3))

Confusion Matrix :
[[599 181]
 [1555 431]]

Classification Report :

      precision    recall  f1-score   support

    0      0.81      0.95      0.87      6154
    1      1.00      0.67      0.79      1395

 accuracy      0.74      0.82      0.78      8149
 macro avg      0.74      0.82      0.78      8149
weighted avg      0.78      0.80      0.78      8149

Validating the Result for Precision, Recall, F1-Score and Accuracy.

In [116]: prec = 599/(599+181)
pre = 431/(431+181)
rec = 599/(599+181)
print('Precision 0 : ', prec)
print('\n')
print('Precision 1 : ', pre)
print('\n')
print('Recall 0 : ', rec)
print('\n')
print('Recall 1 : ', pre)
print('\n')
F0 = 2*prec*rec/(prec + rec)
F1 = 2*pre*pre/(pre + rec)
print('F1-score 0 : ', F0)
print('\n')
print('F1-score 1 : ', F1)
print('\n')
acc = (599+431)/(599+181+431+1555)
print('Accuracy : ', acc)

Precision 0 : 0.8666666666666666
Precision 1 : 0.7055623097453752

Recall 0 : 0.9423889177778056
Recall 1 : 0.2178913893756296

F1-score 0 : 0.8795058139534864
F1-score 1 : 0.3428634286342857

Accuracy : 0.798314863148664

Percentage of Misclassification.
per_mis = (1555+181)/(599+181+431+1555)
print('Percentage of Misclassification : ', per_mis*100)

Percentage of Misclassification : 20.389585889589587

d. Logistic Regression.

In [106]: from sklearn.linear_model import LogisticRegression
from sklearn.metrics import confusion_matrix,classification_report

In [109]: s4 = LogisticRegression(solver = 'liblinear')
s4.fit(x_train,y_train)

LogisticRegression(solver='liblinear')

Out[109]:
print('Training Score', s4.score(x_train,y_train))
print('Testing Score', s4.score(x_test,y_test))

Training Score 0.788885480548528
Testing Score 0.797665878658477

Income Prediction for Test Data.

In [111]: ypredict_s4 = s4.predict(x_test)
print(ypredict_s4)

[0 0 0 ... 0 0 0]

Confusion Matrix and Classification Report.

In [112]: conf_mat = confusion_matrix(y_test,ypredict_s4)
print('Confusion Matrix :\n')
print(conf_mat)
print('\n')
print('Classification Report :\n')
print(classification_report(y_test,ypredict_s4))

Confusion Matrix :
[[599 181]
 [1416 579]]

Classification Report :

      precision    recall  f1-score   support

    0      0.81      0.95      0.87      6154
    1      1.00      0.67      0.79      1395

 accuracy      0.74      0.82      0.78      8149
 macro avg      0.74      0.82      0.78      8149
weighted avg      0.78      0.80      0.78      8149

Validating the Result for Precision, Recall, F1-Score and Accuracy.

In [115]: prec = 592/(592+1391)
pre = 579/(579+1391)
rec = 592/(592+1391)
print('Precision 0 : ', prec)
print('\n')
print('Precision 1 : ', pre)
print('\n')
print('Recall 0 : ', rec)
print('\n')
print('Recall 1 : ', pre)
print('\n')
F0 = 2*prec*rec/(prec + rec)
F1 = 2*pre*pre/(pre + rec)
print('F1-score 0 : ', F0)
print('\n')
print('F1-score 1 : ', F1)
print('\n')
acc = (592+579)/(592+1391+579+1416)
print('Accuracy : ', acc)

Precision 0 : 0.808916791743464
Precision 1 : 0.6682686767676767

Recall 0 : 0.9423889177778056
Recall 1 : 0.2948978759135355

F1-score 0 : 0.3968688916576364
F1-score 1 : 0.798314863148664

Accuracy : 0.798314863148664

Percentage of Misclassification.
per_mis = (1391+1416)/(592+1391+579+1416)
print('Percentage of Misclassification : ', per_mis*100)

Percentage of Misclassification : 20.454545454545457

Model With Best Accuracy.
From the above Machine Learning Models,
Random Forest Classifier Model has an Accuracy of 0.86' and Misclassification of approximately 14%.
```