

In [1]:

```
import numpy as np
import pandas as pd
```

In [5]:

```
dataset = pd.read_csv("C:/Users/Sampath/Downloads/adult.csv")
dataset.head()
dataset.head()
```

Out[5]:

	39	State-gov	77516	Bachelors	13	Never-married	Adm-clerical	Not-in-family	White	Male	2174	0
0	50	Self-emp-not-inc	83311	Bachelors	13	Married-civ-spouse	Exec-managerial	Husband	White	Male	0	0
1	38	Private	215646	HS-grad	9	Divorced	Handlers-cleaners	Not-in-family	White	Male	0	0
2	53	Private	234721	11th	7	Married-civ-spouse	Handlers-cleaners	Husband	Black	Male	0	0
3	28	Private	338409	Bachelors	13	Married-civ-spouse	Prof-specialty	Wife	Black	Female	0	0
4	37	Private	284582	Masters	14	Married-civ-spouse	Exec-managerial	Wife	White	Female	0	0

In [6]:

```
colnames= ['Age', 'Workclass', 'Fnlwgt', 'Education', 'education_num', 'marital_status', 'occupation', 'relationship']
dataset = pd.read_csv("C:/Users/Sampath/Downloads/adult.csv", names=colnames, header=None)
dataset.head()
```

Out[6]:

	Age	Workclass	Fnlwgt	Education	education_num	marital_status	occupation	relationship
0	39	State-gov	77516	Bachelors	13	Never-married	Adm-clerical	Not-in-family
1	50	Self-emp-not-inc	83311	Bachelors	13	Married-civ-spouse	Exec-managerial	Husband
2	38	Private	215646	HS-grad	9	Divorced	Handlers-cleaners	Not-in-family
3	53	Private	234721	11th	7	Married-civ-spouse	Handlers-cleaners	Husband
4	28	Private	338409	Bachelors	13	Married-civ-spouse	Prof-specialty	Wife

In [7]:

```
print("Original features:\n", list(dataset.columns), "\n")
data_dummies = pd.get_dummies(dataset)
print("Features after get_dummies:\n", list(data_dummies.columns))
data_dummies.head()
```

Original features:

```
['Age', 'Workclass', 'Fnlwgt', 'Education', 'education_num', 'marital_status', 'occupation', 'relationship', 'race', 'sex', 'capital_gain', 'capital_loss', 'hours_per_week', 'native_country', 'income']
```

Features after get_dummies:

```
['Age', 'Fnlwgt', 'education_num', 'capital_gain', 'capital_loss', 'hours_per_week', 'Workclass_?', 'Workclass_Federal-gov', 'Workclass_Local-gov', 'Workclass_Never-worked', 'Workclass_Private', 'Workclass_Self-emp-inc', 'Workclass_Self-emp-not-inc', 'Workclass_State-gov', 'Workclass_Without-pay', 'Education_10th', 'Education_11th', 'Education_12th', 'Education_1st-4th', 'Education_5th-6th', 'Education_7th-8th', 'Education_9th', 'Education_Assoc-acdm', 'Education_Assoc-voc', 'Education_Bachelors', 'Education_Doctorate', 'Education_HS-grad', 'Education_Masters', 'Education_Preschool', 'Education_Prof-school', 'Education_Some-college', 'marital_status_Divorced', 'marital_status_Married-AF-spouse', 'marital_status_Married-civ-spouse', 'marital_status_Married-spouse-absent', 'marital_status_Never-married', 'marital_status_Separated', 'marital_status_Widowed', 'occupation_?', 'occupation_Adm-clerical', 'occupation_Armed-Forces', 'occupation_Craft-repair', 'occupation_Exec-managerial', 'occupation_Farming-fishing', 'occupation_Handlers-cleaners', 'occupation_Machine-op-inspct', 'occupation_Other-service', 'occupation_Priv-house-serv', 'occupation_Prof-specialty', 'occupation_Protective-serv', 'occupation_Sales', 'occupation_Tech-support', 'occupation_Transport-moving', 'relationship_Husband', 'relationship_Not-in-family', 'relationship_Other-relative', 'relationship_Own-child', 'relationship_Unmarried', 'relationship_Wife', 'race_Amer-Indian-Eskimo', 'race_Asian-Pac-Islander', 'race_Black', 'race_Other', 'race_White', 'sex_Female', 'sex_Male', 'native_country_?', 'native_country_Cambodia', 'native_country_Canada', 'native_country_China', 'native_country_Columbia', 'native_country_Cuba', 'native_country_Dominican-Republic', 'native_country_Ecuador', 'native_country_El-Salvador', 'native_country_England', 'native_country_France', 'native_country_Germany', 'native_country_Greece', 'native_country_Guatemala', 'native_country_Haiti', 'native_country_Holand-Netherlands', 'native_country_Honduras', 'native_country_Hong', 'native_country_Hungary', 'native_country_India', 'native_country_Iran', 'native_country_Ireland', 'native_country_Italy', 'native_country_Jamaica', 'native_country_Japan', 'native_country_Laos', 'native_country_Mexico', 'native_country_Nicaragua', 'native_country_Outlying-US(Guam-USVI-etc)', 'native_country_Peru', 'native_country_Philippines', 'native_country_Poland', 'native_country_Portugal', 'native_country_Puerto-Rico', 'native_country_Scotland', 'native_country_South', 'native_country_Taiwan', 'native_country_Thailand', 'native_country_Trinidad&Tobago', 'native_country_United-States', 'native_country_Vietnam', 'native_country_Yugoslavia', 'income_<=50K', 'income_>50K']
```

Out[7]:

	Age	Fnlwgt	education_num	capital_gain	capital_loss	hours_per_week	Workclass_?	W
0	39	77516	13	2174	0	40	0	
1	50	83311	13	0	0	13	0	

	Age	Fnlwgt	education_num	capital_gain	capital_loss	hours_per_week	Workclass	W
2	38	215646	9	0	0	40	0	
3	53	234721	7	0	0	40	0	
4	28	338409	13	0	0	40	0	

5 rows × 110 columns

In [8]:

```
dataset.info('income')
```

```
<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
```

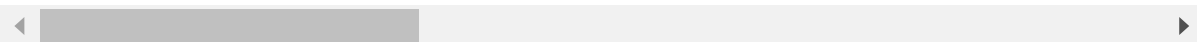
In [9]:

```
X = data_dummies.iloc[:, :-2].values
y = data_dummies.iloc[:, -1].values
data_dummies.head()
```

Out[9]:

	Age	Fnlwgt	education_num	capital_gain	capital_loss	hours_per_week	Workclass_?	Workc Fe
0	39	77516	13	2174	0	40	0	
1	50	83311	13	0	0	13	0	
2	38	215646	9	0	0	40	0	
3	53	234721	7	0	0	40	0	
4	28	338409	13	0	0	40	0	

5 rows × 110 columns



In [11]:

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25, random_state =
```

In [12]:

```
print(y_train)
```

```
[0 0 1 ... 0 1 0]
```

In [13]:

```
print(y_train)
```

```
[0 0 1 ... 0 1 0]
```

In [14]:

```
print(X_test)
```

```
[[ 27 177119 10 ... 1 0 0]
 [ 27 216481 13 ... 1 0 0]
 [ 25 256263 12 ... 1 0 0]
 ...
 [ 38 229236 9 ... 0 0 0]
 [ 63 213095 13 ... 1 0 0]
 [ 47 186009 10 ... 0 0 0]]
```

In [17]:

```
from sklearn.preprocessing import StandardScaler
sc = StandardScaler()
X_train = sc.fit_transform(X_train)
X_test = sc.transform(X_test)
```

In [18]:

```
print(X_train)
```

```
[[ 1.49894077 -1.2079589 -0.80338299 ...  0.33936755 -0.04296689
 -0.02307887]
 [ 2.38045121 -0.08366135  1.53114709 ...  0.33936755 -0.04296689
 -0.02307887]
 [ 0.2501343 -0.77396423  1.53114709 ...  0.33936755 -0.04296689
 -0.02307887]
 ...
 [-1.14559057 -1.36469155 -0.0252063 ...  0.33936755 -0.04296689
 -0.02307887]
 [ 0.47051192  0.24865765  2.30932378 ...  0.33936755 -0.04296689
 -0.02307887]
 [-0.99867216 -0.02508527 -0.0252063 ...  0.33936755 -0.04296689
 -0.02307887]]
```

In [19]:

```
print(X_test)
```

```
[[ -0.85175375 -0.11784965 -0.0252063 ...  0.33936755 -0.04296689
 -0.02307887]
 [ -0.85175375  0.25451337  1.14205874 ...  0.33936755 -0.04296689
 -0.02307887]
 [ -0.99867216  0.63084957  0.7529704 ...  0.33936755 -0.04296689
 -0.02307887]
 ...
 [ -0.04370251  0.37517518 -0.41429464 ... -2.94665767 -0.04296689
 -0.02307887]
 [ 1.79277758  0.22248194  1.14205874 ...  0.33936755 -0.04296689
 -0.02307887]
 [ 0.61743032 -0.03375059 -0.0252063 ... -2.94665767 -0.04296689
 -0.02307887]]
```

DECISION TREE

In [20]:

```
from sklearn.tree import DecisionTreeClassifier
classifier = DecisionTreeClassifier(criterion = 'entropy', random_state = 0)
classifier.fit(X_train, y_train)
```

Out[20]:

```
DecisionTreeClassifier(criterion='entropy', random_state=0)
```

In [21]:

```
print(classifier.predict(sc.transform(X_train)))
```

```
[0 0 1 ... 0 1 0]
```

In [22]:

```
y_pred = classifier.predict(X_test)
print(np.concatenate((y_pred.reshape(len(y_pred),1), y_test.reshape(len(y_test),1)),1))
```

```
[[0 0]
 [0 0]
 [0 0]
 ...
 [0 0]
 [1 1]
 [0 0]]
```

In [23]:

```
from sklearn.metrics import confusion_matrix, accuracy_score, classification_report
cm = confusion_matrix(y_test, y_pred)
print(cm)
DT_accu=accuracy_score(y_test, y_pred)
print(classification_report(y_test, y_pred))
```

```
[[5412  747]
 [ 756 1226]]
```

	precision	recall	f1-score	support
0	0.88	0.88	0.88	6159
1	0.62	0.62	0.62	1982
accuracy			0.82	8141
macro avg	0.75	0.75	0.75	8141
weighted avg	0.82	0.82	0.82	8141

RANDOM FOREST CLASSIFIER

In [25]:

```
from sklearn.ensemble import RandomForestClassifier
classifier = RandomForestClassifier(n_estimators = 10, criterion = 'entropy', random_state
classifier.fit(X_train, y_train)
```

Out[25]:

```
RandomForestClassifier(criterion='entropy', n_estimators=10, random_state=0)
```

In [26]:

```
print(classifier.predict(sc.transform(X_train)))
```

```
[0 0 1 ... 0 1 0]
```

In [27]:

```
y_pred = classifier.predict(X_test)
print(np.concatenate((y_pred.reshape(len(y_pred),1), y_test.reshape(len(y_test),1)),1))
```

```
[[0 0]
 [0 0]
 [0 0]
 ...
 [0 0]
 [1 1]
 [0 0]]
```

In [28]:

```
from sklearn.metrics import confusion_matrix, accuracy_score, classification_report
cm = confusion_matrix(y_test, y_pred)
print(cm)
RF_accu=accuracy_score(y_test, y_pred)
print(classification_report(y_test, y_pred))
```

```
[[5741  418]
 [ 852 1130]]
```

	precision	recall	f1-score	support
0	0.87	0.93	0.90	6159
1	0.73	0.57	0.64	1982
accuracy			0.84	8141
macro avg	0.80	0.75	0.77	8141
weighted avg	0.84	0.84	0.84	8141

LOGISTIC REGRESSION

In [29]:

```
from sklearn.linear_model import LogisticRegression
classifier = LogisticRegression(random_state = 0)
classifier.fit(X_train, y_train)
```

Out[29]:

```
LogisticRegression(random_state=0)
```

In [30]:

```
print(classifier.predict(sc.transform(X_train)))
```

```
[0 0 1 ... 0 1 0]
```

In [31]:

```
y_pred = classifier.predict(X_test)
print(np.concatenate((y_pred.reshape(len(y_pred),1), y_test.reshape(len(y_test),1)),1))
```

```
[[0 0]
 [0 0]
 [0 0]
 ...
 [0 0]
 [1 1]
 [0 0]]
```

In [32]:

```
from sklearn.metrics import confusion_matrix, accuracy_score, classification_report
cm = confusion_matrix(y_test, y_pred)
print(cm)
LR_accu=accuracy_score(y_test, y_pred)
print(classification_report(y_test, y_pred))
```

```
[[5719  440]
 [ 779 1203]]
```

	precision	recall	f1-score	support
0	0.88	0.93	0.90	6159
1	0.73	0.61	0.66	1982
accuracy			0.85	8141
macro avg	0.81	0.77	0.78	8141
weighted avg	0.84	0.85	0.85	8141

KNN CLASSIFIER

In [33]:

```
from sklearn.neighbors import KNeighborsClassifier
classifier = KNeighborsClassifier(n_neighbors = 5, metric = 'minkowski', p = 2)
classifier.fit(X_train, y_train)
```

Out[33]:

```
KNeighborsClassifier()
```

In [36]:

```
print(classifier.predict(sc.transform(X_train)))
```

```
[0 0 1 ... 0 1 0]
```


In [38]:

```
y_pred = classifier.predict(X_test)
print(np.concatenate((y_pred.reshape(len(y_pred),1), y_test.reshape(len(y_test),1)),1))

[[0 0]
 [0 0]
 [1 0]
 ...
 [0 0]
 [1 1]
 [0 0]]
```

In [39]:

```
from sklearn.metrics import confusion_matrix, accuracy_score, classification_report
cm = confusion_matrix(y_test, y_pred)
print(cm)
KNN_accu = accuracy_score(y_test, y_pred)
print(classification_report(y_test, y_pred))
```

```
[[5562  597]
 [ 854 1128]]
```

	precision	recall	f1-score	support
0	0.87	0.90	0.88	6159
1	0.65	0.57	0.61	1982
accuracy			0.82	8141
macro avg	0.76	0.74	0.75	8141
weighted avg	0.82	0.82	0.82	8141

SVC CLASSIFIER(with linear kernel)

In [40]:

```
# Importing SVM from sklearn
from sklearn import svm
```

In [41]:

```
# Creating the SVM model and training the model with X_train and y_train
SVM_m = svm.SVC()
SVM_m.fit(X_train, y_train)
```

Out[41]:

SVC()

In [42]:

```
# Using the X_test to predict the values
y_pred_SVM = SVM_m.predict(X_test)
```

In [43]:

```
# Calculating the Confusion Matrix and Accuracy_Score of the model
from sklearn.metrics import confusion_matrix, accuracy_score, classification_report
cm = confusion_matrix(y_test, y_pred_SVM)
print("Confusion Matrix = ", cm)
print(""*80)
# Classification Report
print(classification_report(y_test, y_pred_SVM))
print(""*80)
SVM_accu = accuracy_score(y_test, y_pred_SVM)
print("SVM Classifier Accuracy = ", SVM_accu)
```

```
Confusion Matrix = [[5789  370]
 [ 867 1115]]
```

```
*****
```

```
****
```

	precision	recall	f1-score	support
0	0.87	0.94	0.90	6159
1	0.75	0.56	0.64	1982
accuracy			0.85	8141
macro avg	0.81	0.75	0.77	8141
weighted avg	0.84	0.85	0.84	8141

```
*****
```

```
****
```

```
SVM Classifier Accuracy = 0.8480530647340622
```

In [44]:

```
print("Decision Tree Classifier Accuracy = ", DT_accu)
print("Random Forest Classifier Accuracy = ", RF_accu)
print("KNN Classifier Accuracy = ", KNN_accu)
print("Logistic Regression Accuracy = ", LR_accu)
print("SVM Classifier Accuracy = ", SVM_accu)
```

```
Decision Tree Classifier Accuracy = 0.8153789460754207
```

```
Random Forest Classifier Accuracy = 0.8439995086598698
```

```
KNN Classifier Accuracy = 0.8217663677680874
```

```
Logistic Regression Accuracy = 0.8502640953199853
```

```
SVM Classifier Accuracy = 0.8480530647340622
```

Logistic Regression has the best accuracy

In []: