## Abstract :

The main aim is to Create a classification model to predict whether a person makes over $50k a year

## Data Understanding :

```python
import pandas as pd
df=pd.read_csv("C:/Users/asus/Downloads/adult.csv",names=
["age","Workclass","Fnlwgt","Education","education_num","marital_status","occupation","relationship","race","sex","capital_gain","capital_loss","hours_per_week","native_country","income"])
data
```

| | Age | Workclass | Fnlwgt | Education | education_num | marital_status | occupation | relationship | race | sex | capital_gain | capital_loss | hours_per_week | native_country | income |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 30 | State-gov | 77516 | Bachelors | 13 | Never-married | Adm-clerical | Not-in-family | White | Male | 2174 | 0 | 40 | United-States | <=50K |
| 1 | 50 | Self-emp-not-inc | 83311 | Bachelors | 13 | Married-civ-spouse | Exec-managerial | Husband | White | Male | 0 | 0 | 13 | United-States | <=50K |
| 2 | 38 | Private | 215646 | HS-grad | 9 | Divorced | Handlers-cleaners | Not-in-family | White | Male | 0 | 0 | 40 | United-States | <=50K |
| 3 | 53 | Private | 234721 | 11th | 7 | Married-civ-spouse | Handlers-cleaners | Husband | Black | Male | 0 | 0 | 40 | United-States | <=50K |
| 4 | 28 | Private | 338409 | Bachelors | 13 | Married-civ-spouse | Prof-specialty | Wife | Black | Female | 0 | 0 | 40 | Cuba | <=50K |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 32556 | 27 | Private | 257302 | Assoc-acdm | 12 | Married-civ-spouse | Tech-support | Wife | White | Female | 0 | 0 | 38 | United-States | <=50K |
| 32557 | 40 | Private | 154374 | HS-grad | 9 | Married-civ-spouse | Machine-op-inspct | Husband | White | Male | 0 | 0 | 40 | United-States | >50K |
| 32558 | 58 | Private | 151910 | HS-grad | 9 | Widowed | Adm-clerical | Unmarried | White | Female | 0 | 0 | 40 | United-States | <=50K |
| 32559 | 22 | Private | 201490 | HS-grad | 9 | Never-married | Adm-clerical | Own-child | White | Male | 0 | 0 | 20 | United-States | <=50K |
| 32560 | 52 | Self-emp-inc | 287927 | HS-grad | 9 | Married-civ-spouse | Exec-managerial | Wife | White | Female | 15024 | 0 | 40 | United-States | >50K |

32561 rows × 15 columns

```python
data.isnull().sum() #Checking if there are any null values in data set
```

```
age               0
Workclass         0
Fnlwgt            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
```

## Data Cleaning :

```python
# Checking if there are any unwanted values in data set
(data["Age"]=="?").sum()
(data["Workclass"]=="?").sum()
(data["Fnlwgt"]=="?").sum()
(data["Education"]=="?").sum()
(data["education_num"]=="?").sum()
(data["marital_status"]=="?").sum()
(data["occupation"]=="?").sum()
(data["relationship"]=="?").sum()
(data["race"]=="?").sum()
(data["sex"]=="?").sum()
(data["capital_gain"]=="?").sum()
(data["capital_loss"]=="?").sum()
(data["hours_per_week"]=="?").sum()
(data["native_country"]=="?").sum()
(data["income"]=="?").sum()
```

```
0
```

Now checking count of values in each row to replace unwanted value with this :

```python
data["Workclass"].value_counts()
```

```
Private             22696
Self-emp-not-inc     2541
Local-gov            2093
?                    1836
State-gov            1298
Self-emp-inc         1116
Federal-gov           960
Without-pay            14
Never-worked            7
Name: Workclass, dtype: int64
```

```python
data["occupation"].value_counts()
```

```
Prof-specialty       4140
Craft-repair         4099
Exec-managerial      4066
Adm-clerical         3770
Sales                3650
Other-service        3295
Machine-op-inspct    2002
?                    1843
Transport-moving     1597
Handlers-cleaners    1370
Farming-fishing       994
Tech-support          928
Protective-serv       649
Priv-house-serv       149
Armed-Forces            9
Name: occupation, dtype: int64
```

```python
data["native_country"].value_counts()
```

```
United-States                 29170
Mexico                          643
?                               583
Philippines                     198
Germany                         137
Canada                          121
Puerto-Rico                     114
El-Salvador                     106
India                           100
Cuba                             95
England                          90
Jamaica                          81
South                            80
China                            75
Italy                            73
Dominican-Republic               70
Vietnam                          67
Guatemala                        64
Japan                            62
Poland                           60
Columbia                         59
Taiwan                           51
Haiti                            44
Iran                             43
Portugal                         37
Nicaragua                        34
Peru                             31
Greece                           29
France                           29
Ecuador                          28
Ireland                          24
Hong                             20
Trinadad&Tobago                  19
Cambodia                         19
Thailand                         18
Laos                             18
Yugoslavia                       16
Outlying-US(Guam-USVI-etc)       14
Honduras                         13
Hungary                          13
Scotland                         12
Holand-Netherlands                1
Name: native_country, dtype: int64
```

```python
#Replacing the unwanted values
data["Workclass"]=data["Workclass"].replace(to_replace=" ?",value=" Private")
data["occupation"]=data["occupation"].replace(to_replace=" ?",value=" Prof-specialty")
data["occupation"]=data["occupation"].replace(to_replace=" ?",value=" Prof-specialty")
data["native_country"]=data["native_country"].replace(to_replace=" ?",value=" United-States")
```

```python
data.columns
```

As all the values should be in integer datatype , string values has to be converted into integer values by using LabelEncoder

```python
from sklearn.preprocessing import LabelEncoder
l1=LabelEncoder()
data["Workclass"]=l1.fit_transform(data["Workclass"])
data["Education"]=l1.fit_transform(data["Education"])
data["marital_status"]=l1.fit_transform(data["marital_status"])
data["occupation"]=l1.fit_transform(data["occupation"])
data["relationship"]=l1.fit_transform(data["relationship"])
data["race"]=l1.fit_transform(data["race"])
data["sex"]=l1.fit_transform(data["sex"])
data["native_country"]=l1.fit_transform(data["native_country"])
data["income"]=l1.fit_transform(data["income"])
```

```python
data.shape
```

```
(32561, 15)
```

## Spliting dataset into train and test data :

```python
x=data.iloc[:,:14].values
y=data.iloc[:,14].values
x.shape
```

```
(32561, 14)
```

```python
y.shape
```

```
(32561,)
```

## Data Preprocessing :

```python
from sklearn.model_selection import train_test_split
xtrain,xtest,ytrain,ytest=train_test_split(x,y,test_size=0.2,random_state=0)
```

```python
xtrain.shape
```

```
(26048, 14)
```

```python
ytrain.shape
```

```
(26048,)
```

## Model Builiding :

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

```python
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.neighbors import KNeighborsClassifier
from sklearn.svm import SVC
```

```python
#creating function for all models
def apply_model(xtrain,xtest,ytrain,ytest,model):
    model.fit(xtrain,ytrain)
    ypred=model.predict(xtest)
    print("prediction of income for test data :",ypred)
    print("Training score:",model.score(xtrain,ytrain))
    print("Testing score:",model.score(xtest,ytest))
    print("Accuracy score:",accuracy_score(ytest,ypred)*100)
    cm=confusion_matrix(ytest,ypred)
    print("confusion matrix:\n",cm)
    print("classification report :",classification_report(ytest,ypred))
    ps=cm[0][0]/(cm[0][0]+cm[0][1])
    rs=cm[0][0]/(cm[0][0]+cm[1][0])
    print("precision:",ps)
    print("recall:",rs)
    print("f1-score: %.2f%%(ps*rs)/(ps+rs))
    print("Accuracy:",acc)
    print(" ")
    Accuracy[(((cm[0][0]+cm[1][1])/(cm[0][0]+cm[0][1]+cm[1][0]+cm[1][1])))]
    print("Accuracy",acc)
    print(" ")
    print("percentage of misclassification :\n",mss)
```

Model 1 : Decision Tree Classifier

```python
Model1=DecisionTreeClassifier(criterion='gini')
apply_model(xtrain,xtest,ytrain,ytest,Model1)
```

```
prediction of income for test data : [' <=50K' ' <=50K' ' <=50K' ... ' <=50K' ' <=50K' ' <=50K']
Training score: 0.9999616093058994
Testing score: 0.8118894519879779
Accuracy score: 81.18894519879779
confusion matrix:
 [[4566  661]
 [ 557  882]]
classification report :
              precision    recall  f1-score   support
       <=50K       0.89      0.87      0.88      5026
        >50K       0.58      0.60      0.61      1487

    accuracy                           0.76      6513
   macro avg       0.74      0.75      0.74      6513
weighted avg       0.82      0.82      0.82      6513

precision: 0.8868562007141
recall: 0.8949338983616
f1-score: 0.8778233928284264
Accuracy: 0.8118894519879
percentage of misclassification :
 337.4194491292043
```

Model 2 : Random Forest Classifier

```python
Model2=RandomForestClassifier()
apply_model(xtrain,xtest,ytrain,ytest,Model2)
```

```
prediction of income for test data : [' <=50K' ' <=50K' ' <=50K' ... ' <=50K' ' <=50K' ' <=50K']
Training score: 0.9996160930589936
Testing score: 0.84156294333546566
Accuracy score: 84.156294333546566
confusion matrix:
 [[4647  379]
 [ 647  840]]
classification report :
              precision    recall  f1-score   support
       <=50K       0.90      0.92      0.91      5026
        >50K       0.72      0.56      0.68      1487

    accuracy                           0.86      6513
   macro avg       0.81      0.78      0.79      6513
weighted avg       0.86      0.86      0.86      6513

precision: 0.9075450606254588
recall: 0.8563812809709591
f1-score: 0.85932992809085
Accuracy: 0.9996160930589936
percentage of misclassification :
 535.0501910298984
```

Model 3 : Logistic Regression

```python
Model3=LogisticRegression(solver='liblinear')
apply_model(xtrain,xtest,ytrain,ytest,Model3)
```

```
prediction of income for test data : [' <=50K' ' <=50K' ' <=50K' ... ' <=50K' ' <=50K' ' <=50K']
Training score: 0.7919589377866666
Testing score: 0.8020881314409569
Accuracy score: 80.20881314409569
confusion matrix:
 [[4696  330]
 [ 958  529]]
classification report :
              precision    recall  f1-score   support
       <=50K       0.83      0.93      0.88      5026
        >50K       0.62      0.36      0.45      1487

    accuracy                           0.80      6513
   macro avg       0.72      0.62      0.66      6513
weighted avg       0.78      0.80      0.77      6513

precision: 0.9343659619181262
recall: 0.9291563825510789
f1-score: 0.9775577407044838
Accuracy: 0.7919589377866666
percentage of misclassification :
 1209.0592699367137
```

Model 4 : KNN Classifier

```python
Model4=KNeighborsClassifier(n_neighbors=5)
apply_model(xtrain,xtest,ytrain,ytest,Model4)
```

```
prediction of income for test data : [' <=50K' ' <=50K' ' <=50K' ... ' <=50K' ' <=50K' ' <=50K']
Training score: 0.8429007992633546
Testing score: 0.77291570704430796
Accuracy score: 77.291570704430796
confusion matrix:
 [[4692  584]
 [ 932  555]]
classification report :
              precision    recall  f1-score   support
       <=50K       0.83      0.90      0.86      5026
        >50K       0.55      0.40      0.45      1487

    accuracy                           0.77      6513
   macro avg       0.67      0.64      0.65      6513
weighted avg       0.76      0.77      0.76      6513

precision: 0.8323200904931584
recall: 0.8464430933007088
f1-score: 0.84564563574387628
Accuracy: 0.7705727870424438
percentage of misclassification :
 1472.26310838371
```

Model 5 : SVC Classifier (with Linear Kernel)

```python
Model5=SVC(kernel='linear')
apply_model(xtrain,xtest,ytrain,ytest,Model5)
```

```
prediction of income for test data : [' <=50K' ' <=50K' ' <=50K' ... ' <=50K' ' <=50K' ' <=50K']
Training score: 0.7929612262521948886
Testing score: 0.80512450269385
Accuracy score: 80.512450269385
confusion matrix:
 [[4646  327]
 [1266  273]]
classification report :
              precision    recall  f1-score   support
       <=50K       0.80      1.00      0.89      5026
        >50K       0.07      0.15      0.13      1487

    accuracy                           0.80      6513
   macro avg       0.44      0.58      0.51      6513
weighted avg       0.64      0.80      0.71      6513

precision: 0.7924991988518451
recall: 0.6985639729558396
f1-score: 0.8766740751751075257
Accuracy: 0.80512450269385098
percentage of misclassification :
 1242.2621565915929
```

## Conclusion:

Random Forest is the best model in terms of accuracy