# EDA, FE and Classification Model (Census Income Dataset)

## Aryan Shriva

**Linkedin:** https://www.linkedin.com/in/aryan-shriva-0811/ (https://www.linkedin.com/in/aryan-shriva-0811/)

**GitHub** https://github.com/AryanShriva/Machine-Learning (https://github.com/AryanShriva/Machine-Learning)

**For the code please check out my Machine Learning repository on GitHub**

## 1. EDA and FE

1. Data Profiling
2. Stastical analysis
3. Graphical Analysis
4. Data Cleaning
5. Data Scaling

## 2. Logistic Regression Model

1. Linear Regression Model
2. Performance metrics for above model
3. Hyper-Parameter Tuning for above model

## 3. Support Vector Classifier Model

1. Support Vector Classifier Model
2. Performance metrics for above model
3. Hyper-Parameter Tuning for above model

**Dataset:** https://archive.ics.uci.edu/ml/datasets/Census+Income (https://archive.ics.uci.edu/ml/datasets/Census+Income)

## 1.0 Importing required libraries

```
In [ ]:    1  ### Pandas and Numpy
           2  import pandas as pd
           3  import numpy as np
           4
           5  ### MongoDB Library
           6  import pymongo
           7
           8  ### Machine Learning libraries
           9  from sklearn.model_selection import train_test_split
          10  from sklearn.preprocessing import OneHotEncoder
          11  from sklearn.preprocessing import StandardScaler
          12  from sklearn.compose import make_column_transformer
          13  from sklearn.linear_model import LogisticRegression
          14  from sklearn.svm import SVC
          15  from sklearn.model_selection import GridSearchCV
          16  from sklearn.metrics import confusion_matrix, accuracy_score, classification
          17
          18  ### To ignore warnings
          19  import warnings
          20  warnings.filterwarnings('ignore')
```

# 2.0 Retrieving data from MongoDB

```
In [2]:    1  ### Retriving data from Mongodb
           2  ### creating connection with MongoDB
           3
           4  client = pymongo.MongoClient("mongodb+srv://{username}:{password}@clustershu
```

```
In [3]:    1  db=client['Census_income']
           2  collection=db['Census_income_data']
```

```
In [4]:    1  ### Locating our collection and data in MongoDb using find() method
           2  data_from_mongodb=collection.find()
```

```
In [5]:    1  ### converting data from MongoDb to Dataframe in pandas
           2  data_mongodb=pd.DataFrame(data_from_mongodb)
```

```
In [6]:    1   ### first 5 records in dataset
           2   data_mongodb.head()
```

Out[6]:

|   | _id | index | age | workclass | fnlwgt | education | education_num | n |
|---|-----|-------|-----|-----------|--------|-----------|---------------|---|
| 0 | 63635b0506652a035edadc4c | 31 | 20 | Private | 266015 | Some_college | 10 | M |
| 1 | 63635b0506652a035edadc2d | 0 | 39 | other | 77516 | Bachelors | 13 | M |
| 2 | 63635b0506652a035edadc3a | 13 | 32 | Private | 205019 | other | 12 | M |
| 3 | 63635b0506652a035edadc3d | 16 | 25 | other | 176756 | HS_grad | 9 | M |
| 4 | 63635b0506652a035edadc2e | 1 | 50 | other | 83311 | Bachelors | 13 | Marrie |

```
In [7]:    1   ### dropping _id and index feature from dataset imported from MongoDB
           2   data_mongodb.drop(['_id','index'], axis=1, inplace=True)
           3   data_mongodb.head()
```

Out[7]:

|   | age | workclass | fnlwgt | education | education_num | marital_status | occupation | rel |
|---|-----|-----------|--------|-----------|---------------|----------------|------------|-----|
| 0 | 20 | Private | 266015 | Some_college | 10 | Never_married | Sales | C |
| 1 | 39 | other | 77516 | Bachelors | 13 | Never_married | Adm_clerical | Not_ |
| 2 | 32 | Private | 205019 | other | 12 | Never_married | Sales | Not_ |
| 3 | 25 | other | 176756 | HS_grad | 9 | Never_married | other | C |
| 4 | 50 | other | 83311 | Bachelors | 13 | Married_civ_spouse | Exec_managerial | |

# 3.0 Model and Evaluation

## 3.1 Seperating Independent and Dependent features

```
In [8]:    1   ### Splitting data into independent feature dataframe and dependent feature
           2   X=data_mongodb.iloc[:,:-1]
           3   y=data_mongodb.iloc[:,-1]
           4   X.head()
```

Out[8]:

|   | age | workclass | fnlwgt | education | education_num | marital_status | occupation | rel |
|---|-----|-----------|--------|-----------|---------------|----------------|------------|-----|
| 0 | 20 | Private | 266015 | Some_college | 10 | Never_married | Sales | C |
| 1 | 39 | other | 77516 | Bachelors | 13 | Never_married | Adm_clerical | Not_ |
| 2 | 32 | Private | 205019 | other | 12 | Never_married | Sales | Not_ |
| 3 | 25 | other | 176756 | HS_grad | 9 | Never_married | other | C |
| 4 | 50 | other | 83311 | Bachelors | 13 | Married_civ_spouse | Exec_managerial | |

In [9]:
```python
1  y.head()
```

Out[9]: 
```
0    0
1    0
2    0
3    0
4    0
Name: salary, dtype: int64
```

## 3.2 Train Test Split

In [10]:
```python
1  ### random state train test split will be same with all people using random_
2
3  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, ra
4  X_train.head()
```

Out[10]:

| | age | workclass | fnlwgt | education | education_num | marital_status | occupation |
|---|---|---|---|---|---|---|---|
| 34576 | 30 | Private | 23778 | Some_college | 10 | Never_married | Exec_managerial |
| 33148 | 41 | Private | 112763 | Some_college | 10 | Married_civ_spouse | Adm_clerical |
| 2109 | 20 | Private | 241752 | HS_grad | 9 | Married_civ_spouse | other |
| 33501 | 35 | Private | 211494 | Bachelors | 13 | Never_married | Exec_managerial |
| 47110 | 24 | Private | 408585 | other | 4 | Married_civ_spouse | other |

In [11]:
```python
1  y_train.head()
```

Out[11]:
```
34576    0
33148    0
2109     0
33501    0
47110    0
Name: salary, dtype: int64
```

In [12]:
```python
1  X_test.head()
```

Out[12]:

| | age | workclass | fnlwgt | education | education_num | marital_status | occupation | re |
|---|---|---|---|---|---|---|---|---|
| 3436 | 52 | Private | 48925 | Some_college | 10 | Married_civ_spouse | Adm_clerical | |
| 16332 | 27 | Private | 31659 | Bachelors | 13 | Married_civ_spouse | other | |
| 39798 | 67 | Private | 187553 | other | 4 | Divorced | Prof_specialty | Nc |
| 12405 | 45 | other | 255559 | HS_grad | 9 | Never_married | Adm_clerical | Nc |
| 7584 | 32 | Private | 169955 | Some_college | 10 | Married_civ_spouse | Other_service | |

```python
In [13]:    1  y_test.head()
```

```
Out[13]:  3436      0
          16332     1
          39798     0
          12405     0
          7584      0
          Name: salary, dtype: int64
```

```python
In [14]:    1
            2  ### both will have same shape
            3  X_train.shape, y_train.shape
```

```
Out[14]:  ((36609, 14), (36609,))
```

```python
In [15]:    1  ### both will have same shape
            2  X_test.shape, y_test.shape
```

```
Out[15]:  ((12204, 14), (12204,))
```

### 3.3 Feature Encoding

```python
In [16]:    1  column_trans=make_column_transformer(
            2          (OneHotEncoder(), ['workclass','education', 'marital_status', 'occup
            3          remainder='passthrough')
```

```python
In [17]:    1  X_train=column_trans.fit_transform(X_train)
```

```python
In [18]:    1  X_test=column_trans.transform(X_test)
```

### 3.4 Feature Scaling

```python
In [19]:    1  scaler=StandardScaler()
```

```python
In [20]:    1  X_train=scaler.fit_transform(X_train)
```

```python
In [21]:    1  X_test=scaler.transform(X_test)
```

In [22]:
```python
1  X_train_scaled=pd.DataFrame(X_train)
2  X_train_scaled.head()
```

Out[22]:

|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 0.573765 | -0.573765 | -0.444065 | -0.691388 | 1.869548 | -0.638185 | -0.394399 | -0.923205 | 1.426647 |
| 1 | 0.573765 | -0.573765 | -0.444065 | -0.691388 | 1.869548 | -0.638185 | -0.394399 | 1.083183 | -0.700944 |
| 2 | 0.573765 | -0.573765 | -0.444065 | 1.446367 | -0.534889 | -0.638185 | -0.394399 | 1.083183 | -0.700944 |
| 3 | 0.573765 | -0.573765 | 2.251920 | -0.691388 | -0.534889 | -0.638185 | -0.394399 | -0.923205 | 1.426647 |
| 4 | 0.573765 | -0.573765 | -0.444065 | -0.691388 | -0.534889 | 1.566943 | -0.394399 | 1.083183 | -0.700944 |

5 rows × 34 columns

In [23]:
```python
1  X_test_scaled=pd.DataFrame(X_test)
2  X_test_scaled.head()
```

Out[23]:

|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 0.573765 | -0.573765 | -0.444065 | -0.691388 | 1.869548 | -0.638185 | -0.394399 | 1.083183 | -0.700944 |
| 1 | 0.573765 | -0.573765 | 2.251920 | -0.691388 | -0.534889 | -0.638185 | -0.394399 | 1.083183 | -0.700944 |
| 2 | 0.573765 | -0.573765 | -0.444065 | -0.691388 | -0.534889 | 1.566943 | 2.535503 | -0.923205 | -0.700944 |
| 3 | -1.742874 | 1.742874 | -0.444065 | 1.446367 | -0.534889 | -0.638185 | -0.394399 | -0.923205 | 1.426647 |
| 4 | 0.573765 | -0.573765 | -0.444065 | -0.691388 | 1.869548 | -0.638185 | -0.394399 | 1.083183 | -0.700944 |

5 rows × 34 columns

## 3.5 Logestic Regression Model

In [24]:
```python
1  ### model
2  logistic_reg=LogisticRegression()
3  logistic_reg
```

Out[24]: LogisticRegression()

In [25]:
```python
1  logistic_reg.fit(X_train, y_train)
```

Out[25]: LogisticRegression()

In [26]:
```python
1  logistic_reg_pred=logistic_reg.predict(X_test)
2  logistic_reg_pred
```

Out[26]: array([0, 1, 0, ..., 0, 0, 0])

```
In [27]:   1  confusion_mat=confusion_matrix(y_test, logistic_reg_pred)
           2  confusion_mat
```

```
Out[27]:  array([[8618,  644],
                 [1175, 1767]])
```

```
In [28]:   1  truly_positive=confusion_mat[0][0]
           2  falsely_positive=confusion_mat[0][1]
           3  falsely_negative=confusion_mat[1][0]
           4  truly_negative=confusion_mat[1][1]
```

```
In [29]:   1  classification_rep_log_reg=classification_report(y_test, logistic_reg_pred)
           2  print(classification_rep_log_reg)
```

```
              precision    recall  f1-score   support

           0       0.88      0.93      0.90      9262
           1       0.73      0.60      0.66      2942

    accuracy                           0.85     12204
   macro avg       0.81      0.77      0.78     12204
weighted avg       0.84      0.85      0.85     12204
```

## 3.6 Support Vector Classifier Model

```
In [30]:   1  X_train1, X_test1, y_train1, y_test1 = train_test_split(X, y, test_size=0.25
```

```
In [31]:   1  column_trans_svc=make_column_transformer(
           2          (OneHotEncoder(), ['workclass','education', 'marital_status', 'occup
           3          remainder='passthrough')
```

```
In [32]:   1  X_train1=column_trans_svc.fit_transform(X_train1)
```

```
In [33]:   1  X_test1=column_trans_svc.transform(X_test1)
```

```
In [34]:   1  scaler_svc=StandardScaler()
           2  scaler_svc
```

```
Out[34]:  StandardScaler()
```

```
In [35]:   1  X_train1=scaler_svc.fit_transform(X_train1)
```

```
In [36]:   1  X_test1=scaler_svc.transform(X_test1)
```

```
In [ ]:   1 svc=SVC()
          2 svc
```

Out[38]: SVC()

```
In [ ]:   1 svc.fit(X_train1, y_train1)
```

Out[39]: SVC()

```
In [ ]:   1 svc_pred=svc.predict(X_test1)
          2 svc_pred
```

Out[40]: array([0, 1, 0, ..., 0, 0, 0])

```
In [ ]:   1 confusion_mat_svc=confusion_matrix(y_test1, svc_pred)
          2 confusion_mat_svc
```

Out[41]: array([[8718,  544],
                [1256, 1686]])

```
In [ ]:   1 truly_positive=confusion_mat_svc[0][0]
          2 falsely_positive=confusion_mat_svc[0][1]
          3 falsely_negative=confusion_mat_svc[1][0]
          4 truly_negative=confusion_mat_svc[1][1]
```

```
In [ ]:   1 classification_rep_svc=classification_report(y_test1, svc_pred)
          2 print(classification_rep_svc)
```

```
              precision    recall  f1-score   support

           0       0.87      0.94      0.91      9262
           1       0.76      0.57      0.65      2942

    accuracy                           0.85     12204
   macro avg       0.82      0.76      0.78     12204
weighted avg       0.85      0.85      0.85     12204
```

### 3.7 Hyper-Parameter Tuning Logistic Regression Model

```
In [ ]:   1 param_grid = [
          2     {'penalty' : ['l1', 'l2', 'elasticnet', 'none'],
          3     'C' : np.logspace(-4, 4, 5),
          4     'solver' : ['lbfgs','newton-cg','liblinear','sag','saga'],
          5     'max_iter' : [100, 500]
          6     }
          7 ]
```

```
In [ ]:   1 log_reg_hpt=LogisticRegression()
          2 log_reg_hpt
```

Out[45]: LogisticRegression()

```
In [ ]:  1  hpt_log_reg=GridSearchCV(log_reg_hpt, param_grid = param_grid)
```

```
In [ ]:  1  best_hpt_log_reg=hpt_log_reg.fit(X_train, y_train)
         2  best_hpt_log_reg
```

```
Out[47]: GridSearchCV(estimator=LogisticRegression(),
                  param_grid=[{'C': array([1.e-04, 1.e-02, 1.e+00, 1.e+02, 1.e+04]),
                             'max_iter': [100, 500],
                             'penalty': ['l1', 'l2', 'elasticnet', 'none'],
                             'solver': ['lbfgs', 'newton-cg', 'liblinear', 'sag',
                                       'saga']}])
```

```
In [ ]:  1  ### getting best parameters for Logistic Regression model after gridsearchCV
         2  print("Best parameters are {} for optimal accuracy.".format(best_hpt_log_reg
```

Best parameters are LogisticRegression(C=0.01, penalty='l1', solver='liblinea
r') for optimal accuracy.

```
In [ ]:  1  ### getting best accuracy for Logistic Regression model after gridsearchCV
         2  print("Best accuracy is {}".format(best_hpt_log_reg.score(X_test, y_test)))
```

Best accuracy is 0.8504588659455916

## 3.8 Hyper-Parameter Tuning Support Vector Classifier Model

```
In [37]:  1  svc_hpt=SVC()
          2  svc_hpt
```

```
Out[37]: SVC()
```

```
In [38]:  1  #### using gridsearchcv to increase model efficiency by combining above para
          2  param_grid={'C':[1,2,3], 'kernel':['rbf']}
          3  hpt_svc=GridSearchCV(svc_hpt, param_grid=param_grid)
```

```
In [39]:  1  best_hpt_svc=hpt_svc.fit(X_train1, y_train1)
          2  best_hpt_svc
```

```
Out[39]: GridSearchCV(estimator=SVC(),
                  param_grid={'C': [10, 20], 'degree': [2, 3],
                              'kernel': ['linear', 'rbf', 'poly', 'sigmoid']})
```

```
In [40]:  1  ### getting best parameters for Logistic Regression model after gridsearchCV
          2  print("Best parameters are {} for optimal accuracy.".format(best_hpt_svc.bes
```

Best parameters are SVC(C=20, degree=2, kernel='poly') for optimal accuracy.

```
In [41]:  1  ### getting best accuracy for Logistic Regression model after gridsearchCV
          2  print("Best accuracy is {}".format(best_hpt_svc.score(X_test, y_test)))
```

Best accuracy is 0.8491478203867584

In [42]:
```python
1  svc_hpt1=SVC()
2  svc_hpt1
```

Out[42]:  SVC()

In [43]:
```python
1  #### using gridsearchcv to increase model efficiency by combining above para
2  param_grid1={'C':[1,2,3], 'kernel':['rbf']}
3  hpt_svc1=GridSearchCV(svc_hpt1, param_grid=param_grid1)
```

In [44]:
```python
1  best_hpt_svc1=hpt_svc1.fit(X_train1, y_train1)
2  best_hpt_svc1
```

Out[44]:  GridSearchCV(estimator=SVC(), param_grid={'C': [1, 2, 3], 'kernel': ['rbf']})

In [45]:
```python
1  ### getting best parameters for Logistic Regression model after gridsearchCV
2  print("Best parameters are {} for optimal accuracy.".format(best_hpt_svc1.be
```

Best parameters are SVC(C=2) for optimal accuracy.

In [47]:
```python
1  ### getting best accuracy for Logistic Regression model after gridsearchCV
2  print("Best accuracy is {}".format(best_hpt_svc1.score(X_test1, y_test1)))
```

Best accuracy is 0.8517699115044248

**Note: Please refer my github repo for ROC AUC curve implementation**