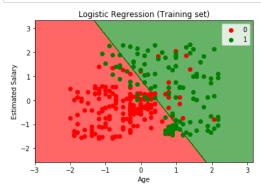
Assignment 5

- 1. Implement logistic regression using Python/R to perform classification on Social Network Ads.csv dataset.
- 2. Compute Confusion matrix to find TP, FP, TN, FN, Accuracy, Error rate, Precision, Recall on the given dataset..

```
In [10]:
          import pandas as pd
          import numpy as np
          import matplotlib.pyplot as plt
In [11]: dataset=pd.read_csv("C:\\Users\\Admin\\Desktop\\Social_Network_Ads.csv")
In [12]: print(dataset)
                User ID
                          Gender
                                   Age
                                        EstimatedSalary
                                                          Purchased
          0
               15624510
                            Male
                                    19
                                                   19000
                                                                   0
          1
               15810944
                            Male
                                    35
                                                   20000
                                                                   a
               15668575
                                                   43000
          2
                          Female
                                    26
                                                                   a
                                                   57000
               15603246
                                    27
                                                                   0
                          Female
               15804002
                                                   76000
          4
                            Male
                                    19
                                                                   0
                                                   58000
               15728773
                                    27
                                                                   0
                            Male
               15598044
                                    27
                                                   84000
                                                                   0
          6
                          Female
               15694829
                          Female
                                    32
                                                  150000
                                                                   1
               15600575
                            Male
                                    25
                                                   33000
                                                                   0
               15727311
                          Female
                                    35
                                                   65000
          10
               15570769
                                                   80000
                                                                   0
                          Female
                                    26
          11
               15606274
                          Female
                                    26
                                                   52000
          12
               15746139
                            Male
                                    20
                                                   86000
                                                                   0
               15704987
                                                   18000
                                                                   0
                            Male
                                    32
          14
               15628972
                            Male
                                    18
                                                   82000
                                                                   0
          15
               15697686
                            Male
                                    29
                                                   80000
                                                                   0
          16
               15733883
                            Male
                                    47
                                                   25000
                                                                   1
          17
               15617482
                            Male
                                    45
                                                   26000
                                                                   1
                                                   28000
          18
               15704583
                            Male
                                    46
                                                                   1
          19
               15621083
                          Female
                                    48
                                                   29000
                                                                   1
          20
               15649487
                            Male
                                    45
                                                   22000
                                                                   1
               15736760
                                                   49000
          21
                          Female
                                    47
                                                                   1
          22
                                                   41000
                                                                   1
               15714658
                            Male
                                    48
               15599081
                                                   22000
          23
                                    45
                                                                   1
                          Female
          24
               15705113
                            Male
                                    46
                                                   23000
                                                                   1
                                                   20000
          25
               15631159
                            Male
                                    47
                                                                   1
          26
               15792818
                            Male
                                    49
                                                   28000
                                                                   1
          27
               15633531
                          Female
                                    47
                                                   30000
                                                                   1
          28
               15744529
                            Male
                                    29
                                                   43000
                                                                   0
          29
               15669656
                                                                   0
                            Male
                                    31
                                                   18000
          370
               15611430
                          Female
                                    60
                                                   46000
                                                                   1
               15774744
                                                   83000
          371
                            Male
                                                                   1
          372
               15629885
                                    39
                                                   73000
                                                                   0
                          Female
          373
               15708791
                            Male
                                    59
                                                  130000
                                                                   1
          374
               15793890
                          Female
                                    37
                                                   80000
                                                                   0
          375
               15646091
                          Female
                                    46
                                                   32000
                                                                   1
          376
               15596984
                          Female
                                    46
                                                   74000
                                                                   0
                                                                   0
          377
               15800215
                          Female
                                    42
                                                   53000
                                                   87000
          378
               15577806
                            Male
                                    41
                                                                   1
          379
               15749381
                          Female
                                    58
                                                   23000
                                                                   1
                                                   64000
                                                                   0
          380
               15683758
                            Male
                                    42
          381
               15670615
                                                   33000
                                                                   1
                            Male
                                    48
                                                  139000
               15715622
                                    44
                                                                   1
          382
                          Female
                                                   28000
          383
               15707634
                            Male
                                    49
                                                                   1
               15806901
                                    57
                                                   33000
                                                                   1
          384
                          Female
          385
               15775335
                                                   60000
                                                                   1
                            Male
                                    56
               15724150
                          Female
                                    49
                                                   39000
                                                                   1
          386
          387
               15627220
                            Male
                                    39
                                                   71000
                                                                   0
          388
               15672330
                            Male
                                    47
                                                   34000
          389
               15668521
                          Female
                                    48
                                                   35000
                                                                   1
          390
               15807837
                            Male
                                    48
                                                   33000
                                                                   1
          391
               15592570
                                    47
                                                   23000
                            Male
                                                                   1
          392
               15748589
                          Female
                                    45
                                                   45000
          393
               15635893
                            Male
                                    60
                                                   42000
          394
               15757632
                          Female
                                    39
                                                   59000
                                                                   0
          395
               15691863
                          Female
                                    46
                                                   41000
                                                                   1
          396
               15706071
                            Male
                                    51
                                                   23000
                                                                   1
                                                   20000
          397
               15654296
                          Female
                                    50
                                                                   1
          398
               15755018
                            Male
                                    36
                                                   33000
                                                                   0
          399
               15594041 Female
                                    49
                                                   36000
          [400 rows x 5 columns]
In [13]: dataset.isnull().sum()
Out[13]: User ID
          Gender
          Age
          EstimatedSalary
          Purchased
                               0
          dtype: int64
```

```
In [14]: X = dataset.iloc[:, [2, 3]].values
         y = dataset.iloc[:, 4].values
In [17]: 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 = 0)
In [18]: print(X_train[:3])
print('-'*15)
         print(y_train[:3])
         print('-'*15)
         print(X_test[:3])
         print('-'*15)
         print(y_test[:3])
               44 390001
         П
               32 120000]
               38 5000011
          [
         [0 1 0]
         [[ 30 87000]
              38 50000]
          [ 35 75000]]
         [0 0 0]
In [19]: from sklearn.preprocessing import StandardScaler
         sc_X = StandardScaler()
         X_train = sc_X.fit_transform(X_train)
         X_test = sc_X.transform(X_test)
         C:\ProgramData\Anaconda3\lib\site-packages\sklearn\utils\validation.py:475: DataConversionWarning: Data with input dtype int6
         4 was converted to float64 by StandardScaler.
           warnings.warn(msg, DataConversionWarning)
In [20]: | print(X_train[:3])
         print('-'*15)
         print(X_test[:3])
         [[ 0.58164944 -0.88670699]
          [-0.60673761 1.46173768]
          [-0.01254409 -0.5677824 ]]
         [[-0.80480212 0.50496393]
          [-0.01254409 -0.5677824 ]
          [-0.30964085 0.1570462 ]]
In [21]: from sklearn.linear_model import LogisticRegression
         classifier = LogisticRegression(random_state = 0, solver='lbfgs') classifier.fit(X_{train}, y_{train})
         y_pred = classifier.predict(X_test)
         print(X_test[:10])
         [[-0.80480212 0.50496393]
          [-0.01254409 -0.5677824 ]
          [-0.30964085 0.1570462
          [-0.80480212 0.27301877]
          [-0.30964085 -0.5677824 ]
          [-1.10189888 -1.43757673]
          [-0.70576986 -1.58254245]
          [-0.21060859 2.15757314]
          [-1.99318916 -0.04590581]
          [ 0.8787462 -0.77073441]]
In [22]: | print('-'*15)
         print(y_pred[:10])
         [0000000101]
In [23]: print(y_pred[:20])
         print(y_test[:20])
         [0 0 0 0 0 0 0 1 0 1 0 0 0 0 0 0 0 0 1 0]
         In [25]: from sklearn.metrics import confusion_matrix
         cm = confusion_matrix(y_test, y_pred)
         print(cm)
         [[65 3]
          [ 8 24]]
```



In []: