```
CO3234 Pattern Recognition Assignment-1
                                                       by Siddhant Verma - 2K18/EC/167
 In [2]: import numpy as np
         import pandas as pd
         import matplotlib.pyplot as plt
         import seaborn as sns
         from sklearn.naive_bayes import GaussianNB
         from sklearn.metrics import confusion_matrix
         from sklearn.metrics import roc_curve, plot_roc_curve, roc_auc_score
         from sklearn.preprocessing import label_binarize
         from sklearn.model_selection import train_test_split
 In [3]: !1s
         BankNote_Authentication.csv q1.ipynb
                                                              t10k-labels.idx1-ubyte
         pattern.ipynb
                                      q2.ipynb
                                                              train-images.idx3-ubyte
                                      t10k-images.idx3-ubyte train-labels.idx1-ubyte
         pattern2.ipynb
         Importing Dataset as Dataframe
 In [4]: df = pd.read_csv('./BankNote_Authentication.csv')
         df
 Out[4]:
               variance skewness curtosis entropy class
            0 3.62160
                       8.66610 -2.8073 -0.44699
            1 4.54590
                       8.16740 -2.4586 -1.46210
            2 3.86600
                       -2.63830
                              1.9242 0.10645
            3 3.45660
                       9.52280 -4.0112 -3.59440
            4 0.32924
                       -4.45520 4.5718 -0.98880
          1367 0.40614
                       1.34920
                              -1.4501 -0.55949
          1368 -1.38870
                       -4.87730 6.4774 0.34179
                      -13.45860 17.5932 -2.77710
          1369 -3.75030
          1370 -3.56370
                       -8.38270 12.3930 -1.28230
          1371 -2.54190
                       -0.65804 2.6842 1.19520
         1372 rows × 5 columns
         Data Visualization and Exploration
 In [5]: sns.pairplot(df)
 Out[5]: <seaborn.axisgrid.PairGrid at 0x7f409c5fce80>
             5.0
             2.5
             0.8
                                                                             -7.5 -5.0 -2.5 0.0 2.5 0.00 0.25 0.50 0.75 1.00
 In [6]: plt.plot(df['class'])
 Out[6]: [<matplotlib.lines.Line2D at 0x7f4054fe60a0>]
          0.8
          0.6
          0.4
                             600
                                  800 1000 1200 1400
         Separating the target variable from the Dataset
 In [7]: X = df.iloc[:, :4]
 Out[7]:
               variance skewness curtosis entropy
            0 3.62160
                       8.66610 -2.8073 -0.44699
                        8.16740 -2.4586 -1.46210
                       -2.63830
            2 3.86600
                              1.9242 0.10645
            3 3.45660
                        9.52280 -4.0112 -3.59440
            4 0.32924
                       -4.45520 4.5718 -0.98880
          1367 0.40614
                       1.34920 -1.4501 -0.55949
          1368 -1.38870
                      -13.45860 17.5932 -2.77710
          1369 -3.75030
                       -8.38270 12.3930 -1.28230
          1371 -2.54190
                       -0.65804 2.6842 1.19520
         1372 rows × 4 columns
 In [8]: Y = df.iloc[:,4].values
 Out[8]: array([0, 0, 0, ..., 1, 1, 1])
 In [9]: plt.plot(Y)
 Out[9]: [<matplotlib.lines.Line2D at 0x7f4054fc83d0>]
          1.0
          0.8
          0.6
          0.4
          0.2
                   200
                        400
                             600
                                  800
                                       1000 1200 1400
         Prior Probability of both classes
In [10]: unique, counts = np.unique(df['class'], return_counts=True)
         counts
         print(f"Dataset:\nPrior Probability of Genuine Class is : {counts[0]/1372}\nPrior Probability of Forged Class is :
          {counts[1]/1372}")
         Dataset:
         Prior Probability of Genuine Class is: 0.5553935860058309
         Prior Probability of Forged Class is: 0.4446064139941691
         Test-Train Splitting
In [11]: X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.5, random_state=40)
In [12]: plt.plot(y_train)
Out[12]: [<matplotlib.lines.Line2D at 0x7f4054f20be0>]
          0.8
          0.6
          0.4
          0.2
          0.0
                   100
                             300
                                        500
                        200
                                   400
                                             600
In [13]: unique, counts = np.unique(y_test, return_counts=True)
         counts
         print(f"After Splitting:\nPrior Probability of Genuine Class is : {counts[0]*2/1372}\nPrior Probability of Forged Cl
         ass is : {counts[1]*2/1372}")
         After Splitting:
         Prior Probability of Genuine Class is : 0.5568513119533528
         Prior Probability of Forged Class is: 0.44314868804664725
         Building the classifier
In [14]: model = GaussianNB()
In [15]: fit = model.fit(X_train, y_train)
In [16]: predictions = fit.predict(X_test)
         con_matrix = confusion_matrix(y_test, predictions)
         print(con_matrix)
         [[340 42]
          [ 61 243]]
In [17]: def draw_confusionmatrix(ytest, yhat):
             plt.figure(figsize=(10,7))
             cm = confusion_matrix(ytest, yhat)
             ax = sns.heatmap(cm, annot=True, fmt="d")
             plt.ylabel('True label')
             plt.xlabel('Predicted label')
         Confusion Matrix
In [18]: draw_confusionmatrix(y_test, predictions)
                                                                         - 300
                                                                         - 150
                                                                         - 100
                                    Predicted label
In [19]: def diag_sum(conmat):
             sum = 0
             for i in range(len(conmat)):
                  for j in range(len(conmat[0])):
                     if i == j:
                         sum += conmat[i,j]
             return sum
         Accuracy
In [20]: | accurate = diag_sum(con_matrix)
         total = np.sum(con_matrix)
         acc = accurate/total
         print(f"accuracy = {accurate}/{total} = {acc*100}%")
         accuracy = 583/686 = 84.98542274052478%
In [21]: probs = fit.predict_proba(X_test)
         preds = probs[:,1]
         fpr, tpr, threshold = roc_curve(y_test, preds)
         ROC Curve
In [24]: plt.figure(figsize=(10,7))
         plt.plot(fpr,tpr)
         plt.ylabel('True Positive Rate')
         plt.xlabel('False Positive Rate')
         plt.title('FAR vs GAR')
Out[24]: Text(0.5, 1.0, 'FAR vs GAR')
                                            FAR vs GAR
            1.0
            0.8
            0.0
                                                                             1.0
                                          False Positive Rate
In [23]: auc = roc_auc_score(y_test, preds)
         print(f"AUC: {auc*100}%")
         AUC: 93.26346789749242%
         Splitting the data such that 90% is genuine (class 0) and 10% is forged (class 1)
In [26]: df
Out[26]:
               variance skewness curtosis entropy class
                        8.66610 -2.8073 -0.44699
            0 3.62160
            1 4.54590
                        8.16740
                               -2.4586
                                     -1.46210
            2 3.86600
                       -2.63830 1.9242 0.10645
                               -4.0112 -3.59440
            3 3.45660
            4 0.32924
                       -4.45520
                              4.5718 -0.98880
          1367 0.40614
                       1.34920 -1.4501 -0.55949
          1368 -1.38870
                       -4.87730
                               6.4774 0.34179
          1369 -3.75030
                      -13.45860 17.5932 -2.77710
          1370 -3.56370
                       -8.38270
                              12.3930 -1.28230
          1371 -2.54190 -0.65804 2.6842 1.19520
         1372 rows × 5 columns
In [29]: uniques, counts = np.unique(df['class'], return_counts=True)
         counts
Out[29]: array([762, 610])
In [34]: df = df[:846]
In [37]: uniques, counts = np.unique(df['class'], return_counts=True)
         ratio = counts[0]/np.sum(counts)
         print(f"After Redistribution\nPrior Probability of Genuine class = {ratio}\nPrior Probability of Forged Class = {1-ratio}\nPrior Probability of Probability = {1-rat
         atio}")
         After Redistribution
         Prior Probability of Genuine class = 0.900709219858156
         Prior Probability of Forged Class = 0.099290780141844
In [39]: plt.plot(df['class'])
Out[39]: [<matplotlib.lines.Line2D at 0x7f40395d3d90>]
          1.0
          0.8
          0.4
          0.2
          0.0
                      200
         Data Splitting
In [40]: X = df.iloc[:, :4]
Out[40]:
               variance skewness curtosis entropy
            0 3.621600
                       8.66610 -2.80730 -0.44699
                       8.16740 -2.45860 -1.46210
            1 4.545900
            2 3.866000
                       -2.63830 1.92420 0.10645
            3 3.456600
                       9.52280 -4.01120 -3.59440
            4 0.329240
                       -4.45520 4.57180 -0.98880
                      -12.84610 12.79570 -3.13530
          841 -3.885800
                       -6.78930 5.27610 -0.32544
          842 -1.896900
          843 -0.526450
                       -0.24832 -0.45613 0.41938
                       3.56120 -4.40700 -4.41030
          844 0.009661
          845 -3.882600
                       4.89800 -0.92311 -5.08010
         846 rows × 4 columns
In [43]: Y = df.iloc[:,4].values
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1,
                1, 1, 1, 1, 1, 1, 1, 1, 1, 1])
In [44]: X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.5, random_state=40)
In [45]: plt.plot(y_train)
Out[45]: [<matplotlib.lines.Line2D at 0x7f40398fffa0>]
          0.8
          0.6
          0.4
          0.2
         Building the bayes classifier
In [47]: model = GaussianNB()
In [48]: fit = model.fit(X_train, y_train)
In [49]: predictions = fit.predict(X_test)
         con_matrix = confusion_matrix(y_test, predictions)
         print(con_matrix)
         [[372 7]
          [ 22 22]]
         Confusion Matrix
In [50]: draw_confusionmatrix(y_test, predictions)
                                                                         - 350
                                                                         - 300
                                                                         - 250
```

1

FAR vs GAR

False Positive Rate

0

accuracy = 394/423 = 93.14420803782507%

Accuracy

In [51]: accurate = diag_sum(con_matrix)
total = np.sum(con_matrix)

acc = accurate/total

In [52]: probs = fit.predict_proba(X_test)

preds = probs[:,1]

ROC Curve

1.0

0.8

9.0 gg

을 0.4

0.2

0.0

In [54]: auc = roc_auc_score(y_test, preds)
print(f"AUC: {auc*100}%")

AUC: 95.53849844087311%

In [53]: plt.figure(figsize=(10,7))
 plt.plot(fpr,tpr)

plt.title('FAR vs GAR')

Out[53]: Text(0.5, 1.0, 'FAR vs GAR')

plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')

Predicted label

print(f"accuracy = {accurate}/{total} = {acc*100}%")

fpr, tpr, threshold = roc_curve(y_test, preds)