





```
In [32]: # For SVM
SVM_model2 = SVC(C = 10, gamma = 0.01)

# Classification based on SVM
SVM_accuracy_list = []
SVM_precision_list = []
SVM_recall_list = []
SVM_f1_list = []

# Iteration for Retraining
for i in random.seed.iterations():
    # Random Dataset Split
    x_train_dataA_SVM, x_test_dataA_SVM, y_train_dataA_SVM, y_test_dataA_SVM = train_test_split(x_standar
ardized_dataA, np_label_dataA.ravel(), test_size = 0.3, random_state = i)

    # Training the Model
    SVM_model2.fit(x_train_dataA_SVM, y_train_dataA_SVM)

    # Test Dataset Prediction
    y_predict_SVM_list = SVM_model2.predict(x_test_dataA_SVM)

    # Classification Accuracy
    accuracy = metrics.accuracy_score(y_test_dataA_SVM, y_predict_SVM_list)
    SVM_accuracy_list.append(accuracy)

    #Precision
    precision = metrics.precision_score(y_test_dataA_SVM, y_predict_SVM_list)
    SVM_precision_list.append(precision)

    #Recall
    recall = metrics.recall_score(y_test_dataA_SVM, y_predict_SVM_list)
    SVM_recall_list.append(recall)

    #F1-score
    f1score = metrics.f1_score(y_test_dataA_SVM, y_predict_SVM_list)
    SVM_f1_list.append(f1score)

print("Accuracy list is :", SVM_accuracy_list)
print("\n")
print("Precision list is :", SVM_precision_list)
print("\n")
print("Recall list is:", SVM_recall_list)
print("\n")
print("F1-Score list is:", SVM_f1_list)
print("\n")

#Average and standard Deviation of Classification values
mean_accuracy_list_SVM = sum(SVM_accuracy_list)/len(SVM_accuracy_list)
variance_accuracy_list_SVM = sum([(x-mean_accuracy_list_SVM)**2] for x in SVM_accuracy_list)/len(SVM_accuracy_list)
sd_accuracy_list_SVM = variance_accuracy_list_SVM**0.5

print("Mean for accuracy in case of SVM is: ", mean_accuracy_list_SVM)
print("\n")
print("Standard Deviation for accuracy in case of SVM is: ", sd_accuracy_list_SVM)
print("\n")
mean_precision_list_SVM=sum(SVM_precision_list)/len(SVM_precision_list)
variance_precision_list_SVM=sum([(x-mean_precision_list_SVM)**2] for x in SVM_precision_list)/len(SVM_precision_list)
sd_precision_list_SVM=variance_precision_list_SVM**0.5

print("Mean for precision list in case of SVM is: ", mean_precision_list_SVM)
print("\n")
print("Standard Deviation for precision list in case of SVM is: ", sd_precision_list_SVM)
print("\n")
mean_recall_score_list_SVM=sum(SVM_recall_list)/len(SVM_recall_list)
variance_recall_score_list_SVM=sum([(x-mean_recall_score_list_SVM)**2] for x in SVM_recall_list)/len(SVM_recall_list)
sd_recall_score_list_SVM=variance_recall_score_list_SVM**0.5

print("Mean for recall score in case of SVM is: ", mean_recall_score_list_SVM)
print("\n")
print("Standard Deviation for recall score in case of SVM is: ", sd_recall_score_list_SVM)
print("\n")
mean_f1_score_list_SVM=sum(SVM_f1_list)/len(SVM_f1_list)
variance_f1_score_list_SVM=sum([(x-mean_f1_score_list_SVM)**2] for x in SVM_f1_list)/len(SVM_f1_list)
sd_f1_score_list_SVM=variance_f1_score_list_SVM**0.5

print("Mean for f1 score in case of SVM is: ", mean_f1_score_list_SVM)
print("\n")
print("Standard Deviation for f1 score in case of SVM is: ", sd_f1_score_list_SVM)
print("\n")

Accuracy list is : [0.9227272727272727, 0.9121212121212121, 0.9121212121212121, 0.9075757575757576, 0.9, 0.8742424242424243, 0.8939393939393939, 0.9181818181818182, 0.9060606060606061, 0.9136363636363637, 0.8893939393939394, 0.8893939393939394, 0.9090909090909091, 0.9166666666666667, 0.9196969696969697, 0.9136363636363637, 0.9075757575757576, 0.9, 0.8954545454545455, 0.9045454545454545, 0.9106060606060606]

Precision list is : [0.93788819875764, 0.8962592063054755, 0.9289940828402367, 0.92879256956594427, 0.9022346368715084, 0.88461584615846, 0.898507462685671, 0.9408284023668639, 0.925595280952381, 0.923461396605059, 0.9303203203203203, 0.913428358208955, 0.9347810108902077, 0.9394904458598726, 0.9371428571428571, 0.9037804878048781, 0.931977713639685, 0.8942857142857142, 0.9118541033434651, 0.9285714285714286, 0.9309309309309309]

Recall list is : [0.9069090909090909, 0.9339339393939394, 0.9022988505747126, 0.8875739644970414, 0.9124933785310734, 0.8717201166180758, 0.8931750741839762, 0.9034090909090909, 0.8936781609194402, 0.9090909090909091, 0.8599439750510365, 0.8742857142857143, 0.8923512747875394, 0.909367231461986, 0.9132400913240091, 0.8937409024745271, 0.893430569343066, 0.9130434782608695, 0.914726821705426, 0.92468265162199, 0.914027149321267, 0.913229018492176, 0.9046242774566474, 0.898680965470891, 0.908296943231411, 0.91310751046554]

mean for accuracy in case of SVM is: 0.9055555555555556

Standard Deviation for accuracy in case of SVM is: 0.011636533989453692

mean for precision list in case of SVM is: 0.9202581656881561

Standard Deviation for precision list in case of SVM is: 0.01620781071136178

mean for recall score in case of SVM is: 0.8971115927716507

Standard Deviation for recall score in case of SVM is: 0.01640915476683593

mean for f1 score in case of SVM is: 0.9083781988965527

Standard Deviation for f1 score in case of SVM is: 0.011053761601686089
```

## Naive Bayes Classifier with 20 Iterations

```
In [33]: # For Naive Bayes
NB_model1 = GaussianNB()

# Classification based on NB
NB_accuracy_list = []
NB_precision_list = []
NB_recall_list = []
NB_f1_list = []

# Iteration for Retraining
for i in random.seed.iterations():
    # Random Dataset Split
    x_train_dataA_NB, x_test_dataA_NB, y_train_dataA_NB, y_test_dataA_NB = train_test_split(x_standar
ardized_dataA, np_label_dataA.ravel(), test_size = 0.3, random_state = i)

    # Training the Model
    NB_model1.fit(x_train_dataA_NB, y_train_dataA_NB)

    #Test Dataset Prediction
    y_predict_NB_list = NB_model1.predict(x_test_dataA_NB)

    #Classification Accuracy
    accuracy = metrics.accuracy_score(y_test_dataA_NB, y_predict_NB_list)
    NB_accuracy_list.append(accuracy)

    #Precision
    precision = metrics.precision_score(y_test_dataA_NB, y_predict_NB_list)
    NB_precision_list.append(precision)

    #Recall
    recall = metrics.recall_score(y_test_dataA_NB, y_predict_NB_list)
    NB_recall_list.append(recall)

    #F1-score
    f1score = metrics.f1_score(y_test_dataA_NB, y_predict_NB_list)
    NB_f1_list.append(f1score)

print("Accuracy list is :", NB_accuracy_list)
print("\n")
print("Precision list is :", NB_precision_list)
print("\n")
print("Recall list is:", NB_recall_list)
print("\n")
print("F1-Score list is:", NB_f1_list)
print("\n")

#Average and standard Deviation of Classification values
mean_accuracy_list_NB = sum(NB_accuracy_list)/len(NB_accuracy_list)
variance_accuracy_list_NB = sum([(x-mean_accuracy_list_NB)**2] for x in NB_accuracy_list)/len(NB_accuracy_list)
sd_accuracy_list_NB = variance_accuracy_list_NB**0.5

print("Mean for accuracy in case of Naive Bayes is: ", mean_accuracy_list_NB)
print("\n")
print("Standard Deviation for accuracy in case of Naive Bayes is: ", sd_accuracy_list_NB)
print("\n")
mean_precision_list_NB=sum(NB_precision_list)/len(NB_precision_list)
variance_precision_list_NB=sum([(x-mean_precision_list_NB)**2] for x in NB_precision_list)/len(NB_precision_list)
sd_precision_list_NB=variance_precision_list_NB**0.5

print("Mean for precision list in case of Naive Bayes is: ", mean_precision_list_NB)
print("\n")
print("Standard Deviation for precision list in case of Naive Bayes is: ", sd_precision_list_NB)
print("\n")
mean_recall_score_list_NB=sum(NB_recall_list)/len(NB_recall_list)
variance_recall_score_list_NB=sum([(x-mean_recall_score_list_NB)**2] for x in NB_recall_list)/len(NB_recall_list)
sd_recall_score_list_NB=variance_recall_score_list_NB**0.5

print("Mean for recall score in case of Naive Bayes is: ", mean_recall_score_list_NB)
print("\n")
print("Standard Deviation for recall score in case of Naive Bayes is: ", sd_recall_score_list_NB)
print("\n")
mean_f1_score_list_NB=sum(NB_f1_list)/len(NB_f1_list)
variance_f1_score_list_NB=sum([(x-mean_f1_score_list_NB)**2] for x in NB_f1_list)/len(NB_f1_list)
sd_f1_score_list_NB=variance_f1_score_list_NB**0.5

print("Mean for f1 score in case of Naive Bayes is: ", mean_f1_score_list_NB)
print("\n")
print("Standard Deviation for f1 score in case of Naive Bayes is: ", sd_f1_score_list_NB)
print("\n")

Accuracy list is : [0.8772727272727273, 0.8590909090909091, 0.8787878787878788, 0.8742424242424243, 0.853030303030303, 0.8684545454545455, 0.8712121212121212, 0.8596969696969697, 0.8878787878787879, 0.8757272727272727, 0.8742424242424243, 0.8742424242424243, 0.8590909090909091, 0.8909090909090909, 0.8727272727272727, 0.8681818181818182, 0.8772727272727273, 0.8575757575757575, 0.8469696969696969, 0.8681818181818182, 0.8696969696969697]

Precision list is : [0.853925842666629, 0.83195530726257, 0.8806818181818182, 0.887379939209727, 0.859943977910365, 0.8518518518518519, 0.8446225081961673, 0.882183080409777, 0.894285714285715, 0.8514285714285714, 0.8705054971751462, 0.8803418803418803, 0.8611111111111112, 0.8888888888888888, 0.86285714285714, 0.879531520728391, 0.8628571428571429, 0.878186968395269, 0.8942598157311178, 0.9108635097493036, 0.8865671641791045, 0.9209039548022598, 0.885965217391304, 0.8915626505602408, 0.875815834522111, 0.87947840464243]

Recall list is : [0.9129129129129129, 0.8978978978978979, 0.8908045977014194, 0.863905325443787, 0.867231033707865168, 0.871201166180758, 0.916913468758791, 0.9195909090909091, 0.896551724137931, 0.9030303030303031, 0.879531520728391, 0.8628571428571429, 0.878186968395269, 0.8942598157311178, 0.9108635097493036, 0.8865671641791045, 0.9209039548022598, 0.885965217391304, 0.8915626505602408, 0.875815834522111, 0.87947840464243]

F1-Score list is : [0.8824381644005805, 0.8634124457308249, 0.8857142857142858, 0.8755622188905347, 0.863574331926863, 0.861671469740634, 0.8790896159371213, 0.8771428571428572, 0.8939826080223227, 0.8742424242424244, 0.935742696187869, 0.9028571428571428, 0.941176470582335, 0.927374301675977, 0.904669746418605, 0.9342630098459039, 0.8815977175463623, 0.869565217391304, 0.8915626505602408, 0.886178617861787, 0.872246960746268, 0.889495251023192, 0.8657142857142858, 0.8542568542568543, 0.875815834522111, 0.87947840464243]

Mean for accuracy in case of Naive Bayes is: 0.8694083694083694

Standard Deviation for accuracy in case of Naive Bayes is: 0.01082154877878943

Mean for precision list in case of Naive Bayes is: 0.8649684906823234

Standard Deviation for precision list in case of Naive Bayes is: 0.017104767454261264

Mean for recall score in case of Naive Bayes is: 0.88899375412459191

Standard Deviation for recall score in case of Naive Bayes is: 0.016478153080216752

Mean for f1 score in case of Naive Bayes is: 0.8765869348899032

Standard Deviation for f1 score in case of Naive Bayes is: 0.010196866726204193
```

## Decision Tree with 20 Iterations

```
In [34]: # Initialize Decision Tree
DT_model1 = DecisionTreeClassifier()

# Classification based on DT
DT_accuracy_list = []
DT_precision_list = []
DT_recall_list = []
DT_f1_list = []

# Iteration for Retraining
for i in random.seed.iterations():
    # Random Dataset Split
    x_train_dataA_DT, x_test_dataA_DT, y_train_dataA_DT, y_test_dataA_DT = train_test_split(x_standar
ardized_dataA, np_label_dataA.ravel(), test_size = 0.3, random_state = i)

    # Training the Model
    DT_model1.fit(x_train_dataA_DT, y_train_dataA_DT)

    # Test Dataset Prediction
    y_predict_DT_list = DT_model1.predict(x_test_dataA_DT)

    # Classification Accuracy
    accuracy = metrics.accuracy_score(y_test_dataA_DT, y_predict_DT_list)
    DT_accuracy_list.append(accuracy)

    #Precision
    precision = metrics.precision_score(y_test_dataA_DT, y_predict_DT_list)
    DT_precision_list.append(precision)

    #Recall
    recall = metrics.recall_score(y_test_dataA_DT, y_predict_DT_list)
    DT_recall_list.append(recall)

    #F1-score
    f1score = metrics.f1_score(y_test_dataA_DT, y_predict_DT_list)
    DT_f1_list.append(f1score)

print("Accuracy list is :", DT_accuracy_list)
print("\n")
print("Precision list is :", DT_precision_list)
print("\n")
print("Recall list is :", DT_recall_list)
print("\n")
print("F1-Score list is :", DT_f1_list)
print("\n")

#Average and standard Deviation of Classification values
mean_accuracy_list_DT=sum(DT_accuracy_list)/len(DT_accuracy_list)
variance_accuracy_list_DT=sum([(x-mean_accuracy_list_DT)**2] for x in DT_accuracy_list)/len(DT_accuracy_list)
sd_accuracy_list_DT=variance_accuracy_list_DT**0.5

print("Mean for accuracy in case of Decision Tree is: ", mean_accuracy_list_DT)
print("\n")
print("Standard Deviation for accuracy in case of Decision Tree is: ", sd_accuracy_list_DT)
print("\n")
mean_precision_list_DT=sum(DT_precision_list)/len(DT_precision_list)
variance_precision_list_DT=sum([(x-mean_precision_list_DT)**2] for x in DT_precision_list)/len(DT_precision_list)
sd_precision_list_DT=variance_precision_list_DT**0.5

print("Mean for precision list in case of Decision Tree is: ", mean_precision_list_DT)
print("\n")
print("Standard Deviation for precision list in case of Decision Tree is: ", sd_precision_list_DT)
print("\n")
mean_recall_score_list_DT=sum(DT_recall_list)/len(DT_recall_list)
variance_recall_score_list_DT=sum([(x-mean_recall_score_list_DT)**2] for x in DT_recall_list)/len(DT_recall_list)
sd_recall_score_list_DT=variance_recall_score_list_DT**0.5

print("Mean for recall score in case of Decision Tree is: ", mean_recall_score_list_DT)
print("\n")
print("Standard Deviation for recall score in case of Decision Tree is: ", sd_recall_score_list_DT)
print("\n")
mean_f1_score_list_DT=sum(DT_f1_list)/len(DT_f1_list)
variance_f1_score_list_DT=sum([(x-mean_f1_score_list_DT)**2] for x in DT_f1_list)/len(DT_f1_list)
sd_f1_score_list_DT=variance_f1_score_list_DT**0.5

print("Mean for f1 score in case of Decision Tree is: ", mean_f1_score_list_DT)
print("\n")
print("Standard Deviation for f1 score in case of Decision Tree is: ", sd_f1_score_list_DT)
print("\n")

Accuracy list is : [0.9318181818181818, 0.9227272727272727, 0.9378787878787879, 0.9424242424242424, 0.9151515151515152, 0.9378787878787879, 0.9166666666666667, 0.9212121212121213, 0.9363636363636364, 0.9212121212121213, 0.9409090909090909, 0.9348848484848485, 0.9333333333333333, 0.9439393939393939, 0.9405454545454545, 0.9318181818181818, 0.9378787878787879, 0.9378787878787879, 0.9257575757575758, 0.9390909090909091, 0.9287878787878788]

Precision list is : [0.9260355203985798, 0.9246987951807228, 0.9449275362318841, 0.9491017964071856, 0.938235394176471, 0.95205808392335, 0.9028571428571428, 0.941176470582335, 0.927374301675977, 0.904669746418605, 0.9342630098459039, 0.8815977175463623, 0.869565217391304, 0.8915626505602408, 0.886178617861787, 0.872246960746268, 0.889495251023192, 0.8657142857142858, 0.8542568542568543, 0.875815834522111, 0.87947840464243]

Recall list is : [0.9398398939893984, 0.9219219219219219, 0.9367816091954022, 0.9378698224852071, 0.9011293435028248, 0.927113702623906, 0.9376854598406528, 0.9090909090909091, 0.954022885057471, 0.9424242424242424, 0.9355724269818768, 0.9228571428571429, 0.9376770538243626, 0.936558912386707, 0.9498607242339833, 0.919402880746268, 0.9322033895835084, 0.9327485380116959, 0.920588252941176, 0.9116801911680191, 0.93063583150289]

F1-Score list is : [0.932935916542474, 0.9233082706766917, 0.9408369408369408, 0.94345232809523809, 0.919383573487031, 0.9394387001477105, 0.919947758369723, 0.9248554913294799, 0.9405099150141643, 0.922848664884273, 0.9448373408769448, 0.937590711756168, 0.9376770538243626, 0.943683409436834, 0.9498607242339833, 0.93192133161876, 0.9415121255349501, 0.939617083946991, 0.9274074074074073, 0.9142857142857144, 0.9319823363863633]

Mean for accuracy in case of Decision Tree is: 0.9310968810966812

Standard Deviation for accuracy in case of Decision Tree is: 0.009936822684379112

Mean for precision list in case of Decision Tree is: 0.9372889098968774

Standard Deviation for precision list in case of Decision Tree is: 0.01503965334025902

Mean for recall score in case of Decision Tree is: 0.9303697847389972

Standard Deviation for recall score in case of Decision Tree is: 0.012847793854870058

Mean for f1 score in case of Decision Tree is: 0.9337053309179472

Standard Deviation for f1 score in case of Decision Tree is: 0.009662905260894883
```

## Question 1.5: Comment on obtained results

### Question 1.2:

If we choose smaller value of k then the smaller value can be noisy while a higher value can be computationally very expensive as it would have smooth boundaries i.e. lower variance and higher bias. So we used cross validation to choose and take small data chunks from the dataset as validation dataset and use it against possible value of k which in our case were given. Then each label was predicted for every option in validation set using  $k = 1, 3, \dots, 37$  so that we can get the best performance and take that value as  $k$  value to minimize error in validation. Usually we take  $k = \sqrt{N}$  to choose best value where  $N$  is no of samples in training dataset and the value of  $k$  is kept odd so that there is no confusion between 2 classes of data. The best  $k$  for k-NN from the given set is 15 as the Accuracy is higher among all other  $k$  after performing Cross Validation.

### Question 1.3:

**Penalty Parameter C** is how much miss classified points you need to avoid or how much you actually care in case of RBF misclassification. So a higher  $C$  value would mean that all the data points need to be classified correctly, as opposed to avoid this issue. The best parameter selection is needed not only to get highest accuracy but also to correctly identify the ideal value of  $C$  as  $C$  is 10. Using this value will prevent isolated regions and over-fitting.

The gamma parameter tells how quickly the Gaussian class boundaries disappear when they move away from support vector. So larger values of gamma would decrease the effect caused by any individual support vector. However, choosing a higher value of gamma results in future datapoints to fall in small clusters. So gamma tuning is needed. Hence we found the best value of gamma to avoid this issue. The best parameter selection is needed not only to get highest accuracy but also to correctly identify the situation that is whether clustering of future information should be together or what should be the margins between classes. The best value of gamma parameter obtained after GridSearchCV is 0.01.

### In Question 1.4:

In the Case of KNN Classifier, for  $k = 15$  the accuracy score, Precision, Recall, F1-Score is calculated. The accuracy score for the model is 78% for both the classes. The Recall and F1-score for Class label '1' is higher than that of class label '1', whereas the Precision is higher for class label '1' than of class label '1'. By evaluating the F1-score from the KNN classification report it is clear that the KNN algorithm is unable to classify the class label '1' as the score is lower than that of class label '1'. After varying the randomness in the dataset for 20 iteration the Mean accuracy score for KNN is 74.423% and the Standard Deviation is 0.021. The Mean F1-Score is 0.683. Mean precision is 0.864 and Mean recall is 0.53 after comparing. This values are nearly the same when compared to the values obtained from that of without iteration, for the KNN classifier.

In the case of SVM Classifier, evaluation metrics recall, f1-score and accuracy score are all very similar. It generalizes much better than k-NN classifier. The f1-score and precision for class label '1' and '1' are very similar and they differ by 0.01 in each case. In addition, the evaluation metrics values are higher than that of k-NN classifier. The average accuracy score for SVM classifier is around 91%. The Result from SVM classifier after randomly varying the Training set for 20 iteration is shown above. In this case the Mean accuracy score is 90.85% and Standard Deviation is 0.012. The Mean Precision is 0.920. Mean Recall is 0.897 and Mean F1-Score is 0.908. The results when compared to without iteration is almost the same. SVM performs better even when there are outliers than that of KNN. KNN performs better than that of SVM when the number of training (m) samples are greater than the number of features (n) (m >> n).

For Gaussian Naive Bayes Classifier, the evaluation metrics have all the same values of 0.88 for both the class labels. The accuracy score is 88%, which is higher than k-NN and lower than that of SVM Classifier. The Result from Naive Bayes classifier after randomly varying the Training set for 20 iteration is shown above. In this case the Mean accuracy score is 86.84% and Standard Deviation is 0.011. The Mean Precision is 0.865. Mean Recall is 0.889 and Mean F1-Score is 0.877. The results when compared to without iteration is almost the same. In this that assumed that the features follow Gaussian Distribution.

The results obtained from Decision Tree classifier is better than that of all the other classifiers mentioned above. The f1-score for both the class labels is 0.93 which is highest among all the algorithm. The precision and recall is very similar in values and differ by only 0.2 and 0.1 between two class labels for precision and recall respectively. This values suggest that Decision Tree is performing better than other mentioned algorithms in this case. The Mean accuracy score after 20 iterations of randomly varying the dataset is 83.174% and the Standard Deviation is 0.011. The Mean Precision is 0.857. Mean Recall is 0.832 and the Mean F1-Score is 0.834. Which is again similar to that of the results obtained without iterations. This values however are higher than other algorithms. Decision Trees understand co-linearly better SVM which makes them have higher performance in categorical data. Decision trees use hyper-rectangles in input space whereas, SVM uses kernel trick which makes decision trees better in SVM.