# Assignment 2: Choosing the Best Parameters to Use for a Binary KNN classifier using on 5-fold cross-validation

Dalia Ibrahim[1] and Carlos Dasaed Salcedo[2]

Studnet ID: [1]201893217, [2]201892008

February 16, 2019

## 1 Introduction

For this assignment, we have implemented a cross-fold validation algorithm from scratch using Python and the following libraries: pandas, numpy, math, random, and sys. To improve the efficiency of the algorithm, unsupervised filtering was also used based on the correlation matrix and variances of the data using sklearn libraries. The other sklearn libraries included in the final program were used to calculate KNN, and related precision metrics. To run the program, the following line must be executed from the command line in Linux:
$python3 A2_t2.py [DataFile.tsv]

## 2 Preliminary Steps - Feature Selection

feature selection
cross correlation
low variance

## 3 Main Pseudo-code

This section only includes the functions that are relevant and currently being used by the algorithm. Functions, such as the ones from sklearn, will only be mentioned in the pseudo-code, but will not be individually described.

### 3.1 FoldSplitter()

---

**Algorithm 1** FoldSplitter Function

---

**Require:** kfolds {kfolds = 5 was used}

1: $data \leftarrow dataframe(DataFile.tsv)$
2: $class0 \leftarrow data.where(class=0).ShuffleRows()$
3: $class0partition \leftarrow RowsInclass0/kfolds$
4: $class1 \leftarrow data.where(class=1).ShuffleRows()$
5: $class1partition \leftarrow RowsInclass1/kfolds$
6: $leftOvers0 \leftarrow$ class0 rows after(class0partition * cvfolds)]
7: $leftOvers1 \leftarrow$ class1 rows after(class1partition * cvfolds)]
8: $leftOvers \leftarrow concatenation(leftOvers0, leftOvers1)$
9: $theFolds \leftarrow newPythonDictionary$
10: **for** i=0 **to** kfolds-1 **do**
11:    $createfold$**i** {i = corresponding iteration in the for cycle}
12:    $class0Range \leftarrow$ class0[**from** i * class0partition **to** (i * class0partition)+class0Partition
13:    $class1Range \leftarrow$ class1[**from** i * class1partition **to** (i * class1partition)+class1Partition
14:    $TempOutput \leftarrow Calculate\_Nearest\_Neighbors(BotKPerTest, K)$
15:    $FinalOutput.append(TempOutput)$
16:    $print(FinalOutput)$
17: **end for**
18: $FinalOutput.csv \leftarrow FinalOutput$

---

## 3.2 EuclideanDistance()

**Algorithm 2** Euclidean Distance Function

**Require:** trainRow , testRow
1: $Summation \leftarrow 0$
2: **for** i **in** range(NumFeatures) **do**
3: $\quad d \leftarrow trainRow[i] - testRow[i]$
4: $\quad Summation \leftarrow Summation + d^2$
5: **end for**
6: $Distance \leftarrow [SQRT(Summation), trainRow[LastValue]]$
7: return Distance

## 3.3 Calculate_Nearest_neighbour()

**Algorithm 3** Calculate Nearest Neighbor Function

**Require:** topMatched , K
1: $TieBreaker \leftarrow NewDataframe()$
2: $TieBreaker \leftarrow UniqueRandomNumbers()$
3: $topMatched.CreateColumn('Count') \leftarrow Count('Class')$
4: $concatenate(topMatched, TieBreaker)$
5: $topMatched = topMatched.Filter('Count' = Count.Max())$
6: $topMatched.sortBy('tieBreaker')$
7: $TopMatch \leftarrow topMatched.row(0)$
8: $Probability \leftarrow \frac{topMatch[Count]}{K}$
9: $FinalOutput \leftarrow [TopMatch[Class], Probability)]$
10: $returnFinalOutput$

## 4 Deciding on Performance

Since there are several different performance metrics that can be used to determine how good an algorithm will perform, we decided to use the AUC of the ROC curve

To prevent possible ties, random numbers without repetition are generated and assigned to each of the selected closest neighbors. In the case of a tie in the number of votes, the class of the neighbor with the smallest random number would be selected.
The following modifications were implemented and tested in an effort to improve the algorithm.

1. Normalization of the training and testing data

2. Addition of weights to the neighbors to improve voting in the selection of the closest neighbor

3. Normalization and Addition of weights together.

## 5 Results

To validate the accuracy of our algorithm, we used crossed validation with the help of the sklearn library. These are the steps we followed to implement Cross Validation and test our code:

1. Shuffle the data using $trainingdata.sample()$ function from sklearn Library

2. Split the data using $trainTestSplit$ function in the sklearn Library and make the testing part equal to 40% of the total training data.

3. Calculate average accuracy by comparing actual data with the predicated data, and then applying the formula in (1)

$$Accuracy = \frac{CorrectPredicted}{AllOutputs} \quad (1)$$

4. Generate a report with Precision and Recall using the $classificationReport$ function from the sklearn Library

5. Repeat all pervious steps for Kfold = 5

## 6 Conclusion

Table 1 shows the results of running KNN with serveral modifications, and K=3.

| Method | Accuracy |
|---|---|
| Classical KNN | 85.00% |
| **Classical KNN with TieBreaker** | **100%** |
| KNN Weighted | 84% |
| Normalized KNN | 91.00% |
| Normalize KNN Weighted | 79.00% |

Table 1: Accuracy calculated using sklearn.

## 7 Conclusion

Since we were using the sklearn to validate our results, we were able to generate a table to compare the accuracy of the results of the algorithm as we made and applied the different modifications. However, to our surprise, Classical KNN with a tie breaker function added yielded accuracy reaches up to 100%. Normalizing the data and adding weights actually reduced our accuracy to 79%.

# 8 Appendix

| K=3 | 1 | 2 | 3 | Class 5 | 6 | 7 | Accurcy |
|---|---|---|---|---|---|---|---|
| Classical KNN | Pre = 79%<br>Recall = 85% | Pre = 88%<br>Recall = 90% | Pre = 67%<br>Recall = 40% | Pre = 100%<br>Recall = 40% | Pre = 80%<br>Recall = 80% | Pre = 100%<br>Recall = 91% | 85.00% |
| KNN Weighted | Pre =75%<br>Recall = 81% | Pre = 93%<br>Recall = 83% | Pre = 57%<br>Recall = 67% | Pre = 100%<br>Recall =100% | Pre = 100%<br>Recall = 100% | Pre = 91%<br>Recall = 91% | 84.00% |
| Normalized KNN | Pre = 92%<br>Recall = 89% | Pre = 91%<br>Recall = 94% | Pre = 60%<br>Recall = 60% | Pre = 100%<br>Recall = 100% | Pre = 100%<br>Recall = 100% | Pre = 100%<br>Recall = 100% | **91.00%** |
| Normalize KNN Weighted | Pre = 71%<br>Recall = 87% | Pre = 89%<br>Recall = 74% | Pre = 60%<br>Recall = 43% | Pre = 75%<br>Recall = 100% | Pre = 50%<br>Recall = 67% | Pre = 100%<br>Recall = 79% | 79.00% |
| **Classical KNN with TieBreaker** | Pre = 100%<br>Recall = 100% | Pre = 100%<br>Recall = 100% | Pre = 100%<br>Recall = 100% | Pre = 100%<br>Recall = 100% | Pre = 100%<br>Recall = 100% | Pre = 100%<br>Recall = 100% | **100.00%** |