
k Nearest Neighbors Python Implementation

Damon Nutter

Cal State San Bernardino

CSE 5160 Machine Learning

nutter.damon.a@gmail.com

Abstract

1 This report is on a k Nearest Neighbors Python implementation without the use of
2 any external python libraries for CSE 5160 Machine Learning. The data set used
3 for this machine learning program was the seeds Data Set from the UCI Machine
4 Learning repository. Many different test splits and k values were chosen for this
5 implementation in order to show the variance in accuracy and to come up with an
6 average overall testing error. Overall the testing accuracy of the program was 89%
7 and the total tests ran were 308.

8 1 Implementation Process

9 The implementation process consisted of a function by function implementation of the overall system.
10 It was initially expected that the overall system would contain a four step process: read the input file,
11 calculate the euclidean distance, get the nearest neighbors, and make the classification prediction.
12 This 4 step process was the goal that was set initially, but it was soon realized that the process would
13 be a 6 step process.
14 The first function that was written was the readData function which simply read the seeds data
15 and wrote it into a dictionary with the key as the classifier value. The second function that was
16 implemented was the normalizeData function which served as a min-max normalization of the data.
17 The data needed to be normalized because there were 7 different attribute values and in order for the
18 data to all be weighted the same, the normalization was necessary. Next, the data was split into a
19 training set and a test set. Then the classifier function was created in order to calculate the euclidean
20 distances of the normalized data and classify each test instance. Lastly, an extra function was added
21 called runClassifier in order to test the accuracy of the program. The main function was created
22 in order to contain a series of runs of different test splits and different k values to test all possible
23 outcomes of the system. Each test split (from 3% to 14%) and k value (from 3 to 25) was tested and
24 the total of classification hits and misses were tallied in order to access the overall accuracy.

25 2 Code

26 2.1 readData Function

```
27 #Read data into a dictionary with number of keys = number of different
28 #classifications
29 def readData(filename):
30     dataSet ={1:[],2:[],3:[]}
31     with open(filename, 'r') as reader:
32         line = reader.readline() #read each line of data
```

```

34     while line != '':
35         attributes = line.split() #split each attribute
36         for i in range(0, len(attributes)):
37             attributes[i] = float(attributes[i])
38         if attributes[len(attributes)-1] == 1:
39             dataSet[1].append(attributes)
40         elif attributes[len(attributes)-1] == 2:
41             dataSet[2].append(attributes)
42         elif attributes[len(attributes)-1] == 3:
43             dataSet[3].append(attributes)
44         else:
45             print("error, classification value should be 1, 2, or 3")
46         line = reader.readline()
47     return dataSet

```

49 2.2 normalizeData Function

```

50 #function to min-max normalize the dataset
51 def normalizeData(dataSet):
52     #find max and minimum values
53     maxVals = [0,0,0,0,0,0]
54     minVals = [0,0,0,0,0,0]
55     for classification in dataSet:
56         for attributes in dataSet[classification]:
57             i = 0
58             while i < len(attributes) - 1:
59                 if maxVals[i] < attributes[i]:
60                     maxVals[i] = attributes[i]
61                 if minVals[i] > attributes[i]:
62                     minVals[i] = attributes[i]
63                 i += 1
64     #normalize the values
65     for classification in dataSet:
66         for attributes in dataSet[classification]:
67             i = 0
68             while i < len(attributes) - 1:
69                 attributes[i] = (attributes[i] - minVals[i])/(maxVals[i] -
70                             minVals[i])
71                 i += 1
72     return dataSet

```

75 2.3 splitSet Function

```

76 #Splits the data into a training set and a test set. Test set has 1/m % of the data.
77 #training set has the remaining 1-1/m % of the data.
78 def splitSet(dataSet, m):
79     count = 0
80     trainingSet = {1:[],2:[],3:[]}
81     testSet = {1:[],2:[],3:[]}
82     percent = round(1/m,2)* 100
83     for classification in dataSet:
84         for attributes in dataSet[classification]:
85             count += 1
86             if count%m != 0:
87                 trainingSet[classification].append(attributes)
88             else:
89                 testSet[classification].append(attributes)
90     return [trainingSet, testSet];

```

93 2.4 classifier Function

```
94 #Function that finds the classification of a test instance.
95 #k is the number of nearest neighbors to check.
96 def classifier(trainingSet, test, k):
97     distance = []
98     for classification in trainingSet:
99         for attributes in trainingSet[classification]:
100             #calc euclidean distance of 7 attributes.
101             i = 0
102             values = []
103             while i < len(attributes) - 1:
104                 values.append((attributes[i] - test[i])**2)
105                 i += 1
106             e_distance = math.sqrt(sum(values))
107             distance.append((e_distance,classification))
108     distance = sorted(distance)[:k] #sort distance and select the first k distances
109
110     count1 = 0; #how many match classification 1
111     count2 = 0; #how many match classification 2
112     count3 = 0; #how many match classification 3
113     #loop over i = k distances to see which classification fits.
114     for i in distance:
115         if i[1] == 1:
116             count1 += 1
117         elif i[1] == 2:
118             count2 += 1
119         elif i[1] == 3:
120             count3 += 1
121
122     if count1>count2 and count1>count3:
123         return 1
124     elif count2>count1 and count2>count3:
125         return 2
126     else:
127         return 3
```

129 2.5 runClassifier Function

```
130 #Function that runs the test set against the training set.
131 #function also returns the classification error
132 def runClassifier(trainingSet, testSet, k):
133     classifications=[]
134     correct = 0
135     wrong = 0
136     for classification in testSet:
137         for attributes in testSet[classification]:
138             cl = classifier(trainingSet,attributes,k)
139             classifications.append(cl)
140             if cl == attributes[len(attributes)-1]:
141                 correct += 1
142             else:
143                 wrong += 1
144     total = correct + wrong
145     percentCorrect = correct/total
146     return percentCorrect
```

149 2.6 main Function

```
150 def main():
151     file = "seeds_dataset.txt" #data file
152
```

```

153     dataSet = readData(file) #data set from data file
154     normData = normalizeData(dataSet) #normalize the dataset
155     #Run the data starting with m = 7 and k = 3
156     avgs = []
157     percentCorrect = []
158     m = 7
159     total = 0
160     while m < 35:
161         percentCorrect.clear()
162         k = 3
163         while k < 25:
164             split = splitSet(normData, m)
165             #split the dataset into training and test set using
166             #m = 1/m % of the data.
167             correct = runClassifier(split[0], split[1], k)
168             #run classifier on 2 data sets.
169             percentCorrect.append(correct)
170             k += 2
171             total += 1
172         if(len(percentCorrect)!= 0):
173             avg = sum(percentCorrect)/len(percentCorrect)
174             avgs.append(avg)
175         else:
176             print("div by zero")
177         m += 1
178     print("For k from 3 to 25.")
179     m = 7
180     i = 0
181     while i < len(avgs):
182         print("m = " + str(m) + ", average correct classification was: " +
183             str(avgs[i]))
184         m += 1
185         k += 2
186         i += 1
187     print("The average accuracy for all testing sizes and k values: " +
188         str(sum(avgs)/len(avgs)))
189     print("Total tests ran: " + str(total))

```

191 3 Output and Results

```

192
193 For k from 3 to 25. Where test set size = 1/m of the total data.
194 m = 7, average correct classification was: 0.9212121212121214
195 m = 8, average correct classification was: 0.8531468531468532
196 m = 9, average correct classification was: 0.8695652173913043
197 m = 10, average correct classification was: 0.7619047619047618
198 m = 11, average correct classification was: 0.8995215311004785
199 m = 12, average correct classification was: 0.946524064171123
200 m = 13, average correct classification was: 0.94318181818182
201 m = 14, average correct classification was: 0.9333333333333335
202 m = 15, average correct classification was: 0.9350649350649352
203 m = 16, average correct classification was: 0.9370629370629371
204 m = 17, average correct classification was: 0.8333333333333333
205 m = 18, average correct classification was: 0.9173553719008264
206 m = 19, average correct classification was: 0.909090909090909
207 m = 20, average correct classification was: 0.5999999999999999
208 m = 21, average correct classification was: 0.9090909090909091
209 m = 22, average correct classification was: 0.8787878787878789
210 m = 23, average correct classification was: 0.8787878787878789
211 m = 24, average correct classification was: 1.0
212 m = 25, average correct classification was: 0.75
213 m = 26, average correct classification was: 0.8863636363636364
214 m = 27, average correct classification was: 0.987012987012987

```

```
215 m = 28, average correct classification was: 1.0
216 m = 29, average correct classification was: 1.0
217 m = 30, average correct classification was: 0.8701298701298701
218 m = 31, average correct classification was: 0.9242424242424243
219 m = 32, average correct classification was: 0.8787878787878789
220 m = 33, average correct classification was: 0.9848484848484849
221 m = 34, average correct classification was: 0.8333333333333333
222 The average accuracy for all testing sizes and k values: 0.894345802438572
223 Total tests ran: 308
```

225 4 Strategy

226 The training set vs test set split function played a pivotal role in creating reliable sized sets for the
227 k-NN classifier to run. The different training sizes varied widely so that there would be a larger test
228 accuracy range and it would be easier to tell which test splits yielded better results. For the k-NN
229 classifier, there was first a split of the set, and then a euclidean distance for a test set value was
230 calculated against each training set value in order to find the k nearest neighbors to the test instance.
231 During these calculations, hits and misses are recorded within the runClassifier function in order to
232 calculate accuracy.

233 5 Conclusion

234 Overall, the python implementation of k-nearest neighbors algorithm yielded acceptable results for
235 the seeds data. I am pleased overall with the performance and I think that perhaps on a larger data set,
236 the program could reach higher accuracy numbers. Furthermore, the program is written such that it
237 can import any data set and run the classifier. Additional test sets should be ran in the future in order
238 to further test the accuracy of the program.