Appendix A

```
 1  import collections
 2  from collections import Counter
 3  import math
 4  from sklearn.metrics import confusion_matrix
 5  import numpy as np
 6  import pandas as pd
 7
 8  # set of all words in docs
 9  vocab = set()
10
11  # list of all words
12  allWords = []
13
14  # count of words in vocab
15  countVocab = 0
16
17  # list of all document labels
18  allDocs = []
19
20  # set of all document labels
21  docTypes = set()
22
23  # dictionary of count of each doc type
24  docCou1nt = {}
25
26  # dictionary of probabilities of each doc type
27  docProbabilities = {}
28
29  # dictionary of concatenated docs of each type
30  docCombined = {}
31
32  # dictionary of word counts for each doc type plus count of words in vocab
33  docWordCount = {}
34
35  # dictionary of dictionaries of word probabilities for each doc type
36  docWordProb = collections.defaultdict(dict)
37
38  # dictionary of minimum word probabilities for each doc type
39  docMinProb = {}
40
41  # total number of documents
42  numDocs = 0
43
44  # label for training classes
45  label = ""
46
47  # list of correct classes for test set
48  testAnswers = []
49
50  # list of assigned classes for test set
51  testGuesses = []
52
53  # open and read training file
54  trainingFile = open(r"C:\Users\jeffp\OneDrive\Documents\GitHub\CIS_678_Project2\forumTraining.data", "r")
55
56  # read data from training file and close
57  for row in trainingFile:
58      words = row.split()
59      label = words[0]
60      allDocs.append(label)
61      docTypes.add(label)
62      if label in docCombined.keys():
63          docCombined[label] += words
64      else:
65          docCombined[label] = words
66      del words[0]
67      allWords += words
68      for w in words:
```

```
69          vocab.add(w)
70 trainingFile.close()
71
72 # count total words in vocab
73 countVocab = len(vocab)
74
75 # count total number of docs and number of each type
76 numDocs = len(allDocs)
77 docCount = Counter(allDocs)
78
79 # calculate probabilities, word counts, and word probabilities for each doc type
80 for doc in docTypes:
81     docProbabilities[doc] = docCount[doc] / numDocs
82     docWordCount[doc] = len(docCombined[doc]) + countVocab
83     wordCount = Counter(docCombined[doc])
84     docMinProb[doc] = 1 / docWordCount[doc]
85     for w in wordCount:
86         docWordProb[doc][w] = (wordCount[w] + 1) / docWordCount[doc]
87
88
89 # classifies a new document given a list of its words
90 def classifyDoc(wordList):
91
92     # bring in global variables used
93     global docTypes
94     global docWordProb
95     global docMinProb
96
97     # set starting maximum probability to negative infinity
98     maxProb = float('-inf')
99
100    # find the most likely class for the given word list
101    guessClass = ""
102    for d in docTypes:
103
104        # add the logs of the probabilities instead of multiplying the base figures to avoid underflow
105        prob = math.log(docProbabilities[d])
106        for word in wordList:
107            if word in docWordProb[d].keys():
108                prob += math.log(docWordProb[d][word])
109            else:
110                prob += math.log(docMinProb[d])
111
112        # no need to convert the log figures back since they are only used in comparisons amongst themselves
113        if prob > maxProb:
114            maxProb = prob
115            guessClass = d
116
117    return guessClass
118
119
120 # sets the guesses and answers for the test set, returns the % correct
121 def naiveBayes(exclusionRate):
122
123    # bring in global variables needed
124    global allWords
125    global vocab
126    global testGuesses
127    global testAnswers
128
129    # clear current guesses and answers
130    testAnswers.clear()
131    testGuesses.clear()
132
133    # open test file and set correct count to 0
134    testFile = open(r"C:\Users\jeffp\OneDrive\Documents\GitHub\CIS_678_Project2\forumTest.data", "r")
135    correctCount = 0
136
```

```
137        # create list of the top n% most common words from the training set
138        mostCommon = [word for word, word_count in Counter(allWords).most_common(int(len(vocab) * exclusionRate))]
139
140        # read in the test docs, remove most common words, and predict correct class
141        for line in testFile:
142            words1 = line.split()
143            testAnswers.append(words1[0])
144            del words1[0]
145            wordsNoCommon = [word for word in words1 if word not in mostCommon]
146            testGuesses.append(classifyDoc(wordsNoCommon))
147
148        testFile.close()
149
150        # check the correct guess rate and return it
151        for i, val in enumerate(testAnswers):
152            if testGuesses[i] == val:
153                correctCount += 1
154        return correctCount / len(testGuesses)
155
156
157 # finds the optimal n% (between .1% and 1%) most common words to exclude for highest correct rate
158 def findOptimalExclusion():
159        exclusionRate = 0.000
160        maxAccuracy = -1
161        optimalExclusion = exclusionRate
162        while exclusionRate < 0.011:
163            accuracy = naiveBayes(exclusionRate)
164            if accuracy > maxAccuracy:
165                maxAccuracy = accuracy
166                optimalExclusion = exclusionRate
167            exclusionRate += 0.001
168        print('optimal exclusion: %r' % "{:.2%}".format(optimalExclusion))
169        return optimalExclusion
170
171
172 # find the optimal exclusion % and use that to create lists of answers and guesses
173 correctRate = naiveBayes(findOptimalExclusion())
174
175 # create labels for confusion matrix and set display settings
176 docTypesList = list(docTypes)
177 pd.set_option('display.max_columns', 500)
178 pd.set_option('display.width', 1000)
179
180 # compare answers and guesses using precision, recall, CM, F1, and misclassification rate
181 # recall and precision method taken from https://stats.stackexchange.com/questions/51296
182 cm = confusion_matrix(testAnswers, testGuesses, labels=docTypesList)
183 recall = np.diag(cm) / np.sum(cm, axis=1)
184 precision = np.diag(cm) / np.sum(cm, axis=0)
185 meanRecall = np.mean(recall)
186 meanPrecision = np.mean(precision)
187 f1 = (meanPrecision * meanRecall) / (meanPrecision + meanRecall)
188 misclassificationRate = 1 - correctRate
189
190 # print results
191 print('Recall: %r' % "{:.2%}".format(meanRecall))
192 print('Precision: %r' % "{:.2%}".format(meanPrecision))
193 print('F1: %r' % "{:.2%}".format(f1))
194 print('Misclassification Rate: %r' % "{:.2%}".format(misclassificationRate))
195 cmDataFrame = pd.DataFrame(cm, index=docTypesList, columns=docTypesList)
196 print(cmDataFrame)
197
198 # write results to files
199 cmDataFrame.to_csv('confusionMatrix.csv')
200 with open("results.txt", 'w') as f:
201        f.write('Recall: %r' % "{:.2%}".format(meanRecall) + '\n')
202        f.write('Precision: %r' % "{:.2%}".format(meanPrecision) + '\n')
203        f.write('F1: %r' % "{:.2%}".format(f1) + '\n')
204        f.write('Misclassification Rate: %r' % "{:.2%}".format(misclassificationRate) + '\n')
```

Appendix A

```
205     f.write('Classes Guesses \n')
206     for i, val in enumerate(testAnswers):
207         f.write(val + ' ' + testGuesses[i] + '\n')
208     f.close()
209
```