# CS532S19: Assignment #8

Due on Tuesday, April 30, 2019

*Nwala, Alexander C*

**Giridharan Ganeshkumar**

## Question 1

**1. Create two datasets the first called Testing, the second called Training. The Training dataset should - consist of 10 text documents for email messages you consider spam, 10 text documents for email messages you consider not spam.Then the Testing dataset should - consist of 10 text documents for email messages you consider spam, consist of 10 text documents for email messages you consider not spam.**

1. This was a manual excercise where we accessed our mail accounts and created 20 mails from the spam folder and 20 mails from the inbox folder.

2. These are created so that we could train and test with naivebayes methods defined in the doc class

3. The below image list the text files or dataset used for this purpose and the same is also uploaded to the git hub.

```
C:\School\ODU\WS\8\A8\Mails>dir /s /b /o:gn
C:\School\ODU\WS\8\A8\Mails\Testing
C:\School\ODU\WS\8\A8\Mails\Training
C:\School\ODU\WS\8\A8\Mails\Testing\NS
C:\School\ODU\WS\8\A8\Mails\Testing\S
C:\School\ODU\WS\8\A8\Mails\Testing\NS\Testing 1.txt
C:\School\ODU\WS\8\A8\Mails\Testing\NS\Testing 10.txt
C:\School\ODU\WS\8\A8\Mails\Testing\NS\Testing 2.txt
C:\School\ODU\WS\8\A8\Mails\Testing\NS\Testing 3.txt
C:\School\ODU\WS\8\A8\Mails\Testing\NS\Testing 4.txt
C:\School\ODU\WS\8\A8\Mails\Testing\NS\Testing 5.txt
C:\School\ODU\WS\8\A8\Mails\Testing\NS\Testing 6.txt
C:\School\ODU\WS\8\A8\Mails\Testing\NS\Testing 7.txt
C:\School\ODU\WS\8\A8\Mails\Testing\NS\Testing 8.txt
C:\School\ODU\WS\8\A8\Mails\Testing\NS\Testing 9.txt
C:\School\ODU\WS\8\A8\Mails\Testing\S\Testing 1.txt
C:\School\ODU\WS\8\A8\Mails\Testing\S\Testing 10.txt
C:\School\ODU\WS\8\A8\Mails\Testing\S\Testing 2.txt
C:\School\ODU\WS\8\A8\Mails\Testing\S\Testing 3.txt
C:\School\ODU\WS\8\A8\Mails\Testing\S\Testing 4.txt
C:\School\ODU\WS\8\A8\Mails\Testing\S\Testing 5.txt
C:\School\ODU\WS\8\A8\Mails\Testing\S\Testing 6.txt
C:\School\ODU\WS\8\A8\Mails\Testing\S\Testing 7.txt
C:\School\ODU\WS\8\A8\Mails\Testing\S\Testing 8.txt
C:\School\ODU\WS\8\A8\Mails\Testing\S\Testing 9.txt
C:\School\ODU\WS\8\A8\Mails\Training\NS
C:\School\ODU\WS\8\A8\Mails\Training\S
C:\School\ODU\WS\8\A8\Mails\Training\NS\Training 1.txt
C:\School\ODU\WS\8\A8\Mails\Training\NS\Training 10.txt
C:\School\ODU\WS\8\A8\Mails\Training\NS\Training 2.txt
C:\School\ODU\WS\8\A8\Mails\Training\NS\Training 3.txt
C:\School\ODU\WS\8\A8\Mails\Training\NS\Training 4.txt
C:\School\ODU\WS\8\A8\Mails\Training\NS\Training 5.txt
C:\School\ODU\WS\8\A8\Mails\Training\NS\Training 6.txt
C:\School\ODU\WS\8\A8\Mails\Training\NS\Training 7.txt
C:\School\ODU\WS\8\A8\Mails\Training\NS\Training 8.txt
C:\School\ODU\WS\8\A8\Mails\Training\NS\Training 9.txt
C:\School\ODU\WS\8\A8\Mails\Training\S\Training 1.txt
C:\School\ODU\WS\8\A8\Mails\Training\S\Training 10.txt
C:\School\ODU\WS\8\A8\Mails\Training\S\Training 2.txt
C:\School\ODU\WS\8\A8\Mails\Training\S\Training 3.txt
C:\School\ODU\WS\8\A8\Mails\Training\S\Training 4.txt
C:\School\ODU\WS\8\A8\Mails\Training\S\Training 5.txt
C:\School\ODU\WS\8\A8\Mails\Training\S\Training 6.txt
```

# Question 2

**2. Using the PCI book modified docclass.py code and test.py Use your Training dataset to train the Naive Bayes classifier. Use your Testing dataset to test the Naive Bayes classifier and report the classification results.**

1. The naivebayes algorithm is used where we can train the dataset to be classified in categories.

2. The implementation is based on SQLite to store the words and the various categories and select statements mentioned in the classifier method returns the classifier object.

3. As listed in the below script we can call the classify method to check which category the current data belongs.

4. Below is the completed class implementation as per the PCI text book, the next script is the main method where call the spam train and not spam train methods.

Listing 1: Python Script

```python
import sqlite3 as sqlite
import re
import math

class docclass:
    def getwords(doc):
        splitter=re.compile('\\W*')
        #print(doc)
        # Split the words by non-alpha characters
        words=[s.lower() for s in splitter.split(doc)
                if len(s)>2 and len(s)<20]

        # Return the unique set of words only
        toreturn = dict([(w,1) for w in words])
        return toreturn
    def sampletrain(cl):
        cl.train('Nobody owns the water.','good')
        cl.train('the quick rabbit jumps fences','good')
        cl.train('buy pharmaceuticals now','bad')
        cl.train('make quick money at the online casino','bad')
        cl.train('the quick brown fox jumps','good')

    def spamTrain(cl):
        cl.train('the the', 'not spam')
        cl.train('cheap cheap cheap banking the', 'spam')
        cl.train('the', 'not spam')
        cl.train('cheap cheap banking banking banking the the', 'spam')
        cl.train('cheap cheap cheap cheap cheap buy buy the', 'spam')
        cl.train('banking the', 'not spam')
        cl.train('buy banking the', 'not spam')
        cl.train('the', 'not spam')
        cl.train('the', 'not spam')
        cl.train('cheap buy dinner the the', 'not spam')

    def spamTrain(cl):
        cl = docclass.naivebayes(docclass.getwords)
        for x in range(11):
```

Page 4 of 9

```python
38              if (x != 0):
39                  fileTemp = "C:\Mails\Training\S\Training " + str(x) +".txt"
40                  with open(fileTemp, "r") as file:
41                      data = file.read()
42                      cl.train(data,"Spam")

44      def notSpamTrain(cl):
45          cl = docclass.naivebayes(docclass.getwords)
46          for x in range(11):
47              if (x != 0):
48                  fileTemp = "C:\\Mails\\Training\\NS\\Training " + str(x) +".
                        txt"
49                  with open(fileTemp, "r") as file:
50                      data = file.read()
51                      cl.train(data,"Not Spam")

53  class classifier:

55      def __init__(self, getfeatures, filename=None):
56          # Counts of feature/category combinations
57          self.fc={}
58          # Counts of documents in each category
59          self.cc={}
60          self.getfeatures=getfeatures
61          self.setdb('autocreated_db_file')

63      def setdb(self, dbfile):
64          self.con=sqlite.connect(dbfile)
65          print('Connected')
66          self.con.execute('create table if not exists fc(feature,
                category, count)')
67          self.con.execute('create table if not exists cc(category, count)
                ')


70      def incf(self, f, cat):
71          count=self.fcount(f, cat)
72          if count==0:
73              self.con.execute("insert into fc values ('%s','%s',1)"
74                              % (f, cat))
75          else:
76              self.con.execute("update fc set count=%d where feature='%
                    s' and category='%s'" % (count+1, f, cat))

78      def fcount(self, f, cat):
79          res=self.con.execute('select count from fc where feature="%s"
                and category="%s"' %(f, cat)).fetchone()
80          if res==None: return 0
81          else: return float(res[0])

83      def incc(self, cat):
84          count=self.catcount(cat)
85          if count==0:
86              self.con.execute("insert into cc values ('%s',1)" % (cat)
                    )
87          else:
```

```python
88                    self.con.execute("update cc set count=%d where category
                        ='%s'" % (count+1,cat))
89
90        def catcount(self,cat):
91            res=self.con.execute('select count from cc where category="%s"'
                %(cat)).fetchone()
92            if res==None: return 0
93            else: return float(res[0])
94
95        def categories(self):
96            cur=self.con.execute('select category from cc');
97            return [d[0] for d in cur]
98
99        def totalcount(self):
100            res=self.con.execute('select sum(count) from cc').fetchone();
101            if res==None: return 0
102            return res[0]
103
104
105        def train(self,item,cat):
106            features=self.getfeatures(item)
107            # Increment the count for every feature with this category
108            for f in features:
109                self.incf(f,cat)
110
111            # Increment the count for this category
112            self.incc(cat)
113            self.con.commit()
114
115        def fprob(self,f,cat):
116            if self.catcount(cat)==0: return 0
117
118            # The total number of times this feature appeared in this
119            # category divided by the total number of items in this
                category
120            return self.fcount(f,cat)/self.catcount(cat)
121
122        def weightedprob(self,f,cat,prf,weight=1.0,ap=0.5):
123            # Calculate current probability
124            basicprob=prf(f,cat)
125
126            # Count the number of times this feature has appeared in
127            # all categories
128            totals=sum([self.fcount(f,c) for c in self.categories()])
129
130            # Calculate the weighted average
131            bp=((weight*ap)+(totals*basicprob))/(weight+totals)
132            return bp
133
134
135
136
137    class naivebayes(classifier):
138        def __init__(self,getfeatures):
139            docclass.classifier.__init__(self,getfeatures)
140            self.thresholds={}
```

```python
141
142
143          def docprob(self,item,cat):
144              features=self.getfeatures(item)
145
146              # Multiply the probabilities of all the features together
147              p=1
148              for f in features: p*=self.weightedprob(f,cat,self.fprob)
149              return p
150
151          def prob(self,item,cat):
152              catprob=self.catcount(cat)/self.totalcount()
153              docprob=self.docprob(item,cat)
154              return docprob*catprob
155
156          def setthreshold(self,cat,t):
157              self.thresholds[cat]=t
158
159          def getthreshold(self,cat):
160              if cat not in self.thresholds: return 1.0
161              return self.thresholds[cat]
162
163          def classify(self,item,default=None):
164              probs={}
165              # Find the category with the highest probability
166              max=0.0
167              for cat in self.categories():
168                  probs[cat]=self.prob(item,cat)
169                  if probs[cat]>max:
170                      max=probs[cat]
171                      best=cat
172
173              # Make sure the probability exceeds threshold*next best
174              for cat in probs:
175                  if cat==best:
176                      continue
177                  if probs[cat]*self.getthreshold(best)>probs[best]:
178                      return default
179              return best
180
181  class fisherclassifier(classifier):
182      def cprob(self,f,cat):
183          # The frequency of this feature in this category
184          clf=self.fprob(f,cat)
185          if clf==0: return 0
186
187          # The frequency of this feature in all the categories
188          freqsum=sum([self.fprob(f,c) for c in self.categories()])
189
190          # The probability is the frequency in this category divided by
191          # the overall frequency
192          p=clf/(freqsum)
193
194          return p
195      def fisherprob(self,item,cat):
196          # Multiply all the probabilities together
```

```python
197              p=1
198              features=self.getfeatures(item)
199              for f in features:
200                      p*=(self.weightedprob(f,cat,self.cprob))
201
202              # Take the natural log and multiply by -2
203              fscore=-2*math.log(p)
204
205              # Use the inverse chi2 function to get a probability
206              return self.invchi2(fscore,len(features)*2)
207
208       def invchi2(self,chi, df):
209              m = chi / 2.0
210              sum = term = math.exp(-m)
211              for i in range(1, df//2):
212                  term *= m / i
213                  sum += term
214              return min(sum, 1.0)
215       def __init__(self,getfeatures):
216              classifier.__init__(self,getfeatures)
217              self.minimums={}
218
219       def setminimum(self,cat,min):
220              self.minimums[cat]=min
221
222       def getminimum(self,cat):
223              if cat not in self.minimums: return 0
224              return self.minimums[cat]
225       def classify(self,item,default=None):
226              # Loop through looking for the best result
227              best=default
228              max=0.0
229              for c in self.categories():
230                  p=self.fisherprob(item,c)
231                  # Make sure it exceeds its minimum
232                  if p>self.getminimum(c) and p>max:
233                      best=c
234                      max=p
235              return best
```

1. Below is the main method where we create the database file

2. The first step it to use the classify, naivebayes classifier is identified for the purpose and we get hold of the classifier.

3. Next step is to call the spam train and not spam train methods

4. Once trained we can call the classify the method that gives results if the mail passed is an spam or not spam.

Listing 2: Python Script

```python
1  import os
2  from subprocess import check_output
3
4  cl = docclass.naivebayes(docclass.getwords)
```

```
5
6   dbfile = 'anwala.db'
7
8   if( os.path.exists(dbfile) ):
9       check_output(['rm', dbfile])
10      print('removed dbfile:', dbfile)
11
12  cl.setdb(dbfile)
13  docclass.spamTrain(cl)
14  docclass.notSpamTrain(cl)
15
16  for x in range(11):
17      if(x != 0):
18          fileTemp = "C:\\Mails\\Testing\\S\\Testing "+ str(x) +".txt"
19          with open(fileTemp, "r") as file:
20              data = file.read()
21              print(cl.classify(data))
22
23  for x in range(11):
24      if(x != 0):
25          fileTemp = "C:\\Mails\\Testing\\NS\Testing "+ str(x) +".txt"
26          with open(fileTemp, "r") as file:
27              data = file.read()
28              print(cl.classify(data))
```