Machine Learning

# ASSIGNMENT 2

Rushikesh Dixit

1001434559

**Initial Setup**

The Naïve Bayes Classifier is implemented using python on the 20 newsgroups dataset using half as the training data as well as half as test data.

```python
5       #Method to read words from a file and clear it using regular expression
6       def read_words(words_file):
7           with open(words_file, 'r') as word:
8               string=word.read().lower()
9               return re.findall(r"[\w']+", string)
```

The above is a read words method which opens the file passed as parameter and reads it in lowercase. Also the file is read using the specified regular expression which helps to get rid of unneeded special characters.

```python
11      #Finding probability of all the words belonging to a class and also avoiding zero probability error
12      def probablityword(myhash, word, denominator):
13          word=word.lower()
14          if word in myhash:
15              return math.log(myhash[word]+1.0)/denominator
16          return math.log(1.0/denominator)
```

Here the above method calculates the probability of each word and also avoids the zero probability error by adding 1 to numerator and adding the total count of words into the denominator.

```python
18      #Finding total probability
19      def probability(filename, myhash, denominator):
20          l = list(map(lambda x: probablityword(myhash, x, denominator), read_words(filename)))
21          return reduce(lambda x,y: x+y, l)
22
```

Here the total probability is calculated by using the previous two methods: read words and probabilityword.

```python
23      #method used to classify the data
24      def classify(filename, mainhash, trainsum, totalcounter):
25          print(filename)
26          mykeys = trainsum.keys()
27          probabilityvalues = list(map(lambda x: probability(filename, mainhash[x], trainsum[x]+totalcounter), mykeys))
28          minval = min(probabilityvalues)
29          maxval = max(probabilityvalues)
30          median = (minval + maxval)/2
31          probabilityvalues = list(map(lambda x: x - maxval, probabilityvalues))
32          denominator = sum(list(map(lambda x: math.exp(x), probabilityvalues)))
33          probabilityvalues = list(map(lambda x: math.exp(x)/denominator, probabilityvalues))
34          maxval = max(probabilityvalues)
35          print(maxval)
36          maxindex = [i for i in range(len(probabilityvalues)) if probabilityvalues[i] == maxval]
37          if(len(maxindex) > 1):
38              print 'Tie'
39              # TODO: Fix here
40          return mykeys[maxindex[0]]
```

The above line of code is used to define the classify method which classifies the test data by calculating max probability values.

```
57      #count the occurences of every unique word
58      count = dict(zip(path, list(map(lambda x: len(filelist[x]), path))))
59      traincount = {}
60      globalcounter = Counter()
61      trainsum = {}
62   ┌ for key in path:
63          print key
64          cwd = os.getcwd()+'/newsgroups/'+key+'/'
65          c = Counter()
66          for file in trainlist[key]:
67              c = c + Counter(read_words(cwd + str(file)))
68          globalcounter = globalcounter + c
69          traincount[key] = dict(c)
70   ⌐     trainsum[key] = sum(traincount[key].values())
71      globalcounter = len(dict(globalcounter).keys())
72
```

This above snippet is used to count the occurrences of unique word from each document for each class using the counter() method from the Collections package.

```
72      #Classify the test data from the train model
73      print(testlist.keys()[0])
74      length=len(testlist.keys())
75   ┌ for i in range(0, length):
76          lengthoffiles=len(testlist[testlist.keys()[i]])
77   ⌐     for j in range(0, lengthoffiles):
78              print("starting classify.....")
79   ⌐         print(classify(testlist[testlist.keys()[i]][j], traincount, trainsum, globalcounter))
```

The last part shown above classifies the test data from the testlist which includes half of the documents from each class.

**Output:**

```
starting classify.....
/Users/Rushi/PycharmProjects/naivebayes/newsgroups/talk.politics.mideast/76089
1.0
talk.politics.mideast
starting classify.....
/Users/Rushi/PycharmProjects/naivebayes/newsgroups/talk.politics.mideast/75877
1.0
talk.politics.mideast
starting classify.....
/Users/Rushi/PycharmProjects/naivebayes/newsgroups/talk.politics.mideast/77246
1.0
talk.politics.mideast

Process finished with exit code 0
```

The above screenshot shows the output of the classifier where it shows the file and its path that is being tested and shows the output of the classifier as talk.politics.mideast.