# CS4025/CS5057: Continuous Assessment 1
# Classifying Emotions

### Philip Hale

### November 25, 2013

## Contents

## 1 NGram models

### 1.1 Six-way Classification.

#### 1.1.1 Ruby

My first instinct was to classify the emotions using an existing Bayesian classification library, to see how this problem has been approached in the past. Using the language I'm most familiar with, I wrote `bayesian_classifier.rb` to perform sentiment analysis on two of the emotions files.

Having confirmed that a simple Bayesian classifier is able to produce reasonable results, I moved on to producing my own system.

The first attempt used Ruby's support for functional and object-oriented styles to create an alternative implementation from the `Sentiment.py` I had worked with previously. Typically, functional and OOP styles of programming give you more maintainable and testable code at the expense of speed and

simplicity. This program is split into various classes and is supported by comprehensive unit tests.

A DataSet class is used to abstract concerns such as reading from disk, combining the emotions files and partitioning the data into different sets. The held-back set of data is split in two (dev- and test-sets) to allow fine-tuning of the algorithm without introducing bias.

Probabilities and word counts given specific sentiments are generated in a more functional approach, and the results are memoized to reduce running-time. For example, `data_set.word_count(word: 'help', sentiment: :fearful)` will give the number of times 'help' occurs in a fearful context.

### 1.1.2 Python

Unfortunately, two major problems led me to abandon this approach and start again in Python. Firstly, the system had trouble consistently classifying the dataset in under 10 seconds, slowing down development. More importantly, it was unable to classify to a satisfactory accuracy. This impeded further development, since without a reliable unigram classifier the rest of the objectives would be hard to meet.

I cut my losses and started afresh using the Sentiment.py script from the practicals.

My approach this time aimed to deviate as little as possible from code which I knew could classify correctly. To this end, I replicated the existing code to support additional sentiments, almost doubling LOC count. At this point, only unigrams were supported.

Some changes had to be made other than just adding extra files. The threshold above which a classification is determined to be correct was increased to 0.16. In fact, both the threshold and the initially assigned probability for each sentiment can be given a starting value of $\frac{1}{6}$ since the classification task involves six sentiments.

In order to aid future development and understanding of the script, I reduced much of the duplication by nesting the various data structures in dictionaries, where the keys are the names of the sentiments. This greatly reduces the amount of repetition in the code and means adding or removing sentiments only requires changing code in one place. For example, to extend the classifier to support a seventh emotion, you would only need to add the datafile and its name to the `SENTIMENTS()` function.

## 1.2 Combining unigram, bigram and trigram models

I was interested in adopting some of the functional patterns from my initial ruby classifier, and so researched functional methods for producing NGrams. Python doesn't quite have the simplicity of `Enumerable#each_cons` for producing NGrams, but a similar result could be achieved with `zip()`.

A few extra lines of code were written to add bigrams and trigrams including the beginning-of-sentence and end-of-sentence tags.

The models were combined by summing the different order NGrams into a single list of 'words'. This has the effect of approximately tripling the number of words in the dictionary. Since the rest of the code is setup to calculate based on the size of this overall list, the calculations still work out.

## 1.3 System performance

The system is limited by the unadjusted plus-one smoothing approach used for probability discounting. Applying a better discounting algorithm such as Good Turing discounting, in combination with a better combination of NGrams such as Katz Backoff, would have resulted in a higher overall level of accuracy. However, combining the NGrams resulted in a higher overall level of accuracy than any individual NGram regardless of order, suggesting that even a simple approach to combining NGram models yields positive results

Here are the results of running the algorithm on both the training and test set of data:

Listing 1: Training Set

| SENTIMENT | ACCURACY | (CORRECT/TOTAL) |
|-----------|----------|-----------------|
| ALL | 0.35 | (749/2161) |
| angry | 0.99 | (144/145) |
| disgusted | 1.00 | (232/232) |
| fearful | 1.00 | (154/154) |
| happy | 0.00 | (3/834) |
| sad | 0.03 | (12/435) |
| surprised | 0.57 | (204/361) |

Listing 2: Test Set

| SENTIMENT | ACCURACY | (CORRECT/TOTAL) |
|-----------|----------|-----------------|
| ALL | 0.25 | (59/237) |
| angry | 0.93 | (14/15) |
| disgusted | 0.43 | (13/30) |
| fearful | 0.77 | (23/30) |
| happy | 0.02 | (2/91) |
| sad | 0.07 | (2/28) |
| surprised | 0.12 | (5/43) |

# 2 Classification with WEKA

## 2.1 Creating the ARFF file

The ARFF file is generated in a separate subroutine to be run after the classification task. It requires pre-processing the words (strictly speaking NGrams) which are most common in the training data, as a proportion of the sentiments they are found in. This function, `mostUseful()`, takes as arguments the pWord values generated in the training step, and the number of 'top' results to return. Higher numbers give a larger number of features.

A separate function is given to write the actual ARFF file, using this mostUseful data structure. Since computing the data lines for the ARFF file was going to require generating NGrams again, the code for turning a sentence into a large set of uni-, bi- and tri-grams was extracted into a separate method.

Writing the ARFF file takes place in two stages. First, the list of attributes is generated and written. These are the $xy$ most 'influential' NGrams in the data set, where $x$ is the specified number of predictors and $y$ the number of sentiments. In order to produce a total of 600 attributes for a 6-way classification task, $x$ is set to 100. The attribute names are surrounded in quotes to allow for punctuation.

Next, the rows of data are written. This requires iterating through the dataset and ascertaining if each of the most useful NGrams occur in that sentence. There is a 'yes' for every useful NGram which appears in the sentence. Finally, the known sentiment is appended to the end of each line, allowing us to perform cross-validation in WEKA.

## 2.2 System Performance

With a starting size of 120 attributes, the ID3 classifier was only able to correctly classify 49% of all instances in the training set. Increasing the attributes count to 600 resulted in an increased successful classification rate of 54%. Critically though, this increased the time taken to build the model in WEKA by more than an order of magnitude, from less than a second to over twenty. The time taken to generate the ARFF file is negligible in either case.

WEKA also provides an open-source version of a better decision-tree algorithm than ID3, namely C4.5 (J48). However, despite taking even longer to run, it is unable to improve on the 54% accuracy given by ID3.

Since the ARFF file is built from data generated by the trainBayes() function, it inherits some of its weaknesses. Improving the manner in which NGrams are combined and discounted would likely improve on these results.

# Appendices

# A Python Source

```python
#!/usr/bin/env python
# -*- coding: utf-8 -*-
import re, random, math, collections, itertools, pdb, operator

# Gives the emotions we are interested in classifying.  To classify an
# additional emotion, add its name here and an accompanying data file with some
# training data.
def SENTIMENTS():
    return ['angry', 'disgusted', 'fearful', 'happy', 'sad', 'surprised']

# The Bayes classifier gives a number between 0 and 1, expressing the
# confidence that the sentence is of a certain sentiment. The threshold
# determines the required certainty for a sentence to be classified as that
# sentiment.
def THRESHOLD():
    return 1.0 / float(len(SENTIMENTS()))

#  Adapted from Scott Triglia. Elegant N-gram Generation in Python. (2013-01-20).
#  URL: http://locallyoptimal.com/blog/2013/01/20/elegant-n-gram-generation-in-python/.
#  Accessed: 2013-11-21.
#  Archived by WebCite at http://www.webcitation.org/6LImC39zN
def makeNgram(wordList, order):
    ngramList = []
    for ngram in zip(*[wordList[i:] for i in range(order)]):
        ngramList.append('_'.join(ngram))
    return ngramList

def makeNGramList(sentence):
```

4

```python
29        wordList = re.findall(r"[\w']+", sentence)#collect all words
30
31        unigramList = []
32        bigramList = []
33        trigramList = []
34
35        unigramList = wordList
36
37        if len(wordList) > 1:
38            bigramList = ["<sen>_"+wordList[0]]
39        bigramList = bigramList + makeNgram(wordList, 2)
40        if len(wordList) > 1:
41            bigramList = bigramList + [wordList[-1]+"_</sen>"]
42
43        if len(wordList) > 2:
44            trigramList = trigramList + ["<sen>_<sen>_"+wordList[0]]
45        if len(wordList) > 1:
46            trigramList = trigramList + ["<sen>_"+wordList[0]+"_"+ wordList[1]]
47        trigramList = trigramList + makeNgram(wordList, 3)
48        if len(wordList) > 2:
49            trigramList = trigramList + [wordList[-2]+"_"+wordList[-1]+"_</sen>"]
50        if len(wordList) > 1:
51            trigramList = trigramList + [wordList[-1]+"_</sen>_</sen>"]
52
53        return (unigramList + bigramList + trigramList)
54
55  def readFiles(sentencesTrain,sentencesTest):
56      for sentiment in SENTIMENTS():
57          txt = open('emotions/' + sentiment + '.txt')
58          for sentence in re.split(r'\n', txt.read()):
59              if random.randint(1,10)<2:
60                  sentencesTest[sentence] = sentiment
61              else:
62                  sentencesTrain[sentence] = sentiment
63
64  # Calculate p(W|Positive), p(W|Negative), p(W) for all words in training data.
65  def trainBayes(sentencesTrain, pWord, freq):
66      for sentiment in SENTIMENTS():
67          freq[sentiment] = {}
68      dictionary = {}
69      wordTotals = { 'all': 0 }
70      for sentiment in SENTIMENTS():
71          wordTotals[sentiment] = 0
72
73      #iterate through each sentence/sentiment pair in the training data
74      for sentence, sentiment in sentencesTrain.iteritems():
75          if sentence == '':
76              continue
77          for word in (makeNGramList(sentence)):
78              wordTotals['all'] += 1
79              if not dictionary.has_key(word):
80                  dictionary[word] = 1
81              wordTotals[sentiment] += 1
82              if not freq[sentiment].has_key(word):
83                  freq[sentiment][word] = 1
84              else:
85                  freq[sentiment][word] += 1
86
87      # smoothing so min count of each word is 1
88      for word in dictionary:
89          for sentiment in SENTIMENTS():
90              if not freq[sentiment].has_key(word):
```

```python
91                  freq[sentiment][word] = 1
92              else:
93                  freq[sentiment][word] += 1
94
95          # divisor for the p(word) calculation
96          freqWordAll = 0
97          for sentiment in SENTIMENTS():
98              freqWordAll += freq[sentiment][word]
99
100         # p(word|sentiment)
101         for sentiment in SENTIMENTS():
102             pWord[sentiment][word] = \
103                 freq[sentiment][word] / float(wordTotals[sentiment])
104
105         #p(word)
106         pWord['all'][word] = freqWordAll / float(wordTotals['all'])
107
108
109 #INPUTS:
110 #  sentences is a dictonary of { sentence: sentiment } for every sentence.
111 #  pWord is dictionary storing p(word) and p(word|sentiment)
112 def testBayes(sentences, pWord):
113     total = {}
114     correct = {}
115     total['all'] = 0
116     correct['all'] = 0
117     for sentiment in SENTIMENTS():
118         total[sentiment] = 0
119         correct[sentiment] = 0
120
121     #for each sentence, sentiment pair in the dataset
122     for sentence, sentiment in sentences.iteritems():
123         if sentence == '':
124             continue
125         p = {}
126         for s in SENTIMENTS():
127             p[s] = THRESHOLD()
128
129         for word in (makeNGramList(sentence)):
130             if pWord['all'].has_key(word):
131                 for s in SENTIMENTS():
132                     if pWord[s][word] > 0.00000001:
133                         p[s] *= pWord[s][word] * 10000
134
135         total['all'] += 1
136         total[sentiment] += 1
137
138         totalProb = float(sum(p.itervalues()))
139         prob = 0
140         prob = p[sentiment] / totalProb
141         if prob >= THRESHOLD():
142             correct['all'] += 1
143             correct[sentiment] += 1
144
145     accuracy = {}
146     accuracy['all'] = correct['all'] / float(total['all'])
147
148     print " (ALL)=%0.2f" % accuracy['all'] + \
149             " (%d" % correct['all'] + "/%d" % total['all'] + ")"
150     for sentiment in SENTIMENTS():
151         accuracy[sentiment] = correct[sentiment] /  float(total[sentiment])
152         print " (" + sentiment + ")=%0.2f" % accuracy[sentiment] + \
```

```python
153                     " (%d" % correct[sentiment] + "/%d" % total[sentiment] + ")"
154
155 def mostUseful(pWord, usefulWords, predictors):
156     proportion = {}
157     for sentiment in SENTIMENTS():
158         proportion[sentiment] = {}
159         for word in pWord[sentiment]:
160             proportion[sentiment][word] = \
161                     pWord[sentiment][word] / pWord['all'][word]
162         topWords = sorted(proportion[sentiment].iteritems(),
163                 key=operator.itemgetter(1), reverse=True)[:predictors]
164         usefulWords[sentiment] = topWords
165         usefulWords['all'] += topWords
166
167 def writeArff(mostUsefulAll, sentences):
168     f = open("emotion.arff", "w")
169
170     f.write("@relation emotion\n\n")
171     for feature in mostUsefulAll:
172         f.write("@attribute \"" + feature[0] + "\" {yes,no}\n")
173     f.write("@attribute emotion {" + ",".join(SENTIMENTS()) + "}\n\n@data\n")
174
175     for sentence,sentiment in sentences.iteritems():
176         if sentence == '':
177             continue
178         nGrams = makeNGramList(sentence)
179         for word in mostUsefulAll:
180             if word[0] in nGrams:
181                 f.write('yes, ')
182             else:
183                 f.write('no, ')
184         f.write(str(sentiment) + "\n")
185     f.close()
186
187 #---------- Main Script -------------------------
188
189 #initialise datasets and dictionaries
190 sentencesTrain={}
191 sentencesTest={}
192 freq = {}
193 pWord={ 'all': {} }
194 for sentiment in SENTIMENTS():
195     pWord[sentiment] = {}
196
197 # split the sentiment files into training and test sets
198 readFiles(sentencesTrain,sentencesTest)
199
200 #build conditional probabilities using training data
201 trainBayes(sentencesTrain, pWord, freq)
202
203 # establish the most predictive patterns
204 usefulWords = { 'all': [] }
205 for sentiment in SENTIMENTS():
206     usefulWords[sentiment] = []
207
208 mostUseful(pWord, usefulWords, 100)
209
210 writeArff(usefulWords['all'], sentencesTrain)
211
212 #run naive bayes classifier on datasets
213 print "Naive Bayes"
214 print "Train Accuracy"
```

```
215  testBayes(sentencesTrain, pWord)
216  print "Test Accuracy"
217  testBayes(sentencesTest, pWord)
```

# B   Functionality of Ruby classifier

```
 1  classification on real data
 2    #classify
 3      should classify known fearful-text as fearful
 4    #classify_dev_set
 5      doesn't crash at runtime
 6
 7  ArrayExtensions#first_percent
 8    rounds up
 9    supports empty arrays
10    with 0 gives the empty array
11    with 10 gives the first 10% of an array
12      should eq [1]
13    with 100 returns the array
14      should eq [:foo, :bar]
15    graceful degradation
16      treats >100 percentages as 100
17      treats <0 percentages as 0
18
19  EmotionClassifier::Classifier
20    should be initialised with some sentiments
21    initialized with data
22      has a training set
23      has a dev-set of data to test against
24      has a held-back set of test data
25    #probability
26      #probability with only word argument gives proportion of that word in the dataset
27      #probability with word and sentiment arguments gives proportion of that word for that sentiment
28
29  EmotionClassifier::DataSet
30    #sentences strips the sentence of punctuation and downcases
31    splits the data into 80/10/10 for training/test/dev
32    can give all the data
33    assigns a sentiment to each sentence
34    #with_sentiment gives the sentences with a given sentiment
35    uses unigrams by default
36    #set_ngram_order(3) makes the data set return trigrams
37    unigrams
38      #words gives an array with every word by default
39      #words with sentiment argument gives words in that sentiment
40    bigrams
41      #words gives an array of bigrams from all sentences
42
43  EmotionClassifier::Emotion
44    #to_string gives the name of the emotion
45    #negate_to_s gives a string representing the opposite emotion
46    #== is true when emotions have the same name
47
48  EmotionClassifier::Ngram
49    #bigrams gives array of bigrams
50    #trigrams gives array of trigrams
51    can give arbitrary numbered n-grams
52    treats word+punctuation as different from word
53    #ngrams works with arrays as well as sentences
```

```
54
55 Finished in 3.16 seconds
56 33 examples, 0 failures
```

# C Simple Bayesian Classifier

```ruby
1  require 'classifier'
2
3  classifier = Classifier::Bayes.new "Angry", "Fearful"
4
5  File.open('emotions/angry.txt').each { |grr| classifier.train_angry grr }
6  File.open('emotions/fearful.txt').each { |eek| classifier.train_fearful eek }
7
8  puts classifier.classify "I hate you!"
9  #=> Angry
10 puts classifier.classify "Please don't hurt me"
11 #=> Fearful
```