# Programming Project Evaluation Form

| Student Name: | Project: | Date: |
|---|---|---|
| **Correctness Criteria:**<br>• Reads normal input data<br>• Handles incorrect input data<br>• Calculates correct results<br>• Outputs results properly<br>• Prints appropriate error messages | | **Score:_____ / 50** |
| **Design Criteria:**<br>• Problem decomposition<br>• Choice of data structures<br>• Choice of algorithms<br>• Program efficiency (space/time) | | **Score:_____ / 20** |
| **Documentation Criteria:**<br>• Comments for classes and methods<br>• Comments describing algorithms<br>• Comments describing data structures<br>• Description of program design<br>• Description of test results<br>• Description of known problems | | **Score:_____ / 10** |
| **Style Criteria:**<br>• Method and variable names<br>• Program indenting<br>• Use of white space<br>• Ordering of methods<br>• Easy to read code | | **Score:_____ / 10** |
| **Testing Criteria:**<br>• Normal input data<br>• Incorrect input data<br>• Special cases for data structures<br>• Special cases for algorithms | | **Score:_____ / 10** |
| **Grader Comments:**<br><br><br><br><br>**Total:_____ / 100** | | |

# 1. Objective

To disambiguate word pairs using a Naive-Bayesian technique and answer to some questions.

# 2. Installation

**Programming language:** `Python`

**Source code location:** `/home/sbillah/nlp2/`

**Corpus Selection:** `I select "`**`senseval"`**` corpus, which is specially designed for WSD. Here is the download link:` **`http://www.senseval.org/.`**

# 3. Design (Naive Bayesian Approach)

## 3.1   Algorithm:

I use two in-memory Hash-tables to store conditional probabilities of contexts associated with the pseudowords.

```
// preprocess corpus files.
for each sense-file f:
      1. replace the word with pseudoword and remember its sense in
         <tag> region.

      2. extract the context words in both side of the pseudoword, and
      store the contexts in two files: training, and testing by 8:2
      ratio.

//training
For each training file f:
      for each line in f:
            1. update C(context-words), C(word), & sense sk in respected
               hash-tables.
      from the counts, compute P(ci|sk), p(sk)and store in hash-tables.

//testing
for each testing file f:
      for each line in f:
            1. apply Laplace smoothing on conditional probabilities.
            2. compute argmax score(sk) using the formula in the book.
            3. compare the predicted value with actual value.

return accuracy in percentage.
```

**Time Complexity:**

- **Preprocessing phase:** O(# of lines containing word1) + O(#of lines containing word2)
- **Training phase:** 2*O(2*context_size * #lines in training file)
- **Testing phase:** O(2*context_size * #lines in testing file)
- **Overall:** O(2*context_size*(#of word1+ #of word2)

**Space Complexity:**

- **Overall:** O(2 * context_size * (#of word1 +  #of word2) )

## 3.2    Corpus Description:

The **Senseval** WSD corpus has total 35 sense-tagged words. Each word has more than 5 senses. But due to the simplified requirement of our homework, I ignore all those senses. Therefore, for a word pair, I consider only two senses (0,1). Here are my selected word-pairs and their individual occurrence in the corpus.

| Pair | Words | Word Counts |
|---|---|---|
| 1 | amaze | 319 |
| | behaviour | 1003 |
| 2 | sack | 296 |
| | sanciton | 101 |
| 3 | knee | 477 |
| | onion | 29 |
| 4 | accident | 1303 |
| | wooden | 370 |

Below is a snapshot of some lines from "***accident.cor***"  file (context for word, 'accident'):

```
800001
Late on Thursday night it was travelling at about three metres a second in
wind blowing at 20 to 25 knots when an empty car fell off just as it reached
the top.
The <tag "532675">accident</> appeared to have little effect on the Christmas
party, except to lengthen it considerably.

800002
An image of earnest Greenery is almost tangible.
Eighteen years ago she lost one of her six children in an <tag
"532675">accident</> on Stratford Road, a tragedy which has become a pawn in
the pitiless point-scoring of small-town vindictiveness.
```

```
800003
It's a sentiment I recommend to you all.
The <tag "532675">accident</> occurred on the Saturday of the annual Popular
Flying Association (PFA) rally at Cranfield.

...
```

## 3.3    Context Selection:

I varied context length from 1 to 19 (on both side) as shown in the figure below:



Different level of accuracy is obtained under different context size. The results are given in the next chapter.

## 3.4    Laplace Smoothing:

During testing phase, some context-words are not seen before in training phase. Instead of assigning zero probability for them, I use Laplace Smoothing. The Laplace Smoothing is given below:
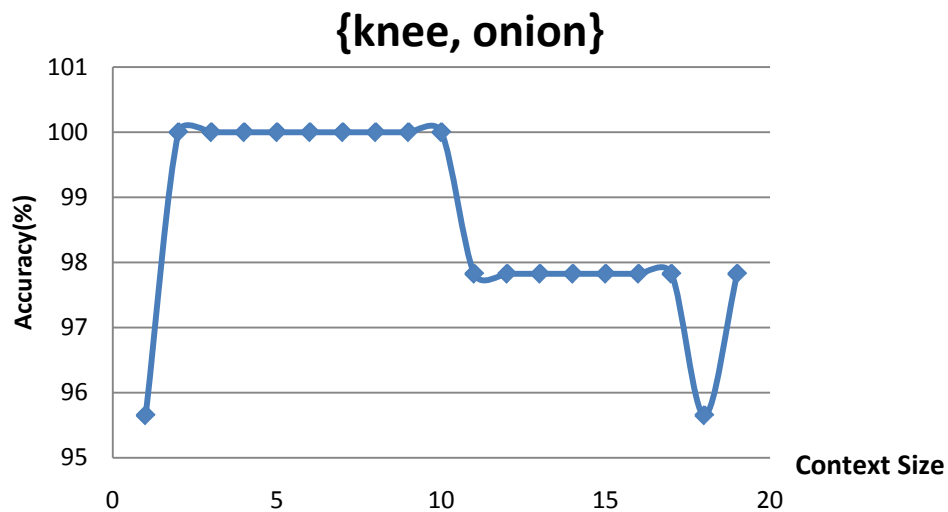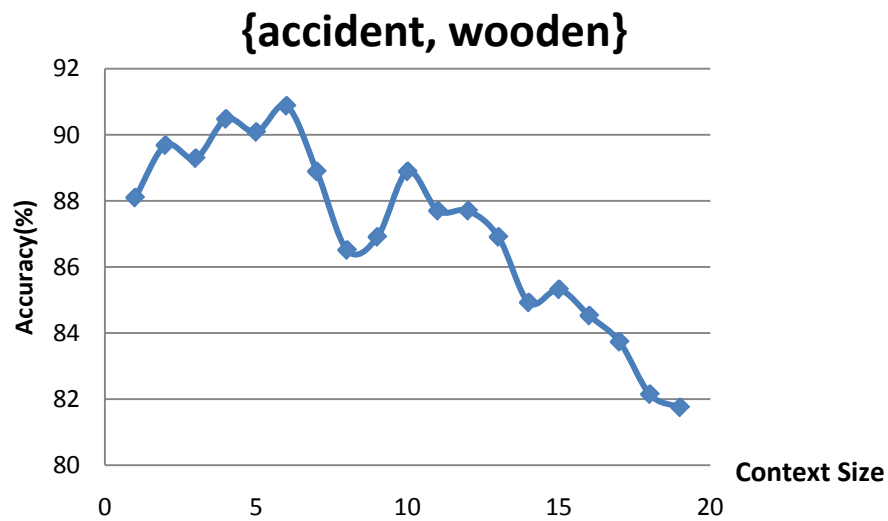
- If $P(ci|s_1) > 0$:
  - $P(ci|s_1)' = (P(ci|s_1)*10000+1)/( 10000 + context\_size)$
- Else:
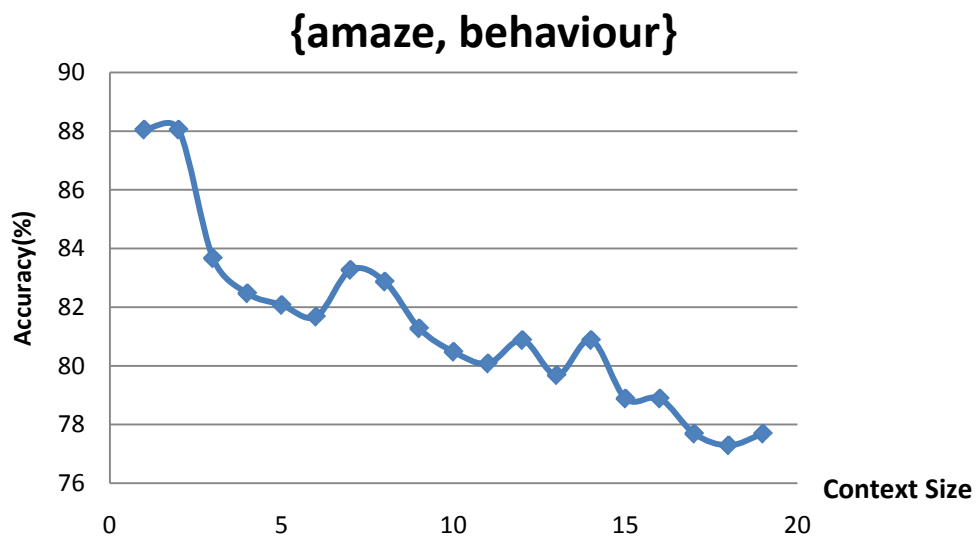  - $P(ci|s_1)' = 1.0/( 10000 + context\_size)$

# 4   Experimental Result

Here, I provide the experimental results for different context size and word-pairs. After the table, I also present these data graphically for better understanding.

| Context Size | Pair 1 accuracy {accident, wooden} | Pair 2 accuracy {knee, onion} | Pair 3 accuracy {sack, sanction} | Pair 4 accuracy {amaze, behaviour} |
|---|---|---|---|---|
| 1 | 88.09524 | 95.65217 | 85.48387 | 88.04781 |
| 2 | 89.68254 | 100 | 85.48387 | 88.04781 |
| 3 | 89.28571 | 100 | 77.41935 | 83.66534 |
| 4 | 90.47619 | 100 | 77.41935 | 82.47012 |
| 5 | 90.07937 | 100 | 79.03226 | 82.07171 |
| 6 | 90.87302 | 100 | 79.03226 | 81.67331 |
| 7 | 88.88889 | 100 | 85.48387 | 83.26693 |
| 8 | 86.50794 | 100 | 87.09677 | 82.86853 |
| 9 | 86.90476 | 100 | 82.25806 | 81.2749 |

| | | | |
|---|---|---|---|
| 10 | 88.88889 | 100 | 82.25806 | 80.47809 |
| 11 | 87.69841 | 97.82609 | 82.25806 | 80.07968 |
| 12 | 87.69841 | 97.82609 | 80.64516 | 80.87649 |
| 13 | 86.90476 | 97.82609 | 82.25806 | 79.68127 |
| 14 | 84.92063 | 97.82609 | 82.25806 | 80.87649 |
| 15 | 85.31746 | 97.82609 | 80.64516 | 78.88446 |
| 16 | 84.52381 | 97.82609 | 80.64516 | 78.88446 |
| 17 | 83.73016 | 97.82609 | 79.03226 | 77.68924 |
| 18 | 82.14286 | 95.65217 | 79.03226 | 77.29084 |
| 19 | 81.74603 | 97.82609 | 75.80645 | 77.68924 |



**{accident, wooden}**



**{knee, onion}**

**{sack, sanction}**



**{amaze, behaviour}**

## Discussion:

1. Accuracy always decrease with the increase of context size.
2. For different word-pairs, maximum accuracy is obtained in different context size.
3. The overall performance of Naive Bayesian disambiguation is around 90%.

# Q&A

2.9 Relative frequency:

**File:** DavidBowie.html
**Entropy:** 3.52728218653

| | | | | | | | | |
|---|---|---|---|---|---|---|---|
| 1 | 2370 | 8 | 917 | h | 7057 | s | 14664 |
| 0 | 1909 | a | 19075 | k | 3306 | r | 12658 |
| 3 | 885 | c | 8595 | j | 423 | u | 4556 |
| 2 | 1502 | b | 4576 | m | 4891 | t | 16047 |
| 5 | 784 | e | 23268 | l | 11101 | w | 5067 |
| 4 | 750 | d | 8301 | o | 11614 | v | 2178 |
| 7 | 881 | g | 3754 | n | 13074 | y | 2469 |
| 6 | 721 | f | 5294 | q | 144 | x | 1038 |
| 9 | 1358 | i | 20024 | p | 7083 | z | 232 |

**File:** Genghis Khan - Wikipedia, the free encyclopedia.html
**Entropy:** 3.55870437955

| | | | | | | | | |
|---|---|---|---|---|---|---|---|
| 1 | 1844 | 8 | 594 | h | 9353 | s | 13442 |
| 0 | 1870 | a | 20479 | k | 4675 | r | 11698 |
| 3 | 873 | c | 7362 | j | 848 | u | 4415 |
| 2 | 1562 | b | 3328 | m | 5568 | t | 15763 |
| 5 | 686 | e | 21399 | l | 11457 | w | 4265 |
| 4 | 566 | d | 7111 | o | 11204 | v | 1788 |
| 7 | 505 | g | 6384 | n | 14745 | y | 2359 |
| 6 | 610 | f | 4749 | q | 263 | x | 1218 |
| 9 | 624 | i | 20236 | p | 6122 | z | 392 |

**File:** Steve Jobs - Wikipedia, the free encyclopedia.html
**Entropy:** 3.61251296863

| | | | | | | | | |
|---|---|---|---|---|---|---|---|
| 1 | 7046 | 8 | 1811 | h | 14120 | s | 31064 |
| 0 | 7798 | a | 36239 | k | 4429 | r | 23907 |
| 3 | 1564 | c | 18155 | j | 2588 | u | 6540 |
| 2 | 6815 | b | 7692 | m | 9194 | t | 32956 |
| 5 | 1808 | e | 47051 | l | 22943 | w | 9403 |
| 4 | 1605 | d | 11058 | o | 23755 | v | 5120 |
| 7 | 1667 | g | 6116 | n | 26050 | y | 4506 |
| 6 | 1966 | f | 10656 | q | 179 | x | 3221 |
| 9 | 2120 | i | 31138 | p | 16747 | z | 504 |

**File:** Winston Churchill - Wikipedia, the free encyclopedia.html
**Entropy:** 3.54645118781

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 6843 | | 8 | 1769 | | H | 21543 | | s | 34129 |
| 0 | 5721 | | a | 45163 | | K | 8332 | | r | 35857 |
| 3 | 1761 | | c | 21142 | | J | 1441 | | u | 10682 |
| 2 | 3863 | | b | 9324 | | M | 12329 | | t | 43826 |
| 5 | 2560 | | e | 56370 | | L | 34877 | | w | 12176 |
| 4 | 1703 | | d | 19968 | | O | 29943 | | v | 5852 |
| 7 | 1273 | | g | 9962 | | N | 34098 | | y | 7013 |
| 6 | 1797 | | f | 13635 | | Q | 445 | | x | 3895 |
| 9 | 2795 | | i | 49316 | | P | 16278 | | z | 654 |

**Finally, the total Corpus frequency:**

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 18103 | | 8 | 5091 | | H | 52073 | | s | 93299 |
| 0 | 17298 | | a | 120956 | | K | 20742 | | r | 84120 |
| 3 | 5083 | | c | 55254 | | J | 5300 | | u | 26193 |
| 2 | 13742 | | b | 24920 | | M | 31982 | | t | 108592 |
| 5 | 5838 | | e | 148088 | | L | 80378 | | w | 30911 |
| 4 | 4624 | | d | 46438 | | O | 76516 | | v | 14938 |
| 7 | 4326 | | g | 26216 | | N | 87967 | | y | 16347 |
| 6 | 5094 | | f | 34334 | | Q | 1031 | | x | 9372 |
| 9 | 6897 | | i | 120714 | | P | 46230 | | z | 1782 |

**2.10 KL Divergence**

| | |
|---|---|
| KL-divergence <1,2> | 0.003688 |
| KL-divergence <2,1> | 0.003688 |
| **KL-divergence <1,3>** | **0.03306** |
| **KL-divergence <3,1>** | **0.03306** |
| KL-divergence <2,3> | 0.028468 |
| KL-divergence <3,2> | 0.028468 |

So, the corpus 1(english1) and corpus 3 (french1) have the highest score. In fact, these two corpus are same and translation of each other, which justifies the result.

```python
1  '''
2  Created on Mar 13, 2013
3
4  @author: Masum
5  '''
6
7  import re
8  from math import log
9  from collections import defaultdict
10 import os
11 import sys
12
13 def replace_tag2(line_text, tag, pseudoword, line_no, context_size=2):  # <tag "532675">
14     header = str(line_no) + ":" + pseudoword + ":" + str(tag) + "\n";
15     line_text = re.sub('<tag "\d+">\s*\w+</>', pseudoword, line_text, 1)
16     line_text = re.sub('[.|"|,|;|:|!|(|)|?|`|\|']', '', line_text);
17     parsed = [tok.lower() for tok in line_text.split(' ') if len(tok)>2];
18     contexts = []
19
20     index = 0;
21     try:
22         index = parsed.index(pseudoword);
23     except:
24         #print line_text+'\n'
25         return None;
26
27     #add left tokens
28     if (index>=context_size):
29         contexts.extend(parsed[index-context_size : index])
30     #add right tokens
31     if len(parsed)>(index+context_size):
32         contexts.extend( parsed[index+1:index+1+context_size])
33     #print contexts
34     line_text = header + (' ').join(contexts) + "\n"
35     return line_text;
36
37 def process_file(file_in, max_lines, file_out, pseudoword, sense, start_line, context_size):
38     mode = 'w+' if (start_line==0) else 'a';
39     f_train = open(file_out+".train.txt", mode);
40     f_test = open(file_out+".test.txt", mode);
41     curr_line = 0;
42     tokenized_line="";
43
44     with open(file_in, 'r') as f:
45         line_text = "";
46         within_a_line = False
47         for line in f:
48             if re.match("\n", line):
49                 # if start_line>3:break;
50                 within_a_line = False;
51                 tokenized_line = replace_tag2(line_text, sense, pseudoword, start_line, context_size)
52                 if tokenized_line:
53                     if curr_line<max_lines*.8:
54                         f_train.write(tokenized_line);
55                     else:
56                         f_test.write(tokenized_line);
```

```python
57                      start_line += 1;
58                  curr_line+=1;
59                  line_text = ""
60                  continue;
61
62              if re.match('\d+', line):
63                  within_a_line = True;
64                  continue;
65
66              if within_a_line:
67                  line_text += line.strip() + " ";
68      f_train.close();
69      f_test.close();
70      return start_line;
71 # end
72
73 def build_corpus(file_in1, line_no1,file_in2, line_no2,file_out, pseudoword, contex_size):
74      # 1st file
75      line_no = process_file(file_in1, line_no1,file_out, pseudoword , 0, 0, contex_size);
76      # 2nd file
77      process_file(file_in2, line_no2,file_out, pseudoword , 1, line_no, contex_size);
78 # end;
79
80
81 def run_traning(file_in, pseudoword):
82      hashes = [defaultdict(float), defaultdict(float)]
83      sense_types= [0,1];
84      count_senses=[0.0, 0.0];
85      current_sense = -1
86
87      with open(file_in, 'r') as f:
88          line_no=0;
89          for line in f:
90              #if line_no>3:break;
91              m = re.match("\d+:"+pseudoword+":(\d)", line) #0:accidentwooden:0
92              if m:
93                  current_sense = int(m.group(1))
94                  count_senses[current_sense]+=1;
95                  line_no += 1;
96              else:
97                  for context in  line.strip().split(" "):
98                      hashes[current_sense][context]+=1;
99          #priors
100         priors = [1.0*count_senses[0]/(count_senses[0]+count_senses[1]), 1.0*count_senses[1]/
    (count_senses[0]+count_senses[1])];
101
102         #conditionals
103         for sense in sense_types:
104             for context in hashes[sense]:
105                 hashes[sense][context] = 1.0*hashes[sense][context]/count_senses[sense];
106
107     #print hashes;
108     return sense_types, priors,hashes ;
109
110
111 def run_disambiguation(file_train, file_test, pseudoword, context_size):
112     sense_types, priors, conditionals =  run_traning(file_train, pseudoword);
```

```
113        actual_sense = -1;
114        predicted_sense = -1;
115        scores  = [0.0, 0.0];
116
117        #print "priors: ", priors;
118        #print "sense types: ",sense_types;
119        #print "training instances: ",count_senses
120        #print conditionals[1];
121
122        factor = 10000;
123        count_corrects=0;
124        count_wrong= 0;
125
126        with open(file_test, 'r') as f:
127            line_no=0;
128            for line in f:
129                m = re.match("\d+:"+pseudoword+":(\d)\n", line) #0:accidentwooden:0
130                if m:
131                    actual_sense = int(m.group(1))
132                    line_no += 1;
133                else:
134                    contexts = line.strip().split(" ");
135                    for sense in sense_types:
136                        scores[sense] = log(priors[sense]);
137
138                        for context in contexts:
139                            if context in conditionals[sense]:
140                                prob = conditionals[sense][context]
141                                scores[sense] += log((prob*factor+1.0)/(factor+context_size));
142                            else:
143                                scores[sense] += log(1.0/factor);
144                            #end inner for
145                        #end sense
146                    #end outer for
147                    if scores[0]>scores[1]:
148                        predicted_sense = 0;
149                    else:
150                        predicted_sense = 1;
151
152                    #calculate accuracy
153                    if predicted_sense==actual_sense:
154                        count_corrects+=1;
155                    else:
156                        count_wrong+=1;
157                #print actual_sense, predicted_sense;
158            #print 'accuracy: ', count_corrects*100.0/(count_corrects+count_wrong), '\n';
159        print context_size, ' ',count_corrects*100.0/(count_corrects+count_wrong);
160
161
162 def run_pair(w1, line_no1, w2, line_no2, context_count):
163     build_corpus('wsd/'+w1+'.cor', line_no1,'wsd/'+w2+'.cor', line_no2,'wsd/'+w1+w2+'.bag',
    w1[2:]+w2[2:], context_count)
164     run_disambiguation('wsd/'+w1+w2+'.bag.train.txt',
    'wsd/'+w1+w2+'.bag.test.txt',w1[2:]+w2[2:], context_count);
165
166
167 def hw3():
```

```python
168     for i in range(1,20):
169         #print 'for run', i,':';
170         run_pair('4.amaze', 319,'4.behaviour', 1003, i);
171         #run_pair('3.sack', 296,'3.sanction', 101, i);
172         #run_pair('2.knee', 477,'2.onion', 29, i);
173         #run_pair('1.accident', 1303,'1.wooden', 370, i);
174 #end
175
176
177 def get_file_chars(fname):
178     alphabet =
    ['a','b','c','d','e','f','g','h','i','j','k','l','m','n','o','p','q','r','s','t','u','v','w','x
    ','y','z','0','1','2','3','4','5','6','7','8','9'];
179     alpha_dict = defaultdict(float);
180
181     total=0
182     with open(fname, 'r') as f:
183         for line in f:
184             for c in line.lower():
185                 total+=1;
186                 if c in alphabet:
187                     alpha_dict[c] += 1.0;
188     return alpha_dict;
189 #end
190
191 def entropy(alpha_dict):
192     info = 0.0
193     total = 0. + sum(alpha_dict.values())
194     for c in alpha_dict:
195         #print c, ' ', alpha_dict[c]
196         p = alpha_dict[c]/total;
197         info += -p*log(p, 2);
198     print 'entropy is: ',info;
199
200 def kldiv(_s, _t):
201     if (len(_s) == 0):return 1e33
202     if (len(_t) == 0):return 1e33
203
204     ssum = 0. + sum(_s.values())
205     tsum = 0. + sum(_t.values())
206
207     vocabdiff = set(_s.keys()).difference(set(_t.keys()))
208     lenvocabdiff = len(vocabdiff)
209
210     """ epsilon """
211     epsilon = min(min(_s.values())/ssum, min(_t.values())/tsum) * 0.001
212
213     """ gamma """
214     gamma = 1 - lenvocabdiff * epsilon
215
216     """ Check if distribution probabilities sum to 1"""
217     sc = sum([v/ssum for v in _s.itervalues()])
218     st = sum([v/tsum for v in _t.itervalues()])
219
220     if sc < 9e-6:
221         print "Sum P: %e, Sum Q: %e" % (sc, st)
222         print "*** ERROR: sc does not sum up to 1. Bailing out .."
```

```
223        sys.exit(2)
224    if st < 9e-6:
225        print "Sum P: %e, Sum Q: %e" % (sc, st)
226        print "*** ERROR: st does not sum up to 1. Bailing out .."
227        sys.exit(2)
228
229    div = 0.
230    for t, v in _s.iteritems():
231        pts = v / ssum
232        ptt = epsilon
233        if t in _t:
234            ptt = gamma * (_t[t] / tsum)
235        ckl = (pts - ptt) * log(pts / ptt)
236        div +=  ckl
237    return div
238 #end of kldiv
239
240
241
242 def get_corpus_count(corpus_path):
243    total_freq = defaultdict(float);
244
245    for filename in os.listdir (corpus_path):
246        #print filename
247        file_count = get_file_chars(os.path.abspath(corpus_path)+os.path.sep+filename);
248        for k in file_count:
249            total_freq[k] += file_count[k];
250        #print '\n'
251    #get_file_chars("kl/English1/test.txt");
252    #for k in total_freq:
253        #print k, ' ', total_freq[k];
254    return total_freq
255
256
257 corpus_path1 = "kl/English1/"
258 corpus_path2 = "kl/English2/"
259 corpus_path3 = "kl/French1/"
260
261 print "KL-divergence <1,2>:",
    kldiv(get_corpus_count(corpus_path1),get_corpus_count(corpus_path2));
262 print "KL-divergence <2,1>:",
    kldiv(get_corpus_count(corpus_path2),get_corpus_count(corpus_path1));
263
264 print "KL-divergence <1,3>:",
    kldiv(get_corpus_count(corpus_path1),get_corpus_count(corpus_path3));
265 print "KL-divergence <3,1>:",
    kldiv(get_corpus_count(corpus_path3),get_corpus_count(corpus_path1));
266
267 print "KL-divergence <2,3>:",
    kldiv(get_corpus_count(corpus_path2),get_corpus_count(corpus_path3));
268 print "KL-divergence <3,2>:",
    kldiv(get_corpus_count(corpus_path3),get_corpus_count(corpus_path2));
269
270
271
```