

Assignment 10

CS532-s16: Web Sciences

Spring 2016

John Berlin

Generated on April 30, 2016

Index

Question 1, 3

Question 2, 5

List of Tables

1	F-Measure knearest nonstemmed	4
2	F-Measure knearest stemmed	4
3	WSDL knearest nonstemmed	4
4	WSDL knearest stemmed	4
5	Genre Ten Fold Non Stemmed Accuracy	5
6	Genre Ten Fold Stemmed Accuracy	6
7	Dr. Nelson Ten Fold Non Stemmed	6
8	Dr. NelsonTen Fold Stemmed	6
9	Genre Ten Fold Non Stemmed	6
10	Genre Ten Fold Stemmed	6
11	Dr. Nelson Ten Fold Non Stemmed	7
12	Dr. NelsonTen Fold Stemmed	7

Listings

1	KNN and SVM Classification	8
2	Proccess Data for A10	14

Question 1

1. Using the data from A8:

- Consider each row in the blog-term matrix as a 500 dimension vector, corresponding to a blog.
- From chapter 8, replace `numpredict.euclidean()` with cosine as the distance metric. In other words, you'll be computing the cosine between vectors of 500 dimensions.
- Use `knestimate()` to compute the nearest neighbors for both:

```
http://f-measure.blogspot.com/  
http://ws-dl.blogspot.com/
```

```
for k={1,2,5,10,20}.
```

Answer

As I did in A8, both the stemmed and non-stemmed words for the blogs were used. The cosign distance used was provided by the [SciPy](#) library. Also from A8 I used the `processData.py` file listing 2 with the relevant methods from `clusters.py` required for this assignment. The code used to generate the results can be seen in listing 1, as usual please consult the comments in the listing for further details.

The k nearest neighbors for F-Measure can be seen in table 1 for non-stemmed and table 2 for stemmed. The results for WSDL can be seen in table 3 for non-stemmed and table 4 for stemmed. Please note that for each k value the previous kth nearest neighbors were included, so I removed them and only include the most unique values for each k value. This means that k=20 only has ten entries, all values for k=10 were apart of the output and can be seen in the k=10 entry.

Just like in A8 the WSDL first nearest neighbor was Dr. Nelsons blog F-Measure but unlike A8 F-Measure had more music blogs in its nearest neighbors.

k nearest neighbors
1 FOLK IS NOT HAPPY
2 Eli Jace — The Mind Is A Terrible Thing To Paste
5 DaveCromwell Writes, The Kids Are Coming Up From Behind, Blog Name Pending
10 The World's First Internet Baby, Did Not Chart, mouxlaloulouda, I Hate The 90s, turnitup!
20 ., I/LOVE/TOTAL/DESTRUCTION, Tremble Under Boom Lights, The Night Mail, Encore, Flatbasset, The Power of Independent Trucking, Steel City Rust, from a voice plantation, Cherry Area

Table 1: F-Measure knearest nonstemmed

k nearest neighbors
1 FOLK IS NOT HAPPY
2 Floorshine Zipper Boots
5 mouxlaloulouda, DaveCromwell Writes, MTJR RANTS & RAVES ON MUSIC
10 Encore, Samtastic! Review, I/LOVE/TOTAL/DESTRUCTION, I Hate The 90s, The Kids Are Coming Up From Behind
20 The World's First Internet Baby, Mile In Mine, Blog Name Pending, from a voice plantation, The Campus Buzz on WSOU, The Night Mail, Did Not Chart, Tremble Under Boom Lights, Interstellar Radio Shower, .

Table 2: F-Measure knearest stemmed

k nearest neighbors
1 F-Measure
2 The Power of Independent Trucking
5 I/LOVE/TOTAL/DESTRUCTION, ., MAGGOT CAVIAR
10 Swinging Singles Club, mouxlaloulouda, Steel City Rust, The World's First Internet Baby, Cherry Area
20 Encore, Flatbasset, DaveCromwell Writes, KiDCHAIR, MTJR RANTS & RAVES ON MUSIC, Blog Name Pending, I Hate The 90s, The Girl at the Rock Show, Eli Jace — The Mind Is A Terrible Thing To Paste, Tremble Under Boom Lights

Table 3: WSDL knearest nonstemmed

k nearest neighbors
1 F-Measure
2 Flatbasset
5 Tremble Under Boom Lights, I Hate The 90s, Riley Haas' blog
10 ., mouxlaloulouda, MTJR RANTS & RAVES ON MUSIC, The Listening Ear, Cherry Area
20 What Am I Doing?, Karl Drinkwater, My Name Is Blue Canary, The World's First Internet Baby, But She's Not Stupid, Samtastic! Review, Pithy Title Here, I/LOVE/TOTAL/DESTRUCTION, Did Not Chart, The Jeopardy of Contentment

Table 4: WSDL knearest stemmed

Question 2

2. Rerun A9, Q2 but this time using LIBSVM. If you have n categories, you'll have to run it n times. For example, if you're classifying music and have the categories:

metal, electronic, ambient, folk, hip-hop, pop

you'll have to classify things as:

metal / not-metal

electronic / not-electronic

ambient / not-ambient

etc.

Use the 500 term vectors describing each blog as the features, and your manually assigned classifications as the true values. Use 10-fold cross-validation (as per slide 46, which shows 4-fold cross-validation) and report the percentage correct for each of your categories.

Answer

As before in A9 the percentage correct is low. Using genre as a label only results in a maximum of a 58% for Pop/Electronic/Hip-Hop accuracy. But when using Dr. Nelsons labeling of his own blog we get a maximum of 66% accuracy for lp review. I have included the metrics requested in A9 along with the accuracy and they had only a marginal increase. I still hold to my conclusion from A9 that using the words alone is not enough to classify these blog posts.

Genre 10 fold accuracy non stemmed can be seen in table 5 and stemmed in table 6. Dr. Nelson label accuracy non stemmed can be seen in table 7 and stemmed in table 8. Genre previous metrics can be seen in table 9 and table 10 for non stemmed and stemmed. Dr. Nelson labels previous metrics can be seen in table 11 and table 12 for stemmed and non stemmed.

label	mean accuracy
Rock, not Rock	49.24%
R&B/Jazz/Mowtown/Country/Other, not R&B/Jazz/Mowtown/Country/Other	56.95%
Pop/Electronic/Hip-Hop, not Pop/Electronic/Hip-Hop	58.39%
Metal/Punk/Hardcore, not Metal/Punk/Hardcore	58.42%
Indie Rock/Alternative, not Indie Rock/Alternative	52.87%

Table 5: Genre Ten Fold Non Stemmed Accuracy

label	mean accuracy
Rock, not Rock	46.29%
R&B/Jazz/Mowtown/Country/Other, not R&B/Jazz/Mowtown/Country/Other	54.84%
Pop/Electronic/Hip-Hop, not Pop/Electronic/Hip-Hop	54.34%
Metal/Punk/Hardcore, not Metal/Punk/Hardcore	55.32%
Indie Rock/Alternative, not Indie Rock/Alternative	50.92%

Table 6: Genre Ten Fold Stemmed Accuracy

label	mean accuracy
the song remains the same, not the song remains the same	56.42%
concert/spotlight, not concert/spotlight	54.37%
lp review, not lp review	66.05%
forgotten song, not forgotten song	46.16%

Table 7: Dr. Nelson Ten Fold Non Stemmed

label	mean accuracy
the song remains the same, not the song remains the same	53.32%
concert/spotlight, not concert/spotlight	49.26%
lp review, not lp review	64.05%
forgotten song, not forgotten song	44.16%

Table 8: Dr. NelsonTen Fold Stemmed

label	precision	recall	f1
Rock	0.50297	0.479219	0.454893
R&B/Jazz/Mowtown/Country/Other	0.508877	0.518228	0.449765
Pop/Electronic/Hip-Hop	0.494767	0.481555	0.446287
Metal/Punk/Hardcore	0.493387	0.492323	0.431816
Indie Rock/Alternative	0.494432	0.48864	0.472867

Table 9: Genre Ten Fold Non Stemmed

label	precision	recall	f1
Rock	0.459143	0.444075	0.43609
R&B/Jazz/Mowtown/Country/Other	0.524861	0.553098	0.458025
Pop/Electronic/Hip-Hop	0.500927	0.508585	0.440113
Metal/Punk/Hardcore	0.515069	0.534176	0.443302
Indie Rock/Alternative	0.497443	0.480556	0.47425

Table 10: Genre Ten Fold Stemmed

label	precision	recall	f1
the song remains the same	0.506409	0.459346	0.434141
concert/spotlight	0.455613	0.431168	0.399722
lp review	0.629662	0.636321	0.61821
forgotten song	0.408316	0.394083	0.382929

Table 11: Dr. Nelson Ten Fold Non Stemmed

label	precision	recall	f1
the song remains the same	0.506409	0.459346	0.434141
concert/spotlight	0.455613	0.431168	0.399722
lp review	0.629662	0.636321	0.61821
forgotten song	0.408316	0.394083	0.382929

Table 12: Dr. NelsonTen Fold Stemmed


```

1 import json
2 import re
3 import statistics
4 from collections import defaultdict
5
6 from sklearn.cross_validation import KFold
7 from sklearn.metrics import f1_score, precision_score, recall_score
8 from tabulate import tabulate
9
10 from numpredict import knnestimate
11 from processData import consume_all
12 from svmutil import svm_problem, svm_predict, svm_parameter,
    svm_train
13 from processData import generate_blogfile, generate_blogfile_stem
14
15 # regex to capture the self labeled topic of the blog post
16 findClass = re.compile("^.+\(((.+))\)$")
17
18 # extract the artist portion. Capture everything until our negative
    look ahead says we have a space - space "
19 artistsExtractor = re.compile("(?!\\s\\-\\s\\")([a-zA-Z0-9.&\\-']+\\s
    (?:[a-zA-Z0-9.&\\-']+\\s)*)"")
20
21 # regex to remove / from my labels for nice file names
22 rem = re.compile("[^a-z]")
23
24 # load the A9 artist to genre file
25 with open("datafiles/artistsToGenre.json", "r") as agi:
26     aTg = json.load(agi)
27
28
29 # read file gotten from clusters.py in A8
30 def readfile(filename):
31     with open(filename, "r") as o:
32         lines = [line.rstrip("\n") for line in o]
33         colnames = lines[0].strip().split('\t')[1:]
34         rownames = []
35         data = []
36         for line in lines[1:]:
37             p = line.strip().split('\t')
38             rownames.append(p[0])
39             data.append({'input': [float(x) for x in p[1:]], 'result':
    p[0]})
40         return rownames, colnames, data
41
42 # using the blogtop500.txt file from A8 do knn on it
43 def nonstemmeda8():
44     # load the data
45     blognames, words, data = readfile("datafiles/blogtop500.txt")
46     # find out where F-Measure is
47     fmeasure_idk = blognames.index('F-Measure')
48     # find out where wsdl is
49     wsdl_idk = blognames.index('Web Science and Digital Libraries
    Research Group')
50     # remove fmeasure from the data as we do not want to include it
51     takeout = data.pop(fmeasure_idk)
52     print("http://f-measure.blogspot.com/")

```

```

53 # hold the results of knn for reporting
54 mout = []
55 klast = None
56 # do knn
57 for k in [1, 2, 5, 10, 20]:
58     # for reporting I am shortening the output to only include
    the new data
59     # successive knn runs return the same data for previous k
    values
60     # ie k=2 [a,b] then k=5 would be [a,b,c,d,e] I only want k
    =2[a,b] and k=5[c,d,e]
61     if k != 1:
62         val = knnestimate(data, data[fmeasure_idk]['input'], k)
        [klast:]
63         else:
64             val = knnestimate(data, data[fmeasure_idk]['input'], k)
65             klast = k
66             mout.append(["%d" % k, ", ".join(val)])
67             print("The %d nearest neighbors are:" % k, knnestimate(data
        , data[fmeasure_idk]['input'], k))
68 # write out a table in latex
69 headers = ['k', 'nearest neighbors']
70 with open("tables/fmeasure-knearest-nonstemmed.text", "w+") as
    out:
71     out.write(tabulate(mout, headers=headers, tablefmt="latex")
    )
72 # do the same for wsdl
73 data.insert(fmeasure_idk, takeout)
74 print("http://ws-dl.blogspot.com/")
75 data.pop(wsdl_idk)
76 mout.clear()
77 for k in [1, 2, 5, 10, 20]:
78     if k != 1:
79         val = knnestimate(data, data[fmeasure_idk]['input'], k)
        [klast:]
80         else:
81             val = knnestimate(data, data[fmeasure_idk]['input'], k)
82             klast = k
83             mout.append(["%d" % k, ", ".join(val)])
84             print("The %d nearest neighbors are:" % k, knnestimate(data
        , data[fmeasure_idk]['input'], k))
85             headers = ['k', 'nearest neighbors']
86             with open("tables/wsdl-knearest-nonstemmed.text", "w+") as out:
87                 out.write(tabulate(mout, headers=headers, tablefmt="latex")
                )
88
89
90 def stemmeda8():
91     # behave exactly as its non-stemmed counterpart except using
    the stemmed file
92     blognames, words, data = readfile("datafiles/blogtop500-stemmed
    .txt")
93     fmeasure_idk = blognames.index('F-Measure')
94     wsdl_idk = blognames.index('Web Science and Digital Libraries
    Research Group')
95     takeout = data.pop(fmeasure_idk)
96     mout = []

```

```

97     klast = None
98     print("http://f-measure.blogspot.com/")
99     for k in [1, 2, 5, 10, 20]:
100         if k != 1:
101             val = knnestimate(data, data[fmeasure_idk]['input'], k)
102             [klast:]
103             else:
104                 val = knnestimate(data, data[fmeasure_idk]['input'], k)
105                 klast = k
106                 mout.append(["%d"%k, " ", " ".join(val)])
107                 print("The %d nearest neighbors are:" % k, knnestimate(data
108 , data[fmeasure_idk]['input'], k))
109             headers = ['k', 'nearest neighbors']
110             with open("tables/fmeasure-knearest-stemmed.text", "w+") as out:
111                 out.write(tabulate(mout, headers=headers, tablefmt="latex")
112 )
113             data.insert(fmeasure_idk, takeout)
114             print("http://ws-dl.blogspot.com/")
115             data.pop(wsd_idk)
116             mout.clear()
117             for k in [1, 2, 5, 10, 20]:
118                 if k != 1:
119                     val = knnestimate(data, data[fmeasure_idk]['input'], k)
120                     [klast:]
121                     else:
122                         val = knnestimate(data, data[fmeasure_idk]['input'], k)
123                         klast = k
124                         mout.append(["%d" % k, " ", " ".join(val)])
125                         print("The %d nearest neighbors are:" % k, knnestimate(data
126 , data[fmeasure_idk]['input'], k))
127                     with open("tables/wsdl-knearest-stemmed.text", "w+") as out:
128                         out.write(tabulate(mout, headers=headers, tablefmt="latex")
129 )
130
131
132
133 # execute question 1
134 def a8():
135     nonstemmeda8()
136     print("_____")
137     stemmeda8()
138
139
140
141 # helper method for classification using SVM
142 # extracts the label for blog structure classification
143 def blog_structure_helper(name):
144     cn = findClass.match(name).group(1).lower().strip()
145     if "concert" in cn or 'spotlight' in cn:
146         cn = 'concert/spotlight'
147     return cn

```

```

147         cn = aTg[m.group(1).rstrip(" -")]
148     else:
149         cn = "R&B/Jazz/Mowtown/Country/Other"
150     return cn
151
152
153 # method to execute the 10 fold SVM classification
154 def tfold_blog(fname, labels, lextactor, outname):
155     """
156     :param fname: the data file to read in
157     :param labels: the labels to be used for classification
158     :param lextactor: the label extractor
159     :param outname: the output file name
160     """
161     # read in the data
162     blognames, words, data = readfile(fname)
163
164     # create lookup table for the blog names
165     blognameidk = {}
166     for i in range(len(blognames)):
167         blognameidk[blognames[i]] = i
168     # reporting data list
169     mout = []
170     mout2 = []
171
172     # do 10 fold validation per label
173     # each label is expanded to label or not label
174     for clazz in labels:
175         print(clazz)
176         # get a fold instance
177         kf = KFold(len(data), 10)
178         # holds the results per fold
179         kfhMetrics = defaultdict(list)
180         ac = []
181         tcIndex = 1
182         # loop through each folds train and text index
183         for train_index, test_index in kf:
184             # holders for data results
185             train = []
186             test = []
187             lab = []
188             tlab = []
189             actualLabels = []
190             # make training data
191             for i in train_index:
192                 clazz1 = lextactor(blognames[i])
193                 if clazz in clazz1:
194                     tlab.append(1)
195                 else:
196                     tlab.append(-1)
197             train.append(data[i]['input'])
198             # make test data
199             for i in test_index:
200                 clazz1 = lextactor(blognames[i])
201                 test.append(data[i]['input'])
202                 if clazz in clazz1:
203                     lab.append(1)

```

```

204         actualLabels.append(clazz)
205     else:
206         lab.append(-1)
207         actualLabels.append("not " + clazz)
208     # create a problem
209     # libsvm 3.21 changed the interface from the example
from Toby Segaran
210     prob = svm_problem(tlab, train)
211     # create parameters: -s 2[one-class svm], -t 0 linear
kernel
212     param = svm_parameter('-s 2 -t 0')
213     # train a model
214     m = svm_train(prob, param)
215     # get the predicted labels, prediction accuracy
216     p_label, p_acc, p_val = svm_predict(lab, test, m)
217     ac.append(p_acc[0])
218     print(p_acc)
219     # transform the labels into reportable format
220     predictedLabels = [clazz if l == 1.0 else "not " +
clazz for l in p_label]
221     # take mean of the values as the scikit learn scores
returns the values for each label
222     kfhMetrics['precision'].append(
223         statistics.mean(precision_score(actualLabels,
predictedLabels, average=None)))
224     kfhMetrics['recall'].append(statistics.mean(
recall_score(actualLabels, predictedLabels, average=None)))
225     kfhMetrics['f1'].append(statistics.mean(f1_score(
actualLabels, predictedLabels, average=None)))
226     ac.append(p_acc[0])
227     tcIndex += 1
228     # create a entry for our report tabel
229     ret1 = [clazz, statistics.mean(kfhMetrics['precision']),
statistics.mean(kfhMetrics['recall']),
230             statistics.mean(kfhMetrics['f1'])]
231     mout.append(ret1)
232     acc = "%.2f" % statistics.mean(ac)
233     mout2.append(["%s, not %s" % (clazz, clazz), acc + "%"])
234
235
236 # write out report
237 headers = ["label", "precision", 'recall', 'f1']
238 with open(outname, "w+") as out:
239     out.write(tabulate(mout, headers=headers, tablefmt="latex")
)
240     headers = ["label", 'mean accuracy']
241     with open(outname+"acc", "w+") as out:
242         out.write(tabulate(mout2, headers=headers, tablefmt="latex"
))
243
244
245 def do_ten_fold():
246     # do ten fold SVM, repeating process done in A9
247     bf_nonstemmed = "datafiles/fmeasure.txt"
248     bf_stemmed = "datafiles/fmeasure_stemmed.txt"
249     blog_structure_labels = ['the song remains the same', 'concert/
spotlight', 'lp review', 'forgotten song']

```

```

250 | blog_genre_labels = ['Rock', 'R&B/Jazz/Mowtown/Country/Other',
251 |                    'Pop/Electronic/Hip-Hop', 'Metal/Punk/Hardcore',
252 |                    'Indie Rock/Alternative']
253 |
254 | tfold_blog(bf_nonstemmed, blog_structure_labels,
255 |           blog_structure_helper, "tables/tenfold-metrics-structure.txt")
256 | tfold_blog(bf_stemmed, blog_structure_labels,
257 |           blog_structure_helper, "tables/tenfold-metrics-structure-
258 |           stemmed.txt")
259 |
260 | tfold_blog(bf_nonstemmed, blog_genre_labels, blog_genre_helper,
261 |           "tables/tenfold-metrics-genre.txt")
262 | tfold_blog(bf_stemmed, blog_genre_labels, blog_genre_helper, "
263 |           tables/tenfold-metrics-genre-stemmed.txt")
264 |
265 | if __name__ == '__main__':
266 |     # read in data-file
267 |     consume_all("datafiles/f-measure.xml")
268 |     generate_blogfile()
269 |     generate_blogfile_stem()
270 |     nonstemmeda8()
271 |     stemmeda8()
272 |     do_ten_fold()

```

Listing 1: KNN and SVM Classification

```

1 import re
2 from collections import defaultdict
3 from operator import add
4 import feedparser
5 import nltk
6 from bs4 import BeautifulSoup
7 from functional import seq
8 from nltk.corpus import stopwords
9 from nltk.stem.snowball import EnglishStemmer
10 import json as jjson
11
12 removeExtra = re.compile('[^a-zA-Z]')
13 stop = stopwords.words('english')
14
15
16 def getwords(doc):
17     splitter = re.compile('[\W*]')
18     # print doc
19     ## Remove all the HTML tags
20     doc = re.compile(r'<[>]+>').sub('', doc)
21     # Split the words by non-alpha characters
22     words = [s.lower().replace('\'', '') for s in nltk.
23               word_tokenize(doc)
24               if len(s) > 2 and len(s) < 20]
25     words = seq(words).map(lambda word: (word, 1)) \
26             .reduce_by_key(add) \
27             .order_by(lambda x: x[1]).to_dict()
28     # Return the unique set of words only
29     return dict([(w, 1) for w in words])
30
31 # modified process text
32 def process_text(text):
33     # clean text
34     t = removeExtra.sub(' ', BeautifulSoup(text.content[0].value, '
35     html5lib').text).lower()
36     ret = []
37     # tokenize according to word and emit word if it is not a
38     # stopword
39     for word in nltk.word_tokenize(t):
40         if len(word) > 2:
41             ret.append(word)
42     ret = seq(ret).map(lambda word: (word, 1)) \
43           .reduce_by_key(add) \
44           .order_by(lambda x: x[1]).to_dict()
45     return {'title': text.title, 'text': ret}
46
47 # modified consume all since
48 def consume_all(fileName):
49     def foldHelp(acum, v):
50         acum[v['title']] = v['text']
51         return acum
52     with open("datafiles/fmeasure.json", "w+") as out:
53         out.write(jjson.dumps(seq(feedparser.parse(fileName).
54                                   entries).map(process_text).fold_left({}, foldHelp), indent=1))
55
56     with open("datafiles/inorder.txt", "w+") as out:
57         seq(feedparser.parse(fileName).entries).map(lambda e: out.

```

```

54     write("%s\n"%e.title)).to_list()
55 # this class represents the word data for a particular feed
56 class feed:
57     # pass flag if we are to do the stemming
58     def __init__(self, fentry, doStem=False):
59         self.title = fentry[0]
60         self.wordCount = fentry[1]
61         self.stemCount = defaultdict(int)
62         if doStem:
63             self._stem_count()
64
65     def _stem_count(self):
66         eng = EnglishStemmer()
67         for word in self.wordCount.keys():
68             self.stemCount[eng.stem(word)] += 1
69
70     def words(self):
71         return set(self.wordCount.keys())
72
73     def __str__(self):
74         return "%s: %s" % (self.title, ' '.join(list(self.wordCount
75 .keys()))))
76
77 # fake tfidf
78 def filter_fun(wc):
79     frac = float(wc[1]) / float(100)
80     return 0.1 < frac < 0.5
81
82
83 # output for the non-stemmed data file
84 def output(f, top):
85     out = [f.title]
86     for wd in top:
87         out.append("%d" % f.wordCount.get(wd, 0))
88     return '\t'.join(out) + "\n"
89
90
91 # output for the stemmed data file
92 def output_stem(f, top):
93     out = [f.title]
94     for wd in top:
95         out.append("%d" % f.stemCount.get(wd, 0))
96     return '\t'.join(out) + "\n"
97
98
99 def generate_blogfile():
100     # read the data in as json and then transform to feeds
101     feedData = seq.json("datafiles/fmeasure.json").map(lambda fe:
102     feed(fe)).to_list() # type: list[feed]
103
104     get the top 500 words by word count over all words:
105     for feed <- feeds, (word,count) <- feed.wordCount.items():
106     yeild (word,count)
107     keep all wordCounts > 10
108     groupby + reduce word:(word1,c1),(word1,c2) -> (word1,c1,c2

```



```

107     ,c3,c4) -> (word1,sumC)
108     keep all wordCounts that meet fake tfidf
109     transform (word,sumC) -> word
110     take top 500
111     transform to list
112     '''
113 top500 = seq(feedData).flat_map(lambda f: list(f.wordCount.
114 items())) \
115     .map(lambda wc: (wc[0], wc[1])) \
116     .reduce_by_key(add) \
117     .filter(filter_fun) \
118     .order_by(lambda wc: -wc[1]) \
119     .map(lambda wc: wc[0]) \
120     .to_list()
121 # sort alphabetically
122 top500 = sorted(top500)
123 print(len(top500))
124 # write resultant to file
125 with open("datafiles/fmeasure.txt", "w+") as out:
126     out.write("Blog\t%s\n" % '\t'.join(top500))
127     for tf in sorted(feedData, key=lambda f: f.title):
128         out.write(output(tf, top500))
129
130 def generate_blogfile_stem():
131     # same as non-stem except use stemmed data
132     feedData = seq.json("datafiles/fmeasure.json").map(lambda fe:
133 feed(fe, True)).to_list() # type: list[feed]
134 top500 = seq(feedData).flat_map(lambda f: list(f.stemCount.
135 items())) \
136     .map(lambda wc: (wc[0], wc[1])) \
137     .reduce_by_key(add) \
138     .filter(filter_fun) \
139     .order_by(lambda wc: -wc[1]) \
140     .map(lambda wc: wc[0]) \
141     .to_list()
142 top500 = sorted(top500)
143 print(len(top500))
144 with open("datafiles/fmeasure_stemmed.txt", "w+") as out:
145     out.write("Blog\t%s\n" % '\t'.join(top500))
146     for tf in sorted(feedData, key=lambda f: f.title):
147         out.write(output_stem(tf, top500))

```

Listing 2: Process Data for A10