

A new ANEW: Evaluation of a word list for sentiment analysis in microblogs

Finn Årup Nielsen

DTU Informatics, Technical University of Denmark, Lyngby, Denmark.
fn@imm.dtu.dk, <http://www.imm.dtu.dk/~fn/>

Abstract. Sentiment analysis of microblogs such as Twitter has recently gained a fair amount of attention. One of the simplest sentiment analysis approaches compares the words of a posting against a labeled word list, where each word has been scored for valence, — a “sentiment lexicon” or “affective word lists”. There exist several affective word lists, e.g., ANEW (Affective Norms for English Words) developed before the advent of microblogging and sentiment analysis. I wanted to examine how well ANEW and other word lists performs for the detection of sentiment strength in microblog posts in comparison with a new word list specifically constructed for microblogs. I used manually labeled postings from Twitter scored for sentiment. Using a simple word matching I show that the new word list may perform better than ANEW, though not as good as the more elaborate approach found in SentiStrength.

1 Introduction

Sentiment analysis has become popular in recent years. Web services, such as socialmention.com, may even score microblog postings on Identi.ca and Twitter for sentiment in real-time. One approach to sentiment analysis starts with labeled texts and uses supervised machine learning trained on the labeled text data to classify the polarity of new texts [1]. Another approach creates a sentiment lexicon and scores the text based on some function that describes how the words and phrases of the text matches the lexicon. This approach is, e.g., at the core of the *SentiStrength* algorithm [2].

It is unclear how the best way is to build a sentiment lexicon. There exist several word lists labeled with emotional valence, e.g., ANEW [3], General Inquirer, OpinionFinder [4], SentiWordNet and WordNet-Affect as well as the word list included in the SentiStrength software [2]. These word lists differ by the words they include, e.g., some do not include strong obscene words and Internet slang acronyms, such as “WTF” and “LOL”. The inclusion of such terms could be important for reaching good performance when working with short informal text found in Internet fora and microblogs. Word lists may also differ in whether the words are scored with sentiment strength or just positive/negative polarity.

I have begun to construct a new word list with sentiment strength and the inclusion of Internet slang and obscene words. Although we have used it for sentiment analysis on Twitter data [5] we have not yet validated it. Data sets with

manually labeled texts can evaluate the performance of the different sentiment analysis methods. Researchers increasingly use Amazon Mechanical Turk (AMT) for creating labeled language data, see, e.g., [6]. Here I take advantage of this approach.

2 Construction of word list

My new word list was initially set up in 2009 for tweets downloaded for on-line sentiment analysis in relation to the United Nation Climate Conference (COP15). Since then it has been extended. The version termed AFINN-96 distributed on the Internet¹ has 1468 different words, including a few phrases. The newest version has 2477 unique words, including 15 phrases that were not used for this study. As SentiStrength² it uses a scoring range from -5 (very negative) to $+5$ (very positive). For ease of labeling I only scored for valence, leaving out, e.g., subjectivity/objectivity, arousal and dominance. The words were scored manually by the author.

The word list initiated from a set of obscene words [7, 8] as well as a few positive words. It was gradually extended by examining Twitter postings collected for COP15 particularly the postings which scored high on sentiment using the list as it grew. I included words from the public domain *Original Balanced Affective Word List*³ by Greg Siegle. Later I added Internet slang by browsing the Urban Dictionary⁴ including acronyms such as WTF, LOL and ROFL. The most recent additions come from the large word list by Steven J. DeRose, *The Compass DeRose Guide to Emotion Words*.⁵ The words of DeRose are categorized but not scored for valence with numerical values. Together with the DeRose words I browsed Wiktionary and the synonyms it provided to further enhance the list. In some cases I used Twitter to determine in which contexts the word appeared. I also used the Microsoft Web n-gram similarity Web service (“Clustering words based on context similarity”⁶) to discover relevant words. I do not distinguish between word categories so to avoid ambiguities I excluded words such as patient, firm, mean, power and frank. Words such as “surprise”—with high arousal but with variable sentiment—were not included in the word list.

Most of the positive words were labeled with $+2$ and most of the negative words with -2 , see the histogram in Figure 1. I typically rated strong obscene words, e.g., as listed in [7], with either -4 or -5 . The word list have a bias towards negative words (1598, corresponding to 65%) compared to positive words (878). A single phrase was labeled with valence 0. The bias corresponds closely to the bias found in the OpinionFinder sentiment lexicon (4911 (64%) negative and 2718 positive words).

¹ http://www2.imm.dtu.dk/pubdb/views/publication_details.php?id=59819

² <http://sentistrength.wlv.ac.uk/>

³ <http://www.sci.sdsu.edu/CAL/wordlist/origwordlist.html>

⁴ <http://www.urbandictionary.com>

⁵ <http://www.deroses.net/steve/resources/emotionwords/ewords.html>

⁶ <http://web-ngram.research.microsoft.com/similarity/>

I compared the score of each word with mean valence of ANEW. Figure 2 shows a scatter plot for this comparison yielding a Spearman’s rank correlation on 0.81 when words are directly matched and including words only in the intersection of the two word lists. I also tried to match entries in ANEW and my word list by applying Porter word stemming (on both word lists) and WordNet lemmatization (on my word list) as implemented in NLTK [9]. The results did not change significantly.

When splitting the ANEW at valence 5 and my list at valence 0 I find a few discrepancies: aggressive, mischief, ennui, hard, silly, alert, mischiefs, noisy. Word stemming generates a few further discrepancies, e.g., alien/alienation, affection/affected, profit/profiteer.

Apart from ANEW I also examined General Inquirer and the OpinionFinder word lists. As these word lists report polarity I associated words with positive sentiment with the valence +1 and negative with -1. I furthermore obtained the sentiment strength from SentiStrength via its Web service⁷ and converted its positive and negative sentiments to one single value by selecting the one with the numerical largest value and zeroing the sentiment if the positive and negative sentiment magnitudes were equal.

3 Twitter data

For evaluating and comparing the word list with ANEW, General Inquirer, OpinionFinder and SentiStrength a data set of 1,000 tweets labeled with AMT was applied. These labeled tweets were collected by Alan Mislove for the *Twittermood*/*“Pulse of a Nation”*⁸ study [10]. Each tweet was rated ten times to get a more reliable estimate of the human-perceived mood, and each rating was a sentiment strength with an integer between 1 (negative) and 9 (positive). The average over the ten values represented the canonical “ground truth” for this study. The tweets were not used during the construction of the word list.

To compute a sentiment score of a tweet I identified words and found the va-

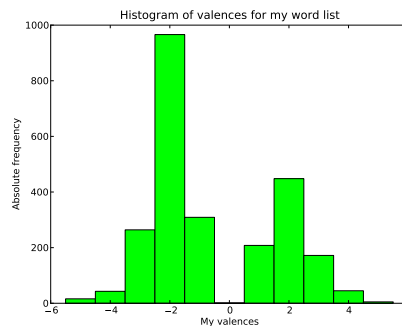


Fig. 1. Histogram of my valences.

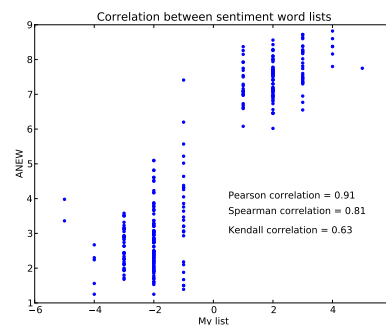


Fig. 2. Correlation between ANEW and my new word list.

⁷ <http://sentistrength.wlv.ac.uk/>

⁸ <http://www.ccs.neu.edu/home/amislove/twittermood/>

Table 1. Example tweet scoring. -5 has been subtracted from the original ANEW score. SentiStrength reported “positive strength 1 and negative strength -2 ”.

Words:	ear	infection	making	it	impossible	2	sleep	headed	2	the	doctors	2	get	new	prescription	so	fucking	early	
My	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	-4	0	-4
ANEW	0	-3.34	0	0	0	0	2.2	0	0	0	0	0	0	0	0	0	0	0	-1.14
GI	0	-1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	-1
OF	0	0	0	0	-1	0	0	0	0	0	0	0	0	0	0	0	0	0	-1
SS																			-2

lence for each word by lookup in the sentiment lexicons. The sum of the valences of the words divided by the number of words represented the combined sentiment strength for a tweet. I also tried a few other weighting schemes: The sum of valence without normalization of words, normalizing the sum with the number of words with non-zero valence, choosing the most extreme valence among the words and quantizing the tweet valences to $+1$, 0 and -1 . For ANEW I also applied a version with match using the NLTK WordNet lemmatizer.

4 Results

My word tokenization identified 15,768 words in total among the 1,000 tweets with 4,095 unique words. 422 of these 4,095 words hit my 2,477 word sized list, while the corresponding number for ANEW was 398 of its 1034 words. Of the 3392 words in General Inquirer I labeled with non-zero sentiment 358 were found in our Twitter corpus and for OpinionFinder this number was 562 from a total of 6442, see Table 1 for a scored example tweet.

I found my list to have a higher correlation (Pearson correlation: 0.564, Spearman’s rank correlation: 0.596, see the scatter plot in Figure 3) with the labeling from the AMT than ANEW had (Pearson: 0.525, Spearman: 0.544). In my application of the General Inquirer word list it did not perform well having a considerable lower AMT correlation than my list and ANEW (Pearson: 0.374, Spearman: 0.422). OpinionFinder with its 90% larger lexicon performed better than General Inquirer but not as good as my list and

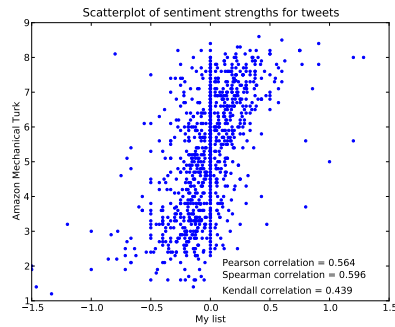


Fig. 3. Scatter plot of sentiment strengths for 1,000 tweets with AMT sentiment plotted against sentiment found by application or my word list.

	My	ANEW	GI	OF	SS
AMT	.564	.525	.374	.458	.610
My		.696	.525	.675	.604
ANEW			.592	.624	.546
GI				.705	.474
OF					.512

Table 2. Pearson correlations between sentiment strength detections methods on 1,000 tweets. AMT: Amazon Mechanical Turk, GI: General Inquirer, OF: OpinionFinder, SS: SentiStrength.

ANEW (Pearson: 0.458, Spearman: 0.491). The SentiStrength analyzer showed superior performance with a Pearson correlation on 0.610 and Spearman on 0.616, see Table 2.

I saw little effect of the different tweet sentiment scoring approaches: For ANEW 4 different Pearson correlations were in the range 0.522–0.526. For my list I observed correlations in the range 0.543–0.581 with the extreme scoring as the lowest and sum scoring without normalization the highest. With quantization of the tweet scores to +1, 0 and −1 the correlation only dropped to 0.548. For the Spearman correlation the sum scoring with normalization for the number of words appeared as the one with the highest value (0.596).

To examine whether the difference in performance between the application of ANEW and my list is due to a different lexicon or a different scoring I looked on the intersection between the two word lists. With a direct match this intersection consisted of 299 words. Building two new sentiment lexicons with these 299 words, one with the valences from my list, the other with valences from ANEW, and applying them on the Twitter data I found that the Pearson correlations were 0.49 and 0.52 to ANEW’s advantage.

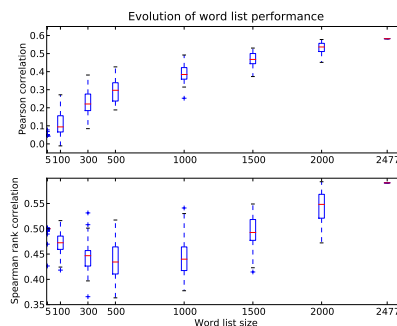


Fig. 4. Performance growth with word list extension from 5 words 2477 words. Upper panel: Pearson, lower: Spearman rank correlation, generated from 50 re-samples among the 2477 words.

5 Discussion

On the simple word list approach for sentiment analysis I found my list performing slightly ahead of ANEW. However the more elaborate sentiment analysis in SentiStrength showed the overall best performance with a correlation to AMT labels on 0.610. This figure is close to the correlations reported in the evaluation of the SentiStrength algorithm on 1,041 MySpace comments (0.60 and 0.56) [2].

Even though General Inquirer and OpinionFinder have the largest word lists I found I could not make them perform as good as SentiStrength, my list and ANEW for sentiment strength detection in microblog posting. The two former lists both score words on polarity rather than strength and it could explain the difference in performance.

Is the difference between my list and ANEW due to better scoring or more words? The analysis of the intersection between the two word list indicated that the ANEW scoring is better. The slightly better performance of my list with the entire lexicon may be due to its inclusion of Internet slang and obscene words.

Newer methods, e.g., as implemented in SentiStrength, use a range of techniques: detection of negation, handling of emoticons and spelling variations [2]. The present application of my list used none of these approaches and might have

benefited. However, the SentiStrength evaluation showed that valence switching at negation and emoticon detection might not necessarily increase the performance of sentiment analyzers (Tables 4 and 5 in [2]).

The evolution of the performance (Figure 4) suggests that the addition of words to my list might still improve its performance slightly.

Although my list comes slightly ahead of ANEW in Twitter sentiment analysis, ANEW is still preferable for scientific psycholinguistic studies as the scoring has been validated across several persons. Also note that ANEW's standard deviation was not used in the scoring. It might have improved its performance.

Acknowledgment I am grateful to Alan Mislove and Sune Lehmann for providing the 1,000 tweets with the Amazon Mechanical Turk labels and to Steven J. DeRose and Greg Siegle for providing their word lists. Mislove, Lehmann and Daniela Balslev also provided input to the article. I thank the Danish Strategic Research Councils for generous support to the 'Responsible Business in the Blogosphere' project.

References

1. Pang, B., Lee, L.: Opinion mining and sentiment analysis. *Foundations and Trends in Information Retrieval* **2**(1-2) (2008) 1–135
2. Thelwall, M., Buckley, K., Paltoglou, G., Cai, D., Kappas, A.: Sentiment strength detection in short informal text. *Journal of the American Society for Information Science and Technology* **61**(12) (2010) 2544–2558
3. Bradley, M.M., Lang, P.J.: Affective norms for English words (ANEW): Instruction manual and affective ratings. Technical Report C-1, The Center for Research in Psychophysiology, University of Florida (1999)
4. Wilson, T., Wiebe, J., Hoffmann, P.: Recognizing contextual polarity in phrase-level sentiment analysis. In: *Proceedings of the conference on Human Language Technology and Empirical Methods in Natural Language Processing*, Stroudsburg, PA, USA, Association for Computational Linguistics (2005)
5. Hansen, L.K., Arvidsson, A., Nielsen, F.Å., Colleoni, E., Etter, M.: Good friends, bad news — affect and virality in Twitter. Accepted for *The 2011 International Workshop on Social Computing, Network, and Services (SocialComNet 2011)* (2011)
6. Akkaya, C., Conrad, A., Wiebe, J., Mihalcea, R.: Amazon Mechanical Turk for subjectivity word sense disambiguation. In: *Proceedings of the NAACL HLT 2010 Workshop on Creating, Speech and Language Data with Amazon's Mechanical Turk*, Association for Computational Linguistics (2010) 195–203
7. Baudhuin, E.S.: Obscene language and evaluative response: an empirical study. *Psychological Reports* **32** (1973)
8. Sapolsky, B.S., Shafer, D.M., Kaye, B.K.: Rating offensive words in three television program contexts. BEA 2008, Research Division (2008)
9. Bird, S., Klein, E., Loper, E.: *Natural Language Processing with Python*. O'Reilly, Sebastopol, California (June 2009)
10. Biever, C.: Twitter mood maps reveal emotional states of America. *The New Scientist* **207**(2771) (July 2010) 14

Supplementary material

A Extra results

AMT	Mislove	ANEW	AFINN-96	AFINN	AFINN(q)	AFINN(s)	AFINN(a)	ANEW(r)	ANEW(l)	ANEW(rs)	ANEW(a)
	0.506		0.546	0.564	0.548	0.581	0.564	0.526	0.527	0.526	0.523
		0.564	0.578	0.588	0.583	0.616	0.780	0.760	0.855	0.950	
			0.970	0.734	0.846	0.802	0.660	0.643	0.580	0.582	
				0.754	0.870	0.824	0.696	0.682	0.599	0.598	
					0.836	0.889	0.531	0.528	0.592	0.606	
						0.884	0.547	0.536	0.628	0.594	
							0.572	0.563	0.606	0.620	
								0.973	0.853	0.819	
									0.829	0.796	
											0.893

Table 3. Pearson correlation matrix. The first data row is the correlation for Amazon Mechanical Turk (AMT) labeling. AFINN is my list. AFINN: sum of word valences with weighting over number of words in text, AFINN(q): (+1, 0, -1)-quantization, AFINN(s): sum of word valences, AFINN(a): average of non-zero word valences, AFINN(x): extreme scoring, ANEW: ANEW with “raw” (direct match) and sum of word valences weighted by number of words, ANEW: with NLTK WordNet lemmatizer word match, ANEW(rs): raw match with sum of word valences, ANEW(a): average of non-zero word valences, GI: General Inquirer, OF: OpinionFinder, SS: SentiStrength.

AMT	Mislove	ANEW	AFINN-96	AFINN	AFINN(q)	AFINN(s)	AFINN(a)	ANEW(r)	ANEW(l)	ANEW(rs)	ANEW(a)
	0.507		0.592	0.596	0.580	0.591	0.581	0.545	0.548	0.537	0.526
		0.630	0.635	0.615	0.625	0.635	0.884	0.857	0.895	0.950	
			0.970	0.912	0.941	0.925	0.656	0.646	0.650	0.647	
				0.941	0.969	0.954	0.668	0.657	0.661	0.656	
					0.944	0.933	0.622	0.620	0.644	0.634	
						0.959	0.630	0.622	0.662	0.644	
							0.641	0.635	0.651	0.644	
								0.972	0.946	0.936	
									0.918	0.905	
											0.943

Table 4. Spearman rank correlation matrix. For explanation of columns see Table 3.

Listings

```
1 #!/usr/bin/env python
2 #
3 # This program go through the list a captures double entries
4 # and ordering problems
5 #
6 # $Id: Nielsen2011New_check.py,v 1.4 2011/03/16 11:34:24 fn Exp $
7
8
9 filebase = '/home/fnielsen/'
10 filename_afinn = filebase + 'fnielsen/data/Nielsen2009Responsible_emotion.csv'
11
12 lines = [ line.split('\t') for line in open(filename_afinn) ]
13 for line in lines:
14     if len(line) != 2:
15         print(line)
16
17 afinn = map(lambda (k,v): (k,int(v)),
18             [ line.split('\t') for line in open(filename_afinn) ])
19 words = [ word for word, valence in afinn ]
20
21
22 swords = sorted(list(set(words)))
23
24 for n in range(len(words)):
25     if words[n] != swords[n]:
26         print(words[n] + " " + swords[n])
27     break
```

```
1 #!/usr/bin/env python
2 #
3 # Construct histogram of AFINN valences.
4 #
5 # $Id: Nielsen2011New_hist.py,v 1.2 2011/03/14 20:53:17 fn Exp $
6
7 import numpy as np
8 import pylab
9 import sys
10 reload(sys)
11 sys.setdefaultencoding('utf-8')
12
13 filebase = '/home/fnielsen/'
14 filename = filebase + 'fnielsen/data/Nielsen2009Responsible_emotion.csv'
15 afinn = dict(map(lambda (k,v): (k,int(v)),
16                [ line.split('\t') for line in open(filename) ]))
17
18 pylab.hist(afinn.values(), bins=11, range=(-5.5, 5.5), facecolor=(0,1,0))
19 pylab.xlabel('My_valences')
20 pylab.ylabel('Absolute_frequency')
21 pylab.title('Histogram_of_valences_for_my_word_list')
22
23 pylab.show()
24
25 # pylab.savefig(filebase + 'fnielsen/eps/' + 'Nielsen2011New_hist.eps')
```

```
1 #!/usr/bin/env python
2 #
3 # $Id: Nielsen2011New_example.py,v 1.2 2011/04/11 17:31:22 fn Exp $
4
5 import csv
```

```

6 import re
7 import numpy as np
8 from nltk.stem.wordnet import WordNetLemmatizer
9 from nltk import sent_tokenize
10 import pylab
11 from scipy.stats.stats import kendalltau, spearmanr
12 import simplejson
13
14 # This variable determines the data set
15 mislove = 2
16
17 # Filenames
18 filebase = '/home/fn/'
19 filename_afinn = filebase + 'fnielsen/data/Nielsen2009Responsible_emotion.csv'
20 filename_afinn96 = 'AFINN-96.txt'
21 filename_mislove1 = "results.txt"
22 filename_mislove2a = "turk.txt"
23 filename_mislove2b = "turk-results.txt"
24 filename_anew = filebase + 'data/ANEW.TXT'
25 filename_gi = filebase + "data/inqtabs.txt"
26 filename_of = filebase + "data/subjclueslen1-HLTEMNLP05.tff"
27
28 # Word splitter pattern
29 pattern_split = re.compile(r"\W+")
30
31
32 afinn = dict(map(lambda (k,v): (k,int(v)),
33                 [ line.split('\t') for line in open(filename_afinn) ]))
34
35
36 afinn96 = dict(map(lambda (k,v): (k,int(v)),
37                  [ line.split('\t') for line in open(filename_afinn96) ]))
38
39
40 # ANEW
41 anew = dict(map(lambda l: (l[0], float(l[2]) - 5) ,
42                [ re.split('\s+', line) for line in open(filename_anew).readlines()[41:1075] ]))
43
44
45 # OpinionFinder
46 of = {}
47 for line in open(filename_of).readlines():
48     elements = re.split('\s|=)', line)
49     if elements[22] == 'positive':
50         of[elements[10]] = +1
51     elif elements[22] == 'negative':
52         of[elements[10]] = -1
53
54
55 # General inquirer
56 csv_reader = csv.reader(open(filename_gi), delimiter='\t')
57 header = []
58 gi = {}
59 previousword = []
60 previousvalence = []
61 for row in csv_reader:
62     if not header:
63         header = row
64     elif len(row) > 2:
65         word = re.search("\w+", row[0].lower()).group()
66         if row[2] == "Positive":
67             valence = 1
68         elif row[3] == "Negativ":
69             valence = -1
70         else:
71             valence = 0
72         if re.search("#", row[0].lower()):
73             if previousword == word:

```

```

74         if previousvalence == []:
75             previousvalence = valence
76         elif previousvalence != valence:
77             previousvalence = 0
78     else:
79         if previousvalence:
80             gi[previousword] = previousvalence
81             previousword = word
82             previousvalence = []
83     elif valence:
84         gi[word] = valence
85
86
87
88 # Lemmatizer for WordNet
89 lemmatizer = WordNetLemmatizer()
90
91
92
93 def words2anewsentiment(words):
94     """
95     Convert words to sentiment based on ANEW via WordNet Lemmatizer
96     """
97     sentiment = 0
98     for word in words:
99         if word in anew:
100             sentiment += anew[word]
101             continue
102         lword = lemmatizer.lemmatize(word)
103         if lword in anew:
104             sentiment += anew[lword]
105             continue
106         lword = lemmatizer.lemmatize(word, pos='v')
107         if lword in anew:
108             sentiment += anew[lword]
109     return sentiment/len(words)
110
111
112 def extremevalence(valences):
113     """
114     Return the most extreme valence. If extremes have different sign then
115     zero is returned
116     """
117     imax = np.argsort(np.abs(valences))[-1]
118     extremes = filter(lambda v: abs(v) == abs(valences[imax]), valences)
119     extremes_samesign = filter(lambda v: v == valences[imax], valences)
120     if extremes == extremes_samesign:
121         return valences[imax]
122     else:
123         return 0
124
125 def corrmatrix2latex(C, columns=None):
126     """
127     s = corrmatrix2latex(C)
128     print(s)
129     """
130     s = '\n\\begin{tabular}{'} + ('r'*(C.shape[0]-1)) + '}\n'
131     if columns:
132         s += "&".join(columns[1:]) + "\\\\n\\hline\n"
133     for n in range(C.shape[0]):
134         row = []
135         for m in range(1, C.shape[1]):
136             if m > n:
137                 row.append("%.3f" % C[n,m])
138             else:
139                 row.append("-----")
140         s += "&".join(row) + "\\\\n"
141     s += '\\end{tabular}\n'

```

```

142     return s
143
144
145 def spearmanmatrix(data):
146     """
147     Spearman r rank correlation matrix
148     """
149     C = np.zeros((data.shape[1], data.shape[1]))
150     for n in range(data.shape[1]):
151         for m in range(data.shape[1]):
152             C[n,m] = spearmanr(data[:,n], data[:,m])[0]
153     return C
154
155
156 def kendallmatrix(data):
157     """
158     Kendall tau rank correlation matrix
159     """
160     C = np.zeros((data.shape[1], data.shape[1]))
161     for n in range(data.shape[1]):
162         for m in range(data.shape[1]):
163             C[n,m] = kendalltau(data[:,n], data[:,m])[0]
164     return C
165
166
167
168 # Read Mislove CSV Twitter data: 'tweets' an array of dictionaries
169 if mislove == 1:
170     csv_reader = csv.reader(open(filename_mislove1), delimiter='\t')
171     header = []
172     tweets = []
173     for row in csv_reader:
174         if not header:
175             header = row
176         else:
177             tweets.append({'id': row[0],
178                          'quant': int(row[1]),
179                          'score_our': float(row[2]),
180                          'score_mean': float(row[3]),
181                          'score_std': float(row[4]),
182                          'text': row[5],
183                          'scores': map(int, row[6:])})
184 elif mislove == 2:
185     if False:
186         tweets = simplejson.load(open('tweets_with_sentistrength.json'))
187     else:
188         csv_reader = csv.reader(open(filename_mislove2a), delimiter='\t')
189         tweets = []
190         for row in csv_reader:
191             tweets.append({'id': row[2],
192                          'score_mislove2': float(row[1]),
193                          'text': row[3]})
194         csv_reader = csv.reader(open(filename_mislove2b), delimiter='\t')
195         tweets_dict = {}
196         header = []
197         for row in csv_reader:
198             if not header:
199                 header = row
200             else:
201                 tweets_dict[row[0]] = {'id': row[0],
202                                       'score_mislove': float(row[1]),
203                                       'score_amt_wrong': float(row[2]),
204                                       'score_amt': np.mean(map(int, re.split("\s+", "_".join(row[3:])))
205                                       )}
206         for n in range(len(tweets)):
207             tweets[n]['score_mislove'] = tweets_dict[tweets[n]['id']]['score_mislove']
208             tweets[n]['score_amt_wrong'] = tweets_dict[tweets[n]['id']]['score_amt_wrong']
209             tweets[n]['score_amt'] = tweets_dict[tweets[n]['id']]['score_amt']

```

```

210
211
212 # Add sentiments to 'tweets'
213 for n in range(len(tweets)):
214     words = pattern_split.split(tweets[n]['text'].lower())
215     tweets[n]['words'] = words
216     afinn_sentiments = map(lambda word: afinn.get(word, 0), words)
217     afinn_sentiment = float(sum(afinn_sentiments))/len(afinn_sentiments)
218     afinn96_sentiments = map(lambda word: afinn96.get(word, 0), words)
219     afinn96_sentiment = float(sum(afinn96_sentiments))/len(afinn96_sentiments)
220     anew_sentiments = map(lambda word: anew.get(word, 0), words)
221     anew_sentiment = float(sum(anew_sentiments))/len(anew_sentiments)
222     gi_sentiments = map(lambda word: gi.get(word, 0), words)
223     gi_sentiment = float(sum(gi_sentiments))/len(gi_sentiments)
224     of_sentiments = map(lambda word: of.get(word, 0), words)
225     of_sentiment = float(sum(of_sentiments))/len(of_sentiments)
226     tweets[n]['sentiment_afinn96'] = afinn96_sentiment
227     tweets[n]['sentiment_afinn'] = afinn_sentiment
228     tweets[n]['sentiment_afinn_quant'] = np.sign(afinn_sentiment)
229     tweets[n]['sentiment_afinn_sum'] = sum(afinn_sentiments)
230     nonzeros = len(filter(lambda nonzero: nonzero, afinn_sentiments))
231     if not nonzeros: nonzeros = 1
232     tweets[n]['sentiment_afinn_nonzero'] = sum(afinn_sentiments)/nonzeros
233     tweets[n]['sentiment_afinn_extreme'] = extremevalence(afinn_sentiments)
234     tweets[n]['sentiment_anew_raw'] = anew_sentiment
235     tweets[n]['sentiment_anew_lemmatize'] = words2anewsentiment(words)
236     tweets[n]['sentiment_anew_raw_sum'] = sum(anew_sentiments)
237     nonzeros = len(filter(lambda nonzero: nonzero, anew_sentiments))
238     if not nonzeros: nonzeros = 1
239     tweets[n]['sentiment_anew_raw_nonzeros'] = sum(anew_sentiments)/nonzeros
240     tweets[n]['sentiment_gi_raw'] = gi_sentiment
241     tweets[n]['sentiment_of_raw'] = of_sentiment
242
243
244 # Index for example tweet
245 words = tweets[index]['words'][:-1]
246 index = 10
247 s = ""
248 s += "Text: &#x2D;\\\multicolumn{%d}{c}{ % len(words) + tweets[index]['text'] + " }-\\\\\\_-[1pt]-\\hline
249 s += "Words: &#x2D; + " &#x2D;".join(words) + " -\\\\\\_-[1pt]-\\n"
250 s += "My&#x2D; + " &#x2D;".join([ str(afinn.get(w,0)) for w in words ]) + " &#x2D;"+ str(sum([ afinn.get(w,0)
251 ])) + " -\\\\\\_-[1pt]-\\n"
252 s += "ANEW&#x2D; + " &#x2D;".join([ str(anew.get(w,0)) for w in words ]) + " &#x2D;"+ str(sum([ anew.get(w,0)
253 ])) + " -\\\\\\_-[1pt]-\\n"
254 s += "GI&#x2D; + " &#x2D;".join([ str(gi.get(w,0)) for w in words ]) + " &#x2D;"+ str(sum([ gi.get(w,0) for
255 ])) + " -\\\\\\_-[1pt]-\\n"
256 s += "OF&#x2D; + " &#x2D;".join([ str(of.get(w,0)) for w in words ]) + " &#x2D;"+ str(sum([ of.get(w,0) for
257 ])) + " -\\\\\\_-[1pt]-\\n"
258 s += "SS_ + " &#x2D; * (len(words)+1) + " -1_ + " \\\_-[1pt]-\\hline_\\n"
259
260 # 'Ear infection making it impossible 2 sleep. headed 2 the doctors 2
261 # get new prescription. so fucking ' has positive strength 1 and negative strength -2
262
263 print(s)
264
265
266 score_amt = np.asarray([ t['score_amt'] for t in tweets ])
267 score_afinn = np.asarray([ t['sentiment_afinn'] for t in tweets ])
268
269 t = ""
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500

```

```

274         sum((1*(score_amt <5)) * (1*(score_afinn >0))),
275         sum((1*(score_amt >5)) * (1*(score_afinn ==0))),
276         sum((1*(score_amt ==5)) * (1*(score_afinn ==0))),
277         sum((1*(score_amt <5)) * (1*(score_afinn ==0))),
278         sum((1*(score_amt >5)) * (1*(score_afinn <0))),
279         sum((1*(score_amt ==5)) * (1*(score_afinn <0))),
280         sum((1*(score_amt <5)) * (1*(score_afinn <0))))
281
282     print(t)
283
284     # 0.1*(277+299+5)

```

```

1  #!/usr/bin/env python
2  #
3  #   The program compares AFINN and ANEW word lists.
4  #
5  #   (Copied from Hansen2010Diffusion and extended.)
6  #
7  # $Id: Hansen2010Diffusion_aneu.py,v 1.3 2010/12/15 15:50:39 fn Exp $
8
9  from nltk.stem.wordnet import WordNetLemmatizer
10 import nltk
11 import numpy as np
12 import pylab
13 import re
14 from scipy.stats.stats import kendalltau, spearmanr
15 import sys
16 reload(sys)
17 sys.setdefaultencoding('utf-8')
18
19
20 filebase = '/home/fn/'
21 filename = filebase + 'fnielsen/data/Nielsen2009Responsible_emotion.csv'
22 afinn = dict(map(lambda (k,v): (k,int(v)),
23                [ line.split('\t') for line in open(filename) ]))
24
25 filename = filebase + 'data/ANEW.TXT'
26 anew = dict(map(lambda l: (l[0], float(l[2])) ,
27                [ re.split('\s+', line) for line in open(filename).readlines()[41:1075] ]))
28
29 lemmatizer = WordNetLemmatizer()
30 stemmer = nltk.PorterStemmer()
31
32 anew_stem = dict([ (stemmer.stem(word), valence) for word, valence in anew.items() ])
33
34
35
36 def word2anewsentiment_raw(word):
37     return anew.get(word, None)
38
39
40 def word2anewsentiment_wordnet(word):
41     sentiment = None
42     if word in anew:
43         sentiment = anew[word]
44     else:
45         lword = lemmatizer.lemmatize(word)
46         if lword in anew:
47             sentiment = anew[lword]
48         else:
49             lword = lemmatizer.lemmatize(word, pos='v')
50             if lword in anew:
51                 sentiment = anew[lword]
52     return sentiment
53
54
55 def word2anewsentiment_stem(word):

```

```

56     return anew_stem.get(stemmer.stem(word), None)
57
58
59
60 sentiments_raw = []
61 for word in afinn.keys():
62     sentiment_anew = word2anewsentiment_raw(word)
63     if sentiment_anew:
64         sentiments_raw.append((afinn[word], sentiment_anew))
65
66
67 sentiments_wordnet = []
68 for word in afinn.keys():
69     sentiment_anew = word2anewsentiment_wordnet(word)
70     if sentiment_anew:
71         sentiments_wordnet.append((afinn[word], sentiment_anew))
72         if (afinn[word] > 0 and sentiment_anew < 5) or \
73             (afinn[word] < 0 and sentiment_anew > 5):
74             print(word)
75
76
77 sentiments_stem = []
78 for word in afinn.keys():
79     sentiment_stem_anew = word2anewsentiment_stem(word)
80     if sentiment_stem_anew:
81         sentiments_stem.append((afinn[word], sentiment_stem_anew))
82         if (afinn[word] > 0 and sentiment_stem_anew < 5) or \
83             (afinn[word] < 0 and sentiment_stem_anew > 5):
84             print(word)
85
86
87 sentiments_raw = np.asarray(sentiments_raw)
88 pylab.figure(1)
89 pylab.plot(sentiments_raw[:,0], sentiments_raw[:,1], '. ')
90 pylab.xlabel('Our_list')
91 pylab.ylabel('ANEW')
92 pylab.title('Correlation between sentiment_word_lists (Direct_match)')
93 pylab.text(1, 3, "Pearson correlation = %.2f" % np.corrcoef(sentiments_raw.T)[1,0])
94 pylab.text(1, 2.5, "Spearman correlation = %.2f" % spearmanr(sentiments_raw[:,0], sentiments_raw[:,1])[0])
95 pylab.text(1, 2, "Kendall correlation = %.2f" % kendalltau(sentiments_raw[:,0], sentiments_raw[:,1])[0])
96 # pylab.savefig(filebase + 'fnielsen/eps/' + 'Nielsen2011New_anew_raw.eps')
97
98 # pylab.show()
99
100
101
102 sentiments = np.asarray(sentiments_wordnet)
103 pylab.figure(2)
104 pylab.plot(sentiments[:,0], sentiments[:,1], '. ')
105 pylab.xlabel('Our_list')
106 pylab.ylabel('ANEW')
107 pylab.title('Correlation between sentiment_word_lists (WordNet_lemmatizer)')
108 pylab.text(1, 3, "Pearson correlation = %.2f" % np.corrcoef(sentiments.T)[1,0])
109 pylab.text(1, 2.5, "Spearman correlation = %.2f" % spearmanr(sentiments[:,0], sentiments[:,1])[0])
110 pylab.text(1, 2, "Kendall correlation = %.2f" % kendalltau(sentiments[:,0], sentiments[:,1])[0])
111 # pylab.savefig(filebase + 'fnielsen/eps/' + 'Nielsen2011New_anew_wordnet.eps')
112
113 # pylab.show()
114
115
116
117 sentiments_stem = np.asarray(sentiments_stem)
118 pylab.figure(3)
119 pylab.plot(sentiments_stem[:,0], sentiments_stem[:,1], '. ')
120 pylab.xlabel('My_list')
121 pylab.ylabel('ANEW')
122 pylab.title('Correlation between sentiment_word_lists (Porter_stemmer)')
123 pylab.text(1, 3, "Correlation = %.2f" % np.corrcoef(sentiments_stem.T)[1,0])

```

```

124 pylab.text(1, 2.5, "Spearman correlation = %.2f" % spearmanr(sentiments_stem[:,0], sentiments_stem[:,0]))
125 pylab.text(1, 2, "Kendall correlation = %.2f" % kendalltau(sentiments_stem[:,0], sentiments_stem[:,0]))
126 # pylab.savefig(filebase + 'fnielsen/eps/' + 'Nielsen2011New_anew_stem.eps')
127
128 # pylab.show()

```

```

1  #!/usr/bin/env python
2  #
3  # $Id: Nielsen2011New.py,v 1.10 2011/03/16 13:41:36 fn Exp $
4
5  import csv
6  import re
7  import numpy as np
8  from nltk.stem.wordnet import WordNetLemmatizer
9  from nltk import sent_tokenize
10 import pylab
11 from scipy.stats.stats import kendalltau, spearmanr
12 import simplejson
13
14 # This variable determines the data set
15 mislove = 2
16
17 # Filenames
18 filebase = '/home/fnielsen/'
19 filename_afinn = filebase + 'fnielsen/data/Nielsen2009Responsible_emotion.csv'
20 filename_afinn96 = 'AFINN-96.txt'
21 filename_mislove1 = "results.txt"
22 filename_mislove2a = "turk.txt"
23 filename_mislove2b = "turk-results.txt"
24 filename_anew = filebase + 'data/ANEW.TXT'
25 filename_gi = filebase + "data/inqtabs.txt"
26 filename_of = filebase + "data/subjclueslen1-HLTEMNLP05.tff"
27
28 # Word splitter pattern
29 pattern_split = re.compile(r"\W+")
30
31
32 afinn = dict(map(lambda (k,v): (k,int(v)),
33                 [ line.split('\t') for line in open(filename_afinn) ]))
34
35
36 afinn96 = dict(map(lambda (k,v): (k,int(v)),
37                  [ line.split('\t') for line in open(filename_afinn96) ]))
38
39
40 # ANEW
41 anew = dict(map(lambda l: (l[0], float(l[2]) - 5),
42                [ re.split('\s+', line) for line in open(filename_anew).readlines()[41:1075] ]))
43
44
45 # OpinionFinder
46 of = {}
47 for line in open(filename_of).readlines():
48     elements = re.split('\s|\=)', line)
49     if elements[22] == 'positive':
50         of[elements[10]] = +1
51     elif elements[22] == 'negative':
52         of[elements[10]] = -1
53
54
55 # General inquirer
56 csv_reader = csv.reader(open(filename_gi), delimiter='\t')
57 header = []
58 gi = {}
59 previousword = []
60 previousvalence = []
61 for row in csv_reader:

```



```

62     if not header:
63         header = row
64     elif len(row) > 2:
65         word = re.search("\w+", row[0].lower()).group()
66         if row[2] == "Positiv":
67             valence = 1
68         elif row[3] == "Negativ":
69             valence = -1
70         else:
71             valence = 0
72         if re.search("#", row[0].lower()):
73             if previousword == word:
74                 if previousvalence == []:
75                     previousvalence = valence
76                 elif previousvalence != valence:
77                     previousvalence = 0
78             else:
79                 if previousvalence:
80                     gi[previousword] = previousvalence
81                     previousword = word
82                     previousvalence = []
83         elif valence:
84             gi[word] = valence
85
86
87
88 # Lemmatizer for WordNet
89 lemmatizer = WordNetLemmatizer()
90
91
92
93 def words2anewsentiment(words):
94     """
95     Convert words to sentiment based on ANEW via WordNet Lemmatizer
96     """
97     sentiment = 0
98     for word in words:
99         if word in anew:
100             sentiment += anew[word]
101             continue
102         lword = lemmatizer.lemmatize(word)
103         if lword in anew:
104             sentiment += anew[lword]
105             continue
106         lword = lemmatizer.lemmatize(word, pos='v')
107         if lword in anew:
108             sentiment += anew[lword]
109     return sentiment/len(words)
110
111
112 def extremevalence(valences):
113     """
114     Return the most extreme valence. If extremes have different sign then
115     zero is returned
116     """
117     imax = np.argsort(np.abs(valences))[-1]
118     extremes = filter(lambda v: abs(v) == abs(valences[imax]), valences)
119     extremes_samesign = filter(lambda v: v == valences[imax], valences)
120     if extremes == extremes_samesign:
121         return valences[imax]
122     else:
123         return 0
124
125 def corrmatrix2latex(C, columns=None):
126     """
127     s = corrmatrix2latex(C)
128     print(s)
129     """

```

```

130     s = '\n\\begin{tabular}{'} + ('r'*(C.shape[0]-1)) + '\\n'
131     if columns:
132         s += "&".join(columns[1:]) + "\\\\n\\hline\\n"
133     for n in range(C.shape[0]):
134         row = []
135         for m in range(1, C.shape[1]):
136             if m > n:
137                 row.append("%.3f" % C[n,m])
138             else:
139                 row.append("-----")
140         s += "&".join(row) + "\\\\n"
141     s += '\\end{tabular}\\n'
142     return s
143
144
145 def spearmanmatrix(data):
146     """
147     Spearman r rank correlation matrix
148     """
149     C = np.zeros((data.shape[1], data.shape[1]))
150     for n in range(data.shape[1]):
151         for m in range(data.shape[1]):
152             C[n,m] = spearmanr(data[:,n], data[:,m])[0]
153     return C
154
155
156 def kendallmatrix(data):
157     """
158     Kendall tau rank correlation matrix
159     """
160     C = np.zeros((data.shape[1], data.shape[1]))
161     for n in range(data.shape[1]):
162         for m in range(data.shape[1]):
163             C[n,m] = kendalltau(data[:,n], data[:,m])[0]
164     return C
165
166
167
168 # Read Mislove CSV Twitter data: 'tweets' an array of dictionaries
169 if mislove == 1:
170     csv_reader = csv.reader(open(filename_mislove1), delimiter='\\t')
171     header = []
172     tweets = []
173     for row in csv_reader:
174         if not header:
175             header = row
176         else:
177             tweets.append({'id': row[0],
178                          'quant': int(row[1]),
179                          'score_our': float(row[2]),
180                          'score_mean': float(row[3]),
181                          'score_std': float(row[4]),
182                          'text': row[5],
183                          'scores': map(int, row[6:])})
184 elif mislove == 2:
185     if True:
186         tweets = simplejson.load(open('tweets_with_sentistrength.json'))
187     else:
188         csv_reader = csv.reader(open(filename_mislove2a), delimiter='\\t')
189         tweets = []
190         for row in csv_reader:
191             tweets.append({'id': row[2],
192                          'score_mislove2': float(row[1]),
193                          'text': row[3]})
194         csv_reader = csv.reader(open(filename_mislove2b), delimiter='\\t')
195         tweets_dict = {}
196         header = []
197         for row in csv_reader:

```



```

266         t['sentiment_afinn_extreme'],
267         t['sentiment_anew_raw'],
268         t['sentiment_anew_lemmatize'],
269         t['sentiment_anew_raw_sum'],
270         t['sentiment_anew_raw_nonzeros'],
271         t['sentiment_gi_raw'],
272         t['sentiment_of_raw'],
273         t['sentistrength']] for t in tweets ])
274
275
276 x = np.asarray(sentiments[:,1]).flatten()
277 y = np.asarray(sentiments[:,0]).flatten()
278 pylab.plot(x, y, '.')
279 pylab.xlabel('My_list')
280 pylab.ylabel('Amazon_Mechanical_Turk')
281 pylab.text(0.1, 2, "Pearson correlation = %.3f" % np.corrcoef(x, y)[1,0])
282 pylab.text(0.1, 1.6, "Spearman correlation = %.3f" % spearmanr(x, y)[0])
283 pylab.text(0.1, 1.2, "Kendall correlation = %.3f" % kendalltau(x, y)[0])
284 pylab.title('Scatterplot of sentiment strengths for tweets')
285 pylab.show()
286 # pylab.savefig(filebase + 'fnielsen/eps/Nielsen2011New_tweetsscatter.eps')
287
288 # Ordinary correlation coefficient
289 C = np.corrcoef(sentiments.transpose())
290 s1 = corrmatrix2latex(C, columns)
291 print(s1)
292
293 f = open(filebase + '/fnielsen/tex/Nielsen2011New_corrmatrix.tex', 'w')
294 f.write(s1)
295 f.close()
296
297
298 # Spearman Rank correlation
299 C2 = spearmanmatrix(sentiments)
300 s2 = corrmatrix2latex(C2, columns)
301 print(s2)
302
303 f = open(filebase + '/fnielsen/tex/Nielsen2011New_spearmanmatrix.tex', 'w')
304 f.write(s2)
305 f.close()
306
307
308 # Kendall rank correlation
309 C3 = kendallmatrix(sentiments)
310 s3 = corrmatrix2latex(C3, columns)
311 print(s3)
312
313 f = open(filebase + '/fnielsen/tex/Nielsen2011New_kendallmatrix.tex', 'w')
314 f.write(s3)
315 f.close()

```

```

1 #!/usr/bin/env python
2 #
3 # This script will call the SentiStrength Web service with the text from
4 # the 1000 tweets and write a JSON file with tweets and the SentiStrength.
5 #
6 # $Id: Nielsen2011New_sentistrength.py,v 1.1 2011/03/13 19:12:46 fn Exp $
7
8 import csv
9 import re
10 import numpy as np
11 import pylab
12 import random
13 from scipy import sparse
14 from scipy.stats.stats import kendalltau, spearmanr
15 import simplejson
16 import sys

```



```

85         'score_amt': np.mean(map(int, re.split("\s+", "\u".join(row[4:]))
86     }
87     for n in range(len(tweets)):
88         tweets[n]['score_mislove'] = tweets_dict[tweets[n]['id']]['score_mislove']
89         tweets[n]['score_amt_wrong'] = tweets_dict[tweets[n]['id']]['score_amt_wrong']
90         tweets[n]['score_amt'] = tweets_dict[tweets[n]['id']]['score_amt']
91
92
93
94
95 for n in range(len(tweets)):
96     url = urlbase + urlencode({'text': tweets[n]['text']})
97     try:
98         html = myopener.open(url).read()
99         positive = int(pattern_positive.findall(html)[0])
100        negative = int(pattern_negative.findall(html)[0])
101        tweets[n]['sentistrength_positive'] = positive
102        tweets[n]['sentistrength_negative'] = negative
103        if positive > abs(negative):
104            tweets[n]['sentistrength'] = positive
105        elif abs(negative) > positive:
106            tweets[n]['sentistrength'] = negative
107        else:
108            tweets[n]['sentistrength'] = 0
109    except Exception, e:
110        error = str(e)
111        tweets[n]['sentistrength_error'] = error
112    print(n)
113
114
115 simplejson.dump(tweets, open("tweets_with_sentistrength.json", "w"))

```

```

1  #!/usr/bin/env python
2  #
3  #   Generates a plot of the evolution of the performance as
4  #   the word list is extended.
5  #
6  # $Id: Nielsen2011New_evolution.py,v 1.2 2011/03/13 23:48:38 fn Exp $
7
8
9  import csv
10 import re
11 import numpy as np
12 import pylab
13 import random
14 from scipy import sparse
15 from scipy.stats.stats import kendalltau, spearmanr
16
17 # This variable determines the data set
18 mislove = 2
19
20 # Filenames
21 filebase = '/home/fnielsen/'
22 filename_afinn = filebase + 'fnielsen/data/Nielsen2009Responsible_emotion.csv'
23 filename_mislove1 = "results.txt"
24 filename_mislove2a = "turk.txt"
25 filename_mislove2b = "turk-results.txt"
26
27
28 # Word splitter pattern
29 pattern_split = re.compile(r"\W+")
30
31
32 afinn = dict(map(lambda (k,v): (k,int(v)),
33                 [ line.split('\t') for line in open(filename_afinn) ]))
34
35

```

```

36
37 # Read Mislove CSV Twitter data: 'tweets' an array of dictionaries
38 if mislove == 1:
39     csv_reader = csv.reader(open(filename_mislove1), delimiter='\t')
40     header = []
41     tweets = []
42     for row in csv_reader:
43         if not header:
44             header = row
45         else:
46             tweets.append({'id': row[0],
47                           'quant': int(row[1]),
48                           'score_our': float(row[2]),
49                           'score_mean': float(row[3]),
50                           'score_std': float(row[4]),
51                           'text': row[5],
52                           'scores': map(int, row[6:])})
53 elif mislove == 2:
54     csv_reader = csv.reader(open(filename_mislove2a), delimiter='\t')
55     tweets = []
56     for row in csv_reader:
57         tweets.append({'id': row[2],
58                       'score_mislove2': float(row[1]),
59                       'text': row[3]})
60     csv_reader = csv.reader(open(filename_mislove2b), delimiter='\t')
61     tweets_dict = {}
62     header = []
63     for row in csv_reader:
64         if not header:
65             header = row
66         else:
67             tweets_dict[row[0]] = {'id': row[0],
68                                   'score_mislove': float(row[1]),
69                                   'score_amt_wrong': float(row[2]),
70                                   'score_amt': np.mean(map(int, re.split("\s+", " ").join(row[4:])))}
71
72     for n in range(len(tweets)):
73         tweets[n]['score_mislove'] = tweets_dict[tweets[n]['id']]['score_mislove']
74         tweets[n]['score_amt_wrong'] = tweets_dict[tweets[n]['id']]['score_amt_wrong']
75         tweets[n]['score_amt'] = tweets_dict[tweets[n]['id']]['score_amt']
76
77
78 allwords = []
79 for n in range(len(tweets)):
80     words = pattern_split.split(tweets[n]['text'].lower())
81     tweets[n]['words'] = words
82     allwords.extend(words)
83
84 print("All words: %d" % len(allwords))
85 print("Number of unique words in Twitter corpus: %d" % len(set(allwords)))
86
87
88 terms = afinn.keys()
89 # terms = list(set(allwords).intersection(afinn.keys()))
90 print("Number of unique words matched to word list: %d" % len(terms))
91
92 term2index = dict(zip(terms, range(len(terms))))
93
94 M = sparse.lil_matrix((len(tweets), len(terms)))
95 M = np.zeros((len(tweets), len(terms))) # Sparse is not necessary
96 for n in range(len(tweets)):
97     for word in tweets[n]['words']:
98         if term2index.has_key(word):
99             M[n, term2index[word]] = afinn[word]
100
101 score_amt = [ t['score_amt'] for t in tweets ]
102
103 # Fix resampling seed for reproducibility

```

```

104 random.seed(a=1729)
105
106 K = [1, 10, 30, 50, 70, 100, 150, 200, 250, 300, 350, 400, len(terms)]
107 K = [5, 100, 300, 500, 1000, 1500, 2000, len(terms)]
108 I = range(len(terms))
109 resamples = 50
110 R = np.zeros((resamples, len(K)))
111 S = np.zeros((resamples, len(K)))
112 for n in range(len(K)):
113     for m in range(resamples):
114         J = random.sample(I, K[n])
115         score_afinn = M[:,J].sum(axis=1)
116         R[m,n] = np.corrcoef(score_amt, score_afinn)[0,1]
117         S[m,n] = spearmanr(score_amt, score_afinn)[0]
118
119
120 # pylab.figure(1)
121 pylab.subplot(2, 1, 1)
122 pylab.boxplot(R, positions=K, widths=40)
123 pylab.ylabel('Pearson correlation')
124 # pylab.xlabel('Word list size')
125 pylab.title('Evolution of word list performance')
126 pylab.axis((0, 2600, -0.05, 0.65))
127 pylab.show()
128
129 # pylab.figure(2)
130 pylab.subplot(2, 1, 2)
131 pylab.boxplot(S, positions=K, widths=40)
132 pylab.ylabel('Spearman rank correlation')
133 pylab.xlabel('Word list size')
134 # pylab.title('Evolution of word list performance')
135 pylab.axis((0, 2600, 0.35, 0.6))
136 pylab.show()
137
138 # pylab.savefig(filebase + 'fnielsen/eps/Nielsen2011New_evolution.eps')

```

```

1 #!/usr/bin/env python
2 #
3 # $Id: Nielsen2011New_anevafinn.py,v 1.1 2011/03/13 17:18:10 fn Exp $
4
5
6
7 import csv
8 import re
9 import numpy as np
10 from nltk.stem.wordnet import WordNetLemmatizer
11 from nltk import sent_tokenize
12 import pylab
13 from scipy.stats.stats import kendalltau, spearmanr
14
15 # This variable determines the data set
16 mislove = 2
17
18 # Filenames
19 filebase = '/home/fnielsen/'
20 filename_afinn = filebase + 'fnielsen/data/Nielsen2009Responsible_emotion.csv'
21 filename_afinn96 = 'AFINN-96.txt'
22 filename_mislove1 = "results.txt"
23 filename_mislove2a = "turk.txt"
24 filename_mislove2b = "turk-results.txt"
25 filename_anew = filebase + 'data/ANEW.TXT'
26 filename_gi = filebase + "data/inqtabs.txt"
27 filename_of = filebase + "data/subjclueslen1-HLTEMNLP05.tff"
28
29 # Word splitter pattern
30 pattern_split = re.compile(r"\W+")
31

```



```

32
33 afinn = dict(map(lambda (k,v): (k,int(v)),
34                 [ line.split('\t') for line in open(filename_afinn) ]))
35
36
37 # ANEW
38 anew = dict(map(lambda l: (l[0], float(l[2]) - 5) ,
39                 [ re.split('\s+', line) for line in open(filename_anew).readlines()[41:1075] ]))
40
41 anewafinn = set(anew.keys()).intersection(afinn.keys())
42
43
44 afinn_intersect = dict([ (k,afinn[k]) for k in anewafinn ])
45 anew_intersect = dict([ (k, anew[k]) for k in anewafinn ])
46
47
48 # Read Mislove CSV Twitter data: 'tweets' an array of dictionaries
49 if mislove == 1:
50     csv_reader = csv.reader(open(filename_mislove1), delimiter='\t')
51     header = []
52     tweets = []
53     for row in csv_reader:
54         if not header:
55             header = row
56         else:
57             tweets.append({'id': row[0],
58                           'quant': int(row[1]),
59                           'score_our': float(row[2]),
60                           'score_mean': float(row[3]),
61                           'score_std': float(row[4]),
62                           'text': row[5],
63                           'scores': map(int, row[6:])})
64 elif mislove == 2:
65     csv_reader = csv.reader(open(filename_mislove2a), delimiter='\t')
66     tweets = []
67     for row in csv_reader:
68         tweets.append({'id': row[2],
69                       'score_mislove2': float(row[1]),
70                       'text': row[3]})
71     csv_reader = csv.reader(open(filename_mislove2b), delimiter='\t')
72     tweets_dict = {}
73     header = []
74     for row in csv_reader:
75         if not header:
76             header = row
77         else:
78             tweets_dict[row[0]] = {'id': row[0],
79                                   'score_mislove': float(row[1]),
80                                   'score_amt_wrong': float(row[2]),
81                                   'score_amt': np.mean(map(int, re.split("\s+", " ").join(row[4:])))}
82
83     for n in range(len(tweets)):
84         tweets[n]['score_mislove'] = tweets_dict[tweets[n]['id']]['score_mislove']
85         tweets[n]['score_amt_wrong'] = tweets_dict[tweets[n]['id']]['score_amt_wrong']
86         tweets[n]['score_amt'] = tweets_dict[tweets[n]['id']]['score_amt']
87
88
89
90 # Computer sentiment for each tweet
91 amt = []
92 afinn_intersect_sentiment = []
93 anew_intersect_sentiment = []
94 for n in range(len(tweets)):
95     amt.append(tweets[n]['score_amt'])
96     words = pattern_split.split(tweets[n]['text'].lower())
97     afinn_sentiments = map(lambda word: afinn_intersect.get(word, 0), words)
98     afinn_intersect_sentiment.append(float(sum(afinn_sentiments))/len(afinn_sentiments))
99     anew_sentiments = map(lambda word: anew_intersect.get(word, 0), words)

```

```
100     anew_intersect_sentiment.append(float(sum(anew_sentiments))/len(anew_sentiments))
101
102
103
104
105 np.corrcoef(amt, afinn_intersect_sentiment)[0,1]
106 np.corrcoef(amt, anew_intersect_sentiment)[0,1]
107
108 spearmanr(amt, afinn_intersect_sentiment)[0]
109 spearmanr(amt, anew_intersect_sentiment)[0]
```
