

# Correlation between weather and sentiment analysis on Twitter

Alex Pellegrini (s132199) & Woody Rousseau (s131642)

**Abstract**—We describe a webservice application which provides several means of visualization of the correlation between the weather and sentiment analysis results, based on data mined tweets which include localization information. We find that only taking into account these two aspects is not enough to obtain a high correlation.

## I. INTRODUCTION

Many studies such as [1] aimed at establishing a correlation between the weather and the mood or other human cognition aspects. It is now well known that many more factors need to be taken into account in order to obtain results corroborating that there is indeed a correlation: the time spent outside, and personal life events which tend to have a much larger impact on the mood than the weather. This is also due to the fact that these studies were unable to acquire enough data to reduce the importance of those personal factors. With that flaw in mind, we focus in our study on data mining techniques to establish such a correlation, by gathering tweets that contain localization information, by finding weather information on that localization and by combining it with a sentiment analysis of the tweet. The data mining tasks are performed using Python and two different APIs: **Twitter**'s and **OpenWeatherMap**'s. The presentation of our results is also made using Python, and more specifically a **Flask** microframework powered webservice. Those two main tasks run in a parallel fashion using threads.

## II. METHODS

### A. Webserver

The webservice relies on **Flask** microframework, defining several routes: one for each visualization methods available, one for a homepage from which all those methods can be accessed, as well as HTTP error handlers (404 and 500). The code for running the webserver is included in the `server.py`, the script that needs to be launched to run the full application.

### B. Visualization Methods

Three visualization methods are available, and on each page, a start button and a stop button are made available to start and stop the data mining of tweets.

- The route `list_data` displays a table containing all acquired tweets as well as their analysis results and updated using websockets powered by **gevent-socketio** as they are data mined. Using websockets allows the server to let the client know that some of the content needs to

be updated (in this case, new acquired tweets), triggering an update.

- The route `plot` displays a scatter plot with all tweets, using **NumPy** and **matplotlib** and our `plotting.py` module. The x-axis correspond to the sentiment value (continuous) and the y-axis to the weather value (discrete), both scaled to be between 0.0 and 1.0. Because we have discrete weather values, a mean sentiment value is computed for each of these weather values. Polynomial (degree 1) curve fitting is then applied to these data points and its result is plotted as well. The image can be updated by clicking on a button, which triggers an Asynchronous JavaScript and XML (AJAX) request, requesting an update from the server.
- The route `map` displays a map powered by **Google Map**'s Javascript API and centered around the United States. A weather layer as well as a heatmap layer based on a simple correlation score<sup>1</sup> is added to the map. Objects added to the layer contain a latitude, longitude and a weight information corresponding to the correlation score. Obviously, only tweets gathered since the start button was first clicked are displayed on the heatmap, in order not to include data coming from a period with a different weather. New tweets are added to the heatmap again using websockets.

### C. Data Mining

Data mining tasks are performed on a separate thread and use the `TweetWeather` class from `tweetweather.py`, in order to keep everything packed in a single Python script while being able to run the webserver and data mine tweets at the same time. The thread is flagged as a "daemon thread" to allow keyboard interruptions to kill the process. For most of the used APIs, app credentials are provided so that the program runs "out of the box", even though they should in principle remain personal. All tweets are saved in a sqlite database.

- **Tweepy** allows an easy Python access to **Twitter**'s API: the initial request only gathers tweets in English, and they are then filtered to remove tweets that do not include localization. When the rate limit is exceeded, the thread is paused and resumes fifteen minutes later, following Twitter's policy.
- **OpenWeatherMap** is used to give a weather score to the acquired localization of a tweet. A manually created dictionary maps weather icons describing the weather to

<sup>1</sup>score = |weatherValue−sentimentValue|

a value between 0.0 and 1.0, ranging from a thunderstorm (0.0) to clear skies (1.0), with snow and night weather (without rain) considered neutral (0.5).

#### D. Sentiment Analysis

For the sentiment analysis we created a large word-values list containing more than 9400 word-value pairs that we load in a single dictionary, so that a value associated to a certain word can be acquired with an  $O(1)$  complexity. The values fit in the  $[-5, 5]$  range. The list is a merge of the **AFINN** word-value list and **Alex Davies's** list which was created by using a learning algorithm and customized to be scaled in the  $[-5, 5]$  interval. The goal is to have words with very high and very low values get a big weight within the overall evaluation, as they are supposed to be quite rare. In order to do that, the "probability" of a word is defined as the probability of its value following a gaussian distribution (properly fitted to the range).

The analysis in itself is performed by splitting the body of each gathered tweet, and evaluating each word with the dictionary mentioned above.

The probability function mentioned above is computed using the likelihood function of a gaussian distribution formula where the mean and deviation are  $\mu = \text{mean}(\text{list.values}())$  and  $\sigma = \sqrt{\frac{\sum_{i=1}^N (x_i - \mu)^2}{N}}$ . The probability is found in the following way:

$$P(w) = \exp\left(-\frac{(\text{value}(w) - \mu)^2}{2\sigma^2}\right) \quad (1)$$

The value of each word is computed by getting the actual value from the list and dividing it by the probability above:

$$\text{value}(w) = \frac{\text{value}(w)}{P(w)} \quad (2)$$

The values in each category (positive, negative and neutral) are summed (*pos*, *neg*, *neu*), and the number of words in each category is also saved (*pos.size()*, *neg.size()*, *neu.size()*). Finally, a global sum of the values of all words (whatever their category may be) is computed (*globalSum*).

With these information, we can "properly" weight each word inside the tweet. The overall value of a tweet is given:

$$\text{value}(\text{tweet}) = \frac{\sum_{ctg \in (\text{pos}, \text{neg}, \text{neu})} ctg \times ctg.size()}{\text{globalSum}} \quad (3)$$

A threshold value was defined in order to limit the values of tweets, which could be very large. This allows the program to scale tweet values in the interval  $[0.0, 1.0]$ :

$$\text{value}(\text{tweet}) = \frac{\text{value}(\text{tweet}) + \text{threshold}}{2 \times \text{threshold}} \quad (4)$$

**NumPy** and the standard **math** modules are used to perform these computations.

#### E. Development process details

The application was developed on two computers running Mac OSX 10.9 with Python 2.7.5 (in order not to run with library compatibility with Python 3 issues). As for text editors, **emacs** and **Sublime Text** were used instead of any IDE. **git** was the chosen revision control system since the beginning of the development process, uploading and maintaining the repository on **Bitbucket**, a popular code cloud service.

The whole module is quite easy to install, as all dependencies can be installed thanks to the `setup.py` file which is included in the project.

### III. RESULTS

#### Data List

ID	Sent.Value	Weather	Weather Info
85	0.5482783166639076	Clouds	few clouds
84	0.9389598525982204	Clouds	few clouds
83	0.5893751940638643	Clouds	overcast clouds
82	0.7647818621818167	Clouds	few clouds
81	0.4807209020871157	Clouds	broken clouds
80	0.725842565328908	Clouds	scattered clouds
79	0.8564392731461345	Clouds	broken clouds
78	0.6748973425970628	Clouds	few clouds
77	0.808392289241573	Clouds	few clouds
76	0.6260259215186186	Clouds	scattered clouds
75	0.7188651554303171	Clear	Sky is Clear
74	1	Clouds	scattered clouds
73	0.6298715181586524	Clouds	overcast clouds
72	0.49820864380843273	Clear	sky is clear
71	0.4991801867445291	Clouds	broken clouds

Fig. 1. Websockets updated table with tweets

#### Heatmap

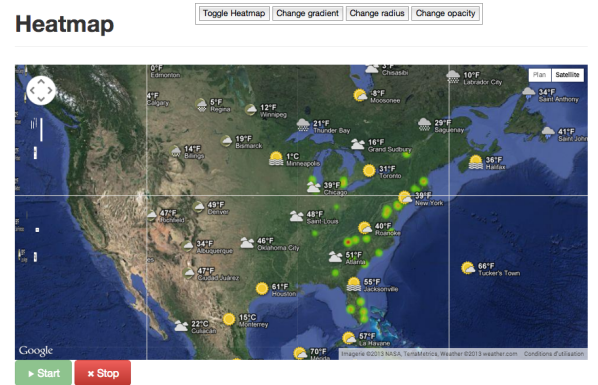


Fig. 2. Websockets updated heatmap with weather layer

Figure 1, Figure 2 and Figure 3 show the three visualization methods. The heatmap does show some points with high weights (high correlation tweets) which have small areas and yet very red centers, and points with low weights (low correlation tweets) which have similar areas but are only colored in green. The scatter plot is useful as it allows a global vision of all data, but also shows that our hypothesis is far from being correct. When many points are acquired, the fitted model is entirely different from the ideal line. This shows some issues with the weather dictionary and the sentiment analysis.

### Scatter Plot

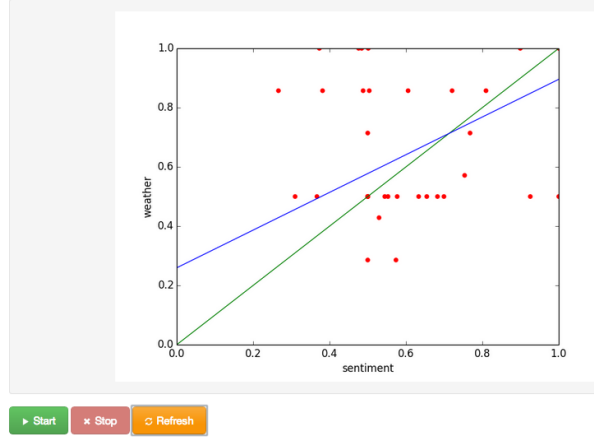


Fig. 3. Scatter plot (red points), ideal line (green), fitted model (blue)

#### A. Code checking

We used `pylint` as well as `pep8` to check our coding quality. Only minor issues (such as a few too long lines) are left after using these tools. These tools made us aware of some “pythonic” conventions, such as using list comprehensions instead of `map`.

#### B. Testing

Testing was performed using `py.test`. Most aspects of our code was not testable as it runs in real time as it gathers tweets. We however focused on the unit testing of our sentiment analysis, on the plotting, as well as on the provided webservice, using Flask’s test client. All tests are runnable from the root directory, through the `python setup.py test` command. In total, 12 tests are created and gathered in three different classes, placed in `tests/test_server.py`, `tests/test_analysis.py` and `tests/test_plotting.py`. Coverage was performed using `pytest-cov`, a `py.test` plugin with the following command:

```
py.test --cov-report term-missing --cov
server/ tests/
```

It returns a coverage of 99% for `analysis.py`, of 89% for `plotting.py` and of 52% for `server.py` (`tweetweather.py` was not tested, for reasons explained above). The low coverage for `server.py` (and to a lesser extent for `plotting.py`) is due to the fact that some methods require an interaction with the data mining thread, with the database or with the javascript part of the application (websockets).

#### C. Profiling

Profiling was performed by timing the data mining, analyzing, and the saving of the data and an overall time of the whole process using the `time.py` module, surrounding each function with the following statement:

```
elapsed = time.time()
method.call()
elapsed = time.time() - elapsed
```

This data was saved in a text file using a simple unix redirection (i.e. `>>`) and analyzed by a new python script (`profile.py`) that for each overall time (i.e. each group of tweets gathered from the twitter stream) computes how much of that time was spent on the above functions. What comes out is that on an overall timing, almost 50% of the time is spent on gathering new tweets using the API:

- gathering: 48.4456012566
- analysis: 1.24803740403
- save: 1.24540493406
- weather: 1.07298472832
- connection: 0.666942267582

`save` is the time spent to save data into the `sqlite` database, `weather` is the time used to gather weather information, and `connection` is the time spent to connect to the database each time.

### IV. DISCUSSION

There are obviously limitations to the application:

- For the data mining, it takes a long time to get a large database (English tweets including localization are scarce). The weather values are also discrete, and they are quite arbitrary. Overall, it would have been a good idea to include more intermediate weather values, and to have many people rate them on a scale from 0.0 to 1.0, taking the means of the results of this study to create the weather dictionary.
- For the Flask application, a good idea would have been to use a cloud service to grow a shared database.
- Obviously, testing and coverage are quite limited by the nature of the application. However, we could have tested the database if we had used Flask’s `sqlalchemy`.
- As for profiling, we did not have the time to implement several versions of our code (for example with several alternative algorithms or data structures) and it would have been interesting to benchmark these different methods.

### V. CONCLUSION

We have built a data mining webservice application. Our original correlation hypothesis was not confirmed, but the application being properly structured, documented and tested, it is thus quite scalable: it could include more complex sentiment analysis or a more complete weather dictionary. The next step would be to establish a state of the art on this correlation, and try to include more factors (time spent outside for instance), although they might be hard to find on Twitter.

### REFERENCES

- [1] M. C. Keller, B. L. Fredrickson, O. Ybarra, S. Côté, K. Johnson, J. Mikels, A. Conway, and T. Wager, “A warm heart and a clear head the contingent effects of weather on mood and cognition,” *Psychological Science*, vol. 16, no. 9, pp. 724–731, 2005.

# APPENDIX A

## CODE LISTINGS

### LISTINGS

server.py . . . . .	4
tweetweather.py . . . . .	6
analysis.py . . . . .	8
plotting.py . . . . .	9
test_server.py . . . . .	10
test_analysis.py . . . . .	11
test_plotting.py . . . . .	12
profile.py . . . . .	12

```

1  """
2  server.py is a python module which implements a flask http
3  server needed by the whole web application.
4  """
5
6  import sys
7  from flask import Flask, render_template, jsonify, Response, request
8  from socketio.server import SocketIOServer
9  from socketio import socketio_manage
10 from socketio.namespace import BaseNamespace
11 from pysqlite2 import dbapi2 as db
12 import gevent
13 import os
14 from analysis import Analyzer
15 from plotting import ScatterPlot
16 import urllib2
17 from tweetweather import TweetWeather
18
19 PATH = os.path.join('.', os.path.dirname(__file__), '../')
20 sys.path.append(PATH)
21
22 PORT = 5000
23
24 app = Flask(__name__)
25
26
27 @app.route('/')
28 def home():
29     """
30     The home page is a static page
31     allowing the choice of the visualization tool
32     """
33     return render_template('hello.html')
34
35
36 @app.route('/list_data')
37 def list_data():
38     """
39     This page displays all objects in a table
40     updated as they are mined.
41     The user can start and stop the data mining
42     using buttons.
43     """
44     data = []
45     if os.path.exists('data.sqlite'):
46         cur = db.connect('data.sqlite').cursor()
47         cur.execute('SELECT id, sentimentValue, weather, infos FROM tweets ORDER BY id DESC')
48         data = cur.fetchall()
49     return render_template('list.html', data=data)
50
51
52 @app.route('/map')
53 def display_map():
54     """
55     This page displays a map centered around the United States
56     containing a weather layer and a heatmap with arriving objects.
57     The user can start and stop the data mining
58     using buttons.
59     """
60     return render_template('map.html')
61
62
63 @app.route("/plot")
64 def plot():
65     """
66     This page displays a scatter plot of all gathered tweets
67     The x-axis is the sentiment value.
68     The y-axis is the weather value.
69     The closer points are to the 'identity' line,

```

```

70     the closer they fit our hypothesis
71     """
72     scatter_plot = ScatterPlot(1)
73     scatter_plot.load_data()
74     img_data = scatter_plot.get_image_data()
75     refresh = request.args.get('refresh', 0, type=int)
76     if refresh:
77         return img_data
78     return render_template('plot.html', data=img_data)
79
80
81 @app.route('/socket.io/<path:remaining>')
82 def socketio(request):
83     """
84     This route configures the WebSocket
85     used to let the client know that new objects
86     were mined
87     """
88     try:
89         socketio_manage(request.environ,
90                         {'/new_posts': BaseNamespace},
91                         request)
92     except:
93         app.logger.error("Exception_while_handling_socketio_connection",
94                         exc_info=True)
95     return Response()
96
97
98 def check_conn():
99     """
100     Checks whether a working Internet connection is available
101     """
102     try:
103         urllib2.urlopen('http://74.125.228.100') # Google IP (no DNS lookup)
104         return True
105     except urllib2.URLError:
106         pass
107     return False
108
109
110 @app.route('/start')
111 def start():
112     """
113     Starts the data mining thread if an internet connection is available
114     """
115     if check_conn():
116         tw_thread.start()
117         return jsonify('true')
118     else:
119         tw_thread.connexion_lost("Absent_Internet_Access")
120         return jsonify('false')
121
122
123 @app.route('/stop')
124 def stop():
125     """
126     Stops the data mining thread
127     """
128     tw_thread.stop()
129     return jsonify('true')
130
131
132 @app.errorhandler(404)
133 def page_not_found(exc):
134     """
135     404 error handler
136     used if a non existant route
137     is requested
138     """
139     return render_template('404.html'), 404
140
141
142 @app.errorhandler(500)
143 def page_not_found(exc):
144     """
145     500 error handler
146     used if there is a server error
147     """
148     return render_template('500.html'), 500
149
150
151 if __name__ == '__main__':
152     analyzer = Analyzer()
153     server = SocketIOServer((' ', PORT), app, resource="socket.io")
154     tw_thread = TweetWeather(server, analyzer, name="Tweet-Weather-Thread")
155     tw_thread.daemon = True

```

```

156 gevent.spawn(tw_thread.new_post, server)
157 gevent.spawn(tw_thread.connexion_lost, server)
158 print "Application_Started:_http://localhost:5000"
159 try:
160     server.serve_forever()
161 except KeyboardInterrupt:
162     tw_thread.stop()
163     server.stop()
164     sys.exit()

```

---

```

1 """Perform the data mining tasks through Twitter and OpenWeatherMap's APIs"""
2
3 import threading
4 from pysqlite2 import dbapi2 as db
5 import simplejson as jsn
6 import urllib
7 import tweepy
8 import os
9 import time
10
11 # secret twitter app credentials
12 CONSUMER_KEY = "Zlfl1aZTxrnydXBMZfeA"
13 CONSUMER_SECRET = "rNSasklWRb8mLbbzTZo6vAHB27EwNRmy4AA5c3G04"
14 ACCESS_KEY = "1889545957-BFTycJVNsAgtlfdKbalVlrwTJqoGGhj0iTxi06k"
15 ACCESS_SECRET = "NmfEez4FykNliGZYUfjYzUvUIksNne2xi6Ovo9Wq00"
16 WEATHER_APPID = "&APPID=4e04cba42b432a01c4226e186f3d23d2"
17
18 '''
19 Dictionary used to map a weather icon to a
20 weather "score" from 0 (worse) to 7 (best).
21 Night weathers (with no rain) as well as snow are considered neutral.
22 The dictionary is normalized to fit in the
23 [0,1] range.
24 '''
25 WEATHER_DICT = {'13d': 3.5, '11d': 0, '09d': 1,
26                 '10d': 2, '50d': 3, '04d': 4,
27                 '03d': 5, '02d': 6, '01d': 7,
28                 '13n': 3.5, '11n': 0, '09n': 1,
29                 '10n': 2, '50n': 3, '04n': 3.5,
30                 '03n': 3.5, '02n': 3.5, '01n': 3.5}
31 WEATHER_DICT = {k: float(v)/7 for (k, v) in WEATHER_DICT.iteritems()}
32
33
34 class TweetWeather(threading.Thread):
35     """
36     TweetWeather inherits the Python Thread class.
37
38     Indeed, data mining tasks need to be performed in a separate thread,
39     to keep the server running.
40
41     Each new object is stored in the database, and sent to the client
42     through a WebSocket.
43     """
44
45     def __init__(self, server, analyzer, name=''):
46         """
47         Checks if there is a working Internet access
48         """
49         threading.Thread.__init__(self)
50         self.name = name
51         self.server = server
52         self.analyzer = analyzer
53         self.root_weather_url = "http://openweathermap.org/data/2.5/weather?lat=%s&lon=%s"
54         auth = tweepy.OAuthHandler(CONSUMER_KEY, CONSUMER_SECRET)
55         auth.set_access_token(ACCESS_KEY, ACCESS_SECRET)
56         self.api = tweepy.API(auth) # Initialization of the tweeter API
57         self.terminated = False
58
59     def new_post(self, *args):
60         """
61         Sends the new object in a packet to the client through a WebSocket
62         """
63         pkt = dict(type="event", name="new_post",
64                   args=args, endpoint="/new_posts")
65         for _, socket in self.server.sockets.iteritems():
66             socket.send_packet(pkt)
67
68     def connexion_lost(self, *args):
69         """
70         Sends a notification packet to let the client know
71         that the internet connexion is lost.
72         """
73         pkt = dict(type="event", name="connexion_lost",
74                   args=args, endpoint="/new_posts")
75         for _, socket in self.server.sockets.iteritems():
76             socket.send_packet(pkt)

```

```

77
78 def run(self):
79     init_database()
80     self.gather_tweets()
81     self.__init__(self.server, self.analyzer, name=self.name)
82
83 def parse_text(self, status):
84     """
85     performs a very basic sentiment analysis on
86     a single tweet by comparing the most
87     significant words found on the afinn
88     word-value list. Furthermore gets the
89     weather conditions from "http://openweathermap.org"
90     of the location where the tweet has been written
91     and saves the result in the database.
92     """
93     conn = db.connect('data.sqlite')
94     cursor = conn.cursor()
95     score = self.analyzer.analyze(status.text)
96     weather_url = self.root_weather_url % tuple(
97         [str(x) for x in status.coordinates['coordinates']])
98     response = urllib.urlopen(weather_url)
99     try:
100         weather = json.load(response)
101     except json.JSONDecodeError:
102         print('Program -> Tweet not saved due to invalid weather json')
103     else:
104         if 'weather' in weather.keys():
105             main = weather['weather'][0]
106             #print(main['main'], status.text, score)
107
108             correlation_score = abs(score-WEATHER_DICT[main['icon']])
109             cursor.execute("INSERT INTO tweets (sentimentValue, "
110                 "weatherValue, correlationScore, weather, "
111                 "latitude, longitude, infos)"
112                 "VALUES(?, ?, ?, ?, ?, ?, ?)",
113                 [score, WEATHER_DICT[main['icon']],
114                 correlation_score, main['main'],
115                 status.coordinates['coordinates'][1],
116                 status.coordinates['coordinates'][0],
117                 main['description']])
118             self.new_post(score, main['main'], main['description'],
119                 status.coordinates['coordinates'][1],
120                 status.coordinates['coordinates'][0],
121                 correlation_score)
122             conn.commit()
123         conn.close()
124
125 def gather_tweets(self):
126     """
127     Performs data mining on tweets which have localization information
128     using the Twitter API (and the tweepy wrapper)
129     """
130     print('Fetching, localizing and analyzing Twitter '
131         'stream data (could take a while due to '
132         'the few geotagged tweets)...')
133     filtered_tweets = []
134     query = 'lang:en'
135     tweet_pages = tweepy.Cursor(self.api.search,
136                                 q=query, lang='en',
137                                 count=100, result_type="recent",
138                                 include_entities=True).pages()
139
140     while True:
141         try:
142             tweets = next(tweet_pages)
143         except tweepy.error.TweepError as exc:
144             if exc.message[0]['code'] == 88: # Rate Limit Exceeded
145                 print "Rate Limit Exceeded. Waiting for 15 minutes."
146                 time.sleep(60*15)
147             tweets = next(tweet_pages)
148         except KeyboardInterrupt:
149             self.stop()
150
151         filtered_tweets = [tweet for tweet in tweets if tweet.coordinates]
152         if not filtered_tweets: # No tweet with coordinates on that page
153             continue
154         for filtered_tweet in filtered_tweets:
155             self.parse_text(filtered_tweet)
156             if self.terminated:
157                 break
158
159 def stop(self):
160     """
161     Stops the thread
162     """
163     self.terminated = True

```



```

163
164
165 def init_database():
166     """
167     Initializes an sqlite database where evaluated tweets
168     will be saved the table created has an id primary key
169     attribute, a main value for the weather and a short
170     description
171     """
172     if not os.path.exists('data.sqlite'):
173         print("Initializing sqlite database for further analysis...")
174         conn = db.connect('data.sqlite')
175         cursor = conn.cursor()
176         cursor.execute("CREATE TABLE tweets(" +
177             "id INTEGER PRIMARY KEY AUTOINCREMENT," +
178             "sentimentValue REAL NOT NULL," +
179             "weatherValue REAL NOT NULL," +
180             "correlationScore REAL NOT NULL," +
181             "weather VARCHAR(255) NOT NULL," +
182             "latitude REAL NOT NULL," +
183             "longitude REAL NOT NULL," +
184             "infos VARCHAR(255) )")
185         conn.commit()
186         conn.close()
187         print(">>>Done<<<")
188     else:
189         print("Connecting to sqlite database")
190         print(">>>Done<<<")

```

---

```

1 """
2 Module needed by the application to analyze a single tweet
3 and give it a sentiment score
4 """
5
6 import numpy as np
7 import math
8 import urllib
9 import os
10 import re
11
12
13 class Analyzer(object):
14     """
15     Analyzer is used to give a real value to the sentiment found
16     in a tweet text.
17     """
18     def __init__(self):
19         """
20         Analyzer constructor, urls and external data
21         management hard coded.
22         """
23         import zipfile
24         self.comp_list = {}
25         self.afinn_url = "http://www2.imm.dtu.dk/pubdb/views/edoc_download.php/6010/zip/imm6010.zip"
26         self.url = "https://dl.dropbox.com/u/3773091/Twitter%20Sentiment/Twitter%20sentiment%20analysis.zip"
27         urllib.urlretrieve(self.afinn_url, 'word_list.zip')
28         word_list_zip = zipfile.ZipFile('word_list.zip')
29         self.comp_list = {unicode(k, 'utf-8'): int(v) # reads AFINN list
30             for (k, v) in [line.split('\t') for line in open(word_list_zip.extract('AFINN/AFINN-111.txt'))]}
31         script_dir = os.path.dirname(__file__)
32         with open(os.path.join(script_dir, 'my_list.txt'), 'r') as comp_file: # reads larger list
33             for line in comp_file:
34                 data = line.split('\t')
35                 self.comp_list[data[0]] = int(float(data[1].strip())) - 5
36
37         with open(os.path.join(script_dir, 'emoticons.csv'), 'r') as smiles: # reads emoticons file
38             for line in smiles:
39                 data = line.split('\t')
40                 self.comp_list[data[0]] = int(data[1].strip())
41
42         # clean temporary files on the fly
43         values = np.array(self.comp_list.values())
44         self.mean = np.mean(values)
45         self.deviation = math.sqrt(sum([pow(x - self.mean, 2) for x in values])
46                                     / float(len(values)))
47         os.remove('word_list.zip')
48         os.remove('AFINN/AFINN-111.txt')
49         os.rmdir('AFINN')
50
51     def analyze(self, tweet):
52         """
53         Analyzes the body of a tweet by comparing words to the
54         AFINN word-value list and using a gaussian distribution
55         to compute the weight of each word
56         """
57         emoticons_groups = re.findall(r"([0-9'\& \- \. \/ \(\) =:; +] | ((?:[;|=](?:-)?(?:\)|D|P)) | (<3))", tweet)

```



```

58     emoticons = [x[0] for x in emoticons_groups if x[0] != ''] # we have three groups in our regexp so we need to check everyone
59     emoticons.extend([x[1] for x in emoticons_groups if x[1] != ''])
60     emoticons.extend([x[2] for x in emoticons_groups if x[2] != ''])
61     data = [self.comp_list.get(word, 0) for word in tweet.split(' ')]
62     data.extend([self.comp_list.get(e, 0) for e in emoticons])
63     ctg_count = {'positive': 0, 'negative': 0, 'neutral': 0} # dict containing the number of positive negative and neutral words i
64     ctg_total = {'positive': 0.0, 'negative': 0.0, 'neutral': 0.0} # dict containing the sum respectively for positive negative and
65     threshold = 22.5
66
67     # computes categories cardinality and global sum of values of each word in tweet
68     vals = self.categories_cardinality(tweet, ctg_count)
69     # weights each category
70     tot_pos, tot_neg, tot_neu = self.weight_categories(data, ctg_total, ctg_count)
71     if vals:
72         total = (sum([tot_pos, tot_neg, tot_neu]) / vals) + threshold
73     else:
74         total = (sum([tot_pos, tot_neg, tot_neu])) + threshold
75     if total > 2*threshold:
76         total = 2*threshold
77     elif total < 0:
78         total = 0
79     else:
80         pass
81     return total / (2*threshold)
82
83 def weight_categories(self, data, ctg_total, ctg_count):
84     """
85     Computes the weight in terms of word value of each category
86     of words ( positive, negative, neutral )
87     """
88     for value in data:
89         if value > 0:
90             ctg_total['positive'] = ctg_total['positive'] + (value / espone(value, self.mean, self.deviation))
91         elif value < 0:
92             ctg_total['negative'] = ctg_total['negative'] + (value / espone(value, self.mean, self.deviation))
93         else:
94             ctg_total['neutral'] = ctg_total['neutral'] + espone(value, self.mean, self.deviation)
95     tot_pos = ctg_total['positive'] * ctg_count['positive']
96     tot_neg = ctg_total['negative'] * ctg_count['negative']
97     tot_neu = ctg_total['neutral'] * ctg_count['neutral']
98     return tot_pos, tot_neg, tot_neu
99
100 def categories_cardinality(self, tweet, ctg_count):
101     """
102     Computes the cardinality in terms of number of words belonging
103     to each category ( positive, negative, neutral ) returns also
104     the sum of the absolute values of each word.
105     """
106     vals = 0
107     for word in (tweet.lower()).split(' '):
108         temp = self.comp_list.get(word, 100)
109         if temp > 0 and temp < 100:
110             ctg_count['positive'] = ctg_count.get('positive') + 1
111             vals = vals + abs(temp)
112         elif temp < 0:
113             ctg_count['negative'] = ctg_count.get('negative') + 1
114             vals = vals + abs(temp)
115         elif temp == 0:
116             ctg_count['neutral'] = ctg_count.get('neutral') + 1
117             vals = vals + abs(temp)
118     return vals
119
120
121 def espone(value, mean, deviation):
122     """
123     The "Likelihood" function using a gaussian probability
124     function.
125     """
126     return math.exp(-(pow((value - mean), 2)) / (2*pow(deviation, 2)))

```

---

```

1  """
2  Module needed by the application to generate the scatter plot
3  of the sentiment value with the weather value
4  """
5
6  import os
7  from pysqlite2 import dbapi2 as db
8  import matplotlib
9  matplotlib.use('Agg')
10 import matplotlib.pyplot as plt
11 import numpy as np
12 from cStringIO import StringIO
13
14
15 class ScatterPlot(object):
16     """

```

```

17     This scatter plots includes the data points ,
18     an ideal correlation line (in green),
19     and a fit of the data available (in blue)
20     """
21
22     def __init__(self, pol_order=1):
23         self.pol_order = pol_order
24         self.x = []
25         self.y = []
26
27     def load_data(self):
28         """
29         Gets all data available
30         from the database
31         """
32         if os.path.exists('data.sqlite'):
33             cur = db.connect('data.sqlite').cursor()
34             cur.execute('SELECT sentimentValue , weatherValue FROM tweets'
35                         ' WHERE sentimentValue >= 0 ORDER BY id DESC')
36             all_fetched = cur.fetchall()
37             self.x = [point[0] for point in all_fetched]
38             self.y = [point[1] for point in all_fetched]
39             return True
40         return False
41
42     def set_data(self, x, y):
43         """
44         Sets x and y from given Lists
45         Used for testing purposes
46         """
47         self.x = list(x)
48         self.y = list(y)
49
50     def get_fit_function(self):
51         """
52         Fits a model following
53         the polynomial order given to the class
54         """
55         mean_x = []
56         mean_y = []
57         for i in range(0, 8):
58             indices = [ind for ind, val in enumerate(self.y) if val == float(i)/7]
59             if indices:
60                 mean_y.append(float(i)/7)
61                 mean_x.append(np.mean([self.x[j] for j in indices]))
62         try:
63             fit_fn = np.poly1d(np.polyfit(mean_x, mean_y, self.pol_order))
64         except TypeError: # Empty Lists
65             return None
66         return fit_fn
67
68     def get_image_data(self):
69         """
70         Prepares the figure and returns the data
71         formatted in a base64 encoded String
72         """
73         fig = plt.figure()
74         axis = fig.add_subplot(1, 1, 1)
75         xs = np.linspace(0, 1, 8)
76         axis.set_xlim([0, 1])
77         axis.set_ylim([0, 1])
78
79         axis.plot(xs, xs,
80                 label='Perfect_Correlation', color='green')
81         fit_fn = self.get_fit_function()
82         if (fit_fn):
83             axis.plot(xs, fit_fn(xs),
84                     label='Observed_Correlation', color='blue')
85         axis.scatter(self.x, self.y,
86                    label='Data_Points', color='red')
87         plt.xlabel('sentiment')
88         plt.ylabel('weather')
89         str_io = StringIO()
90         fig.savefig(str_io, format='png')
91         plt.close("all")
92         return str_io.getvalue().encode('base64')

```

---

```

1 import os
2 import sys
3 MY_PATH = os.path.dirname(os.path.abspath(__file__))
4 sys.path.insert(0, MY_PATH + '/../server')
5 import server
6
7
8 class TestFlaskServer(object):
9     """

```

```

10     Testing a part of the Flask webserver
11     """
12
13     def setup(self):
14         """
15         Called at the beginning of the test module
16         Configures the Flask Test Client
17         """
18         server.app.config['TESTING'] = True
19         self.app = server.app.test_client()
20
21     def test_http_routes(self):
22         """
23         Tests that all routes deliver the
24         pages without errors
25         """
26         response = self.app.get('/')
27         assert response.status_code == 200
28         response = self.app.get('/list_data')
29         assert response.status_code == 200
30         response = self.app.get('/map')
31         assert response.status_code == 200
32         response = self.app.get('/plot')
33         assert response.status_code == 200
34
35     def test_error_handlers(self):
36         """
37         Testing a random URL to
38         catch a 404 HTTP errors
39         """
40         response = self.app.get('/randomurl')
41         assert response.status_code == 404
42
43     def test_plot(self):
44         """
45         Testing the reception of the image
46         data when an ajax request is sent
47         """
48         response = self.app.get('/plot?refresh=1')
49         assert len(response.data) > 0

```

---

```

1 import os
2 import sys
3 MY_PATH = os.path.dirname(os.path.abspath(__file__))
4 sys.path.insert(0, MY_PATH + '/../server')
5 from analysis import Analyzer
6
7
8 class TestAnalysis(object):
9     """
10     Testing the analysis module
11     which deals with sentiment analysis
12     """
13
14     def test_espone(self):
15         """
16         Testing normal distribution
17         value for 0 mean and 1 variance
18         """
19         from analysis import espone
20         assert espone(0, 0, 1) == 1.0
21
22     def test_analyze_empty(self):
23         """
24         Testing empty tweets
25         and tweets including words not
26         in the dictionary
27         """
28         ana = Analyzer()
29         assert ana.analyze("") == 0.5
30         assert ana.analyze("hzoehfsdl") == 0.5
31
32     def test_analyze_bounds(self):
33         """
34         Testing the bounds of the tweets values
35         """
36         ana = Analyzer()
37         assert ana.analyze("this is a test neutral tweet") <= 1.0
38         assert ana.analyze("this is a test neutral tweet") >= 0.0
39
40     def test_analyze_judgement(self):
41         """
42         Testing the proper judgement of the sentiment analysis:
43         * positive and negative
44         * best and worse tweet values
45         """

```

```

46     ana = Analyzer()
47     assert ana.analyze(":)") > 0.5 and ana.analyze(":('") < 0.5
48     assert ana.analyze("yahoo_yahoo_yahoo") == 1.0
49     assert ana.analyze("zzz_zzz_zzz_zzz_zzz") == 0.0
50
51     def test_analyze_judgement_weight(self):
52         """
53         Testing the value order
54         of arbitrary tweets
55         """
56         ana = Analyzer()
57         assert ana.analyze("i_am_so_happy_!_great_day_!D") > ana.analyze("i_am_so_happy_!D")
58         assert ana.analyze("so_sad_!_feeling_depressed_!:'") < ana.analyze("so_depressed_!:'")
59
60     def test_categories_cardinality(self):
61         """
62         Testing the cardinality of the different
63         categorie sums (positive, negative, neutral)
64         """
65         ana = Analyzer()
66         ctg_count = {'positive': 0, 'negative': 0, 'neutral': 0}
67         text = 'great_day_today_lo!;)_but_still_have_to_work'
68         assert ana.categories_cardinality(text, ctg_count) == 15
69         assert ctg_count['positive'] == 4 # great day lol ;)
70         assert ctg_count['neutral'] == 1 # today
71         assert ctg_count['negative'] == 2 # work still
72
73     def test_categories_weight(self):
74         """
75         Testing the weights of the different
76         categorie sums (positive, negative, neutral)
77         """
78         ana = Analyzer()
79         ctg_total = {'positive': 0.0, 'negative': 0.0, 'neutral': 0.0}
80         ctg_count = {'positive': 4, 'negative': 2, 'neutral': 1}
81         data = [2, 3, 0, 2, 2, 0, -4, 0, 0, -2, 2]
82         tot_pos, tot_neg, tot_neu = ana.weight_categories(data, ctg_total, ctg_count)
83         assert (tot_pos, tot_neg, tot_neu) == (99.47646509317096, -49.392885301738836, 3.9750077625545726)

```

---

```

1  import os
2  import sys
3  import pytest
4  import random
5  MY_PATH = os.path.dirname(os.path.abspath(__file__))
6  sys.path.insert(0, MY_PATH + '/../server')
7  from plotting import ScatterPlot
8
9
10 class TestPlotting(object):
11     """
12     Testing the plotting module
13     Since it requires a database,
14     this part is not tested and
15     fake data is inserted instead
16     """
17
18     def test_plotting_empty(self):
19         """
20         Testing the fit function
21         and the image data when no data is available
22         or when weather data are random
23         and thus not among the allowed discrete values
24         """
25         plot = ScatterPlot(1)
26         img_data = plot.get_image_data()
27         assert not(plot.get_fit_function())
28         plot.set_data([random.random(), random.random()], [random.random(), random.random()])
29         assert not(plot.get_fit_function())
30         assert img_data
31
32     def test_plotting(self):
33         """
34         Testing the fit function
35         and the image data when data are available
36         """
37         plot = ScatterPlot(1)
38         plot.set_data([0, 1], [0, 1])
39         img_data = plot.get_image_data()
40         assert plot.get_fit_function()
41         assert img_data

```

---

```

1  """
2  profile.py performs the profiling of the different
3  sections of the application
4  """

```

```

5
6 times = {
7     'gathering': 0.0,
8     'analysis': 0.0,
9     'push': 0.0,
10    'save': 0.0,
11    'weather': 0.0,
12    'connection': 0.0,
13    'overall': 0.0
14 }
15
16 total = {
17     'gathering': 0.0,
18     'analysis': 0.0,
19     'push': 0.0,
20     'save': 0.0,
21     'weather': 0.0,
22     'connection': 0.0,
23     'overall': 0.0
24 }
25
26
27 def run():
28     for key, _ in times.items():
29         times[key] = times[key] / times['overall']
30         total[key] = total[key] + times[key]
31
32
33 def init_times():
34     times = {
35         'gathering': 0.0,
36         'analysis': 0.0,
37         'push': 0.0,
38         'save': 0.0,
39         'weather': 0.0,
40         'connection': 0.0,
41         'overall': 0.0
42     }
43
44
45 def compute_overall():
46     for k in total.keys():
47         if k != 'overall':
48             total[k] = total[k] / total['overall']
49
50
51 def main():
52     with open('times.txt', 'r') as times_file:
53         for line in times_file:
54             data = line.split(',')
55             times[data[0]] = times.get(data[0], 0.0) + float(data[1])
56             if data[0] == 'overall':
57                 run()
58                 init_times()
59             compute_overall()
60     for k in total.keys():
61         if k != 'overall':
62             print(k+" "+str(total[k]*100))
63
64
65 if __name__ == '__main__':
66     main()

```

---