

Healthcare Twitter Analysis

Working Document

George Fisher *

August 20, 2014

Abstract

‘Healthcare Twitter Analysis’ is an Open Source project to investigate ways to improve the quality of medical care with Data Science techniques applied to Twitter. The project began with two handicaps: (1) the data was inadequate and (2) there was little or no understanding on the part of the crowd-sourced team of the analytical tools available. This paper details the efforts of one participant to solve these problems and to provide a base upon which others can build.

The project website is [[Mehta and Saama Technologies, 2013](#)]; the GitHub repository for this paper can be found at [[Fisher, 2014b](#)].

Contents

I	Summary	4
II	Data Acquisition and Management	5
1	Collect Twitter Data	5
2	CSV vs. JSON	5
3	Text, SQL, NoSQL	6

*George Fisher: [GitHub](#), [LinkedIn](#)

4	Supplemental Data	6
4.1	Twitter & Geo Data	6
4.2	Ontologies	7
5	Geo Data	7
5.1	Geo Tagging	7
5.2	Reverse Geo Tagging	8
5.3	Current State of Twitter Geo Tagging	8
III	Exploratory Data Analysis	9
6	Online Twitter Access	9
6.1	Online with Python	9
6.2	Online with R	9
6.3	Analyzing the Static Project Data	9
7	Sentiment	9
7.1	Summary of Sentiment Analyses	9
7.2	Analysis of Breen, AFINN, Score	10
7.3	Digging Deeper into Sentiment Measures	13
8	Analysis of Text	19
8.1	Summary	19
8.2	Word Clouds	19
8.3	Word Frequency	21
8.3.1	Overall Frequencies	22
8.3.2	Co-Occurrence With Top Hashtags	25
8.4	Latent Dirichlet Allocation	25
8.5	1-, 2-, and 3-Grams for each Hashtag	26

CONTENTS

9 Network Analyses	28
10 Plotting Data on a Map	30
11 Time Series Analyses	34
Appendices	36
Appendix A Other Medicine-Related Twitter Projects	36
Appendix B fields_added_to_twitter_json.txt	37
Appendix C Details of the S3 Twitter json Data Files	40
Appendix D Amazon Web Services EC2 & S3	42
Appendix E Sample Programs	45
E.1 R	45
E.2 Python	46
E.3 MongoDB	47

Part I

Summary

The Healthcare Twitter Analysis Project has as its objective finding ways to use Twitter data to help in the provision of health care, whether in the area of medical research, in helping in the daily provision of health care by professionals engaged in the field, in advancing public policy, etc.

The project began with over 6 million tweets relating to medical conditions gathered in the first six months of 2014 and a crowd-sourced team of over one hundred people hoping to do useful analyses of that data.

The project faced two problems initially, which prevented useful progress: (1) the data was incomplete and (2) there was no analytical understanding of the nature of the data. This report details an effort to solve these problems. It does not itself provide any breakthroughs but rather it serves as the baseline of data and analysis upon which further efforts can build.

Most people coming to a project to analyze social media face the same basic problems: they do not have a rich-enough data set to work with and they do not have a solid understanding of the nature of the dataset, especially after it has been filled out with useful supplemental data. This report details a project to remedy these problems.

There were three basics steps:

1. Part II Data Acquisition and Management

Beginning on page [5](#). The data available was very limited; it was voluminous, but fundamental features were missing. Furthermore, there were questions as to the most-appropriate form of data storage to support the analyses. This step solved these problems.

2. Part III Exploratory Data Analyses

Beginning on page [9](#). All of the typical analytical tools were applied to the expanded dataset:

- Online analysis of real-time Twitter data
- Sentiment Analysis
- Textual Analysis
- Network Analysis
- Geographical Plotting
- Time Series

3. Appendix [A](#) Research into Similar Projects

Beginning on page [36](#). Other projects of a similar nature to this one have been done. References have been provided to allow members of this project to benefit from the work of others.

Part II

Data Acquisition and Management

1 Collect Twitter Data

The first step was to add all the Twitter data to the files provided by the project.

The project was provided with 896 csv files by Topsy, a Twitter aggregator [[Topsy, 2014](#)], containing well over 6 million tweets [[Google Drive, 2014](#)]. The files fairly comprehensively cover tweets concerning a wide range of medical conditions, for the first six-months of 2014.

However the data included only the text of the tweet, its originating user and a score calculated by Topsy. While the text might be sufficient for a basic textual analysis, the other data provided by Twitter is clearly of value for more extensive analyses: even simple analyses such as filtering by retweet count or plotting geographic incidence are impossible with the data provided.

My GitHub repo for this project [[Fisher, 2014b](#)] contains a python program that requests the full json from Twitter for each tweet by its id.

I transferred all the project files to Amazon Web Service's EC2 service and ran the python program against them all, producing files that anyone can download for their own use from S3 [[Fisher, 2014a](#)]. See the appendix on page [40](#) for details of the files.

2 CSV vs. JSON

Initially, I focused on creating csv files with this data, and the programs to do so are still on the repo, but after studying Twitter analysis I became convinced that json was more appropriate for two reasons:

1. Every book and paper I have read and every Twitter-analytic program refers to the Twitter data in its json form

2. While MongoDB [[MongoDB, 2014](#)] supports csv files, including a utility for csv loading, MongoDB's native document structure is that of json and since MongoDB seems like a very useful way to store and access Twitter data, it being the one chosen by most other researchers, storing the data in json format seemed to make the most sense.

3 Text, SQL, NoSQL

On the assumption that this project unearths some really useful analytics that can help medical science, it will need to address the question of the best way to store the data. We're using text files at the moment which are easy to use but they get clumsy quickly.

Through 2010 Twitter used MySQL for its data storage. ([[highscalability.com, 2011](#), [Wired, 2014](#), [Quora, 2012](#)]). Subsequent volume growth and the need to serve data from many locations worldwide prompted Twitter to build its own NoSQL database [[Twitter, 2014a](#), [Computerworld, 2010](#)].

Initially I thought we could consider MySQL since Twitter itself had used it but the way it was used was to store the json in a single long text field and to use UDFs to parse it. This is clearly inferior to using MongoDB which does essentially the same thing but is specifically built for indexed json queries.

MongoDB, like all NoSQL datastores, requires map/reduce to perform join operations. To use json, therefore, we must load into MongoDB records that have all the data we need for each record. To the extent that we find useful data in addition to the Twitter data, it will have to be incorporated in the json one way or another.

See the sample program listing in Appendix E on page 47 for an example of loading a MongoDB database with the Twitter json and searching it using Python.

4 Supplemental Data

Finding additional data for the tweets to allow more extensive analyses is clearly a very important area of research but I am aware of only two efforts in the group to do this.

4.1 Twitter & Geo Data

The result of the project detailed in this paper resulted in the following data being added to the Topsy csv files:

- All the Twitter data
- Unix timestamp
- Latitude & longitude
- Country & ISO2 country code
- City
- For country code "US"
 - Zipcode
 - Telephone area code
 - Square miles inside the zipcode
 - 2010 Census population of the zipcode
 - County & FIPS code
 - State name & USPS abbreviation

4.2 Ontologies

Tim Cook [Cook, 2014] has begun work on adding ontology data from BioPortal [BioPortal, 2014]¹.

5 Geo Data

Twitter provides the ability for a user to opt in to tagging their tweets with GPS data; the `tweet["geo"]` and `tweet["place"]` fields contain this data when provided. However, in our dataset fewer than 1.5% of the records have anything but `null` in these fields.

5.1 Geo Tagging

That leaves only the `tweet["user"]["location"]` field which is entered as text by users when they first set up their account and when it is not blank, it often contains either outright junk or, at best, apparently-correct information in a highly-unstructured format

¹Wikipedia: **Ontologies** arise out of metaphysics, which deals with the nature of reality of what exists. The core meaning within computer science is a model for describing the world that consists of a set of types, properties, and relationship types. There is also generally an expectation that the features of the model in an ontology should closely resemble the real world.[Noy and McGuinness,]

that has to be laboriously parsed to provide anything useful to a program attempting to produce a geographic-based analyses.

One example of the problem is the fact that MapQuest [[Mapquest, 2014](#)] returns 47 choices when presented with ‘Pasadena’. As a fan of Jan & Dean, my initial thought was that the city in California near Los Angeles was probably intended most of the time, but it turns out that Pasadena, Texas has a larger population. ‘Swansea’ might make you think of Wales, but there are Swanseas in Australia, Canada and the United States; founded by Welsh emigres, no doubt..

The program `update_geo_data.py` in the repo is my most-recent attempt to crack this code. You are welcome to use it, learn from it or critique it (use the Issues tab of GitHub). The file `HTA_geotagged.json` is the result of running this program on the entire dataset after the full Twitter json had been collected and is publicly available on Amazon’s S3.

5.2 Reverse Geo Tagging

The step following the effort to assign latitude and longitude coordinates is called *reverse geo-tagging* and involves using the coordinates to derive zipcode, city, county, FIPS and state (for the United States). It is unlikely that I will attempt any of this for countries other than the United States, although the data is available and my algorithms may be extensible.

The file `HTA_reversegeo.json` is the output of this process; the culmination of the entire process.

5.3 Current State of Twitter Geo Tagging

I have spot checked a fairly large subset of our dataset and my conclusion is that the results are good enough for exploratory analyses, to test algorithms and to generate hypotheses; surprisingly good given the raw material, but nowhere near good enough to make actual medical-related decisions as the data stands right now.

My location-parsing process can be improved upon but until it is, and really until a much larger number of users opt in for GPS tagging, this data must be considered experimental. What is lacking is the basic lat/lon of the user; the reverse geo-coding process is pretty straightforward and much less subject to criticism –give me a lat/lon and I can very accurately tell you the zipcode– but the fundamental two bits of information, the lat and lon, are simply not reliable at the moment.

Part III

Exploratory Data Analysis

6 Online Twitter Access

6.1 Online with Python

The main section of the repo contains `Instructions for python.pdf` which provides instructions for setting up Python, IPython and installing the prerequisites for online Twitter access.

The `code` folder of the repo contains an IPython notebook `Online Twitter Basics.ipynb` that walks through the process of making online-queries of Twitter and doing simple analyses of the responses. From the notebook you can combine the static project data with real-time queries.

6.2 Online with R

The main section of the repo contains `Instructions for r.pdf` which will get you set up for online Twitter access from RStudio, which is where I did most of the analyses in this document, some of it using the static project data, some it doing real-time Twitter queries.

6.3 Analyzing the Static Project Data

Appendix [E](#) on page [45](#) contains sample programs for processing the static json project data.

7 Sentiment

7.1 Summary of Sentiment Analyses

Sentiment Analysis is a widely-used technique to infer the sentiment of a message from the words and phrases used, including emoticons. The Breen sentiment scoring system seems to be preferable to AFINN and I've run a few samples on subsets of the data. It is not

clear to me that comparing sentiment between disease types is helpful; scoring individual tweets within a diagnosis may be more helpful in identifying candidates for further study.

7.2 Analysis of Breen, AFINN, Score

There are two sentiment measuring systems which popped up in my initial studies of the subject: Jeffrey Breen's [Breen, 2011b, Breen, 2011a] and AFINN [Nielsen, 2011]. In addition, the Topsy data includes a measure called score [Topsy, 2010].

I wondered how the two sentiment measures compared to each other and whether sentiment and score had any relation. Loading the data on Cancer, Cardiovascular and Digestive into R, I had a look:

	breen	afinn	score
min	-6.00	-10.00	6.02
mean	0.00	0.50	8.36
median	0.00	0.00	7.58
stdev	1.23	2.10	1.66
skew	0.00	0.65	1.05
npskew	0.00	0.24	0.47
kurtosis	0.85	2.76	-0.15
max	6.00	16.00	14.62

Table 1: **Statistical Comparison of Sentiments and Score**

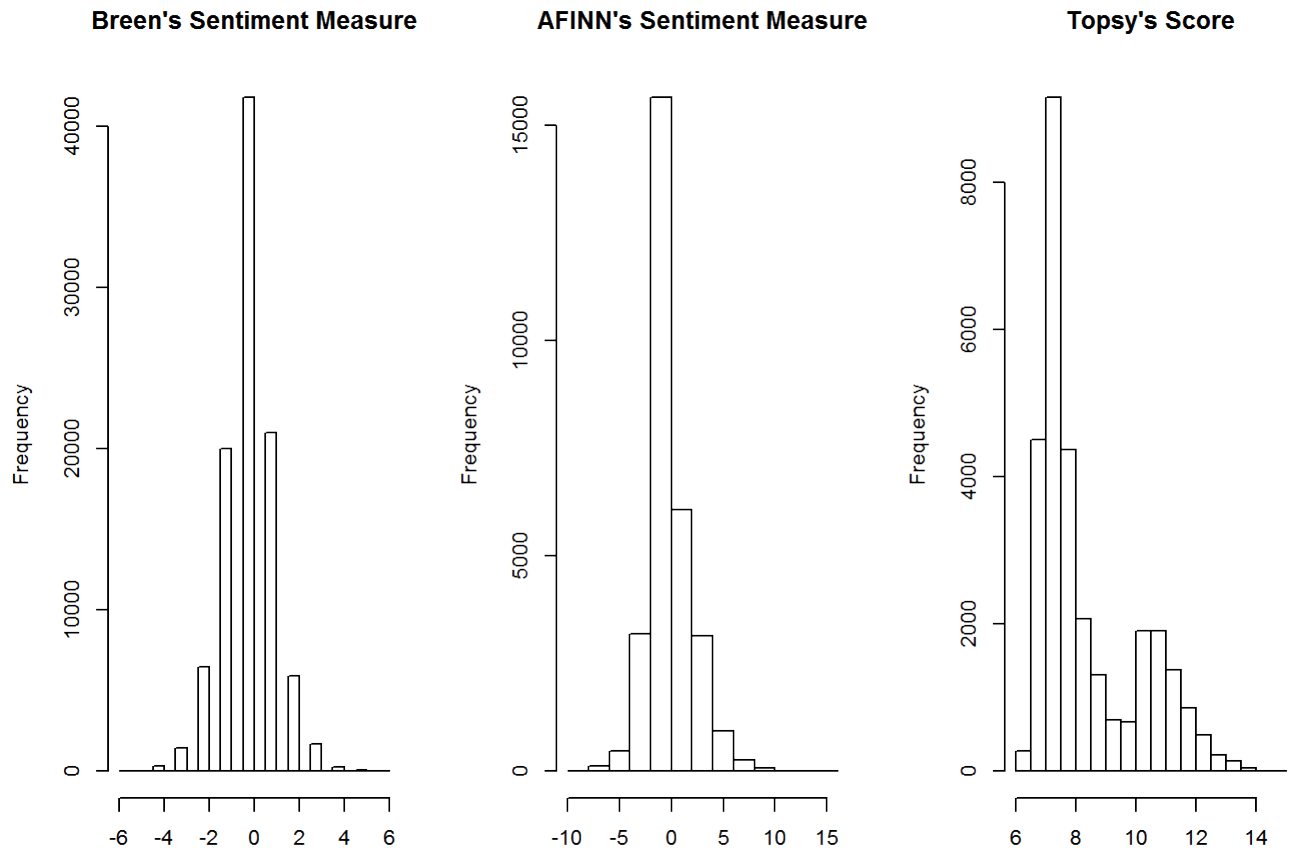


Figure 1: **Distribution of Sentiment Measures and Score** Breen and AFINN are more similar to each other than to score: both have a mean of nearly zero and both are symmetrical around it; but AFINN has a much greater variance and non-normal tail behavior. Score has more of a log or Poisson shape to its distribution, which is bimodal, and is clearly different from the two sentiment scores.

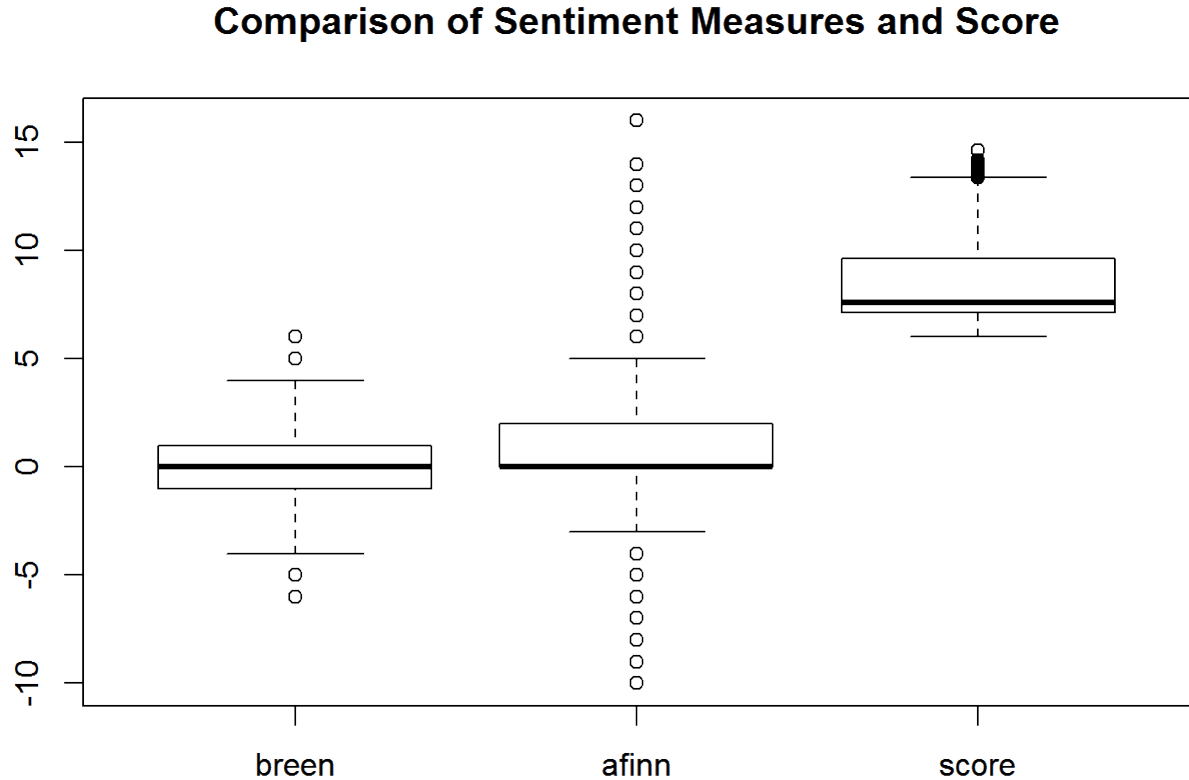


Figure 2: **Distribution of Sentiment Measures and Score** A box plot shows even more starkly the difference in the distributions of these three measures.

First It would seem that score is not created from or predicted by either sentiment measure.

Second The question arises as to which sentiment measure is preferable, if indeed either is adequate: AFINN has a much greater dispersion of its measures, which perhaps is to be expected when dealing with life-destroying diseases; on the other hand, Breen produces a more-nearly-normal distribution and by some accident of Providence, most naturally-occurring phenomena are normally distributed, perhaps including peoples' feelings.

7.3 Digging Deeper into Sentiment Measures

Gaston Sanchez wrote a series in 2012 about Twitter analysis [[Sanchez, 2012](#)]. His work provides an interesting overview of general summary analyses that people do on Twitter data and I have reproduced some of his work here, using R and the Breen sentiment scoring system [[Breen, 2011b](#)], with data from this project in four (randomly-chosen) categories :

1. Blood Disorders
2. Cancer
3. Cardiovascular Diseases
4. Digestive Disorders

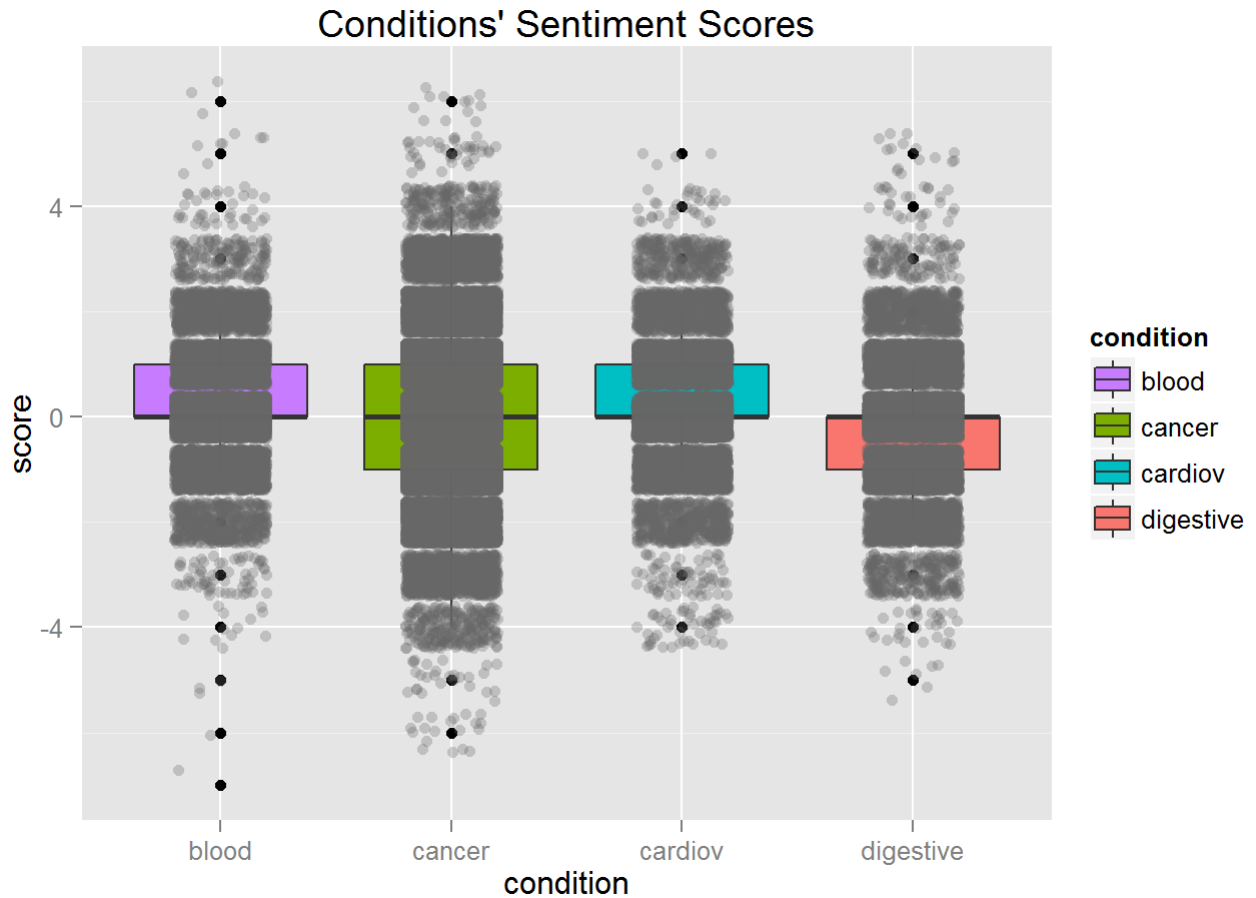


Figure 3: **Distribution of sentiment: Boxplots** The dark gray dots represent the individual data points, roughly 14,000 per condition. The boxes in color represent the inter-quartile distribution of the sentiment for each condition, with bold dots above and below representing outliers beyond the inter-quartile ranges.

They all have their median nearly at zero with a very wide dispersion in both the positive and negative direction. Blood and Cardiovascular disorders seem to be somewhat skewed toward positive overall sentiment while Digestive disorders are skewed toward the negative ranges.

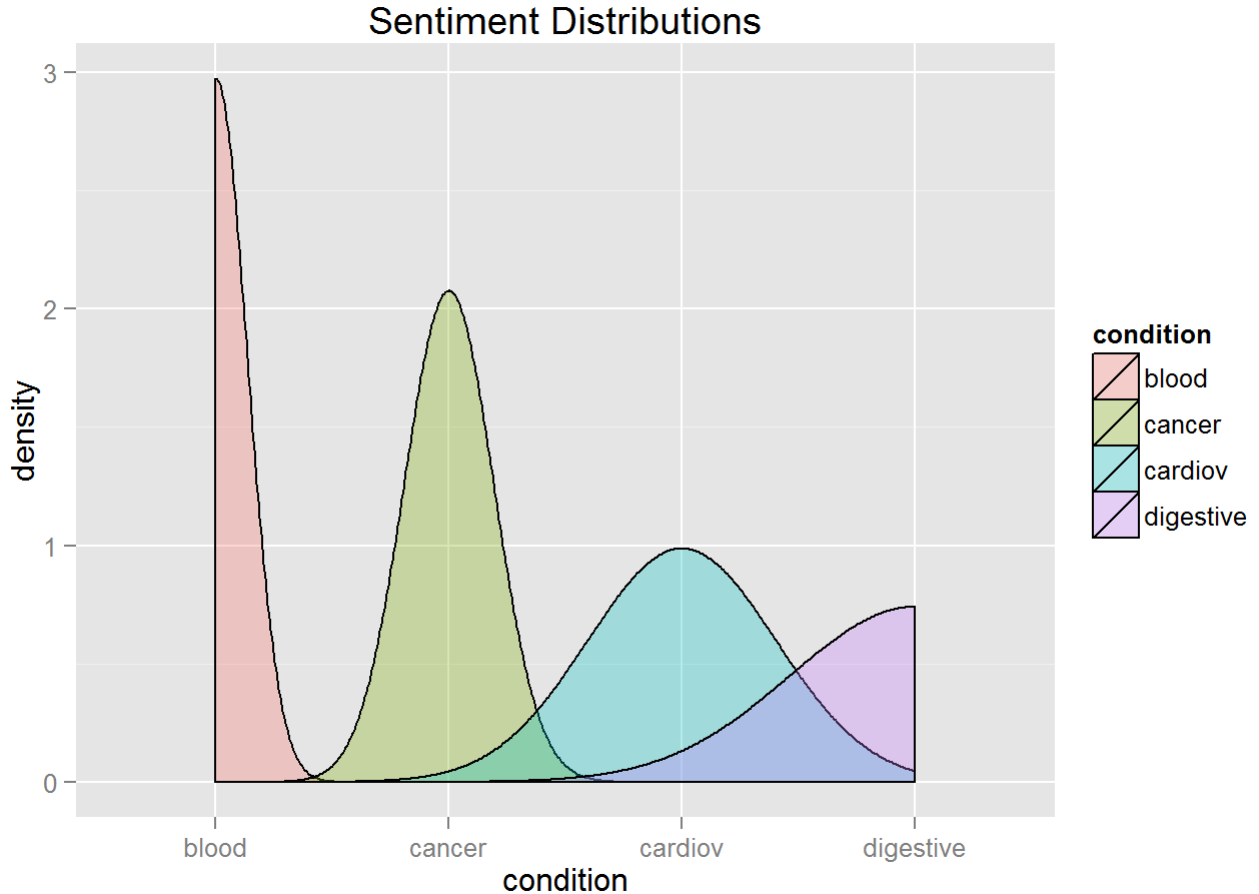


Figure 4: **Distribution of sentiment: Histograms** Another way to look at the distribution of sentiment is to show a smoothed histogram. For each condition, the vertical white line over the label is plotted over the average for that category and the plot shows the distribution around the mean although the left-tail of Blood and the right tail of Digestive are not plotted due to size constraints but they are roughly symmetric. In the study of sentiment measures in section 7.2 beginning on page 10, it was shown that the Breen sentiment measure is symmetric in general and the measures for these specific conditions reflect that.

Blood is in a tight range around its mean, while Digestive has the greatest dispersion.

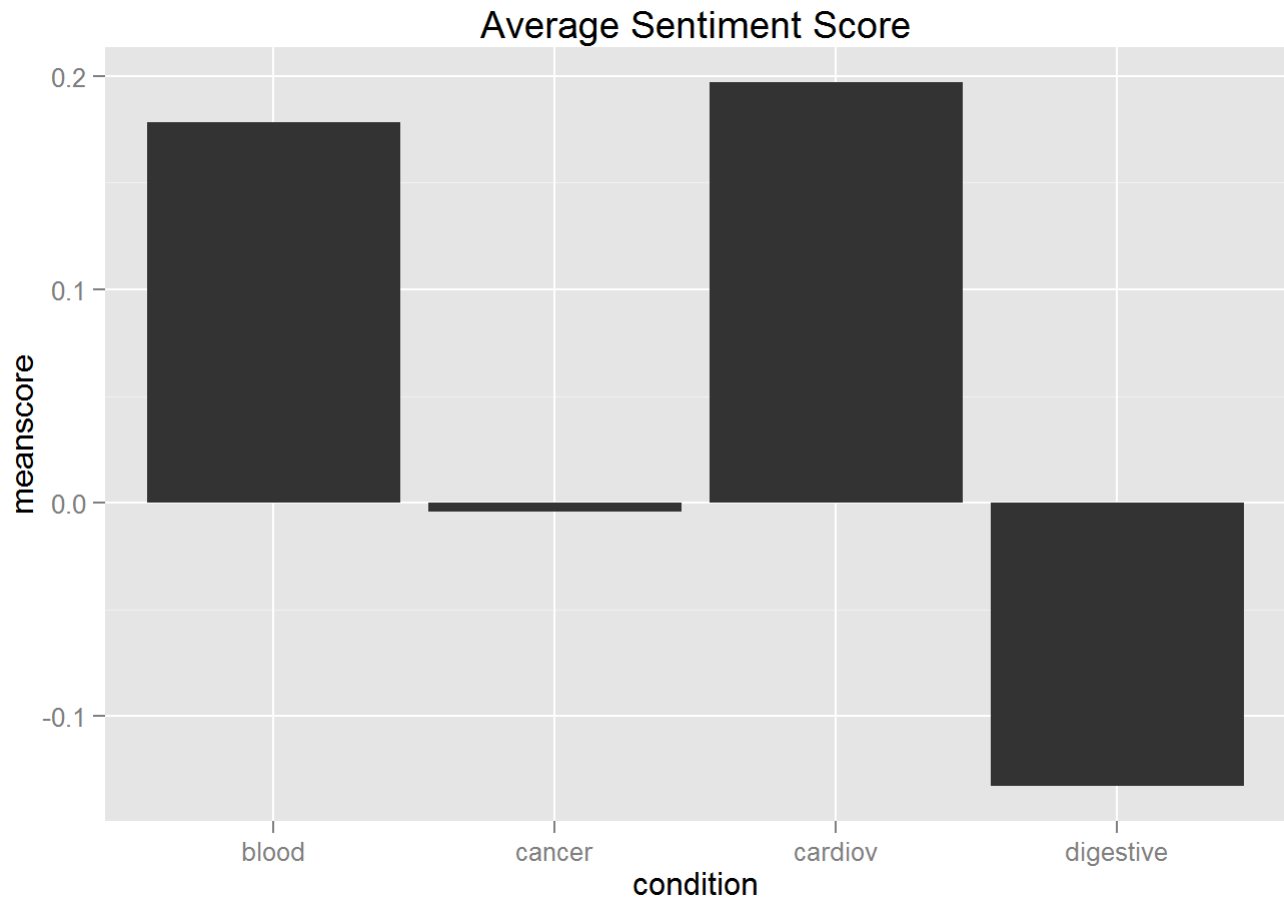


Figure 5: **Average Scores** The averages show us very starkly what we saw in the distributions: Digestive disorders seem to have by far the most negative effect on their sufferers and/or those who tweet about them.

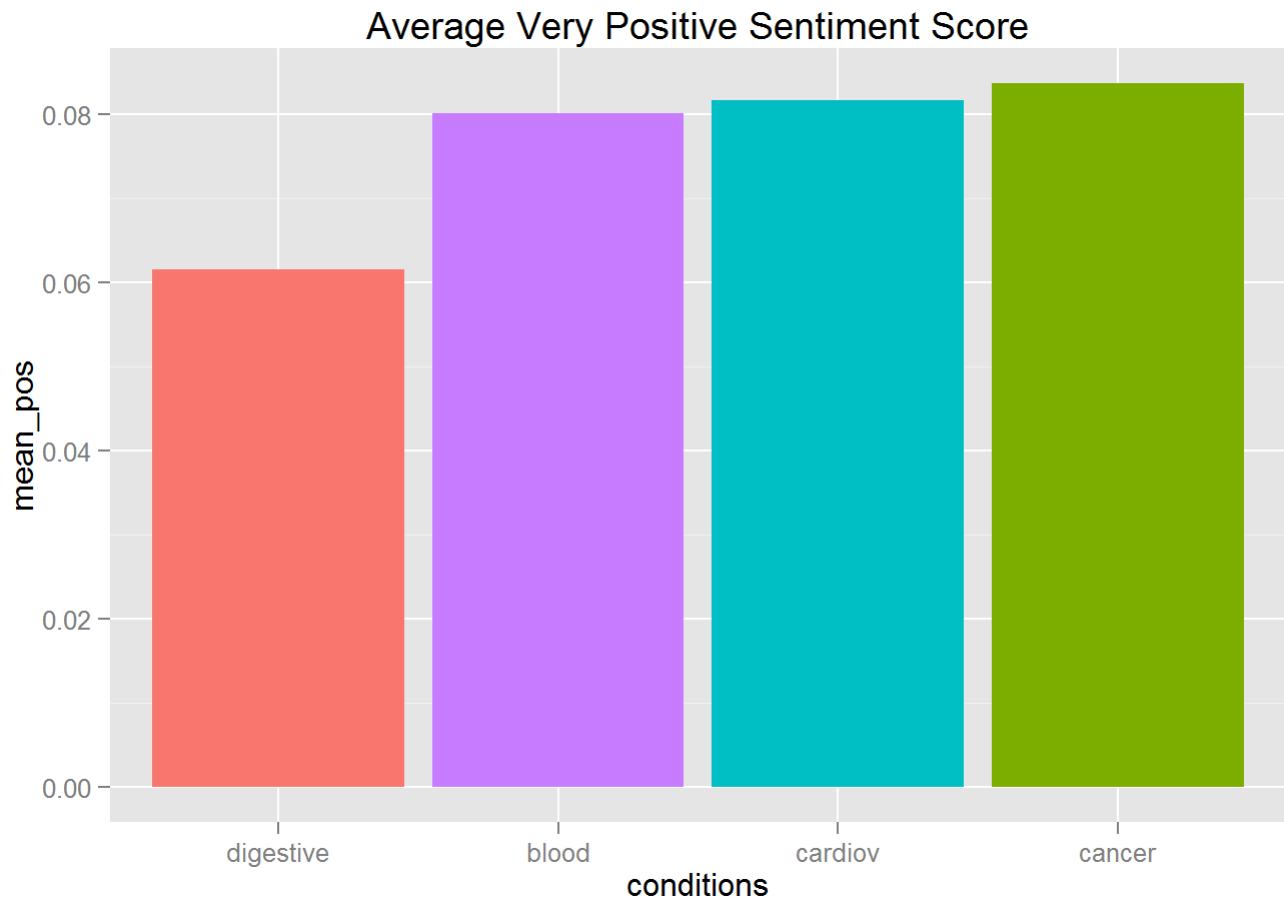


Figure 6: **Average Positive Scores** Looking at the mean scores for only those with a positive sentiment provides more reinforcement for what we have already seen: digestive disorders have a negative psychological effect to the extent of having the lowest mean positive scores.

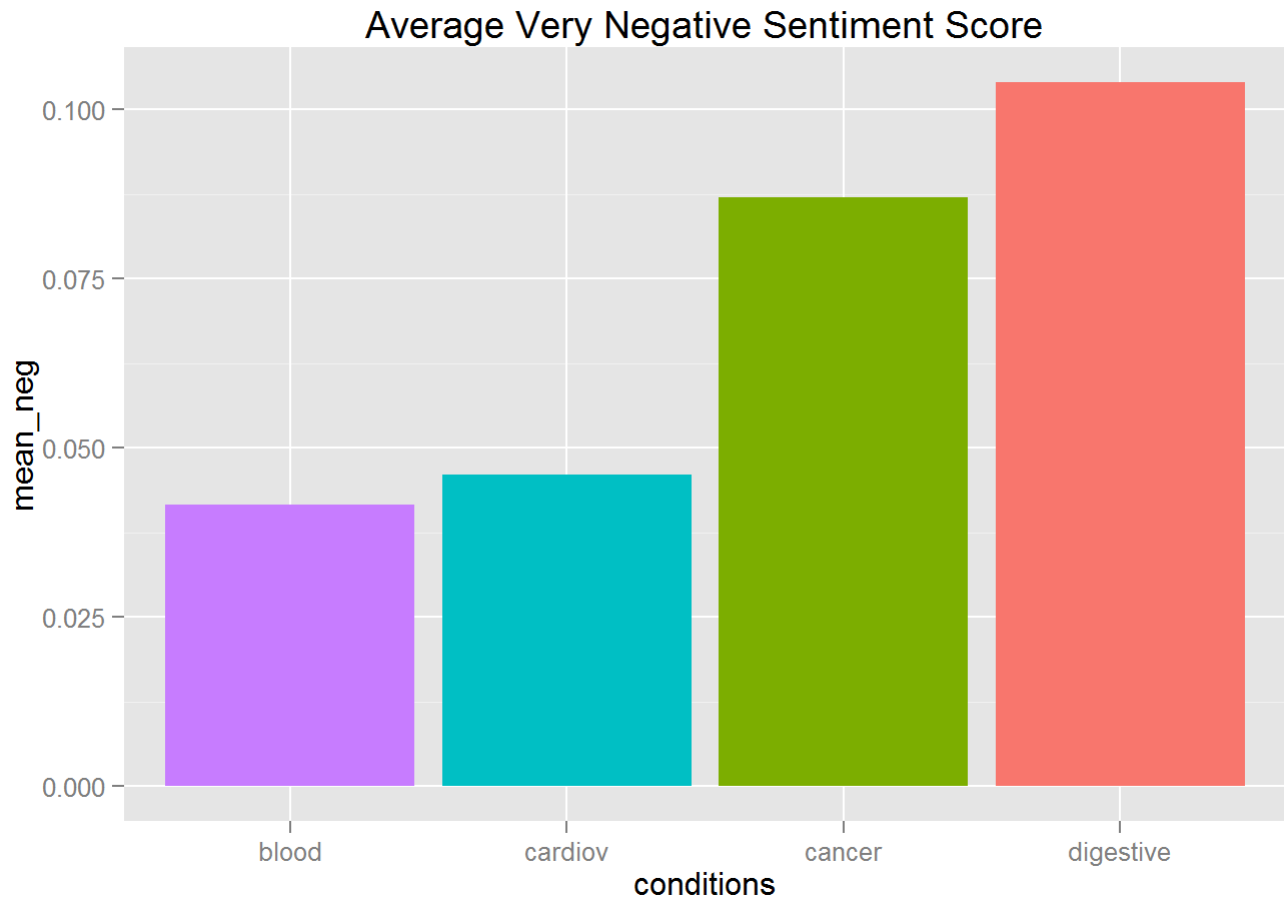


Figure 7: **Average Negative Scores** Looking at the mean scores for only those with a negative sentiment tells the same story: none are good, but of these four, tweets about Digestive Disorders show the greatest tendency toward negativity.

8 Analysis of Text

8.1 Summary

Textual Analysis is another popular analytical technique. By itself, it does not appear to add much value but it is possible that by including additional tags found outside the tweets themselves or else by using more sophisticated techniques we would augment the results sufficiently to provide useful insights.

8.2 Word Clouds

Word Clouds are a very popular EDA technique for text and again with help from Gaston Sanchez [Sanchez, 2012] I have produced a sampling with R and datasets created using the technique described in section 1 starting on page 5.

The corpus was restricted to the first 10,000 tweets in the database for each condition and then further reduced to include only those that had been retweeted more than three times; without these filters the pictures were an incomprehensible mess.



Figure 8: Simple Word Clouds for Four Medical Conditions



Figure 9: **Comparative Word Clouds** show the words specific to the individual conditions. **Commonality Word Clouds** show the words that tweets about the four conditions have in common.

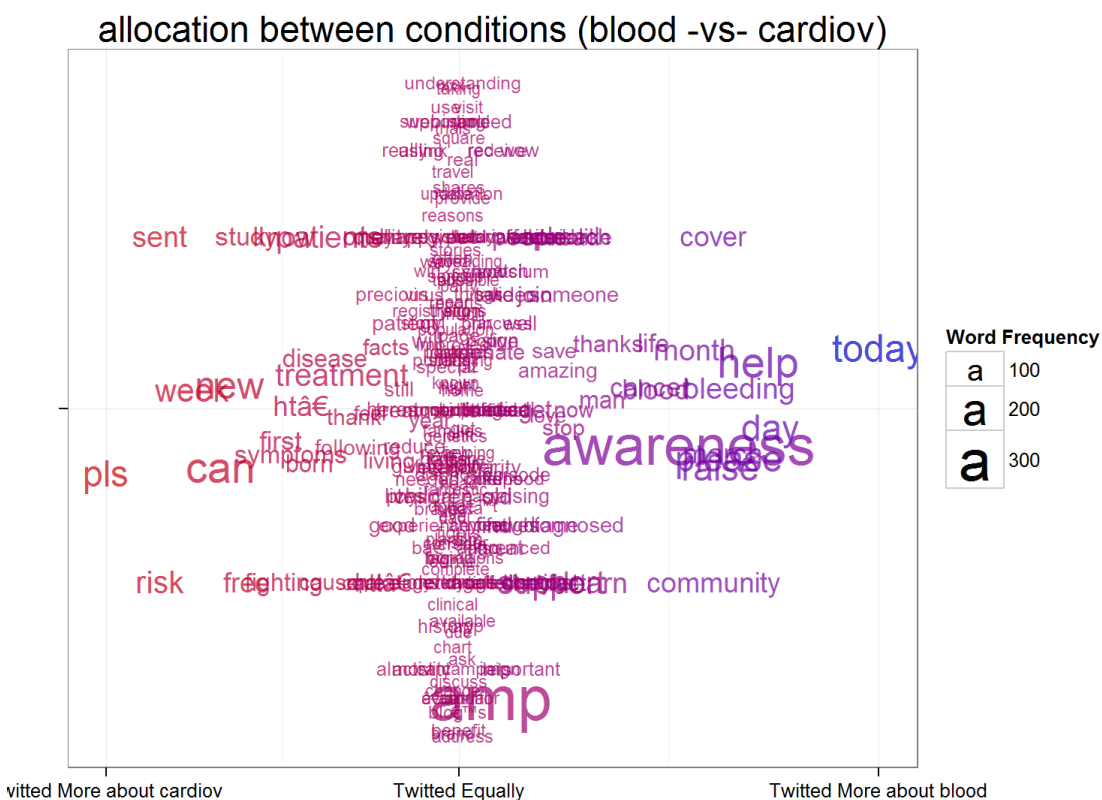


Figure 10: **Conway Comparative Word Cloud of Two Medical Conditions** Comparative word clouds compare all categories together. Conway word clouds show how two categories allocate words between them.

8.3 Word Frequency

The text field of a tweet has four kinds of ‘tokens’:

- Hashtags, beginning with ‘#’, indicating a topic which is propagated across other social media to which the user belongs.
- User Mentions, beginning with ‘@’, indicating a message to/about about a particular user
- URLs, links to other pages or media
- Words, including some emoticons

I have parsed every text field in the database into these four token types, removing stop-words and nuisance strings such as ‘rt’ in the case of words, looking at the various frequencies of tokens:

8.3.1 Overall Frequencies

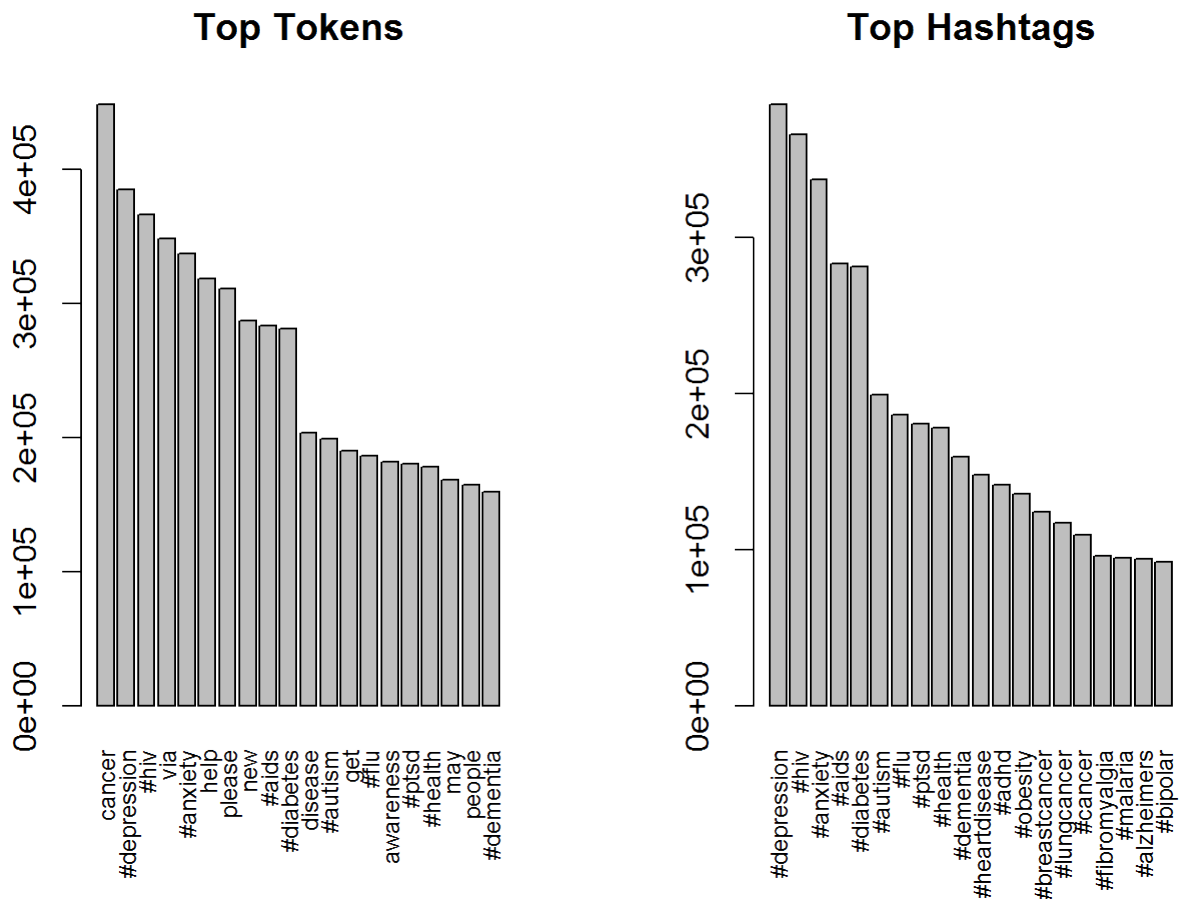


Figure 11: **Most Common Tokens Overall and Most Common Hashtags** The following hashtags are among the top hashtags mentioned in the entire dataset but are not on the list provided by the project:

- #health
- #cancer
- #mentalhealth
- #fibro
- #love
- #pain
- #awareness
- #veterans
- #asd
- #spoonie
- #disability
- #sex
- #weightloss
- #stress
- #glutenfree
- #advice
- #dating

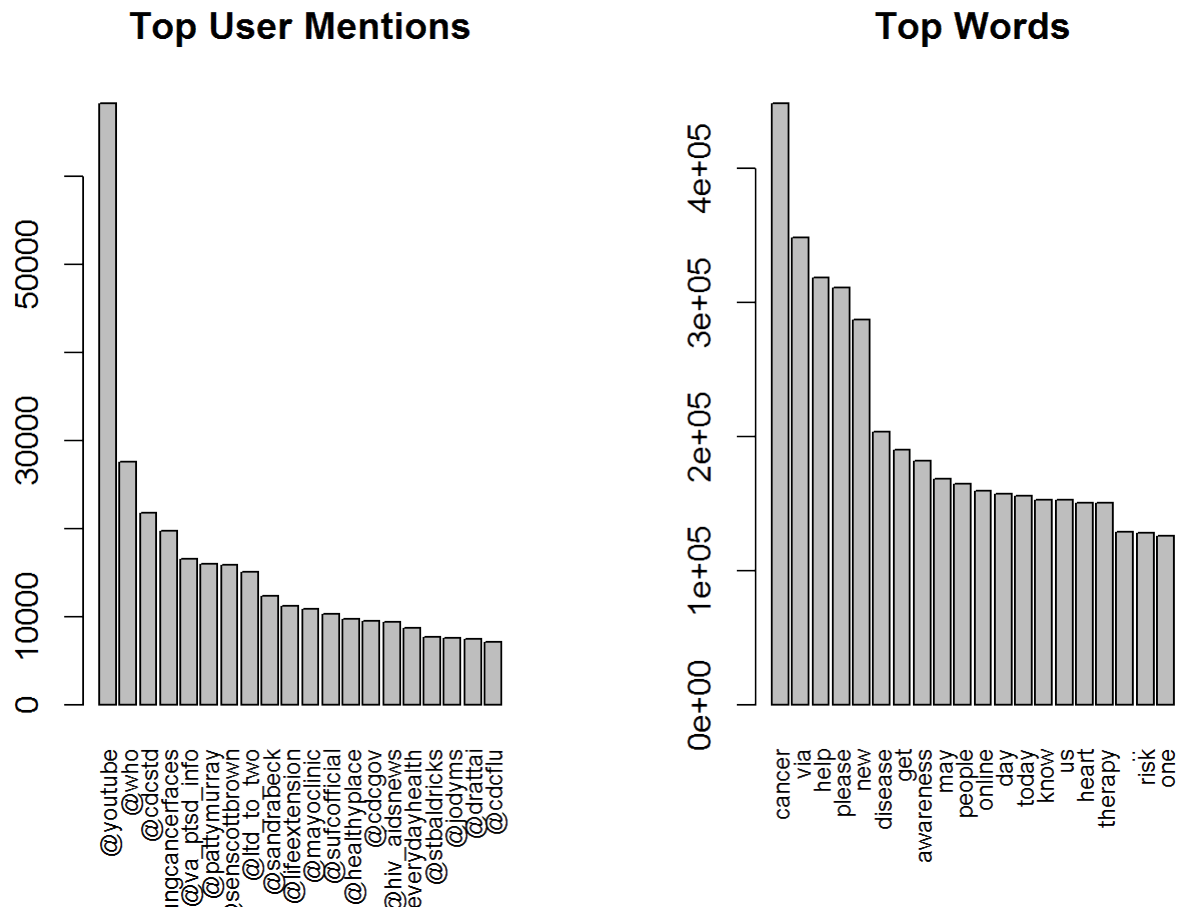


Figure 12: Most Common Users Mentioned and Most Common Words Used

Screen Name	Twitter Description
@youtube	Tweets on news, music and trends from all your favorite channels.
@who	Official Twitter account of the World Health Organization
@cdcstd	Helping people to be safer and healthier by the prevention of STDs
@lungcancerfaces	Faces of Lung Cancer
@va_ptsd_info	National Center for PTSD
@pattymurray	Senator Patty Murray
@senscottbrown	Scott P. Brown
@ltd.to.two	Multiple Sclerosis (PRMS), Fibro may have limited me but it can't destroy me.
@sandrabeck	Sandra Beck #TalkRadio Host #divorce #death #illness #recovery #faith #spiritu
@lifeextension	The latest research on health, wellness, nutrition, & aging.
@mayoclinic	The Mayo Clinic
@sufcofficial	Official Twitter site of Scunthorpe United football club.
@healthyplace	Trusted information on psychological disorders and treatments,
@cdcgov	Centers for Disease Control & Prevention
@hiv_aidsnews	News and developments in the global fight against HIV and AIDS.
@everydayhealth	Powerful weight-loss tools, expert advice & health news and information.
@stbaldricks	Charity funding the world's most promising research to #ConquerKidsCancer.
@jodyms	Writer, blogger. Optimist. Cancer Advocate.
@drattai	Breast Surgeon, President-Elect of @ASBrS
@cdcflu	Flu-related updates from the Centers for Disease Control & Prevention.
@icombat_stress	Motivational Mentor. Hope. Help. Healing. You CAN Turn Your Life Around.
@pozmagazine	The premier HIV/AIDS advocacy
@mndassoc	The Motor Neurone Disease Association.
@clevelandclinic	The Cleveland Clinic
@alldiabetesnews	The Most Comprehensive Diabetes News Aggregator on the Web.

Table 2: **Top Users Mentioned** One current and one ex Senator make the top users mentioned? A soccer club? ...must have been gathered during the World Cup.

8.3.2 Co-Occurrence With Top Hashtags

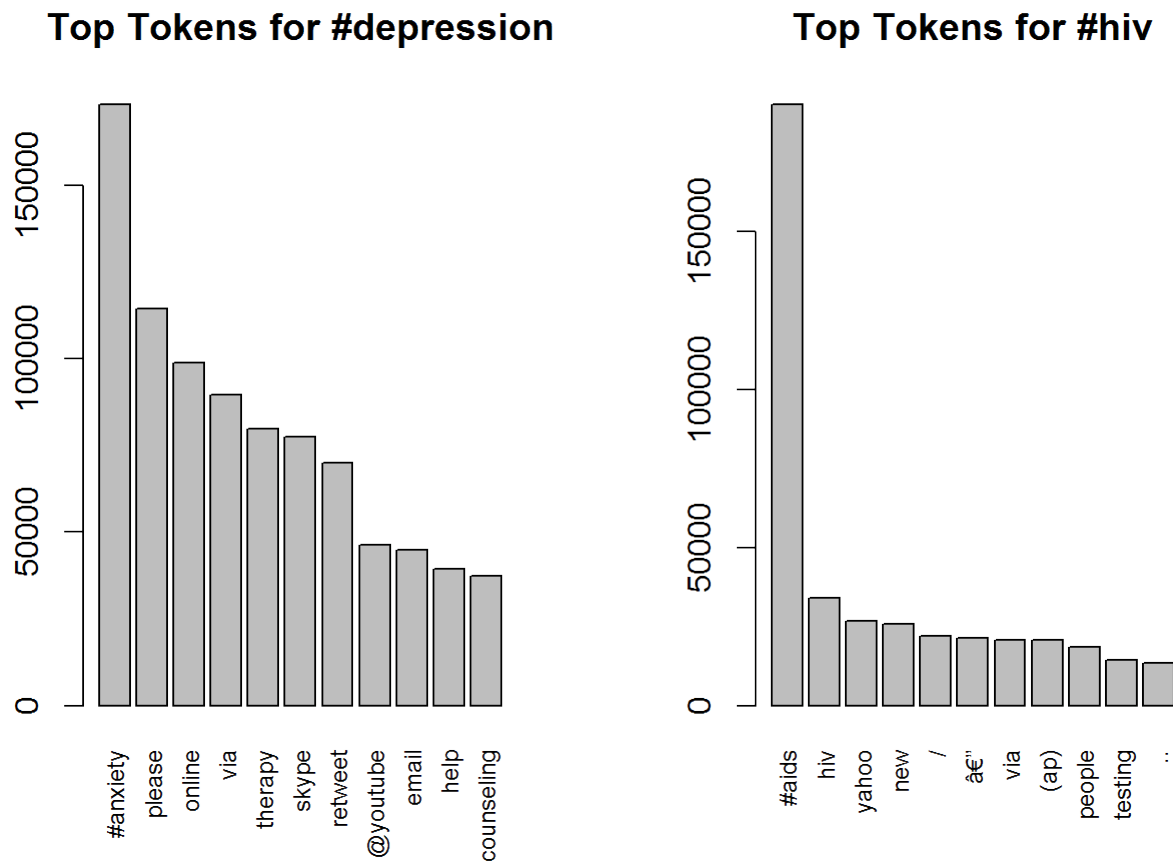


Figure 13: Tokens Most-Commonly Co-Occurring With Two Top Hashtags

8.4 Latent Dirichlet Allocation

I loaded 40,000 tweets of the Blood category into a matrix in R and asked it to tell me the topics; it did a pretty good job: it said there were three:

- sepsis
- myeloma
- hemophilia

8.5 1-, 2-, and 3-Grams for each Hashtag

The following are samples of three csv files in the repo that contain the most-common 1-, 2- and 3-grams associated with each of the hashtags for this project:

hashtag	1-gram	count
rettsyndrome	help	724
influenza	flu	5826
caudaequina	syndrome	20
schizofrenie	van	7
bedwetting	child	95
epilepsy	help	3913
dysautonomia	sharing	915
ppd	postpartum	1089
eds	awareness	1402
sarcoidosis	via	406
trichotillomania	hair	362
afib	atrial	1036
gallbladder	pain	599
testicularcancer	awareness	861
hernia	surgery	123

hashtag	2-gram	count
rettsyndrome	awareness for	250
influenza	out stories	990
caudaequina	please watch	7
bedwetting	your child	59
epilepsy	check out	878
dysautonomia	for sharing	843
ppd	postpartum depression	508
eds	ehlers-danlos syndrome	562
sarcoidosis	news daily	334
trichotillomania	check out	106
afib	atrial fibrillation	931
gallbladder	can cause	274
testicularcancer	to check	225
hernia	detailed general	30

hashtag	3-gram	count
rettsyndrome	\$ awareness for	220
influenza	is out stories	990
caudaequina	please watch share	7
schizofrenie	dialog finse blijkt	3
bedwetting	fitted mattress cover	23
epilepsy	thanks for the	649
dysautonomia	thanks for sharing	760
ppd	should feel ashamed	165
eds	info 085251378519 atau	352
sarcoidosis	news daily review	330
trichotillomania	support this eye-opening	90
afib	with atrial fibrillation	156
gallbladder	can cause severe	273
testicularcancer	about going through	156
hernia	general surgery videos	30
incontinence	disposable pads shaped	211

9 Network Analyses

NodeXL [CodePlex, 2014] is a software package in the form of an Excel template that provides network analysis and visualization. The Python package `networkx` provides a complete programmer's interface for network development and analysis.



Figure 14: Followers of the World Health Organization

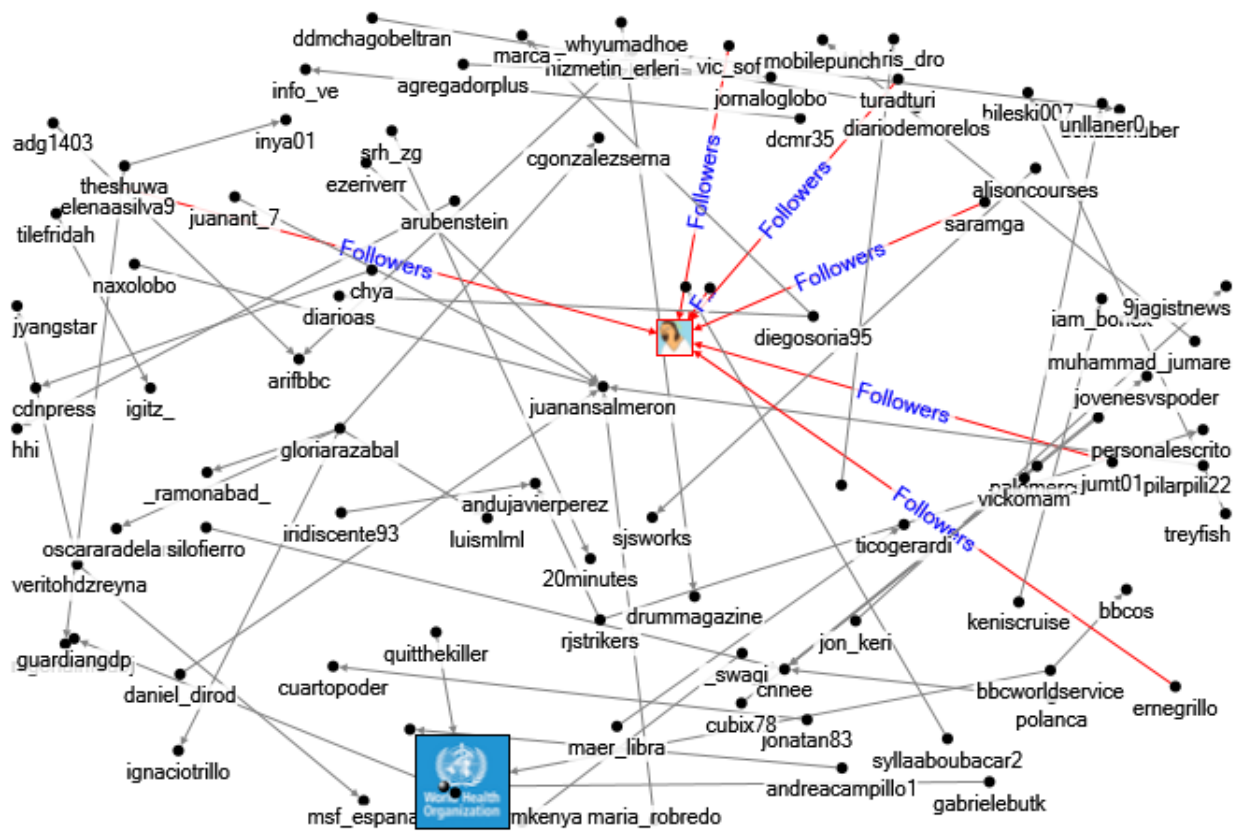


Figure 15: People Tweeting About Ebola

10 Plotting Data on a Map

Worldwide Distribution of All Tweets in the Dataset

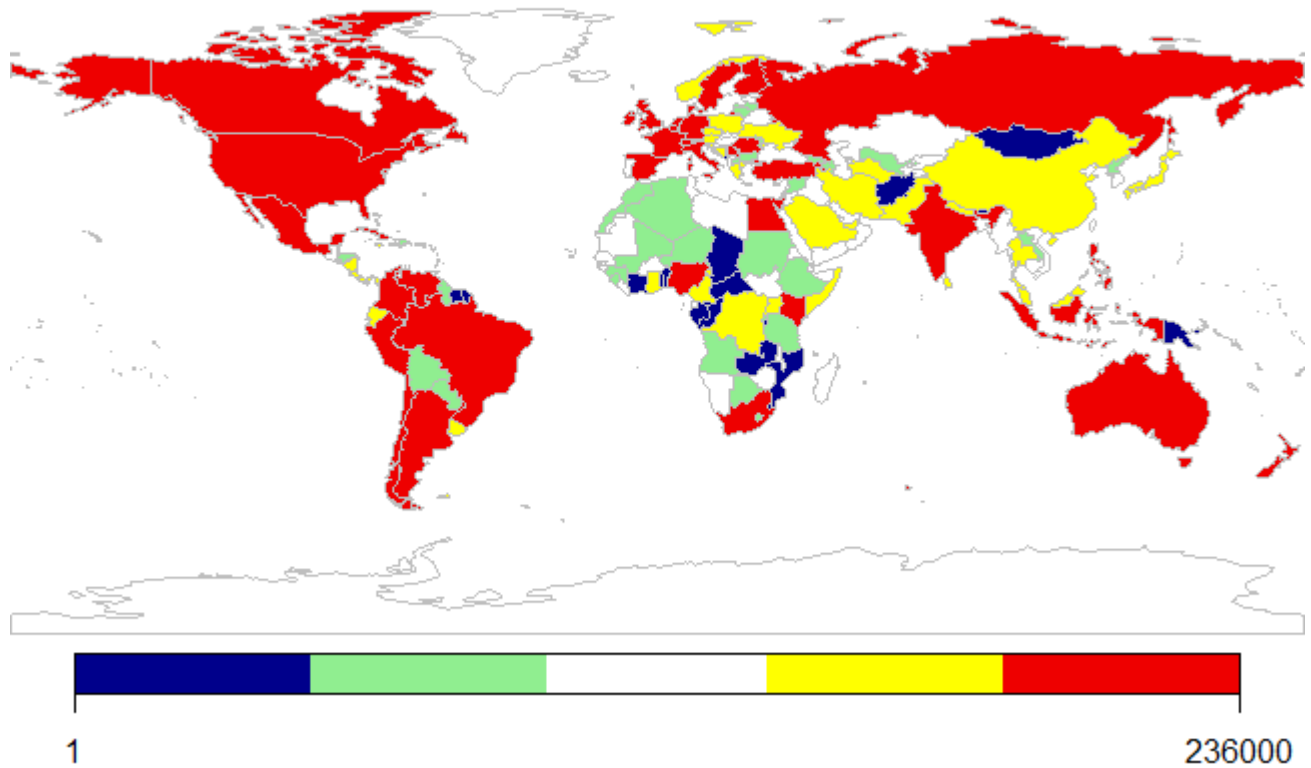


Figure 16: Red: the heaviest users, the top 20%; Yellow: the runners up

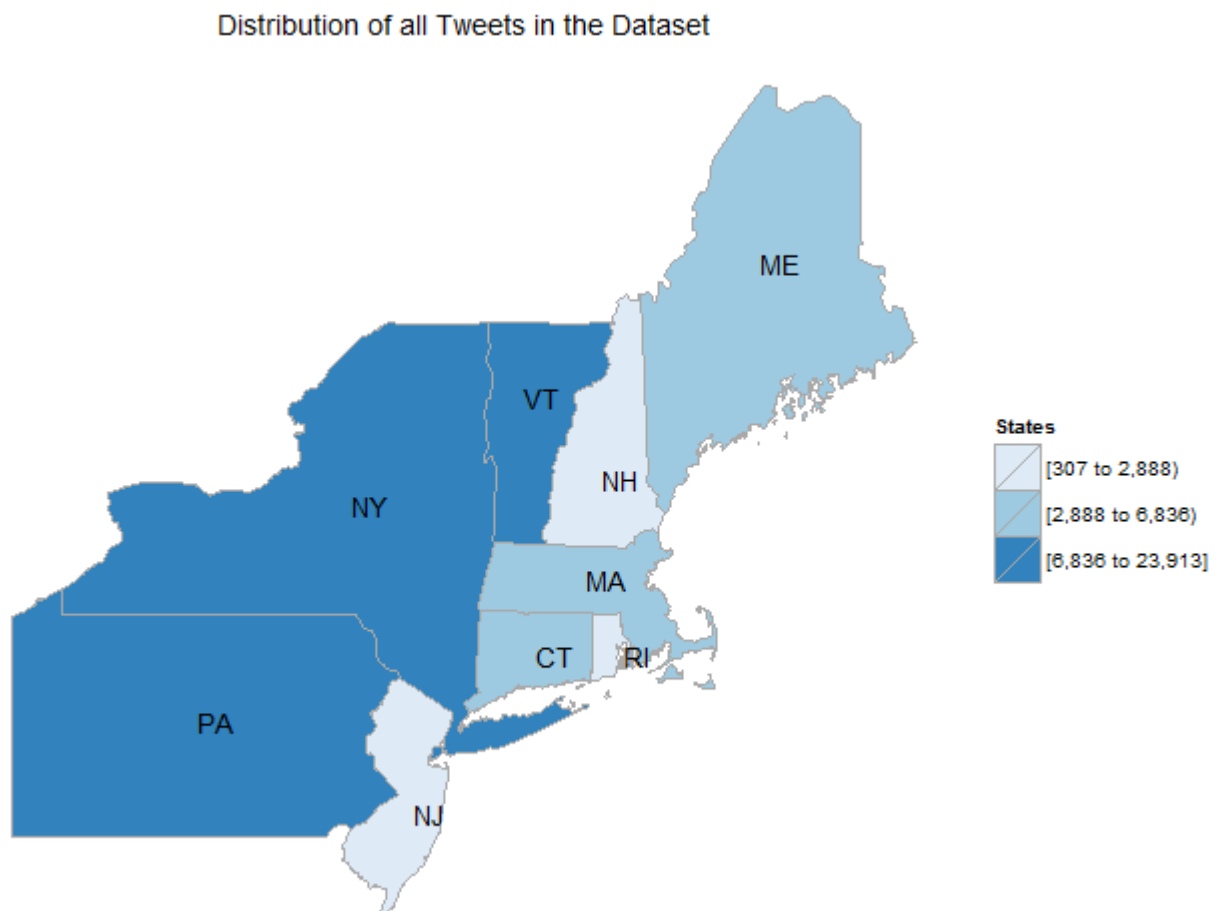


Figure 17: Distribution by US States in the North East

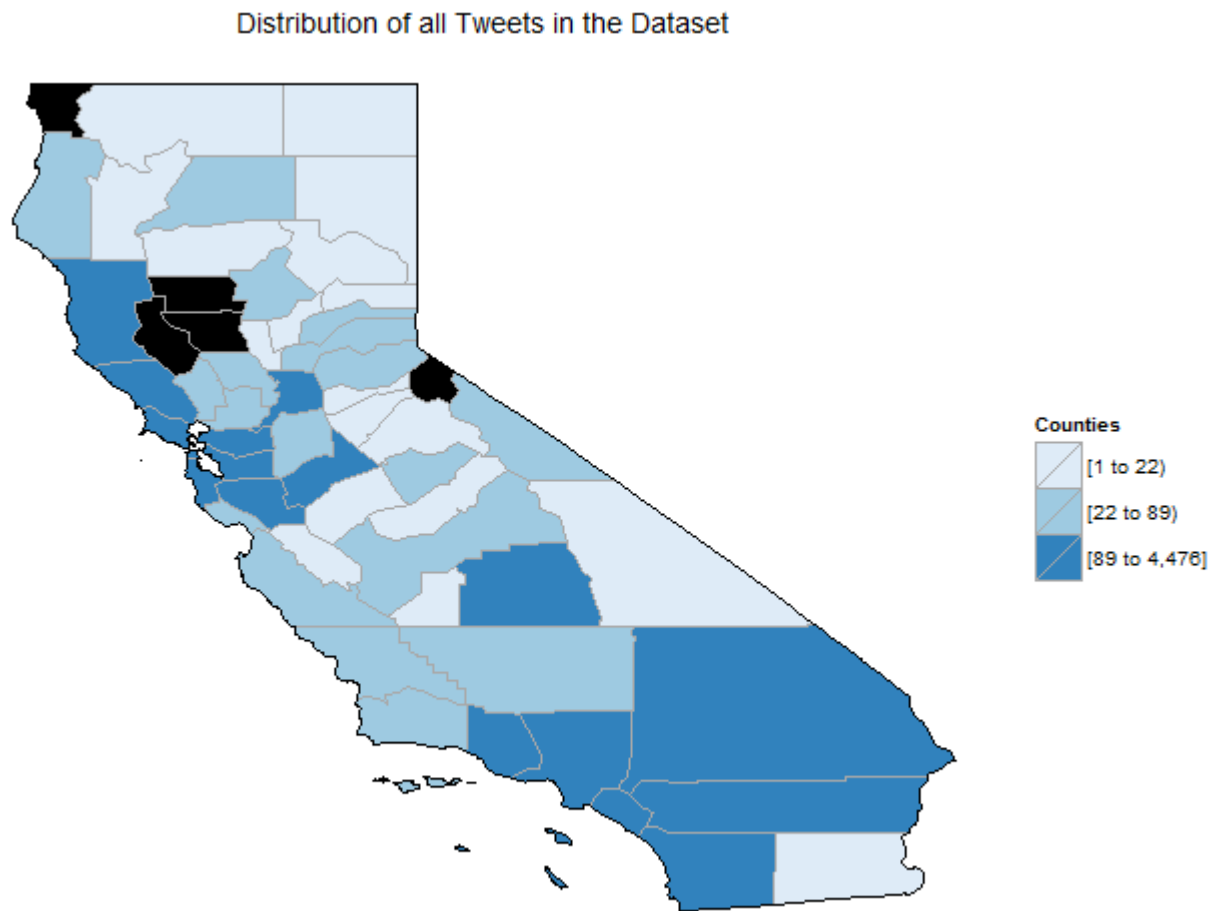


Figure 18: Distribution by County in California

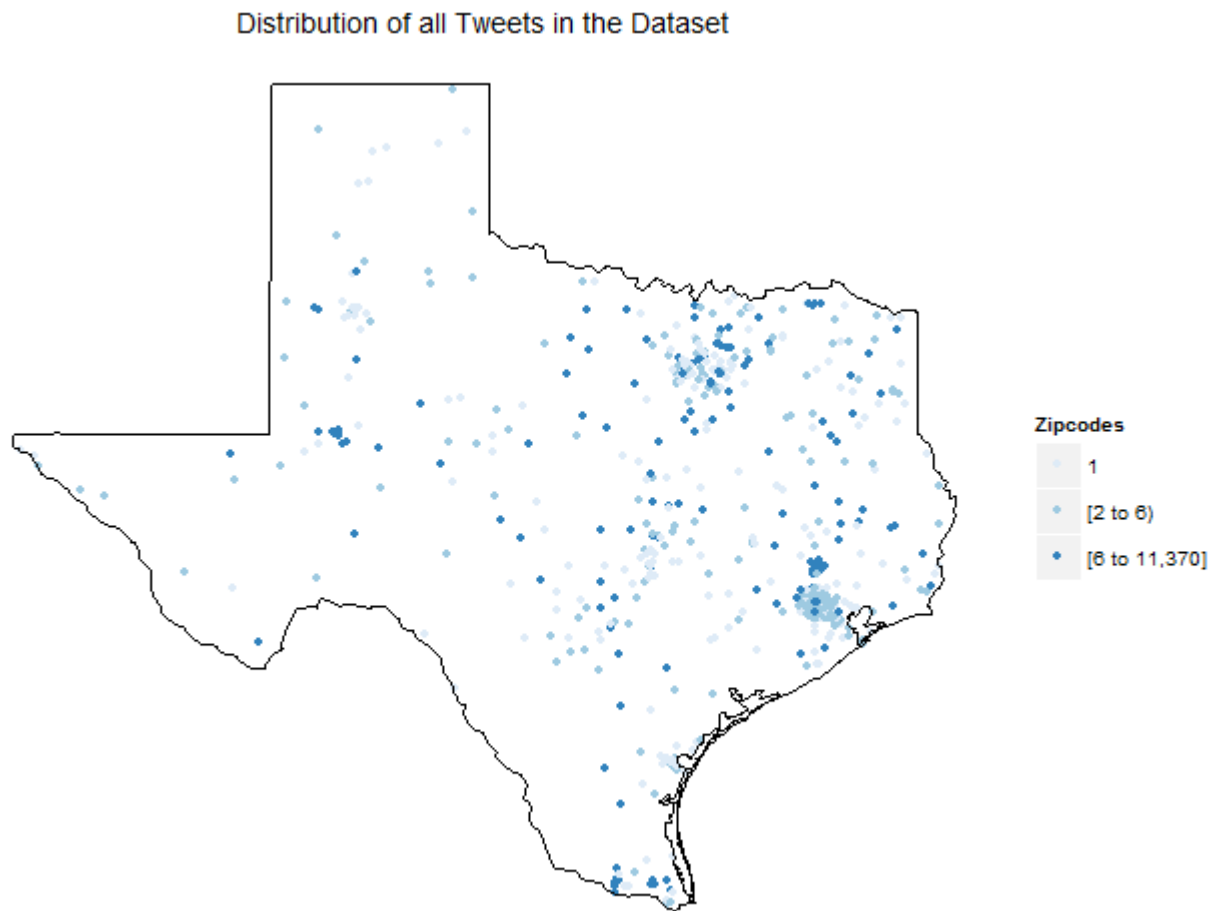


Figure 19: Distribution by zipcode in Texas

11 Time Series Analyses

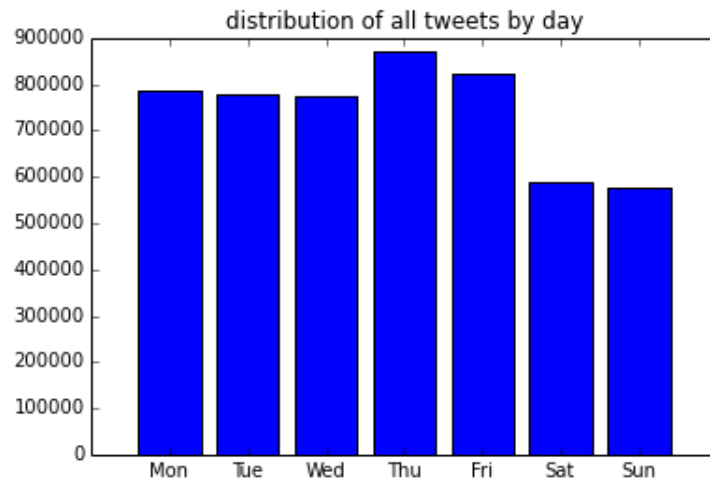


Figure 20: Distribution of all tweets in the dataset by day of the week

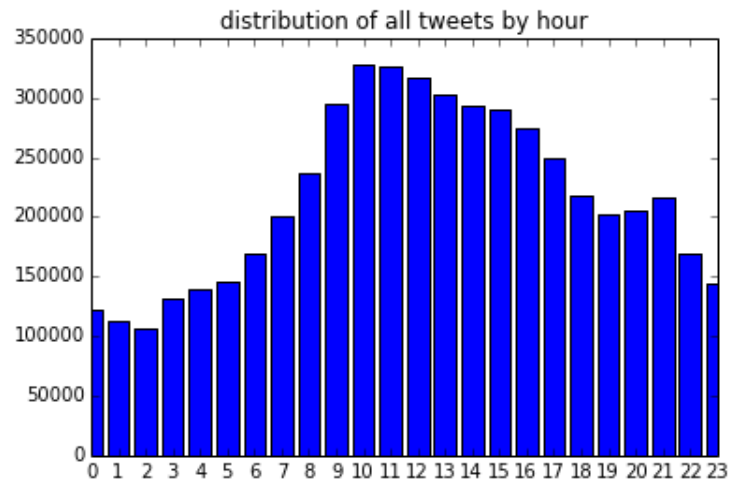


Figure 21: Distribution of all tweets in the dataset by hour of the day

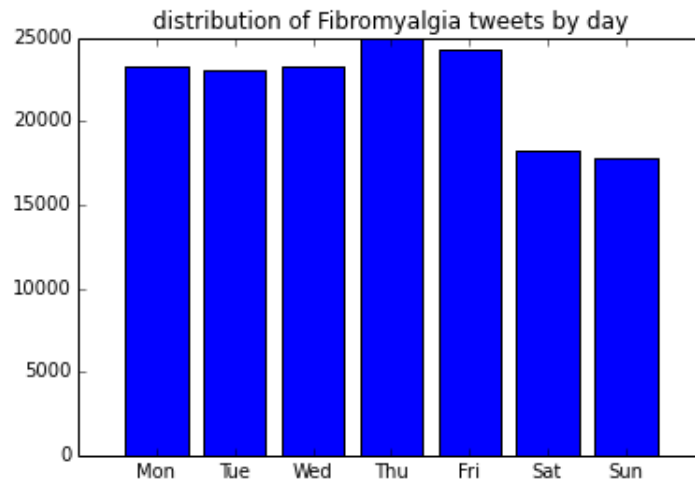


Figure 22: Distribution of Fibromyalgia tweets by day of the week

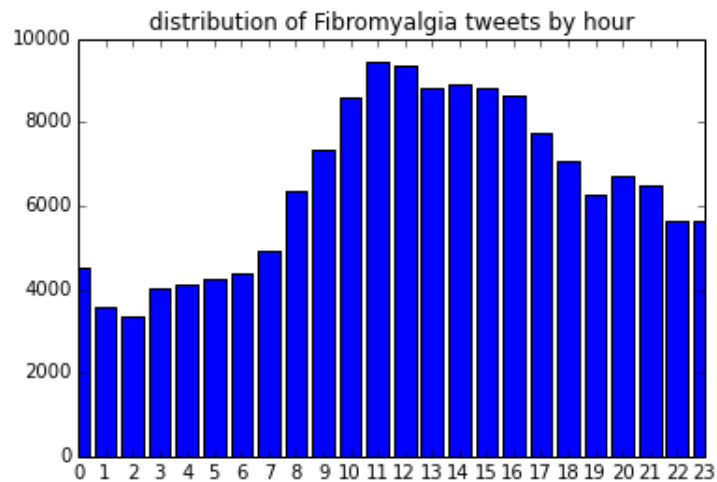


Figure 23: Distribution of Fibromyalgia tweets by hour of the day

Appendices

A Other Medicine-Related Twitter Projects

- How Twitter Is Studied in the Medical Professions:
A Classification of Twitter Papers Indexed in PubMed
[[Williams et al., 2013a](#)]
- What do people study when they study Twitter?
[[Williams et al., 2013b](#)]
- Pandemics in the Age of Twitter:
Content Analysis of Tweets during the 2009 H1N1 Outbreak
[[Chew and Eysenbach, 2010](#)]
- The Use of Twitter to Track Levels of Disease Activity and Public Concern in the U.S. during the Influenza A H1N1 Pandemic
[[Signorini et al., 2011](#)]
- The potential of social networks for early warning and outbreak detection systems:
The swine flu Twitter study
[[Kostkova et al., 2010](#)]
- Using Twitter and other social media platforms to provide situational awareness during an incident
[[Tobias, 2011](#)]
- The other Twitter revolution:
How social media are helping to monitor the NHS reforms
[[McKee et al., 2011](#)]
- A visual backchannel for large-scale events
[[Dork et al., 2010](#)]
- Dissemination of health information through social networks:
Twitter and antibiotics
[[DScanfled et al., 2010](#)]
- Twitter as a communication tool for orthopedic surgery.
[[Franko, 2011](#)]
- Machine intelligence for health information:
Capturing concepts and trends in social media via query expansion
[[Su et al., 2011](#)]

- Social Internet sites as a source of public health information
[[Vance et al., 2009](#)]
- Hospitals are finding ways to use the social media revolution to raise money, engage patients and connect with their communities
[[Galloro, 2011](#)]
- Twitter mining for fine-grained syndromic surveillance
[[syn, 2014](#)]
- Now Trending #health In My Community
[[Department of Health and Human Services, 2012, US Dept. of Health & Human Services, 2012](#)]
- Physicians On Twitter
[[Sabine Tejpar et al., 2011](#)]
- Agencies Use Social Media to Track Food-born Illnesses
[[BM, 2014](#)]
- Social media in vascular surgery
[[Indes et al., 2013](#)]
- Decoding Twitter:
Surveillance and trends for cardiac arrest and resuscitation communication
[[Bosley et al., 2012](#)]
- Twitter as a tool for ophthalmologists
[[Micieli and Micieli, 2012](#)]
- Dissemination of health information through social networks
Twitter and antibiotics
[[Scanfeld et al., 2010](#)]
- All Atwitter About Radiation Oncology:
A Content Analysis of Radiation Oncology-related Traffic on Twitter
[[Jhawar et al., 2012](#)]

B fields_added_to_twitter_json.txt

This file, found in the `files` folder of the repo, shows examples of the json fields added to the full Twitter data by the python programs `get_twitter_json.py`, `update_geo_data.py` and `reverse_geocoding.py`, which can be found in the `code` folder on GitHub [[Fisher, 2014b](#)].

The official guide to Twitter's json structures is here: [[Twitter, 2014b](#)].

The reason for adding the Unix timestamps is for efficient searching in MongoDB; I expect these dates will be part of an eventual index structure and dates in text format are useless for this. To access this field: `tweet["timestamp"]`.

The `twitter["geo"]` field is provided by Twitter and will be filled in from the user's GPS if they opt in, but 98.5% of the records in our dataset had a `null` in this field. The purpose of the program `update_geo_data.py` was to take the `tweet["location"]` field and try to derive the coordinates. `reverse_geocoding.py` is the last step, taking the lat/lon and deriving the city, country and country-code for all geo-tagged tweets as well as zipcode, county, FIPS code and state for US-based tweets.

To access the file name from which this tweet was drawn: `tweet["topsy"]["short_file_name"]`.
To access the originator's screen name as given by Topsy: `tweet["topsy"]["trackback_author_nick"]`.

```
1 Additional Fields in Twitter json
2 =====
3
4 "timestamp": 1389010334.0           # unix timestamp for Twitter's ['created_at'] field
5
6 "user": {                          # provided by Twitter user
7     "location": "New York City",
8 },
9
10 "geo": {                          # derived from ["user"]["location"]
11     "type": "Point",
12     "coordinates": [40.730599,
13                     -73.986581]
14 },
15
16 "geo_reverse": {                  # derived from ["geo"]; data from 2010 US census
17     "areacode": "212",
18     "Land_Sq_Mi": 0.576,
19     "county": "New York",
20     "FIPS": "36061",
21     "state_abbrev": "NY",
22     "country_code": "US",
23     "Type": "",
24     "city": "New York",
25     "country": "United States",
26     "zipcode": "10003",
27     "state": "New York",
28     "Pop_2010": 56024.0
29 },
30
31 "topsy": {                        # fields from original dataset
32
33     "firstpost_date": "01/06/14",
34     "timestamp": 1388984400.0,    # unix timestamp for ["topsy"]["firstpost_date"] fie
35
```

```
36     "url": "http://twitter.com/primary-immune/status/420090415086198784",
37     "score": 7.2846317,
38     "trackback_author_nick": "primary-immune",
39     "trackback_author_url": "http://twitter.com/primary-immune",
40     "trackback_permalink": "http://twitter.com/Primary-Immune/status/420090415086198784",
41
42
43     "file_counter": 2,                # original dataset number and name
44     "short_file_name": "Jan to May\\Blood\\Tweets.BloodCancer.csv"
45 }
```

C Details of the S3 Twitter json Data Files

Original csv Files

The original dataset consisted of 896 csv files with 6,543,272 lines, located on Google Drive https://drive.google.com/folderview?id=0B2io9_E3C0quYWdlWjdU3ozbzg&usp=sharing

Twitter json Files

Because of two network failures during the process which forced restarts, three files were produced by the process to add the full Twitter json to the original data. All three are stored in gzip format on S3. In addition to the original json for each tweet, a new field was added ["topsy"] containing information about the original file from which the record was drawn.

- https://s3-us-west-2.amazonaws.com/healthcare-twitter-analysis/bigtweet_file001.gz
1.48 Gb/12.7 Gb, 3,040,986 lines
- https://s3-us-west-2.amazonaws.com/healthcare-twitter-analysis/bigtweet_file329.gz
340 Mb/2.71 Gb, 651,317 lines
- https://s3-us-west-2.amazonaws.com/healthcare-twitter-analysis/bigtweet_file361.gz
793 Mb/6.55 Gb, 1,509,351 lines

The total number of json records is 5,201,654, which is 1,341,618 fewer than the original files, a loss of 20.50%, because either the original record had no id or Twitter rejected it. Note: I did not search for duplicate records or spam.

geo-tagged files

After the files had all the data from the original tweet, the ["geo"] ["coordinates"] field was filled in with the latitude and longitude.

reverse geo-tagged files

Using the latitude and longitude, each record had a new ["geo_reverse"] field added with city & country information for all records plus state, county, FIPS and zipcode information

added for US records.

Process Summary

To derive the Twitter data, the original data was processed in batches of 100 tweets, a limit imposed by Twitter for automated requests, and a text file was appended and saved after each batch was processed.

The process of running the program is fairly fast: on a Dell Windows 7 laptop it processed about 1,700 tweets per minute; however, the elapsed time was much longer because Twitter imposes a limit of around 14,000 tweets per 15-minute interval, so the program goes to sleep every 13,500.

On Amazon's EC2, the program ran to completion in 8 days, 2 hours. The program terminated several times because of network errors and I lost several hours each time before I noticed and restarted the process.

The process to add the ["geo"] coordinates is similar: if the field ["geo"] exists, it is accepted unchanged; otherwise the ["user"]["location"] field is parsed and MapQuest queried for lat/lon which are put into the ["geo"]["coordinates"] field. This process was also run on AWS, taking ??? days, ??? hours to complete.

Finally, the ["geo"] field of every record was queried and an algorithm used to locate the record in a country; if the country was the United States, into a city, zipcode, county and state. Information about the zipcode was added from the 2010 US census, as well as the county FIPS code which is required for mapping at the county level. The elapsed time for this final process on AWS was ??? days, ??? hours.

The programs involved in this process are located on the GitHub repo.

D Amazon Web Services EC2 & S3

AWS EC2 and S3 have rather obscure documentation and operate in basic command-line mode, however once you've mastered them they are quite useful since you can get essentially as much computing power, storage and Internet access as you could possibly need on demand.

EC2 is the name for the service that provides either Unix or Windows servers on demand. S3 is the name of bit-bucket data storage.

On top of the base operating system you have to build your own programming environment. I used IPython, see pages [43](#) and [44](#).

In addition to being quite useful, it is also inexpensive: even with numerous false starts my total bill for this project was less than \$30.00.

Amazon Web Services for background Python

I assume you have an AWS account and an access Key pair for SSH access. On Windows I used Putty as my SSH terminal and WinSCP for FTP; on my iPhone I used Server Auditor.

Setup:

1. I started an EC2 Ubuntu Server 14.04 LTS (HVM), SSD Volume Type - ami-e7b8c0d7 on an x86_64 t2.micro configuration. The SSH logon for such an image is ubuntu@Public IP.
2. I struggled for an entire day trying to figure out how to access files on S3 from EC2; I gave up and FTP'd the entire bunch from Google Drive to the image I had just started. There is a nice tool at <http://timkay.com/aws/> that is helpful but not for 897 files in recursive folder structures.
3. I also had to download
 1. `get_twitter_json.py` and edit it a little for Ubuntu file formats
 2. `mapquest_key.txt`
 3. `twitter_credentials.py`
 4. `twitter_functions.py`
 5. `filename_list.csv` had to be re-created for the Ubuntu file names and locations

Install Python:

1. I downloaded the Anaconda distro:

```
wget http://09c8d0b2229f813c1b93-c95ac804525aac4b6dba79b00b39d1d3.r79.cf1.rackcdn.com/Anaconda-2.0.1-Linux-x86_64.sh
```
2. ... and installed it

```
bash Anaconda-2.0.1-Linux-x86_64.sh
```

Note: 'q' gets you out of the license agreement
3. Reloaded the `.bashrc` ...

```
source .bashrc
```
4. ... and issued the following commands:

```
sudo -i
apt-get update
apt-get install python-pip
pip install oauth2
apt-get install ipython
```

Then I started up the python program in the background

```
nohup python get_twitter_json.py "filename_list.csv" 1 0 &
```

... and exited the shell

```
exit
```

As the program churned through the files I was able to sign on and monitor progress via the `nohup.out` file. I could also watch system statistics through the AWS Management Console and on the iPhone AWS app. I probably could have used `boto` but I didn't try it.

Create S3 zip file

The first step is to compress it:

```
infilename = 'HTA_geotagged.json'
outfilename = 'HTA_geotagged.gz'

import gzip
f_in = open(infilename, 'rb')
f_out = gzip.open(outfilename, 'wb')
f_out.writelines(f_in)
f_out.close()
f_in.close()
```

... and then to move it to S3

Install utilities from <http://timkay.com/aws/>

```
* sudo -i
* apt-get install curl
* curl https://raw.githubusercontent.com/timkay/aws/master/aws -o aws
* vi ~/.awssecret # AWS credentials Ctrl+o :w <enter> Ctrl-o :q <enter>
* perl aws --install
* chmod +x aws
* cd /home/ubuntu
```

Then you can enter `s3put <S3 bucket name> <local file to be transferred into S3>`

SQLITE3 is required for my reverse geo-tagging functions: `sudo apt-get install sqlite3 libsqlite3-dev`

E Sample Programs

The individual tweets in the `HTA.reversegeo.json` file can be accessed as follows:

E.1 R

Listing 1: Read a json file with R

```
1 #
2 # Note: in all my analyses I used Python to read the json file
3 #       and created csv files with the specific data I needed for
4 #       the analyses I wanted to do in R.
5 #
6 #       Loading over 5 million records into an R data.frame
7 #       is not a good idea: either it won't work at all because
8 #       of the configuration of your machine or else it will be
9 #       horribly slow.
10 #
11 #       Nonetheless, on a reasonable subset, it can be done
12 #       and many packages are available to work with the
13 #       Twitter json record.
14 #
15 library(rjson)
16 file_path = ("HTA.reversegeo.json")
17 tweet_list = fromJSON(sprintf("[%s]", paste(readLines(file_path), collapse=",")))
18
19 retweets = tweet_list[[i]]$retweet_count
20 user_name = tweet_list[[i]]$user$name
21 text      = tweet_list[[i]]$text
22
23 for (i in 1:length(tweet_list)){
24   if (retweets > 100){
25     cat(sprintf("\n\n%d %s\n%s", retweets, user_name, text))
26   }
27 }
28 ## convert to twitterR structure
29 library(twitterR)
30 tweets = import_statuses(raw_data=tweet_list)
```

E.2 Python

Listing 2: Read a json file with Python

```
1 import json
2
3 with open("HTA_reversegeo.json", "r") as tweet_file:
4     for line in tweet_file:
5         tweet = json.loads(line)
6
7         retweets = tweet['retweet_count']
8         user_name = tweet['user']['name']
9         text = tweet['text']
10        location = tweet['user']['location']
11
12        if retweets > 100:
13            print "\n\n%d %s\n%s"%(retweets, user_name, text)
14            print location
```

E.3 MongoDB

Listing 3: MongoDB from Python:
load from json file

```
1 import json
2 from pymongo import MongoClient
3
4 # start up MongoDB
5 # =====
6 client = MongoClient() # assuming you have the MongoDB server running ...
7
8 # list the databases in this MongoDB instance
9 client.database_names()
10
11 # start over fresh
12 db = client['HTA']
13 posts = db.posts
14 db.posts.remove( { } ) # remove the documents
15 #client.drop_database('HTA') # delete the database
16
17 db = client['HTA'] # create/reference the database
18 posts = db.posts
19
20 # read in the tweets and store those you're interested in
21 # =====
22 with open("HTA_reversegeo.json", "r") as tweet_file:
23     for line in tweet_file:
24         tweet = json.loads(line)
25         if " " in tweet['text']: # or whatever, if anything
26             posts.insert(tweet)
```

Listing 4: MongoDB from Python:
read and process the data

```
1 import json
2 from pymongo import MongoClient
3
4 # start up MongoDB
5 # =====
6 client = MongoClient() # assuming you have the MongoDB server running ...
7
8 db = client['HTA'] # reference the database
9 posts = db.posts
10
11 # list the text, location, coordinates and reverse-geo information
12 # =====
13 for tweet in posts.find():
14     if tweet['geo']:
```

```
15     print tweet['text']
16     print tweet['user']['location']
17     print json.dumps(tweet['geo'], indent=4)
18     print json.dumps(tweet['geo_reverse'], indent=4)
```

Listing 5: Sample result of MongoDB program

```
1 Look who was showing off at the Dr's office today! #hemophilia
2 #bleedingdisorders #raredisease http://t.co/TTSDj23FLl
3 Fort Mill, SC
4 {
5     "type": "Point",
6     "coordinates": [
7         35.00737,
8         -80.945076
9     ]
10 }
11 {
12     "city": "Fort Mill",
13     "areacode": "803",
14     "country": "United States",
15     "zipcode": "29708",
16     "Land_Sq_Mi": 19.093,
17     "county": "York",
18     "state": "South Carolina",
19     "FIPS": "45091",
20     "state_abbr": "SC",
21     "country_code": "US",
22     "Pop_2010": 25035.0,
23     "Type": ""
24 }
```

This example does not take advantage of MongoDB's indexing facility which I am sure would improve search performance over the simple Python search I have illustrated.

There is a post in Stack Overflow on how to assign MongoDB databases to different disk drives: [[StackOverflow, 2014](#)]. I haven't investigated this but it seems very useful since this project's data is so large; I have the raw data on an external terabyte drive and collocating the MongoDB database there seems like a good idea. By default, on a Windows 7 system, the MongoDB files are kept on `C:/data/db`.

References

- [syn, 2014] (2014). Twitter mining for fine-grained syndromic surveillance. *Artificial Intelligence in Medicine* 61 (2014) 153-163.
- [BioPortal, 2014] BioPortal (2014). Bioportal, the worlds most comprehensive repository of biomedical ontologies. <http://bioportal.bioontology.org/>.
- [BM, 2014] BM, K. (2014). Agencies use social media to track foodborne illness. *JAMA*, 312(2):117–118.
- [Bosley et al., 2012] Bosley, J. C., Zhao, N. W., Hill, S., Shofer, F. S., Asch, D. A., Becker, L. B., and Merchant, R. M. (2012). Decoding twitter: Surveillance and trends for cardiac arrest and resuscitation communication. [http://www.resuscitationjournal.com/article/S0300-9572\(12\)00871-4/abstract](http://www.resuscitationjournal.com/article/S0300-9572(12)00871-4/abstract).
- [Breen, 2011a] Breen, J. (2011a). Github twitter-sentiment-analysis-tutorial-201107. <https://github.com/jeffreybreen/twitter-sentiment-analysis-tutorial-201107>. Code provided in conjunction with tutorial slides.
- [Breen, 2011b] Breen, J. (2011b). slides from my r tutorial on twitter text mining #rstats. <http://jeffreybreen.wordpress.com/2011/07/04/twitter-text-mining-r-slides/>.
- [Chew and Eysenbach, 2010] Chew, C. and Eysenbach, G. (2010). Pandemics in the age of twitter: Content analysis of tweets during the 2009 h1n1 outbreak. <http://www.plosone.org/article/info%3Adoi%2F10.1371%2Fjournal.pone.0014118>.
- [CodePlex, 2014] CodePlex (2014). Nodexl: Network overview, discovery and exploration for excel. <https://nodexl.codeplex.com/>. An Excel template that provides for network analysis and visualization.
- [Computerworld, 2010] Computerworld (2010). Twitter growth prompts switch from mysql to 'nosql' database. http://www.computerworld.com/s/article/9161078/Twitter_growth_prompts_switch_from_MySQL_to_NoSQL_database.
- [Cook, 2014] Cook, T. (2014). Tim Cook healthcare twitter analysis repo. <https://github.com/twcook/TweetMapping>.
- [Department of Health and Human Services, 2012] Department of Health and Human Services (2012). Now Trending: #Health in My Community. <http://nowtrending.hhs.gov/>. This contest challenged entrants to create a web-based application that searched open source Twitter data for health topics and delivered analyses of that data for both a specified geographic area and the national level.

- [Dork et al., 2010] Dork, M., Gruen, D., Williamson, C., and Carpendale, S. (2010). A visual backchannel for large-scale events.
<http://ieeexplore.ieee.org/xpl/articleDetails.jsp?arnumber=5613451>.
- [DScanfeld et al., 2010] DScanfeld, Scanfeld, V., and Larson, E. (2010). Dissemination of health information through social networks: Twitter and antibiotics.
[http://www.ajicjournal.org/article/S0196-6553\(10\)00034-9/abstract](http://www.ajicjournal.org/article/S0196-6553(10)00034-9/abstract).
- [Fisher, 2014a] Fisher, G. (2014a). Files with twitter json added.
https://s3-us-west-2.amazonaws.com/healthcare-twitter-analysis/bigtweet_file001.gz https://s3-us-west-2.amazonaws.com/healthcare-twitter-analysis/bigtweet_file329.gz https://s3-us-west-2.amazonaws.com/healthcare-twitter-analysis/bigtweet_file361.gz. These files contain the full Twitter json for the records provided by Topsy.
- [Fisher, 2014b] Fisher, G. (2014b). George Fisher healthcare twitter analysis github repo. https://github.com/grfiv/healthcare_twitter_analysis.
- [Franko, 2011] Franko, O. (2011). Twitter as a communication tool for orthopedic surgery. <http://www.ncbi.nlm.nih.gov/pubmed/22050252?dopt=Abstract>.
- [Galloro, 2011] Galloro, V. (2011). Hospitals are finding ways to use the social media revolution to raise money, engage patients and connect with their communities.
<http://www.ncbi.nlm.nih.gov/pubmed/21513035?dopt=Abstract>.
- [Google Drive, 2014] Google Drive (2014). Healthcare twitter analysis data files.
https://drive.google.com/folderview?id=0B2io9_E3C0quYWdlWjdU3ozbZg&usp=sharing.
- [highscalability.com, 2011] highscalability.com (2011). How twitter stores 250 million tweets a day using mysql. <http://highscalability.com/blog/2011/12/19/how-twitter-stores-250-million-tweets-a-day-using-mysql.html>.
- [Indes et al., 2013] Indes, J. E., Gates, L., Mitchell, E. L., and Muhs, B. E. (2013). Social media in vascular surgery.
[http://www.jvascsurg.org/article/S0741-5214\(12\)02104-0/abstract](http://www.jvascsurg.org/article/S0741-5214(12)02104-0/abstract).
- [Jhawar et al., 2012] Jhawar, S., Sethi, R., Yuhas, C., and Schiff, P. (2012). All atwitter about radiation oncology: A content analysis of radiation oncology-related traffic on twitter. [http://www.redjournal.org/article/S0360-3016\(12\)02712-5/abstract](http://www.redjournal.org/article/S0360-3016(12)02712-5/abstract).
- [Kostkova et al., 2010] Kostkova, P., de Quincey, E., and Jawaheer, G. (2010). The potential of social networks for early warning nad outbreak detection systems: the swine flu twitter study.
[http://ijidonline.com/article/S1201-9712\(10\)00507-2/abstract](http://ijidonline.com/article/S1201-9712(10)00507-2/abstract).

- [Mapquest, 2014] Mapquest (2014). Mapquest developer api. <http://developer.mapquest.com/>.
- [McKee et al., 2011] McKee, M., Cole, K., Hurst, L., Aldridge, R., and Horton, R. (2011). The other twitter revolution: how social media are helping to monitor the nhs reforms. <http://www.ncbi.nlm.nih.gov/pubmed/21325389?dopt=Abstract>.
- [Mehta and Saama Technologies, 2013] Mehta, P. and Saama Technologies (2013). Healthcare twitter analysis website. <https://www.coursolve.org/need/184>.
- [Micieli and Micieli, 2012] Micieli, R. and Micieli, J. A. (2012). Twitter as a tool for ophthalmologists. [http://www.canadianjournalofophthalmology.ca/article/S0008-4182\(12\)00294-3/abstract](http://www.canadianjournalofophthalmology.ca/article/S0008-4182(12)00294-3/abstract).
- [MongoDB, 2014] MongoDB (2014). Mongoddb website. <http://www.mongodb.org/>.
- [Nielsen, 2011] Nielsen, F. Å. (2011). Afiin. <http://www2.imm.dtu.dk/pubdb/views/bibtex.php?id=6010>. Informatics and Mathematical Modelling, Technical University of Denmark.
- [Noy and McGuinness,] Noy, N. F. and McGuinness, D. L. Ontology development 101: A guide to creating your first ontology. http://protege.stanford.edu/publications/ontology_development/ontology101-noy-mcguinness.html.
- [Quora, 2012] Quora (2012). Twitter: Which database system(s) does twitter use? <https://www.quora.com/Twitter-1/Which-database-system-s-does-Twitter-use>.
- [Sabine Tejpar et al., 2011] Sabine Tejpar, MD, P., Wendy De Roock, MD, P., and Derek Jonker, M. (2011). Physicians on twitter. *JAMA*, 305(6):566–568.
- [Sanchez, 2012] Sanchez, G. (2012). Mining twitter. https://github.com/gastonstat/Mining_Twitter. An interesting collection of R programs to do analyses of Twitter data. His use of ggplot was out of date but aside from fixing that, the R code worked pretty well in the context of the data for this project.
- [Scanfeld et al., 2010] Scanfeld, D., Scanfeld, V., and Larson, E. L. (2010). Dissemination of health information through social networks: Twitter and antibiotics. [http://www.ajicjournal.org/article/S0196-6553\(10\)00034-9/abstract](http://www.ajicjournal.org/article/S0196-6553(10)00034-9/abstract).
- [Signorini et al., 2011] Signorini, A., Segre, A. M., and Polgreen, P. M. (2011). The use of twitter to track levels of disease activity and public concern in the u.s. during the influenza a h1n1 pandemic. <http://www.plosone.org/article/info%3Adoi%2F10.1371%2Fjournal.pone.0019467>.

- [StackOverflow, 2014] StackOverflow (2014). How to configure mongo to use different volumes for databases? <https://stackoverflow.com/questions/24811046/how-to-configure-mongo-to-use-different-volumes-for-databases>.
- [Su et al., 2011] Su, X., Suominen, H., and Hanlen, L. (2011). Machine intelligence for health information: capturing concepts and trends in social media via query expansion. <http://www.ncbi.nlm.nih.gov/pubmed/21893923?dopt=Abstract>.
- [Tobias, 2011] Tobias, E. (2011). Using twitter and other social media platforms to provide situational awareness during an incident. <http://www.ncbi.nlm.nih.gov/pubmed/22130339?dopt=Abstract>.
- [Topsy, 2010] Topsy (2010). Cool tool: Topsy finds most influential tweeters on any topic. <http://gigaom.com/2010/07/15/cool-tool-topsy-finds-most-influential-tweeters-on-any-topic/>.
- [Topsy, 2014] Topsy (2014). Topsy website. <http://topsy.com/>.
- [Twitter, 2014a] Twitter (2014a). Manhattan, our real-time, multi-tenant distributed database for twitter scale. <https://blog.twitter.com/2014/manhattan-our-real-time-multi-tenant-distributed-database-for-twitter-scale>.
- [Twitter, 2014b] Twitter (2014b). Twitter json documentations. <https://dev.twitter.com/docs>.
- [US Dept. of Health & Human Services, 2012] US Dept. of Health & Human Services (2012). A partnership between the public and the government to solve important challenges. <https://challenge.gov/?q=334-now-trending-health-in-my-community>.
- [Vance et al., 2009] Vance, K., Howe, W., and Dellavalle, R. (2009). Social internet sites as a source of public health information. <http://www.ncbi.nlm.nih.gov/pubmed/19254656?dopt=Abstract>.
- [Williams et al., 2013a] Williams, S. A., Terras, M., and Warwick, C. (2013a). How twitter is studied in the medical professions: A classification of twitter papers indexed in pubmed. <http://www.medicine20.com/2013/2/e2/>. Medicine 2.0: Social Media, Mobile Apps, and Internet/Web 2.0 in Health, Medicine and Biomedical Research.
- [Williams et al., 2013b] Williams, S. A., Terras, M. M., and Warwick, C. (2013b). What do people study when they study twitter? classifying twitter related academic papers. <https://www.emeraldinsight.com/journals.htm?articleid=17088387>.
- [Wired, 2014] Wired (2014). This is what you build to juggle 6,000 tweets a second. <http://www.wired.com/2014/04/twitter-manhattan/>.