

Big Data: Data Wrangling Boot Camp

R Sentiment Analysis

Chuck Cartledge, PhD

24 February 2018

Table of contents (1 of 1)

- 1 Intro.
- 2 Preview
 - Things that will be happening today
 - How we'll get there
 - And implementation
- 3 Sent. analysis
 - What is it, and why should I care?
 - A visualization
- 4 System req.
- 5 Hands-on
- 6 Q & A
- 7 Conclusion
- 8 References
- 9 Files

What are we going to cover?

- Look to the future
- Talk briefly about sentiment analysis
- Address the polyglot of computer languages
- Talk about our sentiment analysis system
- **Data wrangle** tweets using R

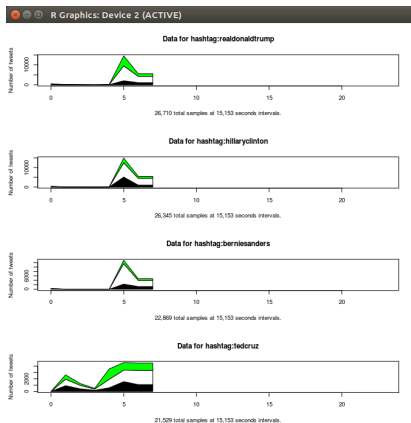




Things that will be happening today

Things that we will be doing.

- 1 **Data wrangle** tweets using R
- 2 Conduct sentiment analysis on tweets
- 3 Look at the sentiments in different ways

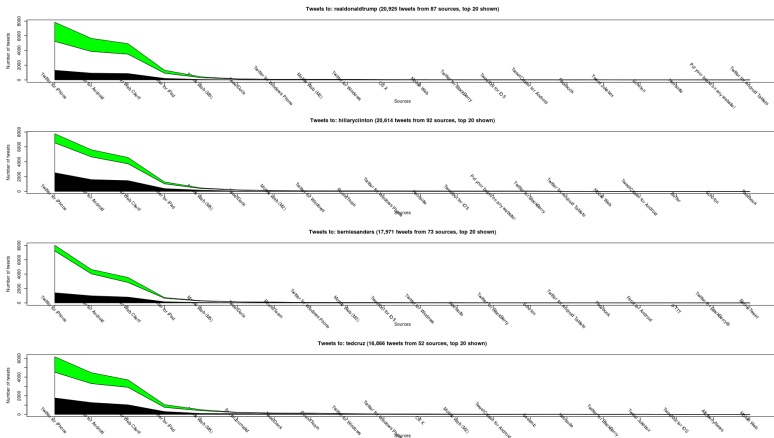


Sentiments over time.



Things that will be happening today

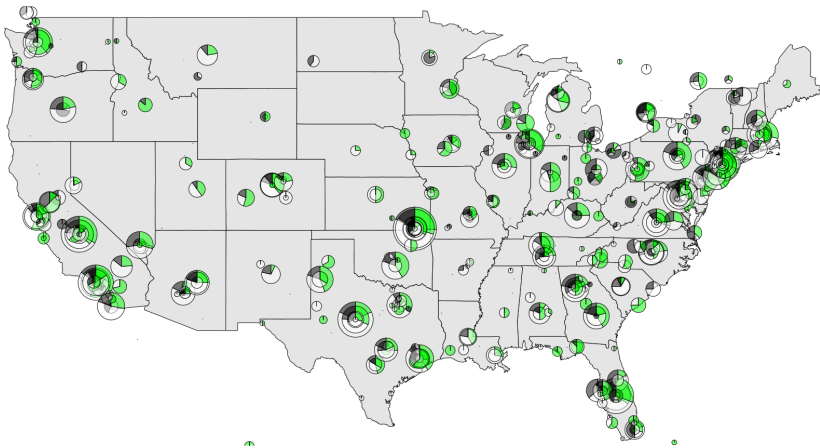
Sentiment by sending device





Things that will be happening today

Sentiment by geographic location

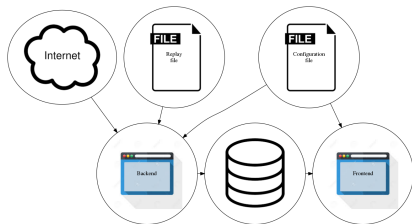


How we'll get to the images

We'll walk before we run.

- Start with a replay file
- **Data wrangle** using the library file
- Go live and download live tweets

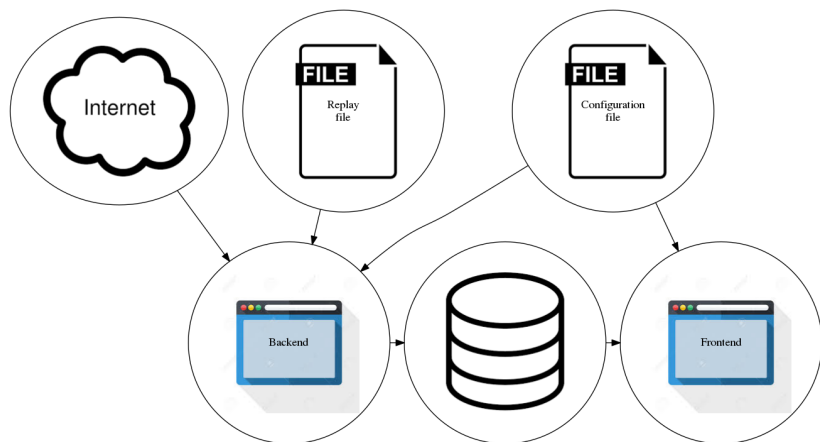
Data flows through the backend, into the database, out the frontend.



The software design document (attached) contains lots of details.

How we'll get there

Same image.

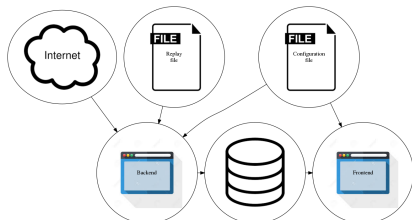


The software design document (attached) contains lots of details.



What some of the files do:

- Replay file – previously recorded tweets
- Configuration file – directives used by both backend and frontend script files
- Shared library file – routines that are common to backend and frontend script files
- Backend file – script file that populates the database from the replay file, or the Internet
- Frontend file – script file that extracts data from the database and presents results



A working definition

“Sentiment analysis (also known as opinion mining) refers to the use of natural language processing, text analysis and computational linguistics to identify and extract subjective information in source materials. Sentiment analysis is widely applied to reviews and social media for a variety of applications, ranging from marketing to customer service.”

W. Staff [3]

More formal definitions

“The field of opinion mining and sentiment analysis is well-suited to various types of intelligence applications. Indeed, business intelligence seems to be one of the main factors behind corporate interest in the field.”

Pang and Lee [2]

“Sentiment analysis, also called opinion mining, is the field of study that analyzes peoples opinions, sentiments, evaluations, appraisals, attitudes, and emotions towards entities such as products, services, organizations, individuals, issues, events, topics, and their attributes.”

Liu [1]



What is it, and why should I care?

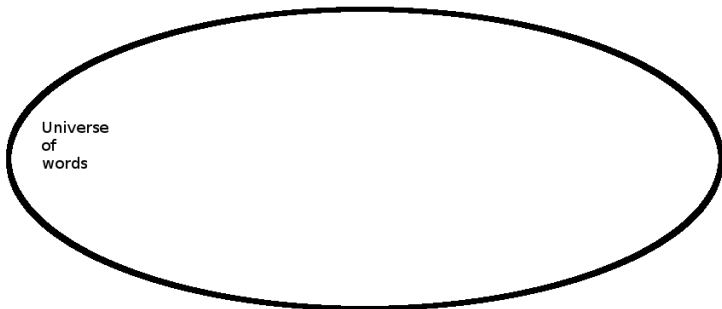
Our approach to sentiment analysis

We will:

- 1 Search the “twitterverse” for tweets using specific hashtags
- 2 Tokenize each tweet
- 3 **Data wrangle** each token
- 4 Remove all stop words from the tokens
- 5 Count number of positive and negative tokens
- 6 Compute the positive, negative, or neutral sentiment for the tokens
- 7 Display the results

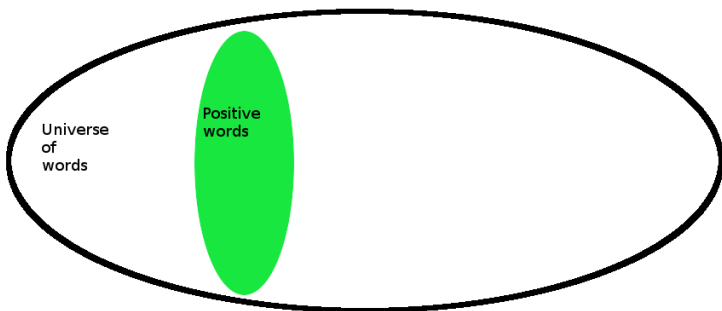
Our approach is language agnostic.

A “Universe” of words





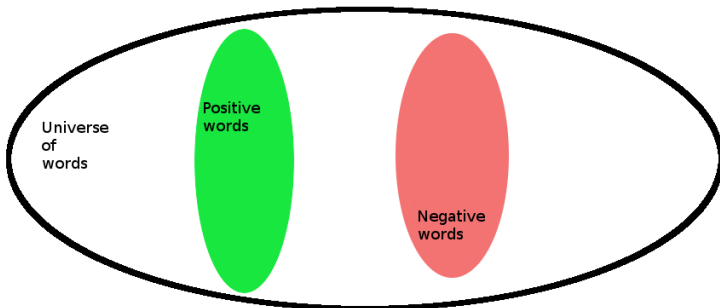
Some words are “Positive”





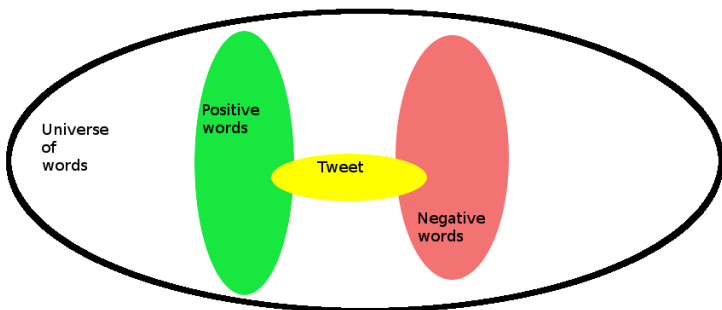
A visualization

Some words are “Negative”



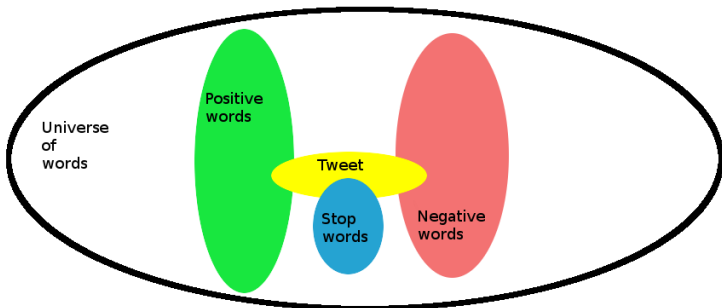


A tweet will/may have positive and negative words





Some words we don't care about



Mechanically this is what we are doing

The steps are:

- 1 Break the tweet into tokens
- 2 Remove stop words from the tokens
- 3 Compute the percentage of remaining tweet tokens that are positive
- 4 Compute the percentage of remaining tweet tokens that are negative
- 5 Classify the tokens as positive, negative, or neutral



Mathematically this is what we are doing

The steps are:

$$tokens = \{\text{words in tweet}\}$$

$$tokensLessStop = tokens - stopWords$$

$$positivePart = positiveWords \cap tokensLessStop$$

$$negativePart = negativeWords \cap tokensLessStop$$

$$classification = \begin{cases} \text{positive,} & \text{if } \frac{positivePart}{tokensLessStop} \leq \text{positiveThreshold} \\ & \text{AND} \\ & \frac{negativePart}{tokensLessStop} < \text{negativeThreshold} \\ \text{negative,} & \text{if } \frac{negativePart}{tokensLessStop} \leq \text{negativeThreshold} \\ & \text{AND} \\ & \frac{positivePart}{tokensLessStop} < \text{positiveThreshold} \\ \text{neutral,} & \text{otherwise} \end{cases}$$

A software design document

The document contains:

- ① Overall system design
- ② Algorithms used through out the system
- ③ Details about the configuration file
- ④ Details about the database tables

softwareOverview.pdf - Adobe Reader
File Edit View Document Tools Window Help

softwareOverview... [X]

ODU Big Data, Data Wrangling Boot Camp
Software Overview and Design
Chuck Carlidge
August 11, 2016

Contents

List of Tables	1
List of Figures	0
1 Introduction	1
2 Software system design	3
2.1 Algorithms used by the software front and back ends	3
2.2 Configuration file	3
3 Database tables	11
4 Notational data structure	12
5 Software on each workstation	13

List of Tables

1 Procedural and functional algorithms, cron scripts	4
2 Configure the server	9
3 Tables to support pattern matching	12
4 System plotting data structure	12
5 System plotting data structure	12

1

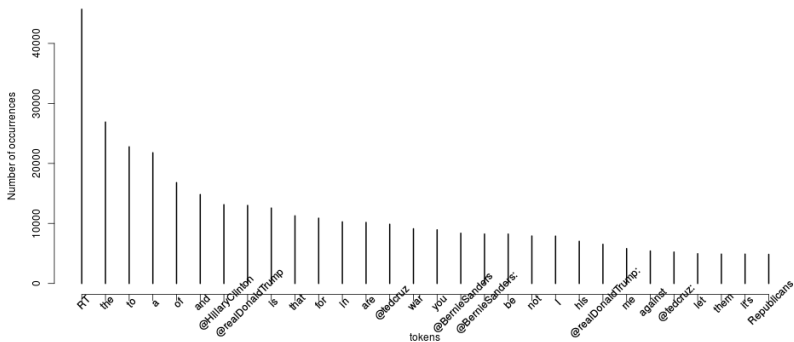
The file is attached.



And implementation

Same image.

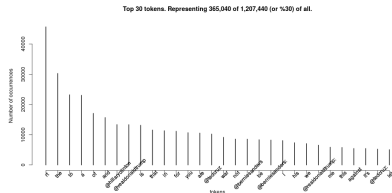
Top 30 tokens. Representing 348,091 of 1,207,440 (or %29) of all.



And implementation

What will happen to the tweets when we make everything the same case

Changing case is easy, unless they are emojis.

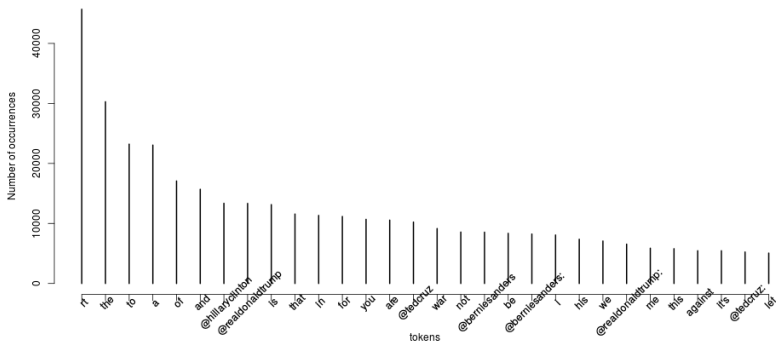




And implementation

Same image.

Top 30 tokens. Representing 365,040 of 1,207,440 (or %30) of all.

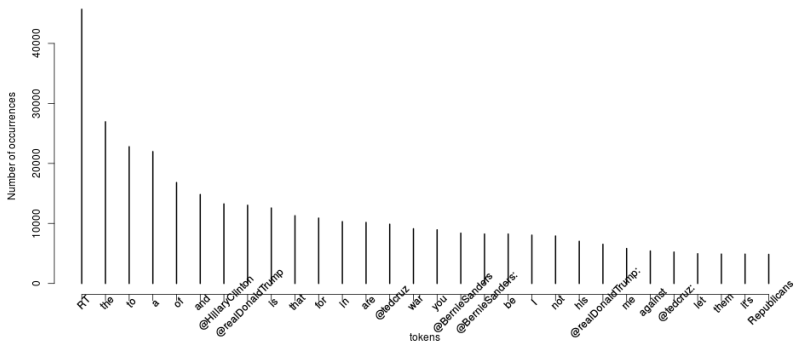




And implementation

Same image.

Top 30 tokens. Representing 348,846 of 1,207,440 (or %29) of all.



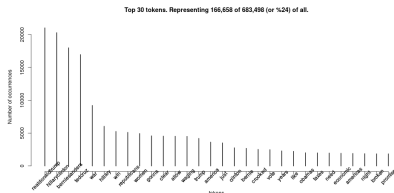


And implementation

What will happen to the tweets when we remove “stop words”

Stop words are words/tokens that have no use in whatever we are doing. Stop words are domain specific.

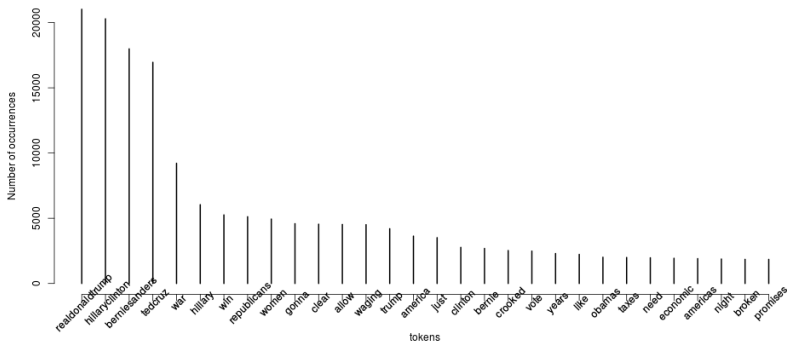
And we change case, remove non-ASCII, remove punctuation, etc.



And implementation

Same image.

Top 30 tokens. Representing 166,658 of 683,498 (or %24) of all.



Using Rstudio

- 1 Set the session working directory
- 2 Load the following files into the editor:
 - 1 `checkPostgres.R`
 - 2 `backEnd.R`
 - 3 `config.txt`
 - 4 `frontEnd.R`
 - 5 `library.R`
- 3 Run the `checkPostgres.R` script to ensure PostGres is configured correctly
- 4 Run the `backEnd.R` script to populate the database with “canned” data
- 5 Run the `frontEnd.R` script to analyze the tweets
- 6 Modify the `config.txt` file, and re-run back and front ends

Q & A time.

Q: How did you get into artificial intelligence?

A: Seemed logical – I didn't have any real intelligence.



What have we covered?

- A preview of today's activities
- An overview of the sentiment analysis



Next: Hands on analysing tweets with R.

References (1 of 1)

- [1] Bing Liu, [Sentiment analysis and opinion mining](#), 2012.
- [2] Bo Pang and Lillian Lee, [Opinion mining and sentiment analysis](#), 2008.
- [3] Wikipedia Staff, [Sentiment analysis](#),
https://en.wikipedia.org/wiki/Sentiment_analysis,
2016.

Files of interest



Software design

document

