# Big Data: Data Wrangling Boot Camp
## Python Sentiment Analysis

# Chuck Cartledge, PhD

# 17 September 2016

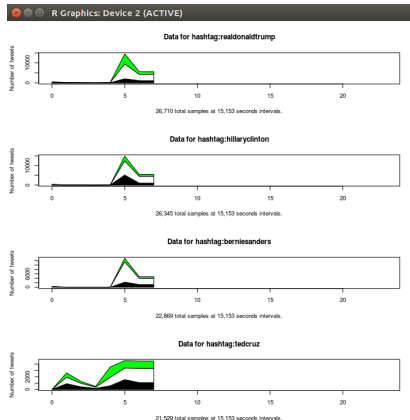## Table of contents I

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 Python

Introduction   **Preview**   Sentiment analysis   Languages   System req.   Q & A   Conclusion   References   Files
○●○         ○○○
            ○○○○○○○

Things that will be happening today

Things that we will be doing.

1. Data wrangle tweets using python
2. Conduct sentiment analysis on tweets
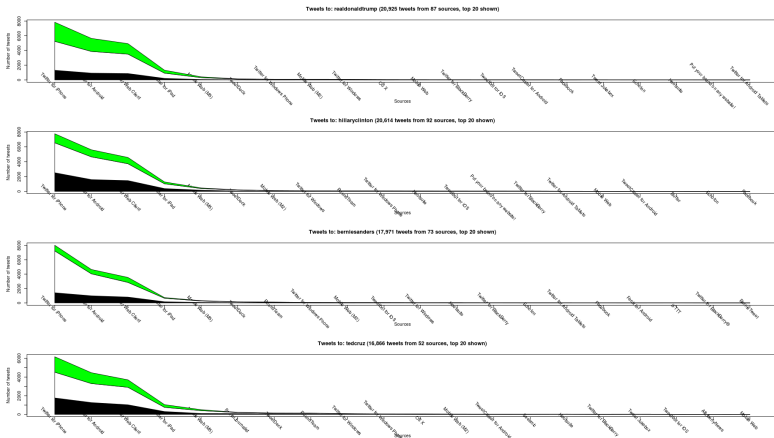3. Look at the sentiments in different ways



Sentiments over time.

Introduction     **Preview**     Sentiment analysis     Languages     System req.     Q & A     Conclusion     References     Files
○●○              ○○○                                    ○○○○○○○○○
                 ○○○○○○○

Things that will be happening today

# Sentiment by sending device

Introduction  **Preview**  Sentiment analysis  Languages  System req.  Q & A  Conclusion  References  Files
○○●            ○○○                             ○○○○○○○○○
              ○○○○○○○

Things that will be happening today

# Sentiment by geographic location

Introduction    Preview    **Sentiment analysis**    Languages    System req.    Q & A    Conclusion    References    Files

What is it, and why should I care?

## 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]

Introduction | Preview | **Sentiment analysis** | Languages | System req. | Q & A | Conclusion | References | Files
000 | 0●0 | | 000000000
0000000

What is it, and why should I care?

## 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]*

Introduction    Preview    **Sentiment analysis**    Languages    System req.    Q & A    Conclusion    References    Files
           ○○○        ○○●
                      ○○○○○○○

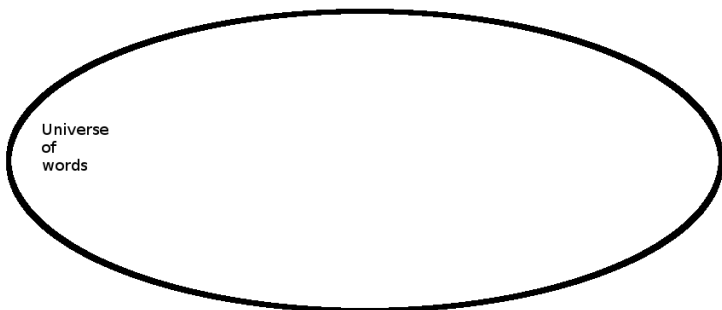What is it, and why should I care?

## 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
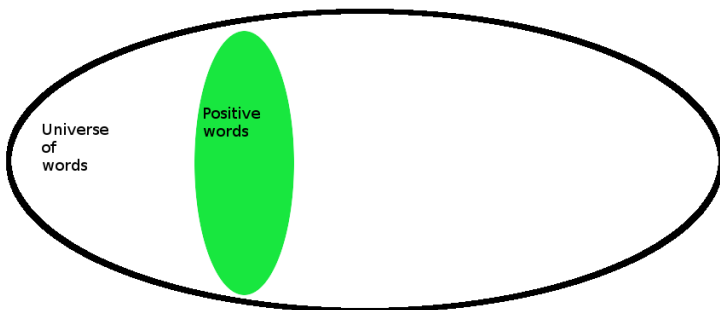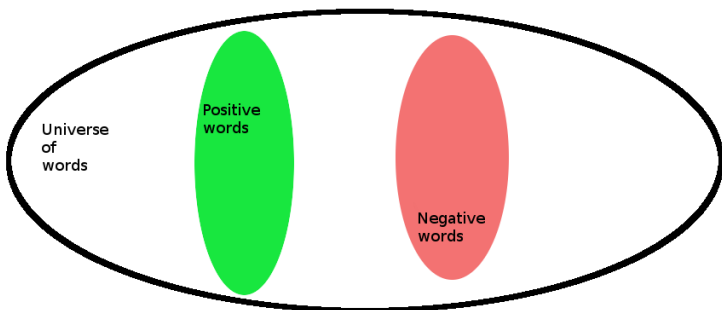7. Display the results

Our approach is language agnostic.

Introduction   Preview   Sentiment analysis   Languages   System req.   Q & A   Conclusion   References   Files
          000          000                               000000000
                               ●000000

A visualization

# A "Universe" of words

Introduction   Preview   Sentiment analysis   Languages   System req.   Q & A   Conclusion   References   Files
     ooo        ooo          ooo                          ooooooooo
                          o●oooooo

A visualization

# Some words are "Positive"



Universe
of
words

Positive
words

Introduction   Preview   **Sentiment analysis**   Languages   System req.   Q & A   Conclusion   References   Files
○○○           ○○○                                          ○○○○○○○○○
              ○○○●○○○○

A visualization

# Some words are "Negative"

Introduction   Preview   Sentiment analysis   Languages   System req.   Q & A   Conclusion   References   Files
   ○○○        ○○○                                ○○○○○○○○○
              ○○○○●○○○

A visualization

# A tweet will/may have positive and negative words

# Some words we don't care about

Introduction    Preview    **Sentiment analysis**    Languages    System req.    Q & A    Conclusion    References    Files
ooo           ooo
           ooooo●o

A visualization

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 token that are positive

4. Compute the percentage of remaining tweet tokens that are negative

5. Classify the tokens as positive, negative, or neutral

Introduction    Preview    **Sentiment analysis**    Languages    System req.    Q & A    Conclusion    References    Files
                 000       000
                           0000000●

A visualization

## Mathematically this is what we are doing

The steps are:

$$
\begin{aligned}
tokens &= \{\text{words in tweet}\} \\
tokensLessStop &= tokens - stopWords \\
positivePart &= positiveWords \cap tokensLessStop \\
negativePart &= negativeWords \cap tokensLessStop
\end{aligned}
$$

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

## Comparing a known language with Python and R



A specific language was not a
requirement for this boot camp.
Python and R used in boot-camp.
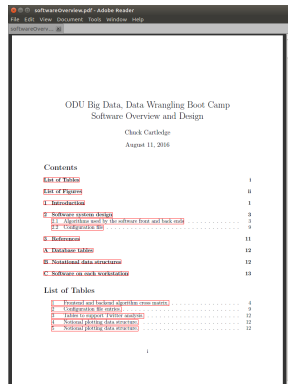Too many languages to compare
to Python and R.

http://hyperpolyglot.org/

Use Hyperployglot.org as a cross reference.

Introduction    Preview    Sentiment analysis    Languages    **System req.**    Q & A    Conclusion    References    Files
 000       000       0000000                         ●00000000

And implementation

## 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



The file is attached.

Introduction   Preview   Sentiment analysis   Languages   **System req.**   Q & A   Conclusion   References   Files
         000      000                          0●00000000
                  0000000

And implementation

What will happen to the tweets when we do nothing.

A histogram of how often a
token occurs.



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

## Same image.



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

Introduction   Preview   Sentiment analysis   Languages   **System req.**   Q & A   Conclusion   References   Files
       ○○○        ○○○                          ○○○●○○○○○
                  ○○○○○○○

And implementation

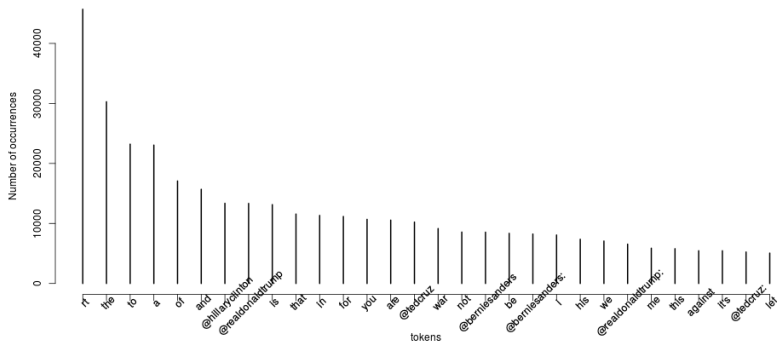## What will happen to the tweets when we make everything the same case
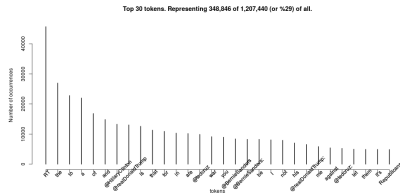
Changing case is easy, unless
they are emojois.



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

## Same image.



Top 30 tokens. Representing 365,040 of 1,207,440 (or %30) of all.
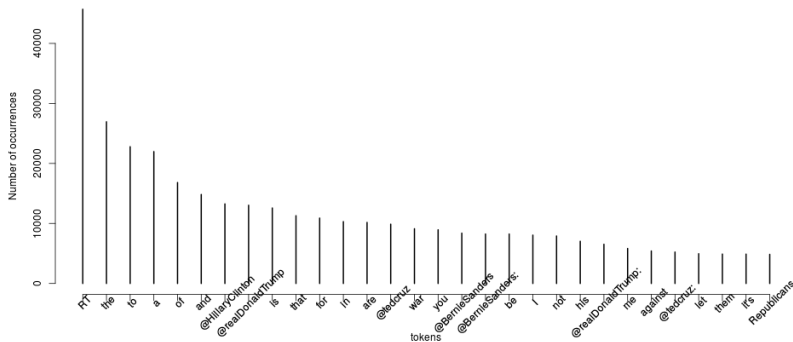
## What will happen to the tweets we remove non-ASCII

We need to talk about ASCII and the rest of the character systems.



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

Introduction   Preview   Sentiment analysis   Languages   System req.   Q & A   Conclusion   References   Files
   ooo       ooo                                ooooooo●oo
       oooooo

And implementation

## Same image.



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

Introduction    Preview    Sentiment analysis    Languages    **System req.**    Q & A    Conclusion    References    Files
           000        000                          0000000        ○○○○○○○●○

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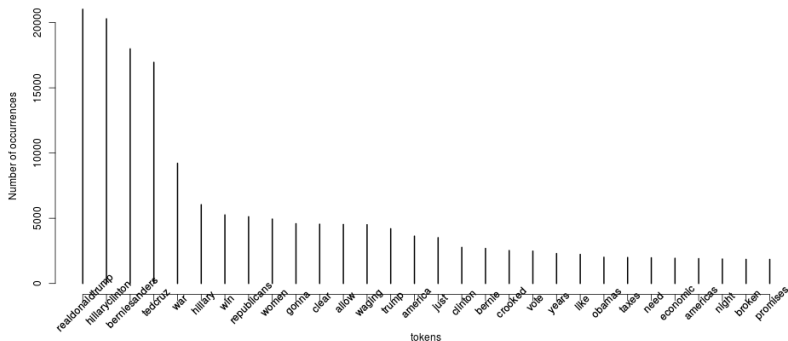And we change case, remove non-ASCII, remove punctuation, etc.



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

And implementation

## Same image.



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

Q & A time.

"'The Answer to the Great
Question . . . Of Life, the Universe
and Everything . . . is
. . . forty-two,' said Deep
Thought, with infinite majesty
and calm."
**Douglas Adams, The
Hitchhiker's Guide to the
Galaxy**

What have we covered?

- A preview of today's activities
- An overview of the sublime IDE



Next: Hands on analysing tweets with Python.

## References I

[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



1. softwareOverview.pdf