

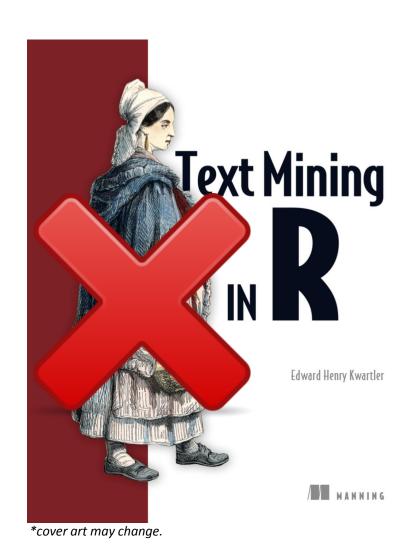
Intro to text mining using tm, openNLP, & topicmodels

www.linkedin.com/in/edwardkwartler

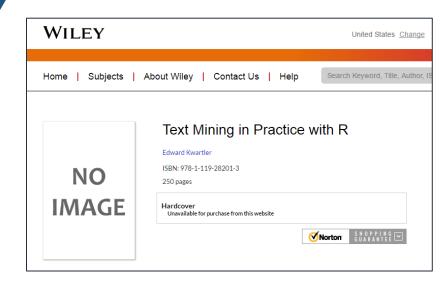


Shameless Plug

Slated for June 2016



Shameless Plug #1



In final editing now!

Shameless Plug #2



Learn Data Science By Doing

Get started for free at DataCamp.com
(You'll be notified of my new course Intro to Text Mining: Bag of Words)

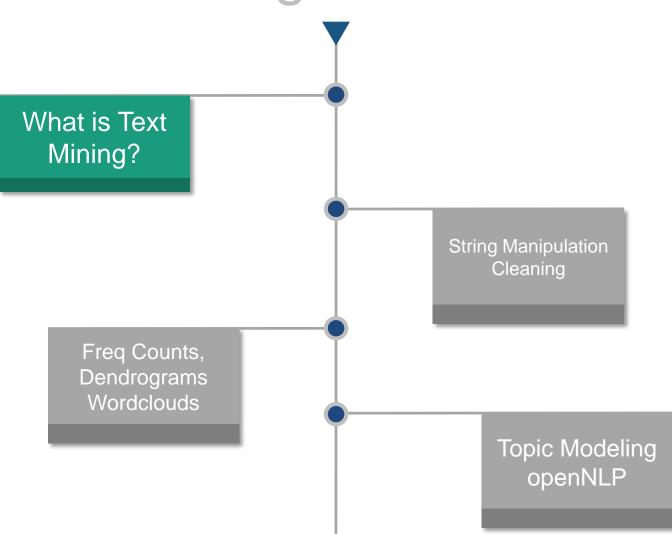
Looking for R or Python training for your team?

Claim your free group account at DataCamp.com/groups

Agenda

What is Text Mining? String Manipulation Cleaning Freq Counts, Dendrograms Wordclouds **Topic Modeling** openNLP

Agenda



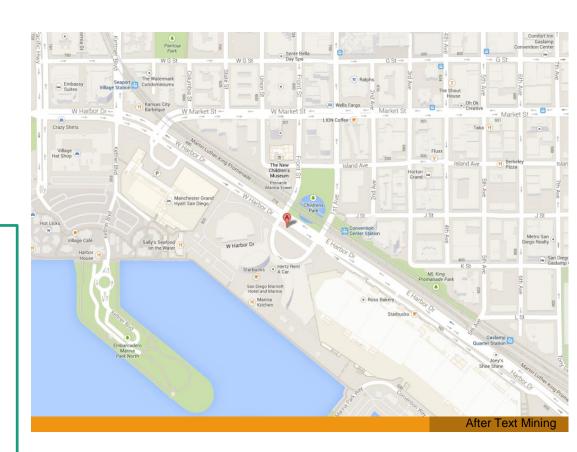
What is text mining?

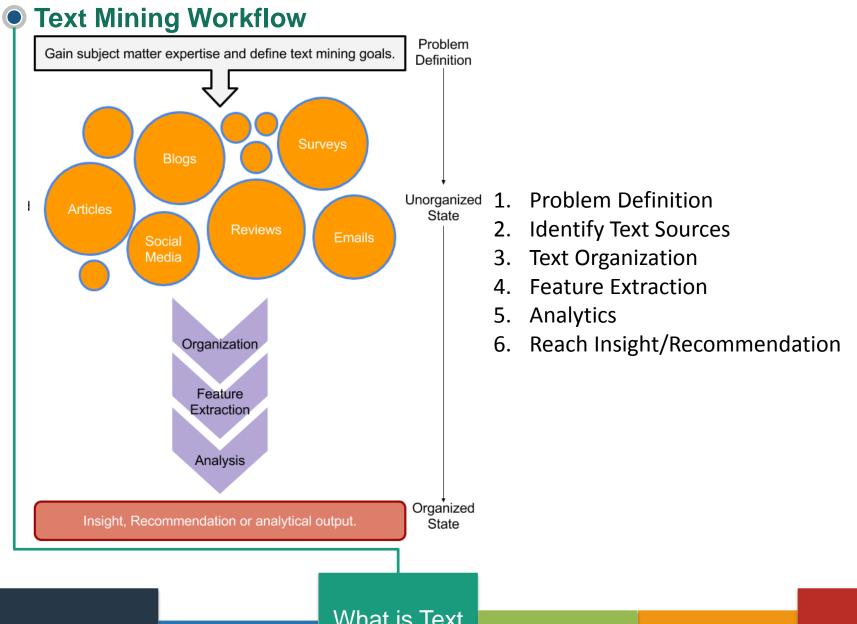
- Extract new insights from text
- Let's you drink from a fire hose of information
- Language is hard; many unsolved problems
 - Unstructured
 - Expression is individualistic
 - Multi-language/cultural implications



What is text mining?

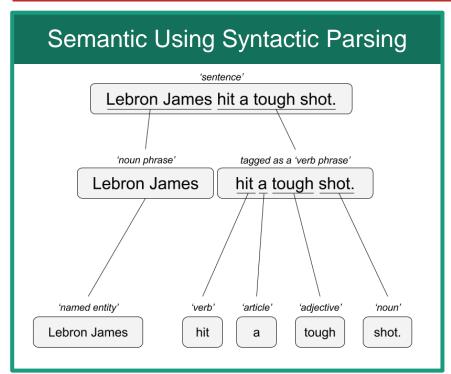
- Extract new insights from text
- Let's you drink from a fire hose of information
- Language is hard; many unsolved problems
 - Unstructured
 - Expression is individualistic
 - Multi-language/cultural implications





Text Mining Approaches

"Lebron James hit a tough shot."





Text Mining Approaches

Some Challenges in Text Mining

- Compound words (tokenization) changes meaning
 - "not bad" versus "bad"
- Disambiguation
- Sarcasm
 - "I like it...NOT!"
- Cultural differences
 - "It's wicked good" (in Boston)

"I made her duck."

- I cooked waterfowl to eat.
- I cooked waterfowl belonging to her.
- I created the (clay?) duck and gave it to her.
- Duck!!

Text Sources

Text can be captured within the enterprise and elsewhere

- Books
- Electronic Docs (PDFs)
- Blogs
- Websites
- Social Media
- Customer Records
- Customer Service Notes
- Notes
- Emails
- Legal Documents

• . . .

The source and context of the medium is important. It will have a lot of impact on difficulty and data integrity.



Enough of me talking...let's do it for real! Scripts in this workshop follow a simple workflow Set the Working Directory **Load Libraries** Make Custom Functions & **Specify Options** Read in Data & Pre-Process Perform Analysis & Save What is Text Mining?

Enough of me talking...let's do it for real!

Setup

Install R/R Studio

- http://cran.us.r-project.org/
- http://www.rstudio.com/products/rstudio/download/

Workshop scripts, corpora (prob best to download at the end) https://goo.gl/2WWkst

Install Packages

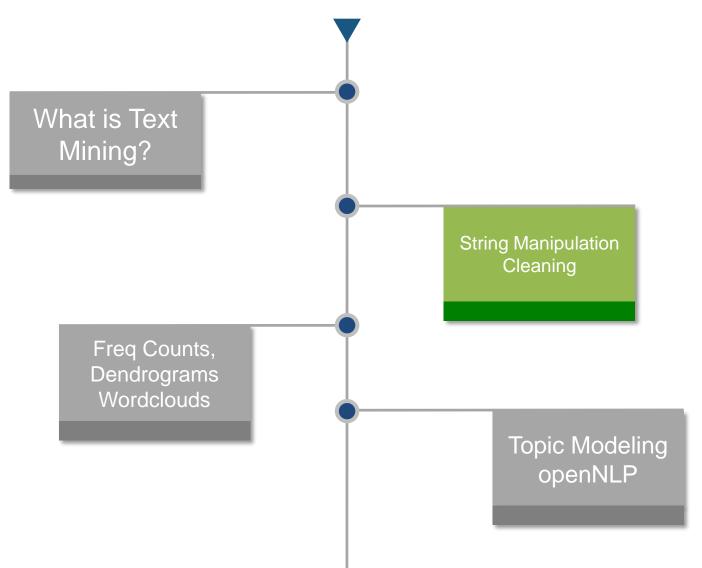
- Run "1_Install_Packages.R" script
 - An error may occur if the Java version doesn't match and depending on the OS. In that case install packages individually.

Warning: Twitter Profanity

- Twitter demographics skew young and as a result have profanity that appear in the examples.
- It's the easiest place to get a lot of messy text fast, if it is offensive feel free to talk to me and I will work to get you other texts for use on your own. No offense is intended.



Keyword Scanning & Text Cleaning



Open the "coffee.csv" to get familiar with the data structure

1000 tweets mentioning "coffee" created truncated replyToSID id replyToUID statusSource screenName retweetCour retweeted 1 @ayyytylerb that is so true drink lots of coffee FALSE 8/9/13 2:43 FALSE FALSE NA 2 RT @bryzy_brib: Senior March tmw morning at 7:25 A.M. in the SENIOR lot. Get up early, make yo coffee/breakfast, cus this will only happen .. <a href="http://carolynicosia FALSE 8/9/13 2:43 FALSE 3.6566E+17 NA FALSE 3 If you believe in #gunsense tomorrow would be a very good day to have your coffee any place BUT @Starbucks Guns+Coffee=#nosense @MomsDemand 8/9/13 2:43 3.6566E+17 NA janeCkay FALSE 4 My cute coffee mug. http://t.co/2udvMU6XIG FALSE 8/9/13 2:43 FALSE 3.6566E+17 NA <a href="htt; AlexandriaO(FALSE NA 5 RT @slaredo21: I wish we had Starbucks here... Cause coffee dates in the morning sound perff FALSE 8/9/13 2:43 FALSE 3.6566E+17 NA <a href="httr Rooosssaaaa FALSE NA NA 6 Does anyone ever get a cup of coffee before a cocktail?? FALSE 8/9/13 2:43 FALSE 3.6566E+17 NA <a href="httr E Z MAC FALSE NΔ NΔ 7 "I like my coffee like I like my women...black, bitter, and preferably fair trade." I love #Archer FALSE 8/9/13 2:43 FALSE 3.6566E+17 NA <a href="http://charlie_3119 FALSE 8 @dreamwwediva ya didn't have coffee did ya? FALSE 8/9/13 2:43 FALSE 3.6566E+17 3.6566E+17 1316942208 <a href="http://descicaSalvat FALSE NA NA 9 RT @iDougherty42: I just want some coffee. 8/9/13 2:43 FALSE 3.6566E+17 NA <a href="http kaytiekirk FALSE 10 RT @Dorkv76: I can't care before coffee. 8/9/13 2:43 3.6566E+17 NA <a href="httplissteria FALSE 11 No lie I wouldn't mind coming home smelling like coffee 8/9/13 2:43 FALSE 3.6566E+17 NA <a href="htt; DOPECROOK FALSE NA FALSE FALSE 3.6566E+17 NA <a href="httrTiffCaruso FALSE NA 12 RT @JonasWorldFeed: Play Ping Pong with Joe, Take a tour of the stage with Nick, Have coffee with Kevin, Charity auction; https://t.co/VTkK. 8/9/13 2:43 NA FALSE FALSE FALSE 13 Have I ever told any of you that Tate Donovan bought my stepmom coffee? 8/9/13 2:43 3.6566E+17 NA web CurlysCrazyN NA NA 14 RT @JonasWorldFeed: Play Ping Pong with Joe. Take a tour of the stage with Nick. Have coffee with Kevin. Charity auction: https://t.co/VTkK... FALSE 8/9/13 2:43 FALSE 3.6566E+17 NA JoeJonasVA FALSE 15 @HeatherWhaley I was about 2 joke it takes 2 hands to hold hot coffee...then I read headline! #Don'tDrinkNShoot FALSE HeatherWha 8/9/13 2:42 FALSE 3.6565E+17 3.6566E+17 26035764 <a href="http://doi.org/10.1001/j.j.gov/2015-10.1001/j.gov/2015-10.1001/j.j.gov/2015-10.1001/j.j.gov/2015-10.1001/j.gov/20 FALSE NA NA 16 RT @MoveTheSticks: Charlie Whitehurst looks like he should be working at a coffee shop in Portland or hosting a renovation show on HGTV. 8/9/13 2:42 FALSE 3.6566E+17 NA <a href="httpmpr4437 FALSE 8/9/13 2:42 FALSE 3.6566E+17 NA sharkshukri FALSE 17 Coffee always makes everything better. web 18 RT @AdelaideReview: Food For Thought: @Annabelleats shares a delicious Venison and Porcini Mushroom Pie Recipe, http://t.co/N807vgFKWN http:// FALSE FALSE 3.6566E+17 NA FALSE NA 8/9/13 2:42 <a href="httpthepaulbake FALSE 3.6566E+17 NA FALSE NA 19 RT @LittleMelss: Imfao!!!"@bryanlaca; nahhh Melanie u is fa sho like an ummm a Coffee table :)) yeeeee Imaoo FALSE 8/9/13 2:42 NA web brvanlaca NA 20 I wonder if Christian Colon will get a cup of coffee once the rosters expand to 40 man in September. Really nothing to lose by doing so FALSE 8/9/13 2:42 FALSE 3.6566E+17 NA <a href="http Shauncore FALSE NA NA

"text\$text" is the vector of tweets that we are interested in.

All other attributes are automatically returned from the twitter API

2_Keyword_Scanning.R

Basic R Unix Commands

grepl returns a vector of T/F if the pattern is present at least once

grepl("pattern", searchable object, ignore.case=TRUE)

grep returns the position of the pattern in the document

grep("pattern", searchable object, ignore.case=TRUE)

[1] 4 214 276 366 479 534 549 620

"library(stringi)" Functions

stri_count counts the number of patterns in a document

stri_count(searchable object, fixed="pattern")

2_String Manipulation.R

Remember This? Problem Gain subject matter expertise and define text mining goals. Definition Unorganized State Organization Feature Extraction Analysis Organized Insight, Recommendation or analytical output. State

R for our Cleaning Steps

Tomorrow I'm going to have a nice glass of Chardonnay and wind down with a good book in the corner of the county :-)



- 1.Remove Punctuation
- 2.Remove extra white space
- 3. Remove Numbers
- 4.Make Lower Case
- 5.Remove "stop" words
- tomorrow going nice glass chardonnay wind down good book corner county

• 3

3_Cleaning and Frequency Count.R

"library(tm)" Functions

VCorpus creates a corpus held in memory.

VCorpus(source)

tm_map applies the transformations for the cleaning

tm_map(corpus, function)

getTransformations() will list all standard tm corpus transformations
We can apply standard R ones too. Sometimes it makes sense to perform all of these or a subset or

even other transformations not listed like "stemming"

tm_map(corpus, removePunctuation) - removes the punctuation from the documents tm_map(corpus, stripWhitespace) - extra spaces, tabs are removed tm_map(corpus, removeNumbers) - removes numbers tm_map(corpus, tolower) - makes all case lower tm_map(corpus, removeWords) - removes specific "stopwords"

New Text Mining Concepts

Corpus- A collection of documents that analysis will be based on.

Stopwords – are common words that provide very little insight, often articles like "a", "the".

Customizing them is sometimes key in order to extract valuable insights.

3_Cleaning and Frequency Counts.R

3_Cleaning and Frequency Count.R

"tryTolower"is poached to account for errors when making lowercase.

```
tryTolower <- function(x){
  # return NA when there is an error
  y = NA
  # tryCatch error
  try_error = tryCatch(tolower(x), error = function(e) e)
  # if not an error
  if (!inherits(try_error, 'error'))
  y = tolower(x)
  return(y)}</pre>
```

"clean.corpus" makes applying all transformations easier.

```
clean.corpus<-function(corpus){
  corpus <- tm_map(corpus, removePunctuation)
  corpus <- tm_map(corpus, stripWhitespace)
  corpus <- tm_map(corpus, removeNumbers)
  corpus <- tm_map(corpus, content_transformer(str_to_lower))
  corpus <- tm_map(corpus,
  content_transformer(replace_contraction))
  corpus <- tm_map(corpus, removeWords, custom.stopwords)
  return(corpus)}</pre>
```

Base: tolower (basic)

Stringr: str_to_lower (wrapper)

Custom: tryTolower (handles errors)

3_Cleaning and Frequency Count.R

"custom.stopwords" combines vectors of words to remove from the corpus

#Create custom stop words custom.stopwords < c(stopwords('english'), 'lol', 'smh')

"custom.reader" keeps the meta data (tweet ID) with the original document

#bring in some text
text<-read.csv('coffee.csv', header=TRUE)

#Keep the meta data, apply the functions to make a clean corpus
custom.reader > readTabular(mapping=list(content="text", id="id"))
corpus <- VCorpus(DataframeSource(text), readerControl=list(reader=custom.reader)))
corpus<-clean.corpus(corpus)

3_Cleaning and Frequency Count.R

Bag of Words means creating a Term Document Matrix or Document Term Matrix*

Term Document Matrix

	Tweet1	Tweet 2	Tweet3	Tweet4		Tweet_n
Term1	0	0	0	0	0	0
Term2	1	1	0	0	0	0
Term3	1	0	0	2	0	0
	0	0	3	0	1	1
Term_n	0	0	0	1	1	0

Document Term Matrix

	Term1	Term2	Term3		Term_n
Tweet1	0	1	1	0	0
Tweet2	0	1	0	0	0
Tweet3	0	0	0	3	0
	0	0	0	1	1
Tweet_n	0	0	0	1	0

"as.matrix" makes the tm's version of a matrix into a simpler version

dtm<-DocumentTermMatrix(corpus)</pre>

tdm<-TermDocumentMatrix(corpus)

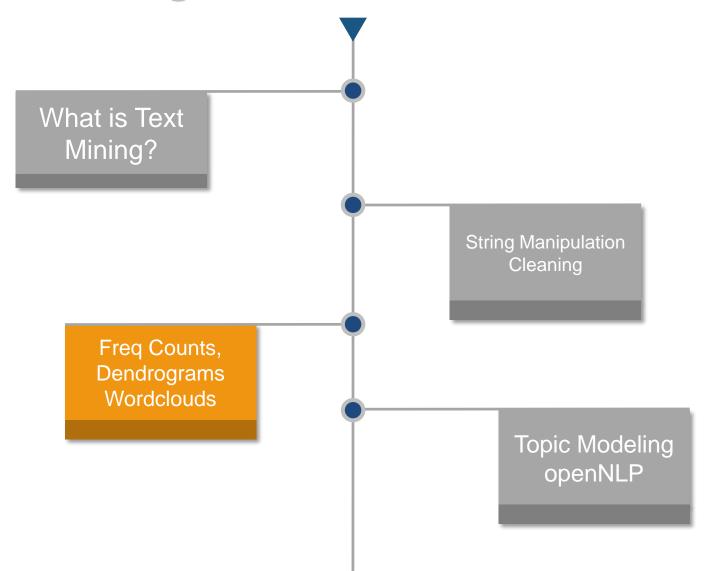
dtm.tweets.m<-as.matrix(dtm)</pre>

tdm.tweets.m<-as.matrix(tdm)

These matrices are often very sparse and large therefore some special steps may be needed and will be covered in subsequent scripts.

*Depends on analysis, both are transpositions of the other

Dendrograms & Word Clouds





4_dendrogram.R script builds on the matrices

First let's explore simple frequencies

#Summed Vector

tdm.m <- as.matrix(tdm)

tdm.v <- sort(rowSums(tdm.m),decreasing=TRUE)

tdm.df <- data.frame(word = names(tdm.v),freq=tdm.v, row.names=NULL)

Term Document Matrix

	Tweet1	Tweet 2	Tweet3	Tweet4		Tweet_n
Term1	0	0	0	0	0	0
Term2	1	1	0	0	0	0
Term3	1	0	0	2	0	0
	0	0	3	0	1	1
Term_n	0	0	0	1	1	0



word	freq
Term1	0
Term2	2
Term3	3
	5
Term_n	2

4_dendrograms.R

4_dendrogram.R script

ggplot2 ggthemes

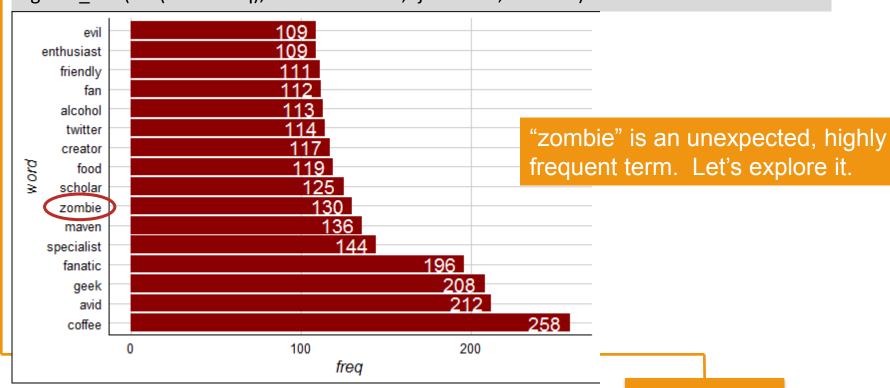
#Make a barplot of the top terms

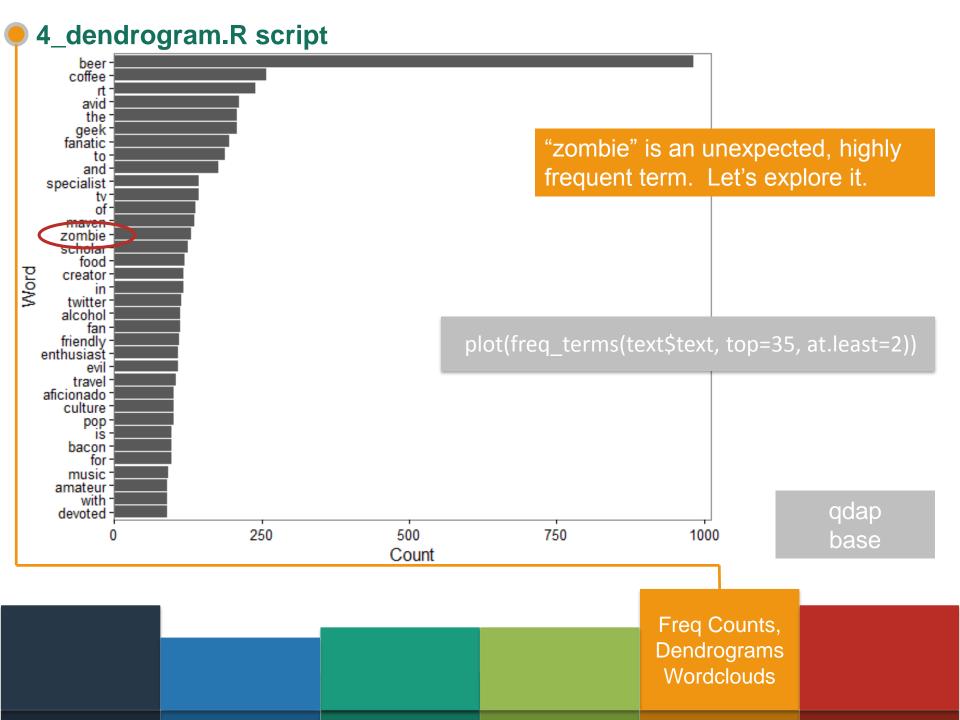
top.words<-tdm.df[which(tdm.df\$freq>=105)]

top.words\$word<-factor(top.words\$word, levels=unique(as.character(top.words\$word)))

ggplot(top.words, aes(x=word, y=freq))+geom_bar(stat="identity", fill='darkred')
+coord_flip()+theme_gdocs()+

geom_text(aes(label=freq), colour="white",hjust=1.25, size=5.0)



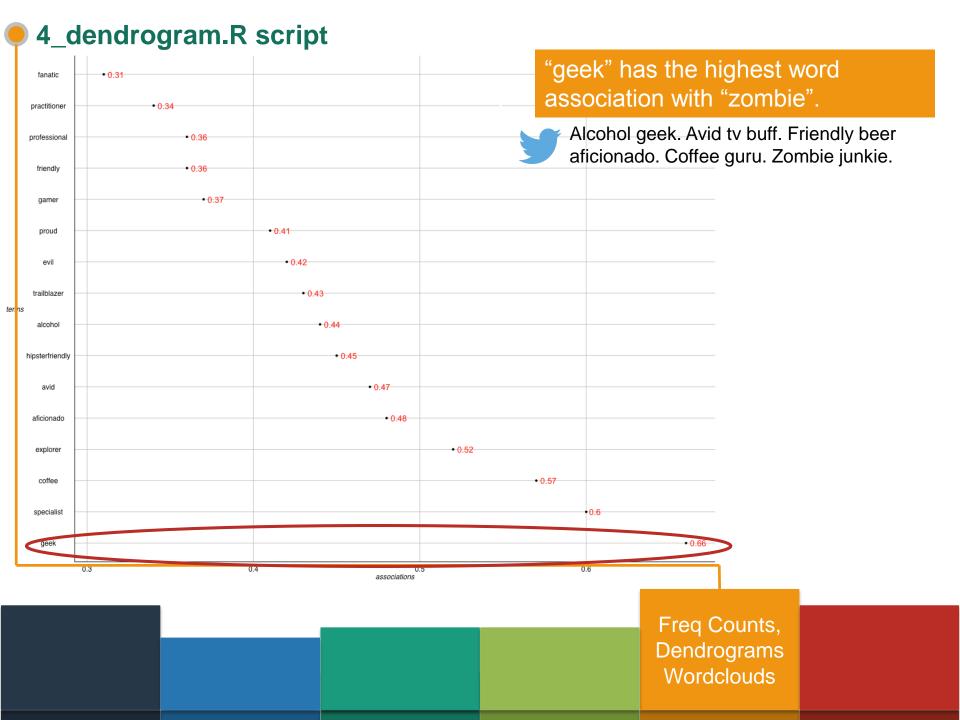


4_dendrogram.R script

Next let's explore word associations, similar to correlation

```
associations<-findAssocs(tdm, 'zombie', 0.30)
a.df<-do.call(cbind,associations)
a.df<-data.frame(terms=row.names(a.df),a.df, row.names=NULL)
a.df$terms<-factor(a.df$terms, levels=a.df$terms)
ggplot(a.df, aes(y=terms)) + geom_point(aes(x=zombie), data=a.df)+
theme_gdocs()+geom_text(aes(x=zombie,label=zombie),
colour="red",hjust=-.25)
```

- Adjust 0.30 to get the terms that are associated .30 or more with the 'zombie' term.
- Treating the terms as factors lets ggplot2 sort them for a cleaner look.



Extracting Meaning using dendrograms

Dendrograms visualize hierarchical clusters based on frequencies.

- Reduces information much like average is a reduction of many observations' values
- Word clusters emerge often showing related terms

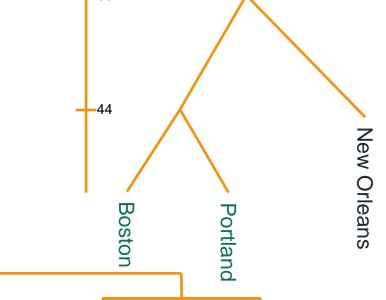
• Term frequency is used to construct the word cluster. Put another way, term A and term B have similar frequencies in the matrix so they are considered a

cluster.

City	Annual Rainfall
Portland	43.5
Boston	43.8
New Orleans	62.7

Boston & Portland are a cluster at height 44. You lose some of the exact rainfall amount

in order to cluster them.





Weird associations! Maybe a dendrogram will help us more

```
#Hierarchical Clustering
tdm2 <- removeSparseTerms(tdm_sparse=0.95) #s oot for ~40 terms
tdm2.df<-as.data.frame(inspect(tdm2))
hc <- hclust(dist(tdm2.df))
hcd <- as.dendrogram(hc)
clusMember <- cutree(hc, 4)
labelColors <- c("#CDB380", "#036564", "#EB6841", "#EDC951")
clusDendro <- dendrapply(hcd, colLab)
plot(clusDendro, main = "Hierarchical Dendrogram", type = "triangle")</pre>
```

% of zeros allowed e.g. higher means more words in TDM/DTM

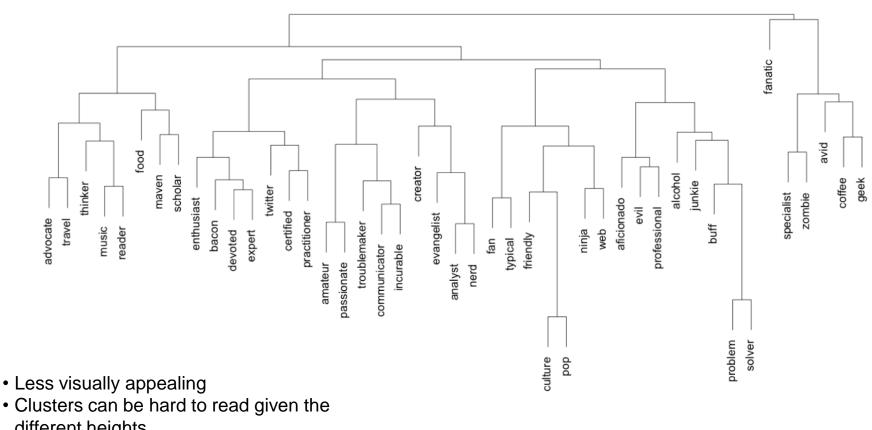
New Text Mining Concept

<u>Sparse</u>- Term Document Matrices are often extremely sparse. This means that any document (column) has mostly zero's. Reducing the dimensions of these matrices is possible by specifying a sparse cutoff parameter. Higher sparse parameter will bring in more terms.

4_dendrogram.R script

Base Plot of a Dendrogram

Cluster Dendrogram



different heights

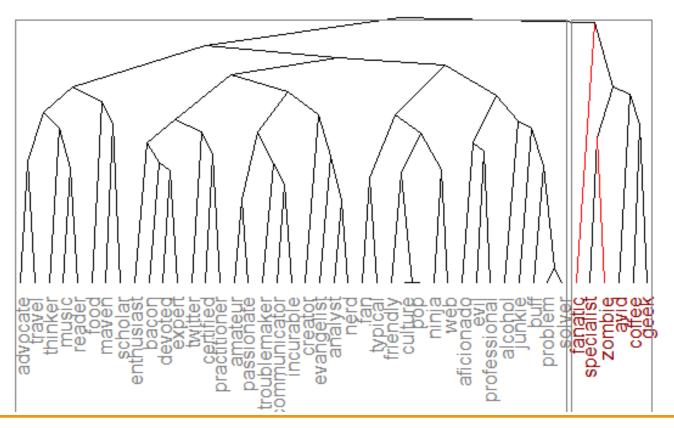
4_dendrogram.R script

Dendextend offers more flexiblity

Better Dendrogram

hcd <- as.dendrogram(hc)
hcd <- branches_attr_by_labels(hcd, c('zombie', 'fanatic'), 'red')
hcd <- color_labels(hcd,2, col = c('grey50', 'darkred'))
plot(hcd, main = "Better Dendrogram", type='triangle',yaxt='n')
rect.dendrogram(hcd, k = 2, border = "grey50")</pre>

 Aesthetically my choice is to have colored clusters and all terms at the bottom.



5_Simple_Wordcloud.R

5_Simple_Wordcloud.R script
Using Rweka Package we create a custom function

(tokenizer) - function(x) NGramTokenizer(x, Weka control(min = 2, max = 2))

It is used as a control parameter when constructing a TDM/DTM

bigram tdm <- TermDocumentMatrix(corpus, control = list(tokenize €tokenizer))

Text Mining is so fun. So do Text Mining!

Unigram

DOCS Terms fun. mining 2 text

Bigram

	Docs
Terms	1
do text	1
fun so	1
is so	1
mining is	1
so do	1
so fun	1
text mining	2

^{*}with common stopwords

New Text Mining Concept

Tokenization- So far we have created single word n-grams. We can create multi word "tokens" like bigrams, or trigrams with this line function. It is applied when making the term document matrix.

5_Simple_Wordcloud.R script

To make a wordcloud we follow the previous steps and create a data frame with the word and the frequency.

#Summed Vector

tdm.m <- as.matrix(bigram tdm)

tdm.v <- sort(rowSums(tdm.m),decreasing=TRUE)

tdm.df <- data.frame(word =

names(tdm.v),freq=tdm.v)

Term Document Matrix

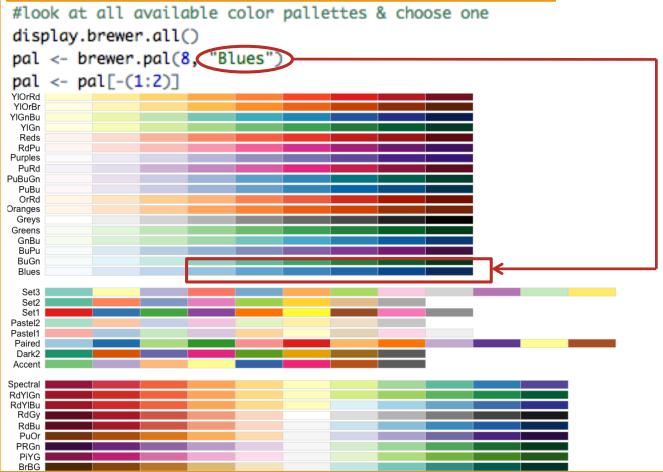
	Tweet1	Tweet 2	Tweet3	Tweet4		Tweet_n
Term1	0	0	0	0	0	0
Term2	1	1	0	0	0	0
Term3	1	0	0	2	0	0
	0	0	3	0	1	1
Term_n	0	0	0	1	1	0



word	freq
Term1	0
Term2	2
Term3	3
•••	5
Term_n	2

5_Simple_Wordcloud.R script

Next we need to select the colors for the wordcloud.



5_Simple_Wordcloud.R script

set.seed(2016) wordcloud(tdm.df\$word,tdm.df\$freq,max.words=50, random.order=FALSE, colors=pal)

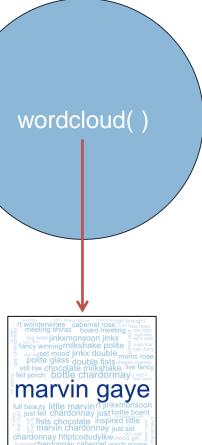
rose https...
rt januaryjames meeting shiraz cabernet
rt januaryjames meeting shiraz love chicken
winning bottle of jinkxmonsoon donjon love
cabernet rose milkshake polite fouryines naked
rose bushes double fists just set
polite glashtle marvilancy winning

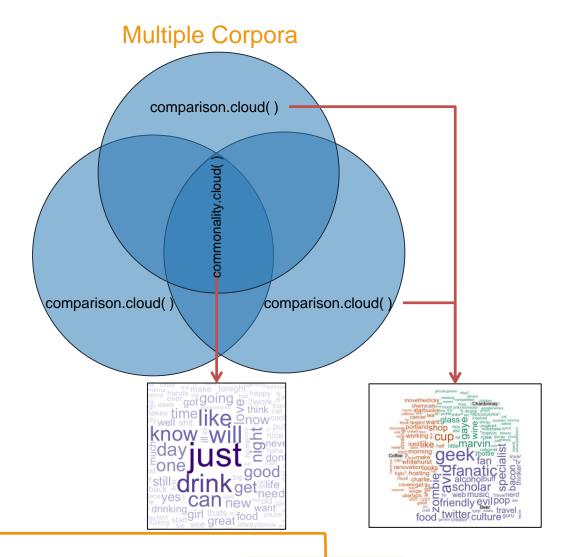
marvin gaye

competition sphocolate milks hakemons rose live fancy jinkx doublenspired littles mood giglass httptcodudylkvjust feller beauty gracefell porch pinot noir full beautyorch moms brought marvingrace just jirl brought wonderwines competition cream mushroomarvin gav

- Bigram Tokenization has captured "marvin gaye"
- A word cloud is a frequency visualization. The larger the term (or bigram here) the more frequent the term.
- You may get warnings if certain tokens are to large to be plotted in the graphics device.

Types of Wordclouds Single Corpus





6_Other_Wordclouds.R

6_Other_Wordcloud.R

Bring in more than one corpora.

```
#bring in some text
text1<-read.csv('chardonnay.csv', header=TRUE)
text2<-read.csv('coffee.csv', header=TRUE)
text3<-read.csv('beer.csv', header=TRUE)

text1<-paste(text1$text,collapse=' ')
text2<-paste(text2$text,collapse=' ')
text3<-paste(text3$text,collapse=' ')
chard.corpus<-clean.corpus(VCorpus(VectorSource(text1)))
coff.corpus<-clean.corpus(VCorpus(VectorSource(text2)))
beer.corpus<-clean.corpus(VCorpus(VectorSource(text3)))</pre>
```

Without the clean.corpus function it is a lot more code!

Extract the clean corpus content and collapse into 3 mega documents

```
all.chardonnay<-unlist(sapply(chard.corpus, `[`, "content"))
all.coffee<-unlist(sapply(coff.corpus, `[`, "content"))
all.beer<-unlist(sapply(beer.corpus, `[`, "content"))

all <- c(all.chardonnay, all.coffee, all.beer)
all.corpus <- VCorpus(VectorSource(all))
```

Commonality Cloud

- The tweets mentioning "chardonnay" "beer", and "coffee" have these words in common.
- Again size is related to frequency.
- Not helpful in this but in diverse corpora it may be more helpful e.g. political speeches.



#Common Words

commonality.cloud(tdm, max.words=300, random.order=FALSE,colors=pal)

Comparison Cloud

- The tweets mentioning "chardonnay" "beer", and "coffee" have these dissimilar words.
- Again size is related to frequency.
- Beer drinkers in this snapshot are passionate (fanatics, geeks, specialists) on various subjects while Chardonnay drinkers mention Marvin Gaye. Coffee mentions up & working.

beauty grace half caused inspired little fell porch portland hosting jinkxmonsoon jinkx Chardonnay grace just shop portland glass httptcodudylkw milkshake polite just set like working fists chocolate set mood chocolate milkshake looks like little marvin just fell working shop jinkx double double fists ≅marvın gay pop culture travel geek single cup tv maven dari jix fanatic evil whitehurst looks fan friendly problem solver avid creator infuriatingly humble Beer web f show hgtv passionate fanatic lover reader

#Comparison Cloud set.seed(2016)

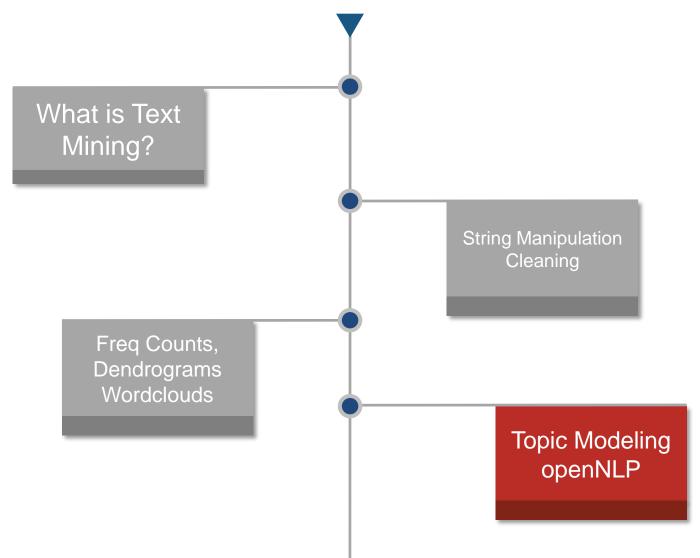
comparison.cloud(tdm, max.words=200, random.order=FALSE,title.size=1.0,colors=brewer.pal(ncol(tdm),"Dark2"))

Both Clouds look very different with different tokenization

```
#Functions
tokenizer <- function(x) NGramTokenizer(x, Weka_control(min = 1, max = 1))</pre>
```

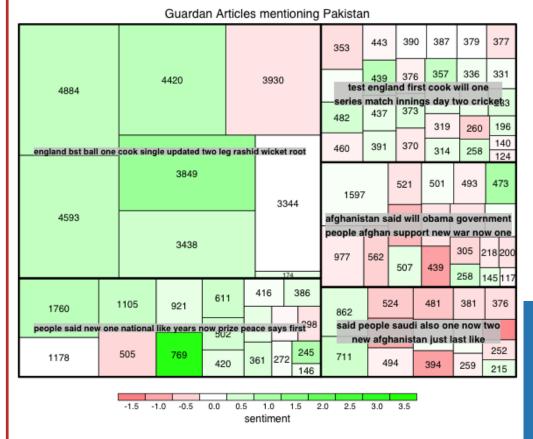
```
see beauty sippchardonnay
                                                                      today winning binot donjon porch
                                                                                                                                                                                                                                                       januaryjames
                                    movethesticks E grace competition live
        caramel prenovation pwcwines naked bushes wonderwines
ever think charles inspired shirazyall board chipotletwins jinkxmonsoonthats moms jinkx httptcodudylkw nice cause table whitehurst cause table white cabernet just milkshake fancy good starbucks show gaye full bottle white chicken the cause of the cause
                                                                                                                                                                                                                         shirazyall board chipotletwins
      chemicals 🧨
                                                                                                                                                                                                                                                 buff dari evil web
                                                                                               alcohol fan total Beer bacon
                                                                                    travel solver writer lifelong Culture never
```

Topic Modeling & openNLP



7_Topic_Modeling_Sentiment.R

TREEMAP: multi dimensional representation of the corpus attributes.



- Color will represent quick, simple polarity sentiment
- Each article is a small square
- The area of the square is related to the number of terms <u>document</u> length
- The larger grouping is based on LDA or CTM <u>Topic Modeling</u>

End result is understanding broad topics, their sentiment and amount of the corpus documents devoted to the identified topic.

Topic Modeling

LDA

Each document is made up of mini topics.

- Probability is assigned to each document for the specific observed topics.
- A document can have varying probabilities of topics simultaneously.

Technical Explanations:

http://en.wikipedia.org/wiki/Latent_Dirichlet_allocati

on

http://cs.brown.edu/courses/csci2950-p/spring2010/lectures/2010-03-03_santhanam.pdf

*Full disclosure I am not an expert in topic modeling

Example

Corpus

- 1.I like to watch basketball and football on TV.
- 2.I watched basketball and Shark Tank yesterday.
- 3. Open source analytics software is the best.
- 4.R is an open source software for analysis.
- 5.I use R for basketball analytics.

LDA Topics

<u>Topic A:</u> 30% basketball, 20% football, 15% watched, 10% TV...something to do with watching sports

<u>Topic B</u>: 40% software, 10% open, 10% source...10% analytics something to do with open source analytics software

Documents

- 1.100% Topic A
- 2.100% Topic A
- 3.100% Topic B
- 4.100% Topic B
- 5.60% Topic B, 40% Topic A

Simple Sentiment Polarity

Scoring

Surprise is a sentiment.

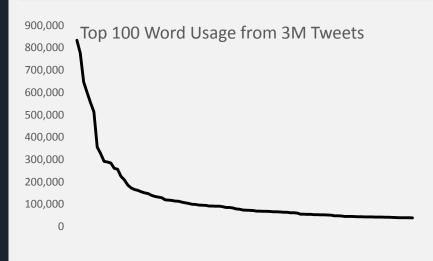
Hit by a bus! – Negative Polarity
Won the lottery!- Positive Polarity

- I loathe BestBuy Service -1
- I <u>love</u> BestBuy Service. They are the <u>best</u>. +2
- I <u>like</u> shopping at BestBuy but <u>hate</u> traffic. 0

R's QDAP polarity function scans for positive words, and negative words as defined by MQPA Academic Lexicon research. It adds positive words and subtracts negative ones along with valence shifters. The final score represents the polarity of the social interaction.

Zipf's Law

Many words in natural language but there is steep decline in everyday usage. Follows a predictable pattern.



Simple Sentiment Polarity

Scoring

```
library(qdap)

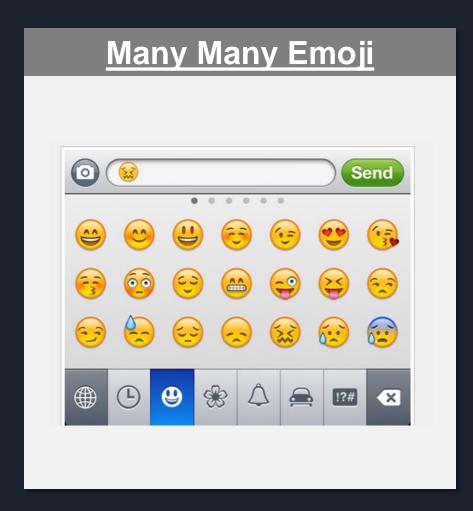
text1<-'i love St Peters University'
text2<-'this lecture is good'
text3<-'this lecture is very good'
text4<-'data science is hard I like it a little'
text5<-'data science is hard'

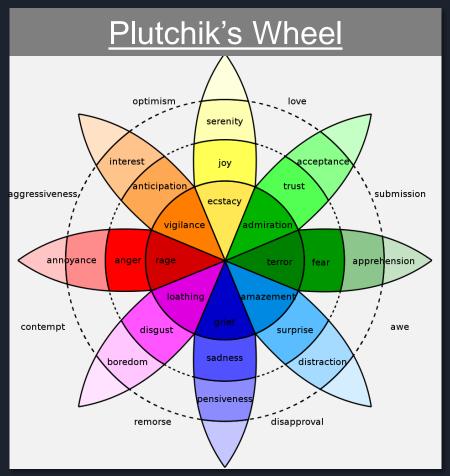
polarity(text1)
polarity(text2)
polarity(text3)
polarity(text4)
polarity(text5)</pre>
```

- <u>Text 1:</u> "love" was identified as positive. The text has 5 words and so 1/sqrt(5) = .447
- <u>Text 2:</u> "good" was identified positively. So 1/sqrt(4)=.5
- <u>Text 3:</u> "good" was found along with the amplifier "very". So (.8+1)/sqrt(5)=.805
- <u>Text 4:</u> hard and like cancel each other out so the polarity is zero. 1-1/sqrt(9)=0
- <u>Text 5:</u> "hard" is -1/sqrt(4)=-.50

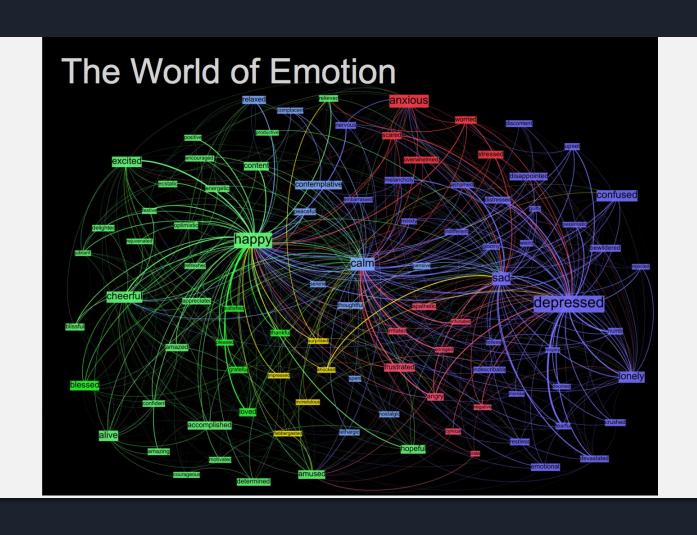
First it looks for the polarized word. Then identifies valence shifters (default 4 words before and two words after) Amplifiers are assigned +.8 and de-amplifiers weight is constrained to -1.

In reality sentiment is more complex.





Kanjoya's Experience Corpus



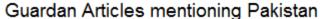


Open the script and let's walk through it line by line because there are multiple additions to the previous scripts

```
1 #Ted Kwartler
2 #Ted@sportsanalytics.org
   #ODSC Workshop: Intro to Text Mining using R
   #11-14-2015
   #v7.3 Topic Modeling, Sentiment and Length Treemap
6
    #Set the working directory
    setwd('/Users/ted/Desktop/ODSC')
   #libraries
11 library(treemap)
12 library(qdap)
13 library(GuardianR)
14 library(topicmodels)
15 library(tm)
    library(SnowballC)
17
18 #options, functions
    options(stringsAsFactors = FALSE) #text strings will not be factors of categories
    Sys.setlocale('LC_ALL','C') #some tweets are in different languages so you may get an error
21
22 - tryTolower <- function(x){
23
      # return NA when there is an error
     y = NA
25
     # tryCatch error
26
     try_error = tryCatch(tolower(x), error = function(e) e)
27
      # if not an error
28
     if (!inherits(try_error, 'error'))
29
        y = tolower(x)
30
      return(y)
31 }
32
33 * clean.corpus<-function(corpus){</pre>
34 corpus <- tm_map(corpus, removePunctuation)</p>
```

7_Topic_Modeling_Sentiment.R

Treemap shows a wider topic sentiment and length range.





"Pakistan" Guardian Mentions 11-1 to 11-8

- English Cricket
 - longer & positive
- Australian & NZ Cricket
 - Long & Negative
- Taliban
 - More numerous & Negative

-1.0 -0.5 0.0 0.5 1.0 1.5 sentiment

> Topic Modeling openNLP

7_Topic_Modeling_Sentiment.R

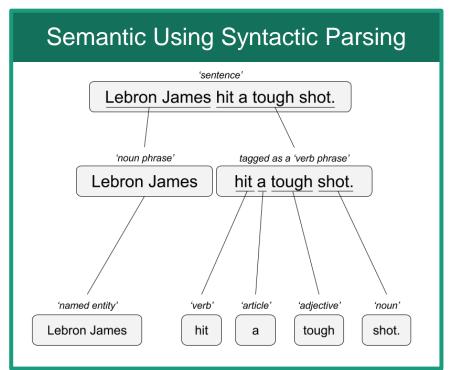
Using the portfolio package is a bit easier but more limited



Topic Modeling openNLP

Remember this? Text Mining Approaches

"Lebron James hit a tough shot."





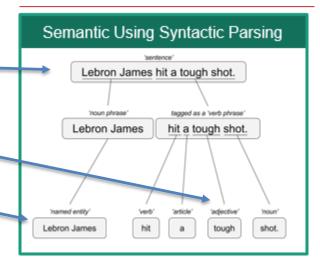
What is Text Mining?

library(openNLP)

- https://rpubs.com/lmullen/nlp-chapter
- Uses Annotators to identify the specific item in the corpus. Then holds all annotations in a plain text document. Think of auto-tagging words in a document and saving the document terms paired with the tag.
- Documentation, examples are hard to come by!

Annotations

- Grammatical or POS (Part of Speech) Tagging
- Sentence Tagging
- Word Tagging —
- Named Entity Recognition
 - Persons
 - Locations
 - Organizations



Annotations Specified

```
#OpenNLP Annotators
persons <- Maxent_Entity_Annotator(kind = 'person')
locations <- Maxent_Entity_Annotator(kind = 'location')
organizations <- Maxent_Entity_Annotator(kind = 'organization')
sent.token.annotator <- Maxent_Sent_Token_Annotator(language = "en")
word.token.annotator <- Maxent_Word_Token_Annotator(language = "en")
pos.tag.annotator <- Maxent_POS_Tag_Annotator(language = "en")</pre>
```

Annotations Applied to Text

```
#annotate text
annotations <- annotate(text.s,list(sent.token.annotator,word.token.annotator,pos
.tag.annotator,persons,locations,organizations))</pre>
```

Topic Modeling openNLP

```
#Extract Entities
entities <- function(doc, kind) {
    s <- doc$content
    a <- annotations(doc)[[1]]
    if(hasArg(kind)) {
        k <- sapply(a$features, `[[`, "kind")
        s[a[k == kind]]
    } else {
        s[a[a$type == "entity"]]
    }
}</pre>
```

The annotated plain text object is large with a complex structure. This function allows us to extract the tokens by tag kind.

```
The function creates a vector of people, locations and orgs that were "recognized"
people<-entities(text.annotations, kind = "person")</pre>
locations<-entities(text.annotations, kind = "location")</pre>
organization<-entities(text.annotations, kind = "organization")
head(people)
head(locations)
head(organization)
> head(people)
[1] "Marvin Gaye"
                   "Marvin"
                                   "LeedsVsSydney" "Marvin Gaye" "Marvin Gaye"
[6] "<U+266B>"
> head(locations)
[1] "Rainbow"
                      "Jerusalem"
                                       "Argentina"
                                                        "France"
[5] "La Petite Ferme" "Blue"
> head(organization)
[1] "RT"
                             "LeedsTour\005\nRT"
                                                       "#radio1xtra"
                             "MT"
                                                       "LeedsTourismBoard\nBlue"
[4] "#radio1xtra"
```

Questions?

https://github.com/kwartler/ODSC West 2016



www.linkedin.com/in/edwardkwartler