



Intro to text mining using tm, openNLP , & topicmodels

www.linkedin.com/in/edwardkwartler

 @tkwartler

Shameless Plug

Slated for June
2016



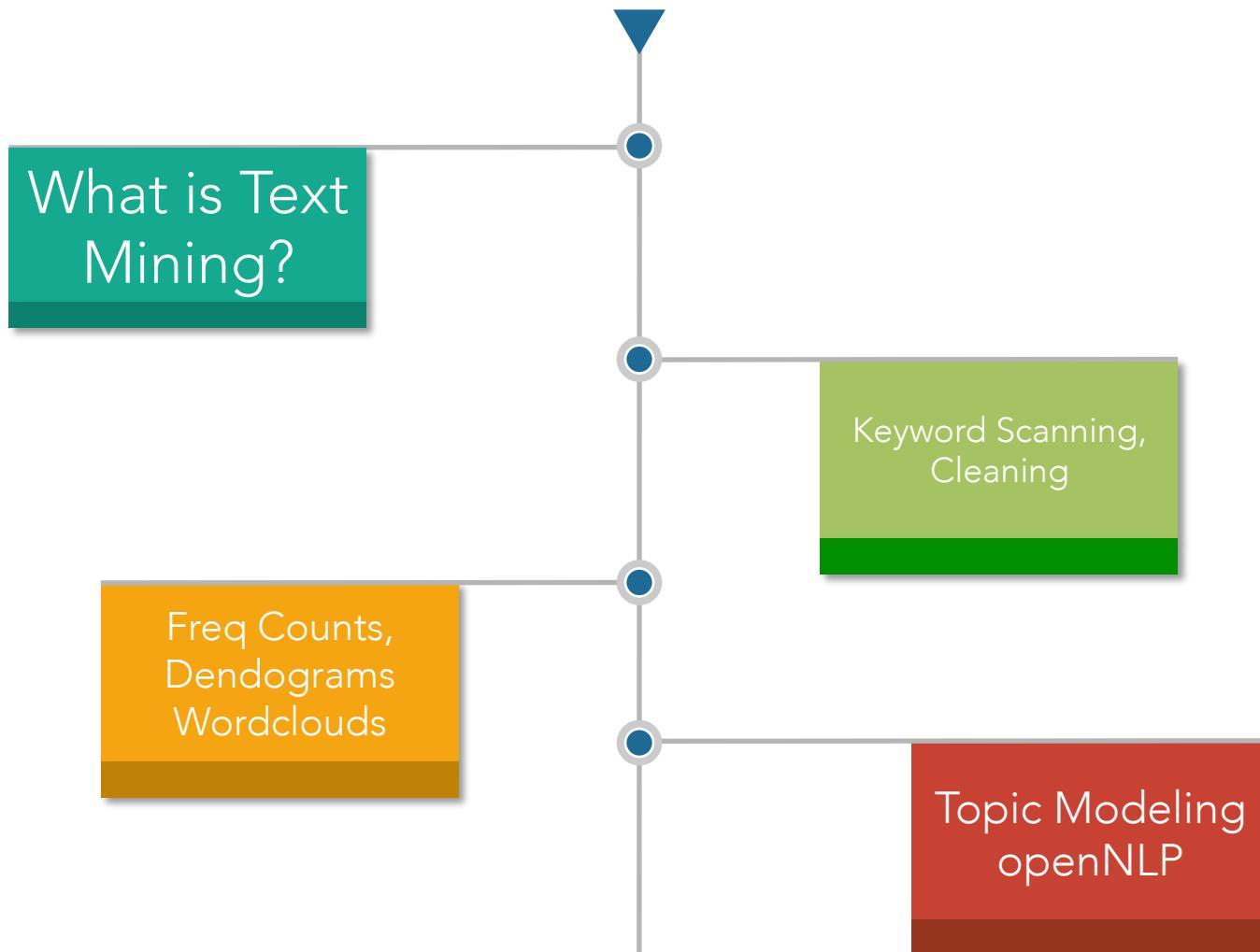
**Text Mining
IN R**

Edward Henry Kwartler

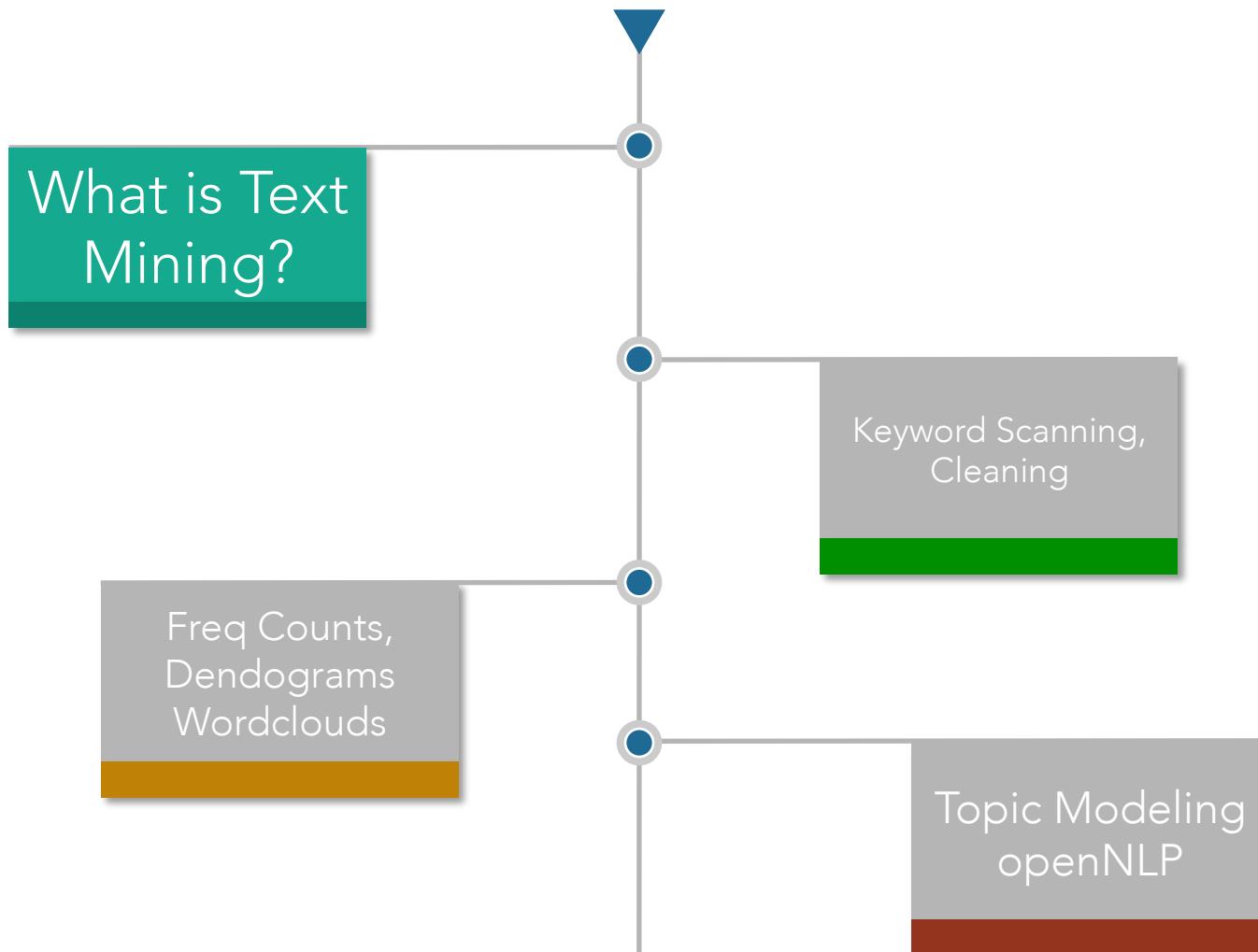
**cover art may change.*

 MANNING

Agenda



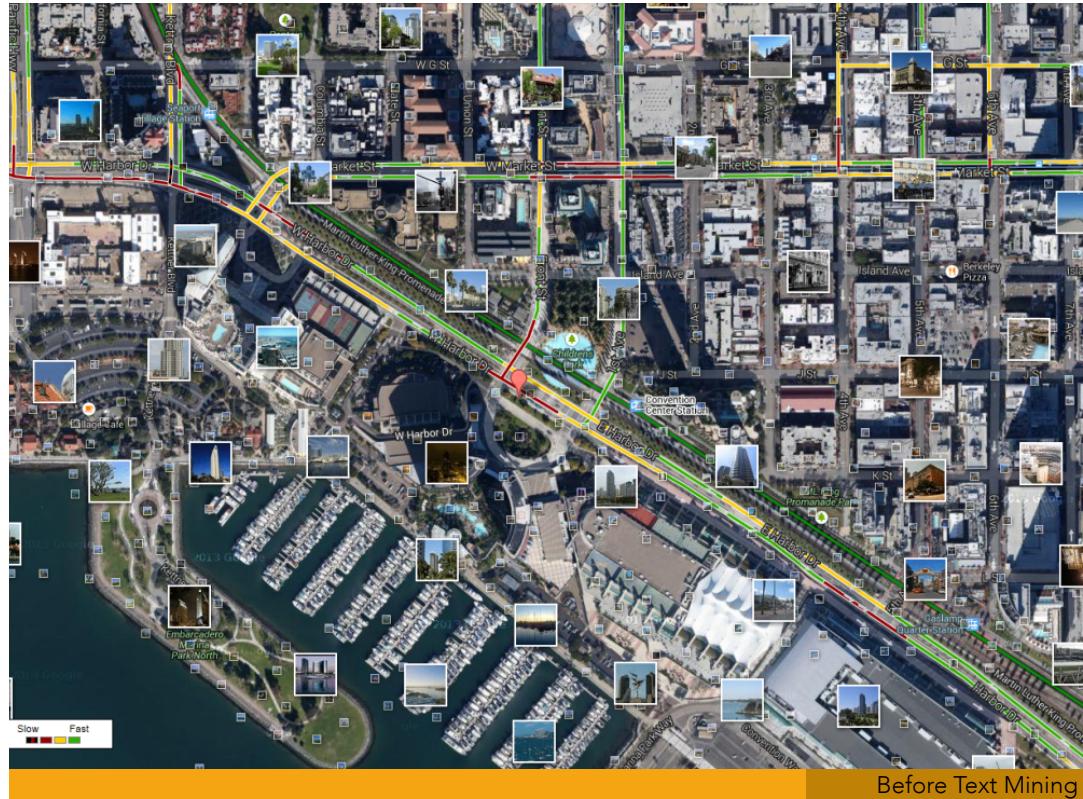
What is Text Mining?



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

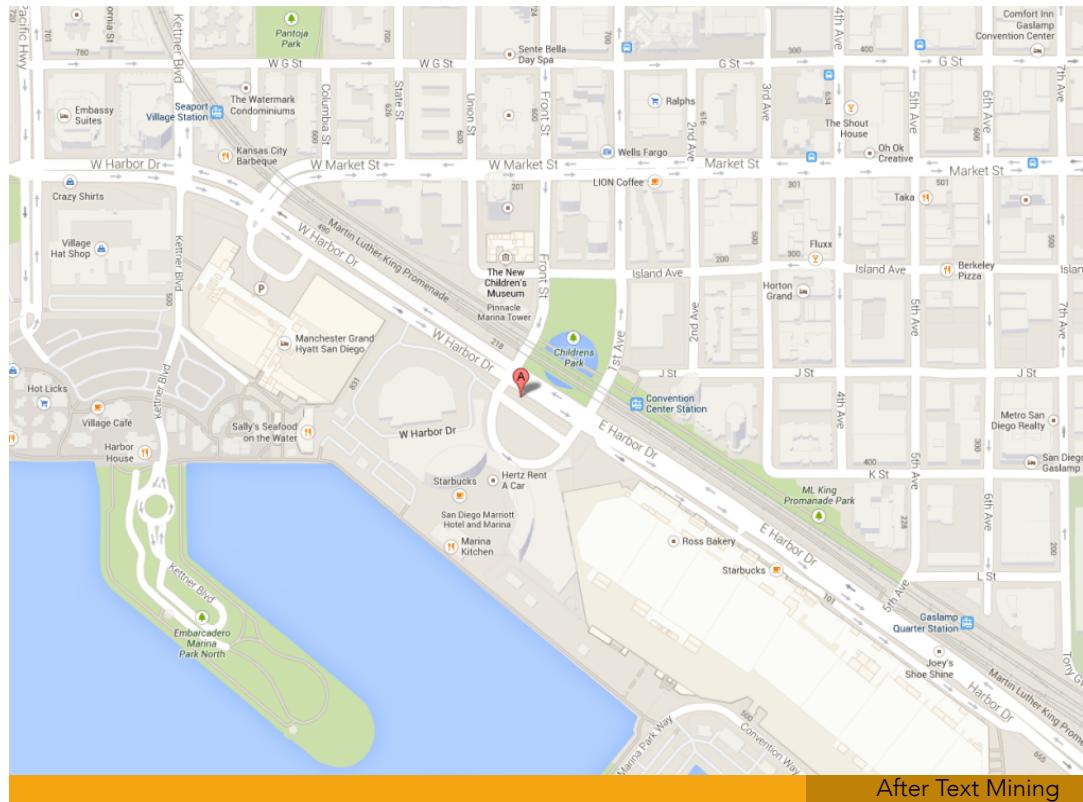
What is Text
Mining?



What is text mining?

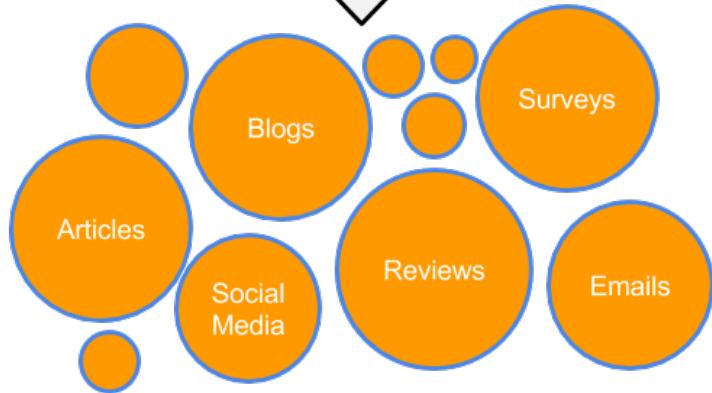
- Extract new insights from text
- Let's you drink from a fire hose of information
- Language is hard; many unsolved problems
 - Unstructured
 - Expression is individualistic
 - Multi-language/cultural implications

What is Text
Mining?



Text Mining Workflow

Gain subject matter expertise and define text mining goals.



Problem Definition

Unorganized State

1. Problem Definition
2. Identify Text Sources
3. Text Organization
4. Feature Extraction
5. Analytics
6. Reach Insight/Recommendation



Insight, Recommendation or analytical output.

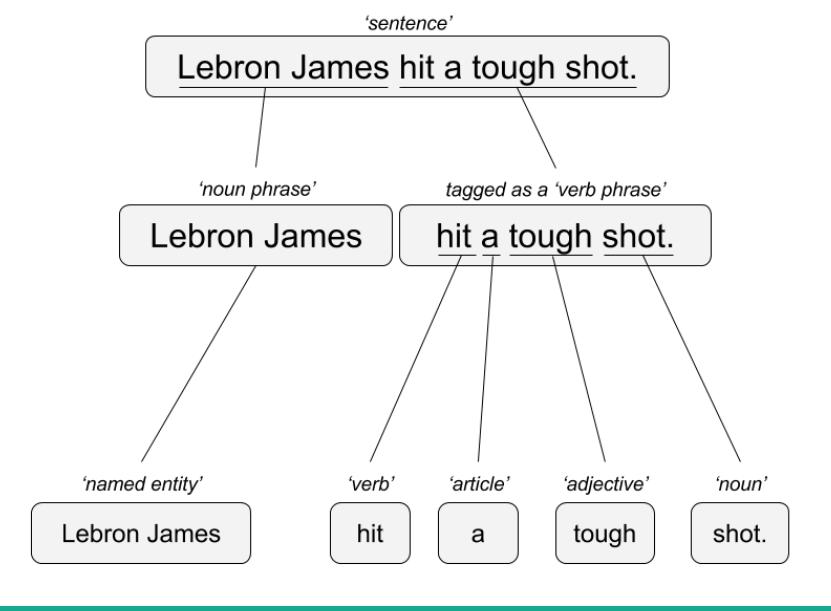
Organized State

What is Text
Mining?

Text Mining Approaches

"Lebron James hit a tough shot."

Semantic Using Syntactic Parsing



Bag of Words



What is Text
Mining?

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!!

What is Text
Mining?

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
- ...

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



What is Text
Mining?

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/>

Download scripts and data, save to a desktop folder & unzip

Install Packages

- Run “1_Install_Packages.R” script
 - An error may occur if the Java version doesn’t match

I have tested these scripts on R 3.1.2 with a Mac.
Some steps may be different in other environments.

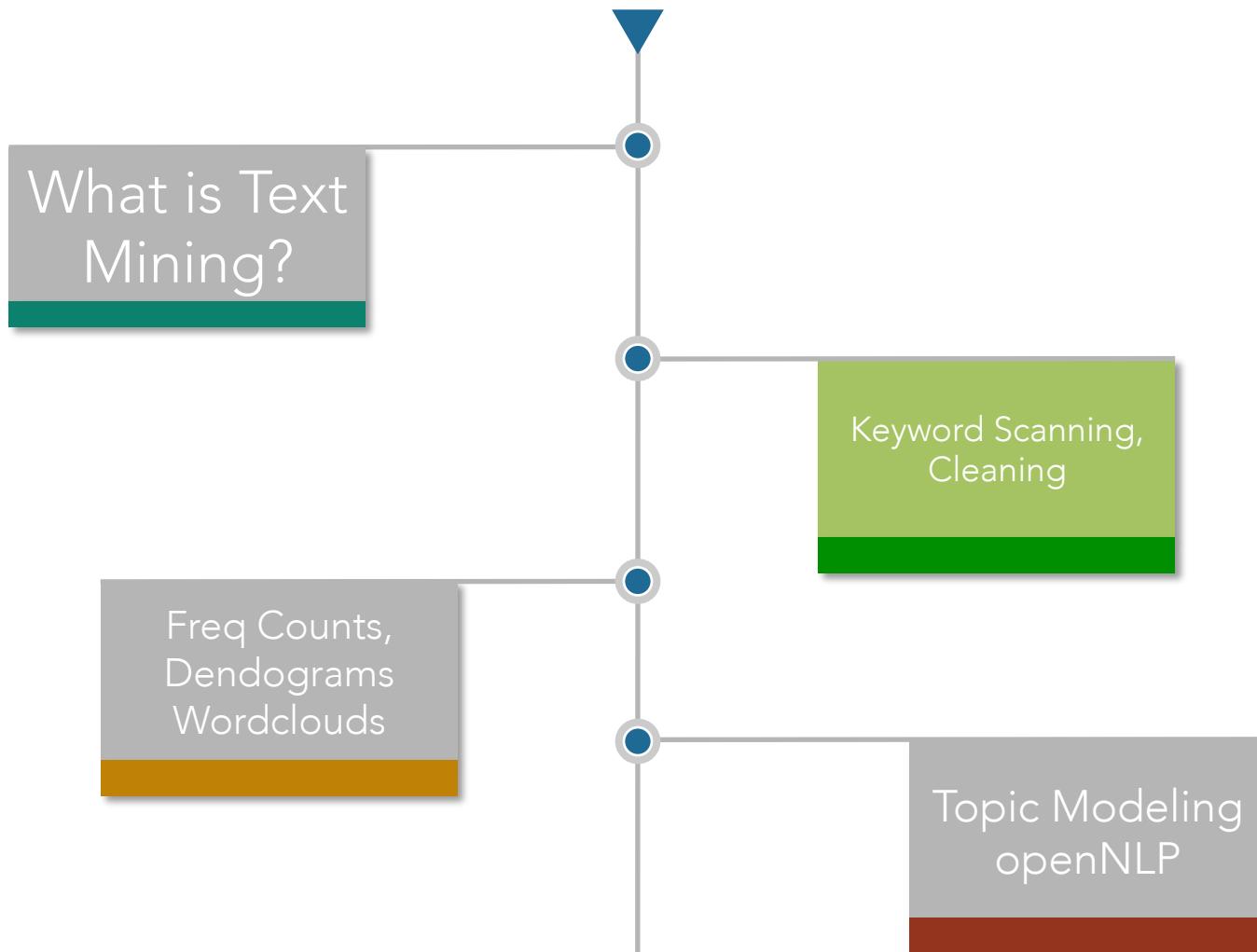
What is Text
Mining?

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, Cleaning & Freq Counts



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

1000 tweets mentioning "coffee"

text	favorited	replyToSN	created	truncated	replyToSID	id	replyToUID	statusSource	screenName	retweetCount	retweeted	longitude	latitude	
1 @ayyytylerb that is so true drink lots of coffee	FALSE	ayyytylerb	8/9/13 2:43	FALSE	3.6566E+17	1637123977	<a href="http://t.co/1hejenagib5	0	FALSE	NA	NA			
2 RT @bryz_brib: Senior March tnw morning at 7:25 A.M. in the SENIOR lot. Get up early, make yo coffee/breakfast, cus this will only happen ...	FALSE	NA	8/9/13 2:43	NA	3.6566E+17	NA	<a href="http://t.co/carolynicosia	1	FALSE	NA	NA			
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	FALSE	NA	8/9/13 2:43	FALSE	NA	3.6566E+17	NA	web	janeCKay	0	FALSE	NA	NA	
4 My cute coffee mug. http://t.co/2udvMU6XIG	FALSE	NA	8/9/13 2:43	FALSE	NA	3.6566E+17	NA	<a href="http://t.co/AlexandriaOr	0	FALSE	NA	NA		
5 RT @slaredo21: I wish we had Starbucks here... Cause coffee dates in the morning sound perfff!	FALSE	NA	8/9/13 2:43	FALSE	NA	3.6566E+17	NA	<a href="http://t.co/Rooosssaaaa	2	FALSE	NA	NA		
6 Does anyone ever get a cup of coffee before a cocktail??	FALSE	NA	8/9/13 2:43	FALSE	NA	3.6566E+17	NA	<a href="http://t.co/_Z_MAC	0	FALSE	NA	NA		
7 "I like my coffee like I like my women...black, bitter, and preferably fair trade." I love #Archer	FALSE	NA	8/9/13 2:43	FALSE	NA	3.6566E+17	NA	<a href="http://t.co/Charlie_3115	0	FALSE	NA	NA		
8 @dreamwwdiva ya didn't have coffee did ya?	FALSE	dreamwwed	8/9/13 2:43	FALSE	3.6566E+17	1316942208	<a href="http://t.co/jessicaSalvat	0	FALSE	NA	NA			
9 RT @Dougherty42: I just want some coffee.	FALSE	NA	8/9/13 2:43	FALSE	NA	3.6566E+17	NA	<a href="http://t.co/kaytieKirk	1	FALSE	NA	NA		
10 RT @Dork76: I can't care before coffee.	FALSE	NA	8/9/13 2:43	FALSE	NA	3.6566E+17	NA	<a href="http://t.co/listeria	2	FALSE	NA	NA		
11 No lie I wouldn't mind coming home smelling like coffee	FALSE	NA	8/9/13 2:43	FALSE	NA	3.6566E+17	NA	<a href="http://t.co/DOPECROOK	0	FALSE	NA	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/VTk...	FALSE	NA	8/9/13 2:43	FALSE	NA	3.6566E+17	NA	<a href="http://t.co/TiffCarusa	6	FALSE	NA	NA		
13 Have I ever told any of you that Tate Donovan bought my stepmom coffee?	FALSE	NA	8/9/13 2:43	FALSE	NA	3.6566E+17	NA	web	Curly'sCrazyN	0	FALSE	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/VTk...	FALSE	NA	8/9/13 2:43	FALSE	NA	3.6566E+17	NA	web	JoeJonasVA	6	FALSE	NA	NA	
15 @HeatherWhaley I was about 2 joint it takes 2 hands to hold hot coffee...then I read headline! #Don'tDrinkNShoot	FALSE	HeatherWha	8/9/13 2:42	FALSE	3.6566E+17	26035764	<a href="http://t.co/AnnaDuleep	0	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.	FALSE	NA	8/9/13 2:42	FALSE	NA	3.6566E+17	NA	<a href="http://t.co/mpr4437	42	FALSE	NA	NA		
17 Coffee always makes everything better.	FALSE	NA	8/9/13 2:42	FALSE	NA	3.6566E+17	NA	web	sharksrukri	0	FALSE	NA	NA	
18 RT @AdelaideReview: Food For Thought: @Annabelleats shares a delicious Venison and Porcini Mushroom Pie Recipe. http://t.co/N8Q7vqFKWN http...	FALSE	NA	8/9/13 2:42	FALSE	NA	3.6566E+17	NA	<a href="http://t.co/theppubake	1	FALSE	NA	NA		
19 RT @LittleMels: lmaoooo!!! @bryanlaca: nahhh Melanie us fa sho like an ummm a Coffee table :)) yeeeeee lmaoo"	FALSE	NA	8/9/13 2:42	FALSE	NA	3.6566E+17	NA	web	bryanlaca	1	FALSE	NA	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	NA	8/9/13 2:42	FALSE	NA	3.6566E+17	NA	<a href="http://t.co/Shauncore	0	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

Keyword
scanning,
Cleaning &
Freq Counts

2_Keyword_Scanning.R

Basic R Unix Commands

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

```
[1] TRUE  
[18] FALSE TRUE TRUE TRUE TRUE FALSE TRUE TRUE
```

grep returns the position of the pattern in the document

```
[1] 4 214 276 366 479 534 549 620
```

“library(stringi)” Functions

`stri_count` counts the number of patterns in a document

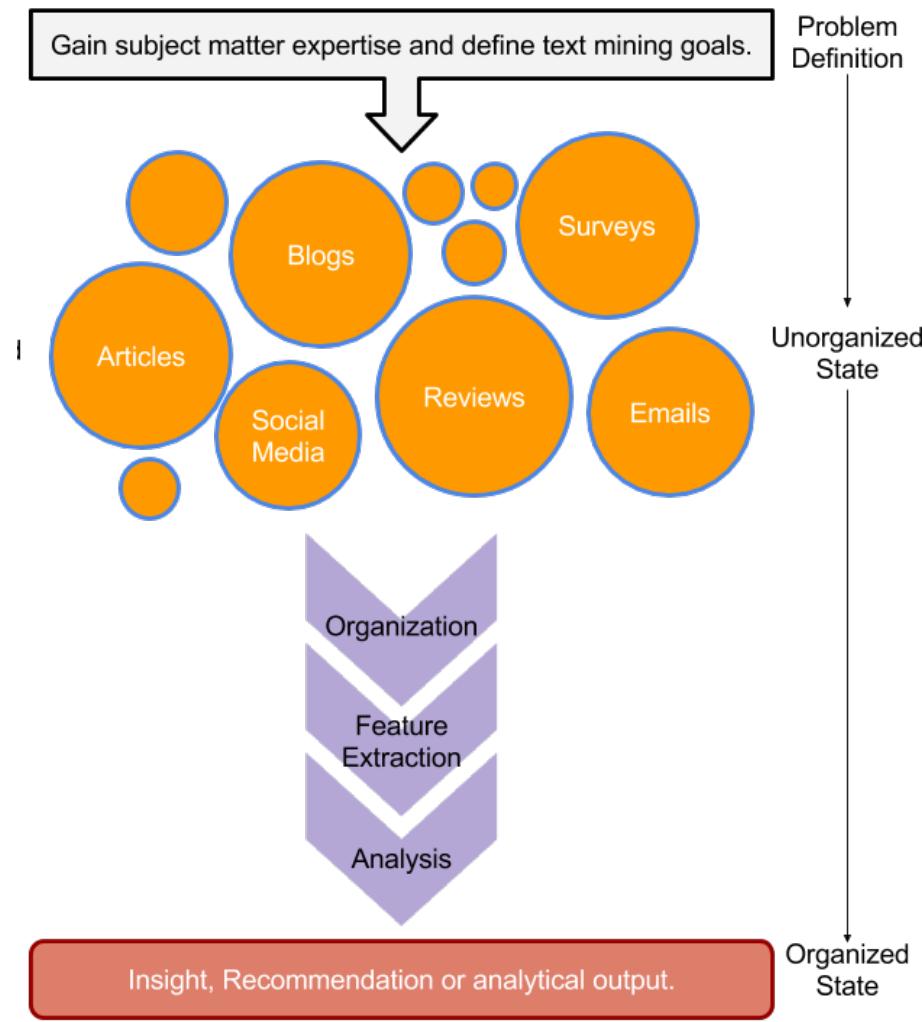
Keyword scanning, Cleaning & Freq Counts

Open 2_Keyword_Scanning.R

SET YOUR WORKING DIRECTORY

- Remember your slashes
- Save it so you won't have to do it again

Remember This?



What is Text
Mining?

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

Keyword
scanning,
Cleaning &
Freq Counts

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.

Keyword
scanning,
Cleaning &
Freq Counts

Open 3_Cleaning and Frequency Counts.R

SET YOUR WORKING DIRECTORY

- Remember your slashes
- Save it so you won't have to do it again

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)  
}
```

"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, tryTolower)  
  corpus <- tm_map(corpus, removeWords, custom.stopwords)  
  return(corpus)
```

"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

```
dd<-data.frame(id=text$id, text=text$text)  
custom.reader <- ReadTabular(mapping=list(content="text", id="id"))  
corpus <- VCorpus(DataframeSource(dd), readerControl=list(reader=custom.reader))
```

Keyword
scanning,
Cleaning &
Freq Counts

3_Cleaning and Frequency Count.R

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

Term Document Matrix

	Tweet1	Tweet2	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

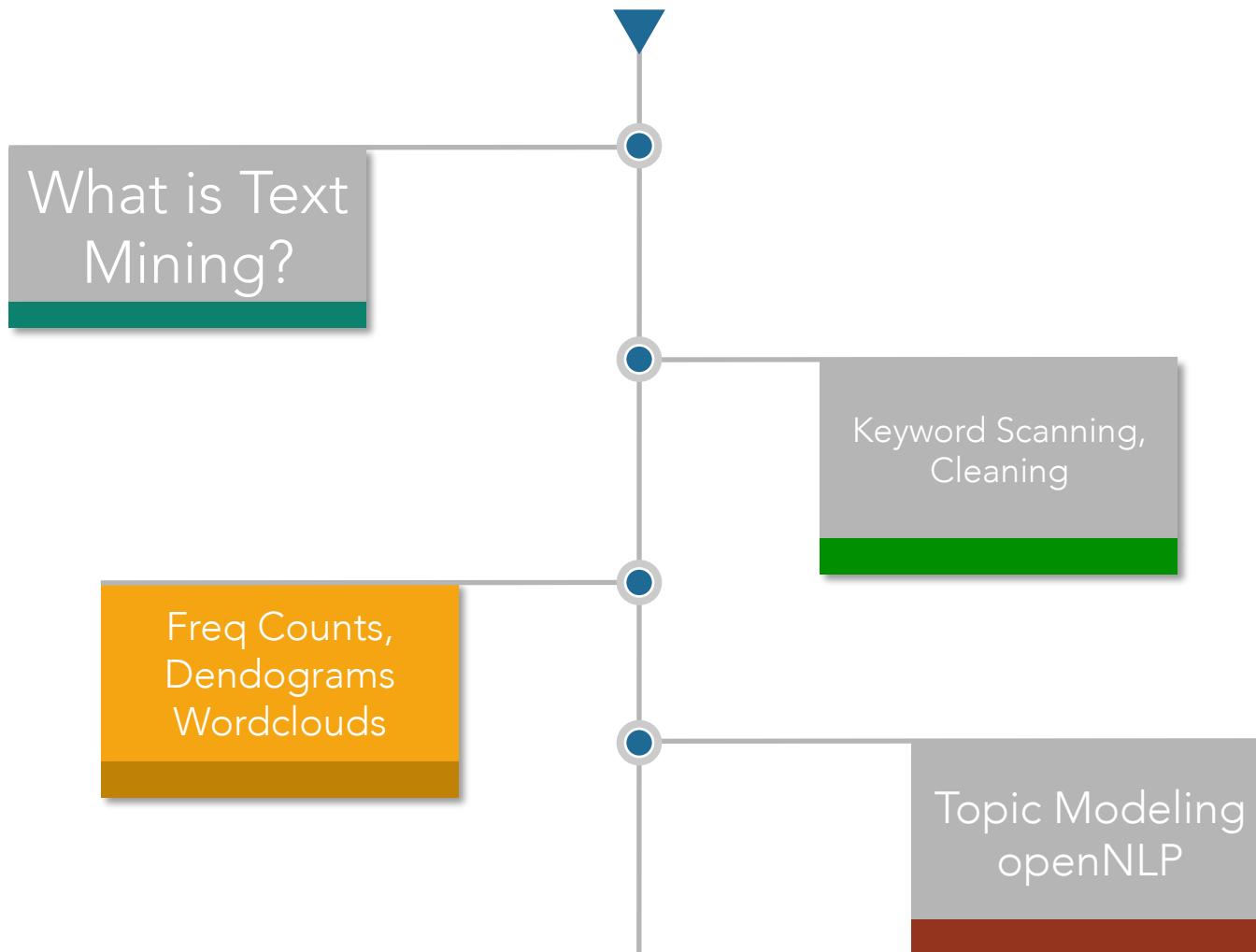
```
dtm<-DocumentTermMatrix(corpus)
tdm<-TermDocumentMatrix(corpus)
dtm.tweets.m<-as.matrix(dtm)
tdm.tweets.m<-as.matrix(tdm)
```

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

Keyword
scanning,
Cleaning &
Freq Counts

Dendograms & Wordclouds



Open 4_Dendograms.R

SET YOUR WORKING DIRECTORY

- Remember your slashes
- Save it so you won't have to do it again

4_Dendogram.R script builds on the matrices

First let's explore simple frequencies

```
#Frequency Data Frame
term.freq<-rowSums(tdm.tweets.m)
freq.df<-data.frame(word=names(term.freq),frequency=term.freq)

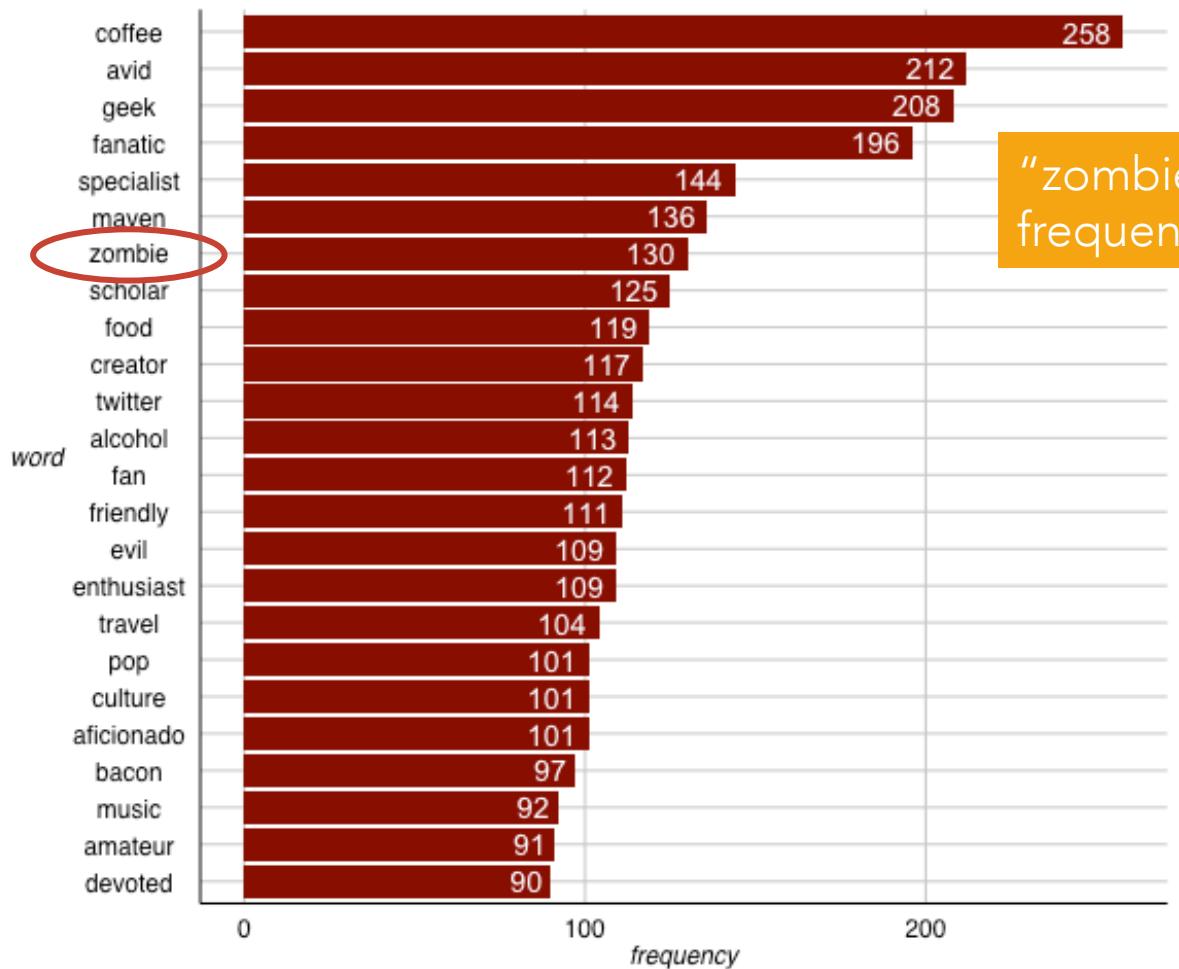
#Save a copy
write.csv(freq.df,'term_frequency.csv', row.names=F)

#Make a barplot of the top terms
top.words<-subset(freq.df, term.freq>=35)
top.words <- top.words[order(top.words$frequency, decreasing=F),]
top.words$word<-factor(top.words$word, levels=unique(as.character(top.words$word)))
ggplot(top.words, aes(x=word, y=frequency))+geom_bar(stat="identity", fill='darkred') +coord_flip() +theme_gdocs()+
  geom_text(aes(label=frequency), colour="white", hjust=1.25, size=5.0)
```

- Adjust $v \geq 35$ to limit the number of words in the dataframe. This will make the visual look better.

Freq Counts,
Dendograms
Wordclouds

4_Dendogram.R script



"zombie" is an unexpected, highly frequent term. Let's explore it.

Freq Counts,
Dendograms
Wordclouds

4_Dendogram.R script

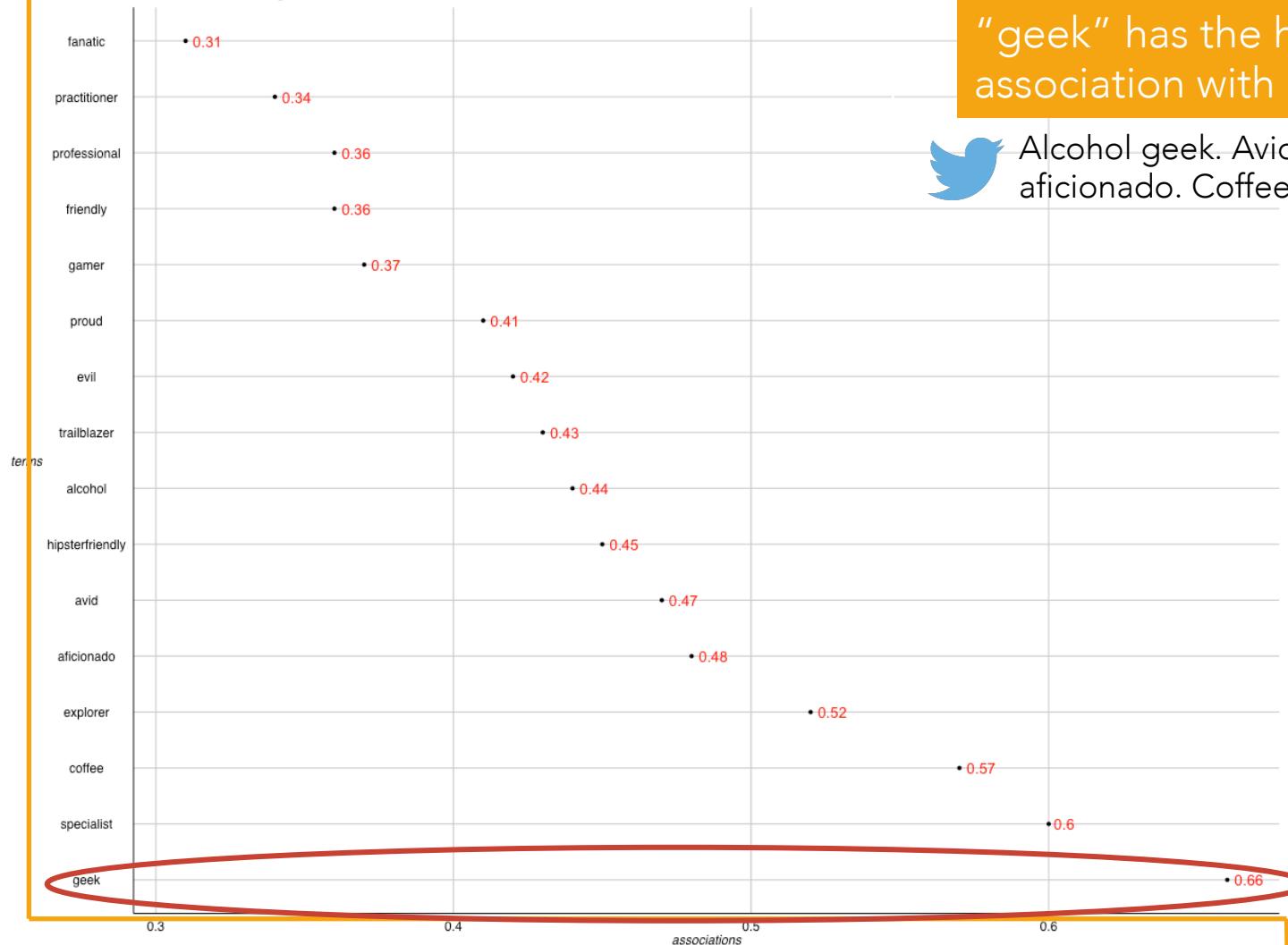
Next let's explore word associations, similar to correlation

```
#inspect word associations
associations<-findAssocs(tdm, 'zombie', 0.30),1]
terms<-row.names(findAssocs(tdm, zombie',0.30))
a.df<-data.frame(associations,terms)
a.df$terms<-factor(a.df$terms, levels=a.df$terms)
ggplot(a.df, aes(y=terms)) + geom_point(aes(x=associations), data=a.df) +
  theme_gdocs() + geom_text(aes(x=associations,label=associations), colour="red", hjust=-.25)
```

- Adjust 0.30 to get the terms that are associated .30 or more with the 'zombie' term.
- Make sure the 0.30s are the same for both in order to make the dataframe which is plotted.
- Treating the terms as factors lets ggplot2 sort them for a cleaner look.

Freq Counts,
Dendograms
Wordclouds

4_Dendogram.R script



"geek" has the highest word association with "zombie".



Alcohol geek. Avid tv buff. Friendly beer aficionado. Coffee guru. Zombie junkie.

Freq Counts,
Dendograms
Wordclouds

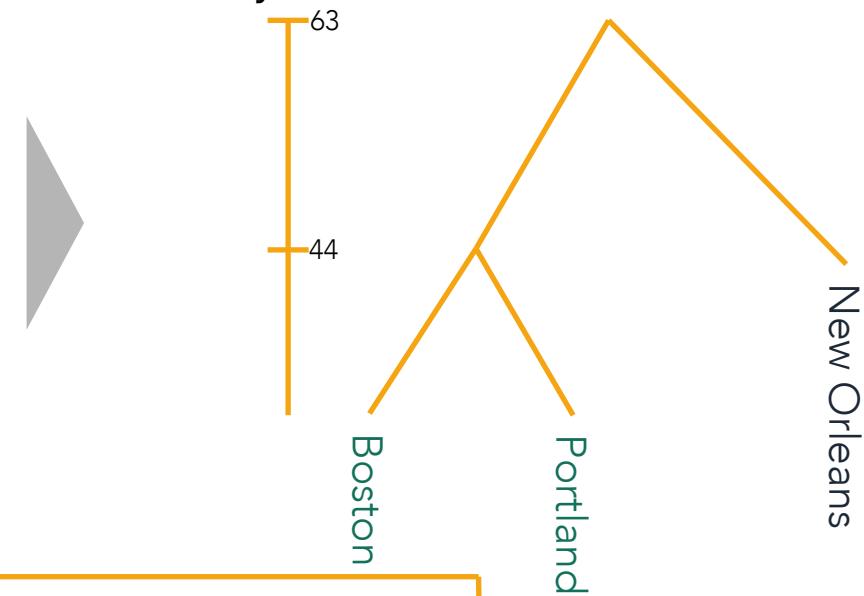
Extracting Meaning using Dendograms

Dendograms 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.



Freq Counts,
Dendograms
Wordclouds

4_Dendogram.R script

Weird associations! Maybe a dendrogram will help us more

```
#Hierarchical Clustering
tdm2 <- removeSparseTerms(tdm, sparse=0.95) #shoot 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")

> tdm
<<TermDocumentMatrix <-->> terms: 2488 documents: 1000>>
Non-/sparse entries: 10266/2477734
Sparsity : 100%
Maximal term length: 91
Weighting : term frequency (tf)

> tdm2
<<TermDocumentMatrix <-->> terms: 45 documents: 1000>>
Non-/sparse entries: 4491/40509
Sparsity : 90%
Maximal term length: 12
Weighting : term frequency (tf)
```

Adjust the sparse parameter to have ~40 to 50 terms in order to have a visually appealing dendrogram.

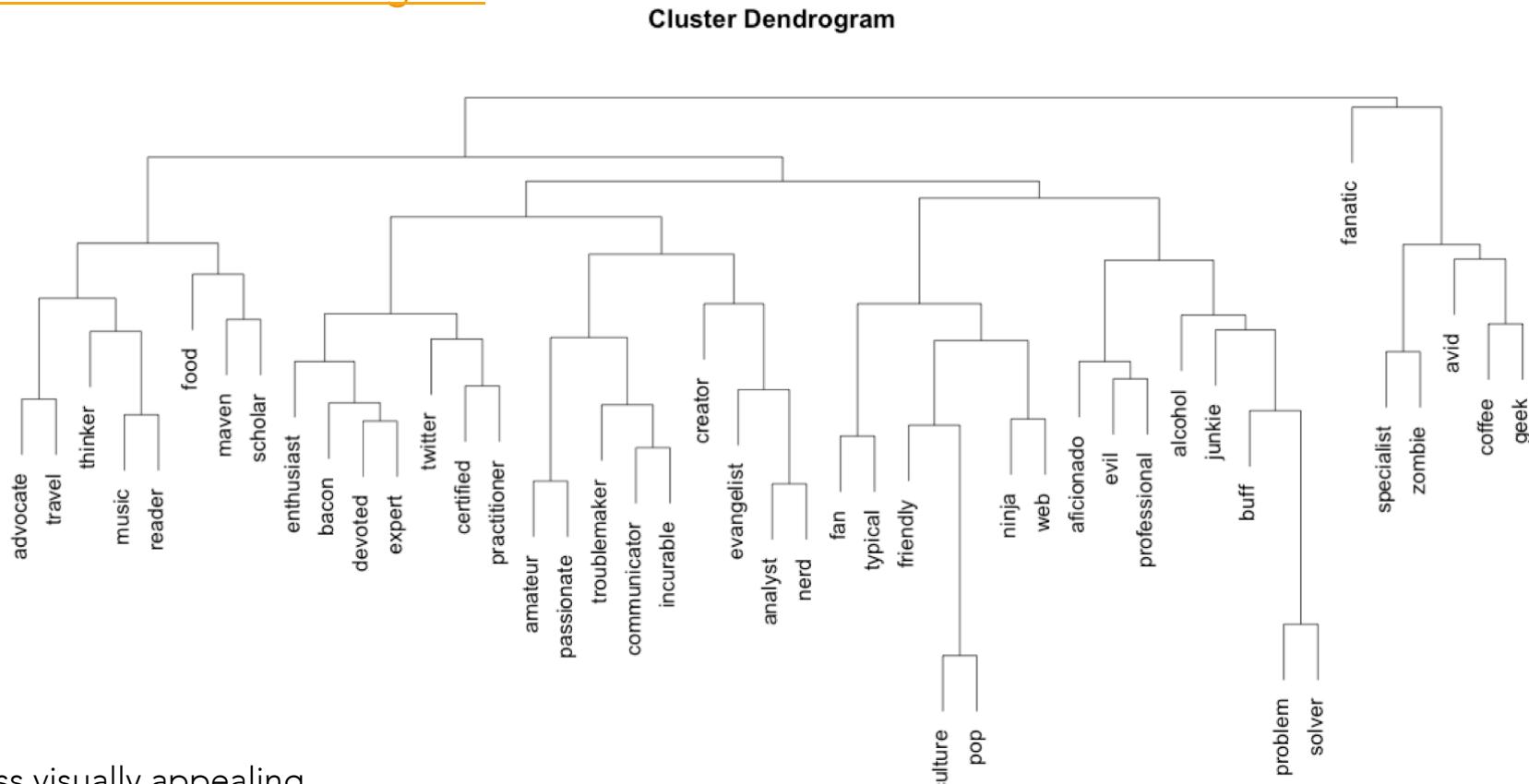
New Text Mining Concept

Sparse- 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.

Freq Counts,
Dendograms
Wordclouds

4_Dendogram.R script

Base Plot of a Dendrogram

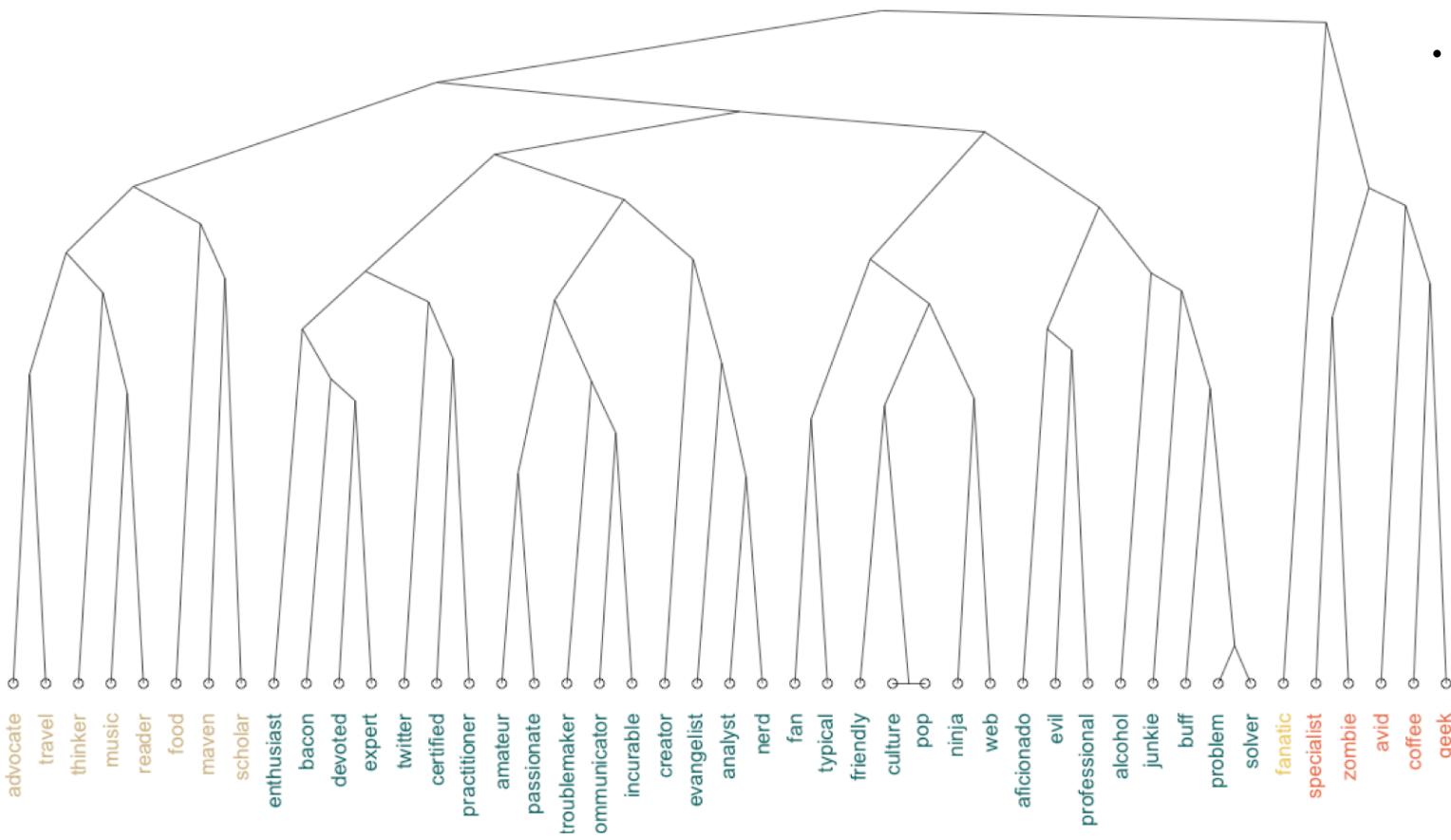


- Less visually appealing
- Clusters can be hard to read given the different heights

Freq Counts,
Dendograms
Wordclouds

4_Dendogram.R script

Hierarchical Dendrogram



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

Freq Counts,
Dendograms
Wordclouds

Open 5_Simple_Wordcloud.R

SET YOUR WORKING DIRECTORY

- Remember your slashes
- Save it so you won't have to do it again

5_Simple_Wordcloud.R script

```
#bigram token maker
bigram.tokenizer<-function(x)
  unlist(lapply(ngrams(words(x), 2), paste, collapse = " "), use.names = FALSE)

#Make a Document Term Matrix or Term Document Matrix depending on analysis
tdm<-TermDocumentMatrix(corpus, control=list(tokenize=bigram.tokenizer))
tdm.m <- as.matrix(tdm)
tdm.v <- sort(rowSums(tdm.m),decreasing=TRUE)
tdm.df <- data.frame(word = names(tdm.v),freq=tdm.v)
```

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.

Freq Counts,
Dendograms
Wordclouds

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.

```
> #Make a Document Term Matrix or Term Document Matrix depending on analysis  
> tdm<-TermDocumentMatrix(corpus, control=list(tokenize=bigram.tokenizer))  
> tdm.m <- as.matrix(tdm)  
> tdm.v <- sort(rowSums(tdm.m),decreasing=TRUE)  
> tdm.df <- data.frame(word = names(tdm.v),freq=tdm.v)  
> head(tdm.df)
```

	word	freq
marvin gaye	marvin gaye	76
gaye chardonnay	gaye chardonnay	69
glass chardonnay	glass chardonnay	54
bottle chardonnay	bottle chardonnay	24
little marvin	little marvin	21
chardonnay just	chardonnay just	19

Freq Counts,
Dendograms
Wordclouds

5_Simple_Wordcloud.R script

Next we need to select the colors for the wordcloud.

```
#look at all available color pallettes & choose one
```

```
display.brewer.all()
```

```
pal <- brewer.pal(8, "Blues")
```

```
pal <- pal[-(1:2)]
```



Freq Counts,
Dendograms
Wordclouds

5_Simple_Wordcloud.R script

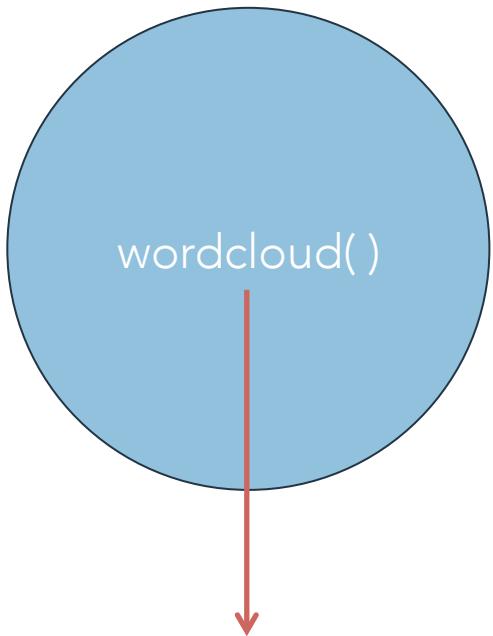
```
#create a simple word cloud  
wordcloud(tdm.df$word, tdm.df$freq, max.words=500, random.order=FALSE, colors=pal)
```



Freq Counts,
Dendograms
Wordclouds

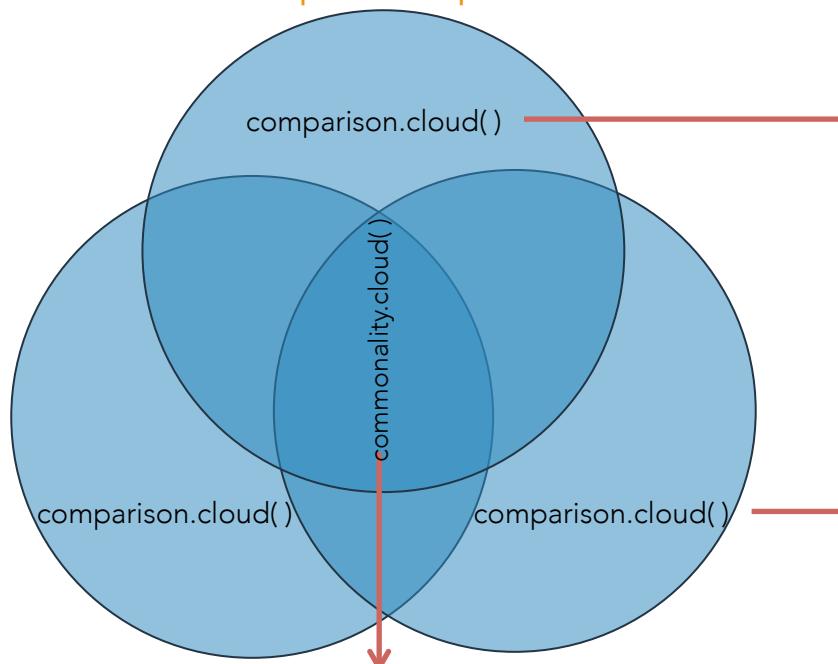
Types of Wordclouds

Single Corpus



lyrics from Marvin Gaye's songs, including:
rt wonderwines cabernet rose
meeting shiraz board meeting
big sean jinkxmonsoon jinkx
fancy winningmilkshake polite
set mood jinkx double moms rose
polite glass double fists mom's chardonnay
still live chocolate milkshake live fancy
fell porch bottle chardonnay
marvin gaye
full beauty little marvin jinkxmonsoon
just fell chardonnay just bottle board
fists chocolate inspired little
marvin chardonnay just set
chardonnay http://codudylkW/mood
it's gonna chardonnay cabernet mom's
don't love beauty grace rose bushes
love chicken grace chardonnay brought marvin

Multiple Corpora



common words from multiple corpora, including:
hot make tonight happy
work think call now
time like love
know will
day just night
one good
big life need
still get
back new
ace yes
drinking food
terlast girl that's want
fucking start see
great always
you're

common words from multiple corpora, including:
jessicajames http://jonon
movethelicks mom's that's Chardonnay
chemicals mom's need jinkxmonsoon
cause starbucks gr... cancer tea
cigarettes...
milkshake port
vitamin working cup lol gay wine
need
icedlike have little marvin wall
Office car's make while
renovation looks hgtv hosting
much charlie caused get
carries simple...
ultra...
205...
jxx craft
food...
friendly...
evil...
travel...
guru...
beer...
fan...
geek...
fanatic...
fan...
specialist...
bacon...
thinker...
http://

Freq Counts,
Dendograms
Wordclouds

Open 6_Other_Wordclouds.R

SET YOUR WORKING DIRECTORY

- Remember your slashes
- Save it so you won't have to do it again

6_Other_Wordcloud.R

Bring in more than one corpora.

```
#read CSV  
data1 <- read.csv(file="chardonnay2.csv", head=TRUE, sep=",")  
data2 <- read.csv(file="coffee.csv", head=TRUE, sep=",")  
data3 <- read.csv(file="beer.csv", head=TRUE, sep=",")
```

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

Apply each transformation separately...more in script.

```
#create a corpus for each group refencing the columns  
item1 <- VCorpus(DataframeSource(data.frame(data1$text)))  
item2 <- VCorpus(DataframeSource(data.frame(data2$text)))  
item3 <- VCorpus(DataframeSource(data.frame(data3$text)))  
  
#remove punctuation  
item1 <- tm_map(item1, removePunctuation)  
item2 <- tm_map(item2, removePunctuation)  
item3 <- tm_map(item3, removePunctuation)
```

Be sure to label the columns in the same order as was brought in.

```
#label the matrix columns in the same order you brought them in  
colnames(tdm) = c("Chardonnay", "Coffee", "Beer")
```

Freq Counts,
Dendograms
Wordclouds

Commonality Cloud

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

#create the word in common word cloud

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

Freq Counts,
Dendograms
Wordclouds

Commonality 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 mention morning cup and work.

```
#create the comparison word cloud
```

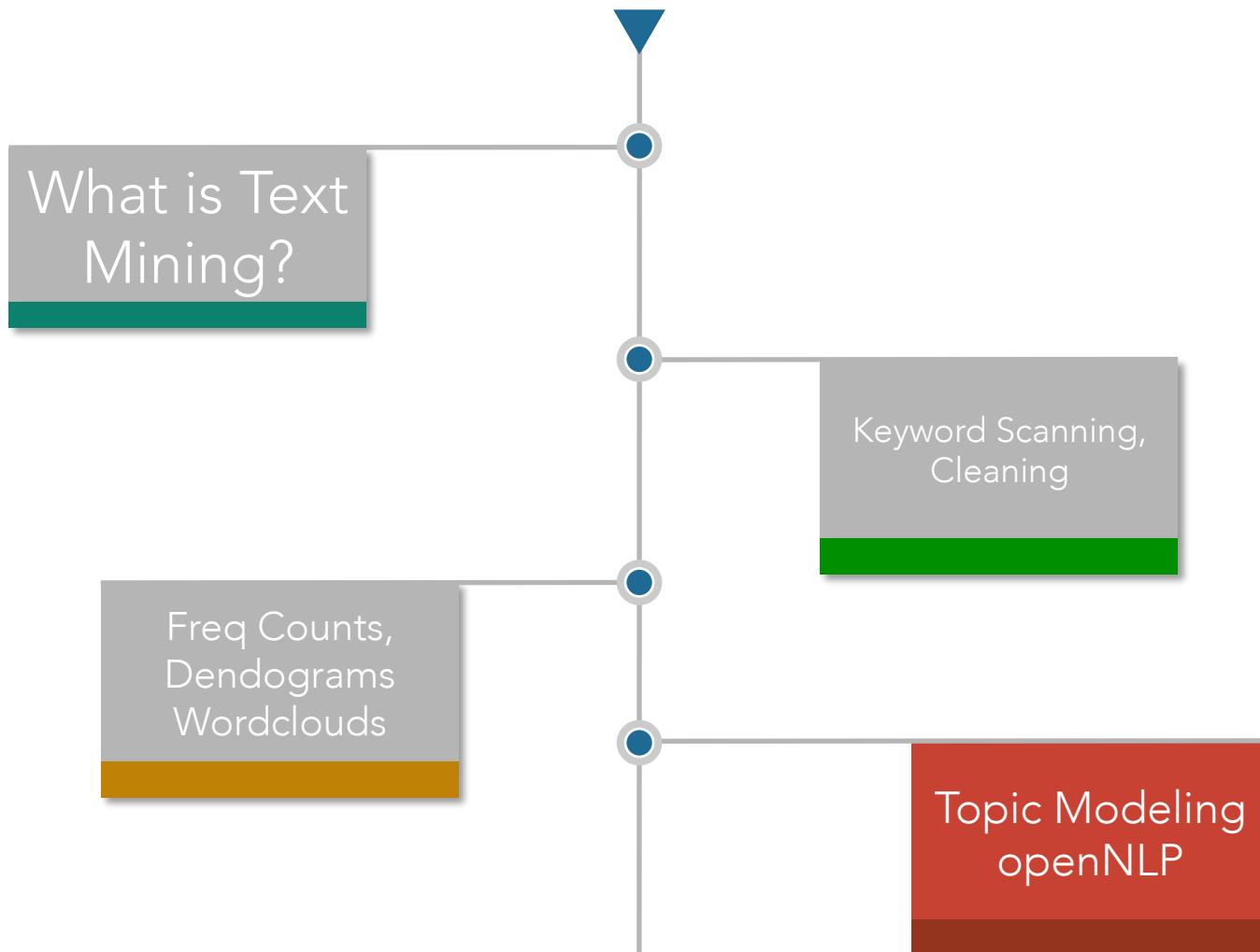
```
set.seed(1237)
```

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



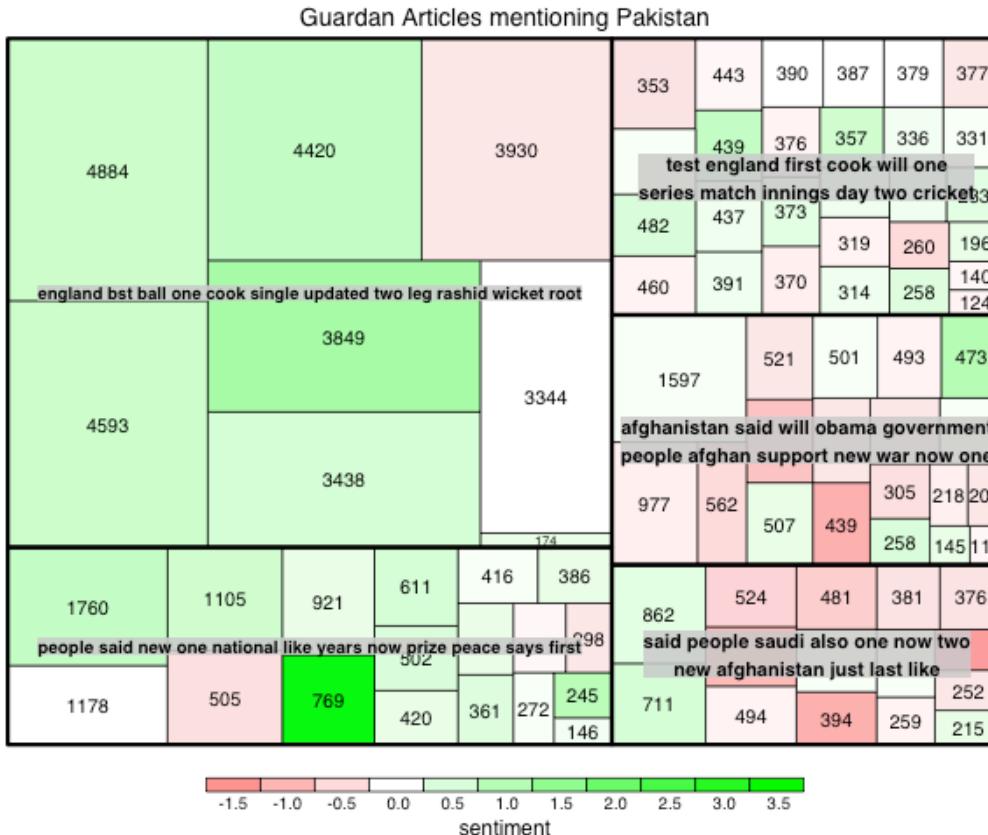
Freq Counts,
Dendograms
Wordclouds

Agenda



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 document length
- The larger grouping is based on LDA or CTM Topic Modeling

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

Second, LDA (Latent Dirichlet Allocation) 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_allocation
http://cs.brown.edu/courses/csci2950-p/spring2010/lectures/2010-03-03_santhanam.pdf

*Full disclosure I am still learning this methodology

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

Topic A: 30% basketball, 20% football, 15% watched, 10% TV...*something to do with watching sports*
Topic B: 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

First, simple Sentiment Polarity

Scoring

Surprise is a sentiment.

Hit by a bus! – Negative Polarity

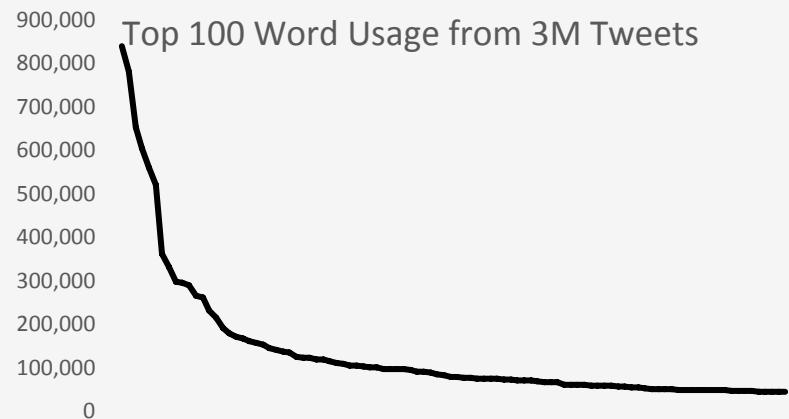
Won the lottery!- Positive Polarity

- I loathe BestBuy Service -1
- I love BestBuy Service. They are the best. +2
- I like shopping at BestBuy but hate 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. 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.



First, 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)
```

- Text 1: “love” was identified as positive. The text has 5 words and so $1/\sqrt{5} = .447$
- Text 2: “good” was identified positively. So $1/\sqrt{4}=0.5$
- Text 3: “good” was found along with the amplifier “very”. So $(.8+1)/\sqrt{5}=0.805$
- Text 4: hard and like cancel each other out so the polarity is zero. $1-1/\sqrt{9}=0$
- Text 5: “hard” is $-1/\sqrt{4}=-0.50$

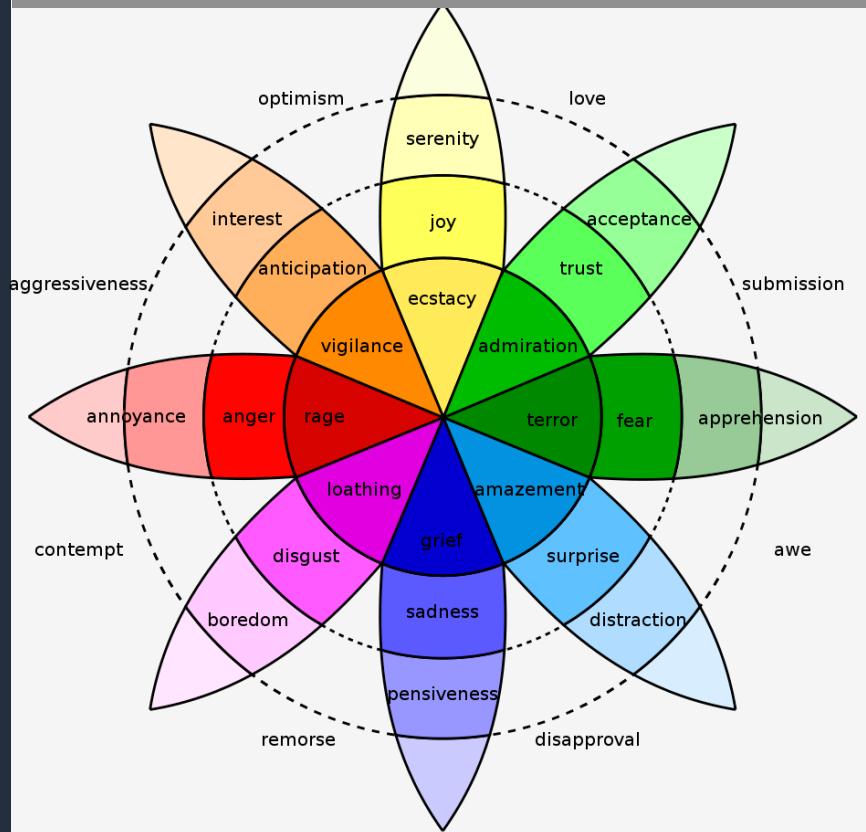
First it looks for the polarized word. Then identifies valence shifters (default 4 words before and two words after) Amplifiers are assigned +.8 and de-amplifiers weight is constrained to -1.

In reality sentiment is more complex.

Countless Emoji

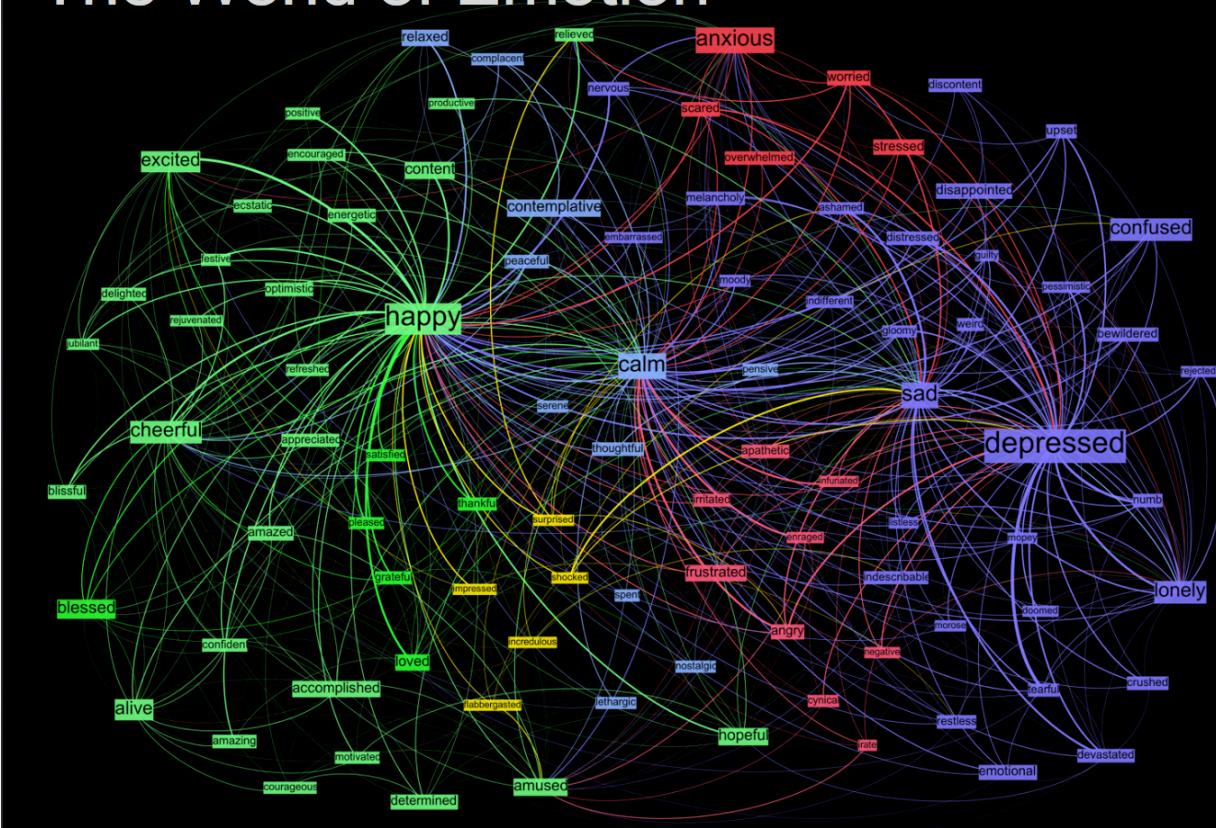


Plutchik's Wheel



Kanjoya's Experience Corpus

The World of Emotion



7_Topic_Modeling_Sentiment.R

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
3 #ODSC Workshop: Intro to Text Mining using R
4 #11-14-2015
5 #v7.3 Topic Modeling, Sentiment and Length Treemap
6
7 #Set the working directory
8 setwd('/Users/ted/Desktop/ODSC')
9
10 #libraries
11 library(treemap)
12 library(qdap)
13 library(GuardianR)
14 library(topicmodels)
15 library(tm)
16 library(SnowballC)
17
18 #options, functions
19 options(stringsAsFactors = FALSE) #text strings will not be factors of categories
20 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
24   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){
34   corpus <- tm_map(corpus, removePunctuation)
```

7_Topic_Modeling_Sentiment.R

Insurance Forum 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
 - Numerous & Negative
- 2 unknown topics needing more analysis

7_Topic_Modeling_Sentiment.R

Using the portfolio package is a bit easier but more limited

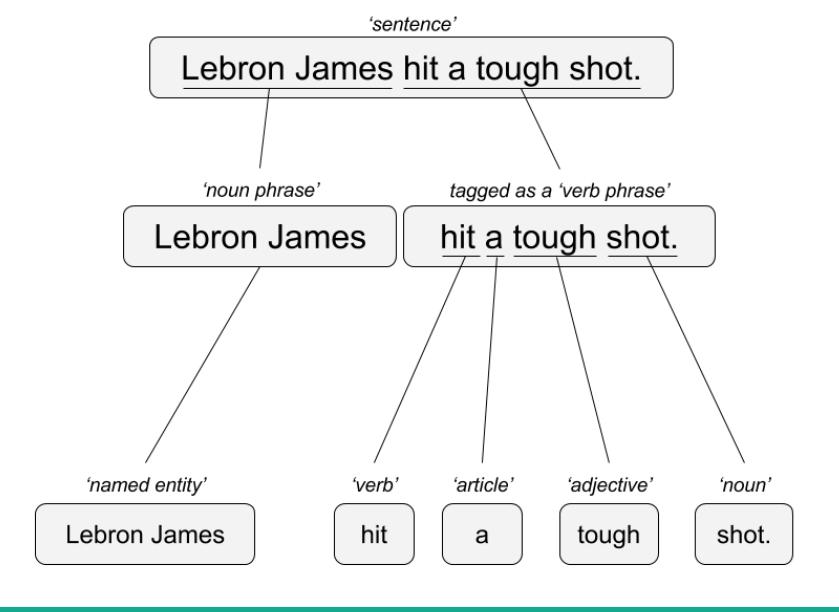
Sentiment/Color, Length/Area, Group/Topic



Remember this? Text Mining Approaches

"Lebron James hit a tough shot."

Semantic Using Syntactic Parsing



Bag of Words



What is Text
Mining?

8_Open_Language_Processing.R

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

Topic
Modeling
openNLP

8_Open_Language_Processing.R

```
#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"]
  }
}
```

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

```
people<-entities(text.annotations, kind = "person")
locations<-entities(text.annotations, kind = "location")
organization<-entities(text.annotations, kind = "organization")
```

8_Open_Language_Processing.R

```
people<-entities(text.annotations, kind = "person")
locations<-entities(text.annotations, kind = "location")
organization<-entities(text.annotations, kind = "organization")
```

```
> head(people)
[1] "Marvin Gaye"    "Marvin"          "LeedsVsSydney" "Marvin Gaye"    "Marvin Gaye"    "Marvin"
> head(locations)
[1] "Rainbow"        "Jerusalem"       "La Petite Ferme" "Blue"           "Sterling"      "Sterling"
> head(organization)
[1] "#radio1xtra"   "#radio1xtra"   "Blue"           "Hamilton"     "Chardonnay"   "Guess"
```

Questions?

 @tkwartler

www.linkedin.com/in/edwardkwartler