OPEN #C DATA SCIENCE CONFERENCE

#ODS©©

@ODSC

Boston | May 2018
Ted Kwartler
Intro to Text Mining in R

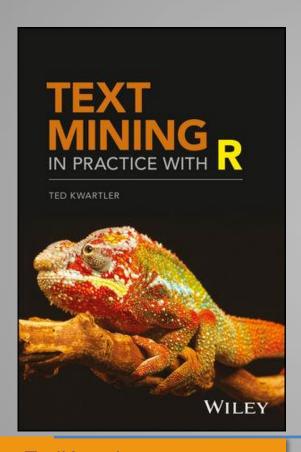


Intro to Natural Language Processing

www.linkedin.com/in/edwardkwartler



Shameless Plug #1



www.TedKwartler.com or www.amazon.com or www.wiley.com Discount Code: VBP65



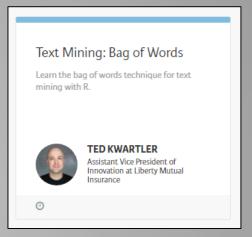


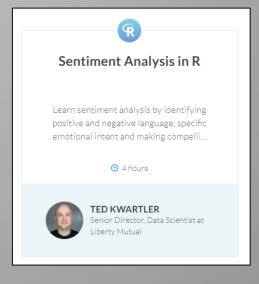


Shameless Plug #2



www.datacamp.com





Workshop Agenda

9-10:45 Intro to TM w/R 10:45-11 Break 11 – 11:45 Gather Text 11:45 – 12 Break 12-12:40 HTML Widgets & Flexdashboard 12:40-12:50 Optional Break

12:50-1 Scheduling

Goals:

- Learn the basics of text mining
- Apply methods to real data
- Follow a business workflow for stakeholders needing text mining insights

Intro to TM w/R Deep Dive

9-10:45 Intro to TM w/R 10:45-11 Break 11 – 11:45 Gather Text 11:45 – 12 Break 12-12:40 HTML Widgets & Flexdashboard

12:40-12:50 Optional Break

12:50-1 Scheduling

- •What is Text Mining?
- String Manipulation & Cleaning
- Frequency Counts, Dendrograms & Wordclouds
- Sentiment Analysis
- Topic Modeling
- OpenNLP/Named Entity Recognition

Goals:

- Learn the basics of text mining
- Apply methods to real data and uses
- Follow a business workflow for stakeholders needing text mining insights

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



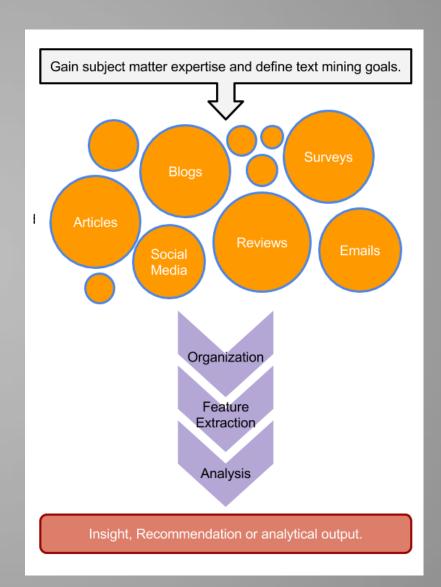
What is Text Mining?

- 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



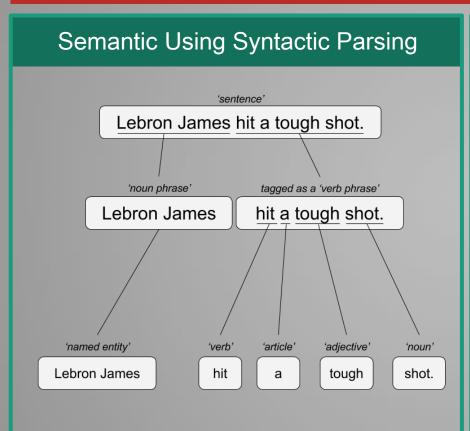
Text Mining Workflow

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



Two Popular 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.



Intro to TM w/R Deep Dive

- •What is Text Mining?
- String Manipulation & Cleaning
- Frequency Counts, Dendrograms & Wordclouds
- Sentiment Analysis
- Topic Modeling
- OpenNLP/Named Entity Recognition

Goals:

- Learn the basics of text mining
- Apply methods to real data and uses
- Follow a business workflow for stakeholders needing text mining insights

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

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

Setup

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

Setup Continued

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.



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 ayyytylerb 8/9/13 2:43 FALSE 3.6566E+17 3.6566E+17 1637123977 http://discount.com/abs/10.6566E+17 163712397 http://discount.com/abs/10.6566E+ 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 FALSE 3.6566E+17 NA janeCkay FALSE FALSE 4 My cute coffee mug. http://t.co/2udvMU6XIG 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 NA 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 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. FALSE FALSE 3.6566E+17 NA <a href="http:TiffCaruso FALSE NA 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 NΔ 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 HeatherWh 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. FALSE 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 <a href="htt; thepaulbake FALSE NA 8/9/13 2:42 FALSE NA 19 RT @LittleMelss: ImfaoIII"@bryanlaca; nahhh Melanie u is fa sho like an ummm a Coffee table :)) yeeeee Imaoo" FALSE 8/9/13 2:42 FALSE NA 3.6566E+17 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="htt; 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

1_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")
```

1_Keyword_Scanning

Intro to TM w/R Deep Dive

- •What is Text Mining?
- String Manipulation & Cleaning
- Frequency Counts, Dendrograms & Wordclouds
- Sentiment Analysis
- Topic Modeling
- OpenNLP/Named Entity Recognition

Goals:

- Learn the basics of text mining
- Apply methods to real data and uses
- Follow a business workflow for stakeholders needing text mining insights

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



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

"library(tm)" Functions

VCorpus creates a corpus held in memory.

VCorpus(source)

tm_map applies the transformations for the cleaning

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"

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.

"library(qdap)" Functions

Multiple Global Substitution

mgsub("search pattern", "replacement pattern", text object)

Family of Replace Functions

```
replace_abbreviation()- Replace Abbreviations
replace_contraction()- Replace Contractions
replace_number()- Replace Numbers With Text Representation
replace_ordinal()- Replace Mixed Ordinal Numbers With Text Representation
replace_symbol()- Replace Symbols with Word Equivalents
```

To use on a corpus you need to apply content_transformer

```
tm_map(corpus, content_transformer(replace_abbreviation))
```

New Text Mining Concepts

<u>Lemmatization</u> in linguistics, is the process of grouping together the different inflected forms of a word so they can be analyzed as a single item.

Poor Man's lemmatization

library(lexicon)

> hash_lemmas

```
token
                    lemma
       furtherst further
         skilled
                    skill.
    3: 'cause
                  because
    4:
                    b Luow
              'em
    5:
                     them
41529:
             Z005
                      Z00
41530:
       zoospores zoospore
       zucchinis zucchini
41531:
           zulus
41532:
                     zulu
41533:
         zygotes
                  zygote
```

Qdap's mgsub() lets you easily aggregate words

```
#Poor Man's Lemmatization
data(hash_lemmas)
text$text <- mgsub(hash_lemmas$token,hash_lemmas$lemma,text$text)</pre>
```

Warning,: Not done in the workshop because it takes a long time.

"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.

```
cleanCorpus<-function(corpus){
  corpus <- tm_map(corpus,
  content_transformer(qdapRegex::rm_url))
  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, customStopwords)
  return(corpus)
}</pre>
```

Base: tolower (basic)

Stringr: str_to_lower (wrapper)

Custom: tryTolower (handles errors)

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

```
#Create custom stop words
custom.stopwords <- c(stopwords('english') 'lol', 'smh')</pre>
```

Add channel specific stop words. E.g. Twitter abbreviations

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

```
# Data
text<-read.csv('coffee.csv', header=TRUE)

# As of tm version 0.7-3 tabular was deprecated
names(text)[1]<-'doc_id' #first 2 columns must be 'doc_id' &
'text'

txtCorpus <- VCorpus(DataframeSource(text))
txtCorpus<-cleanCorpus(txtCorpus)</pre>
```

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

txtDtm<-DocumentTermMatrix(txtCorpus)
txtTdm<-TermDocumentMatrix(txtCorpus)</pre>

txtDtmM<-as.matrix(txtDtm)</pre>

txtTdmM<-as.matrix(txtTdm)</pre>

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

Intro to TM w/R Deep Dive

- •What is Text Mining?
- String Manipulation & Cleaning
- Frequency Counts, Dendrograms & Wordclouds
- Sentiment Analysis
- Topic Modeling
- OpenNLP/Named Entity Recognition

Goals:

- Learn the basics of text mining
- Apply methods to real data and uses
- Follow a business workflow for stakeholders needing text mining insights

3_Dendrogram.R

First let's explore simple frequencies

beerFreq<-rowSums(beerTDMm)</pre>

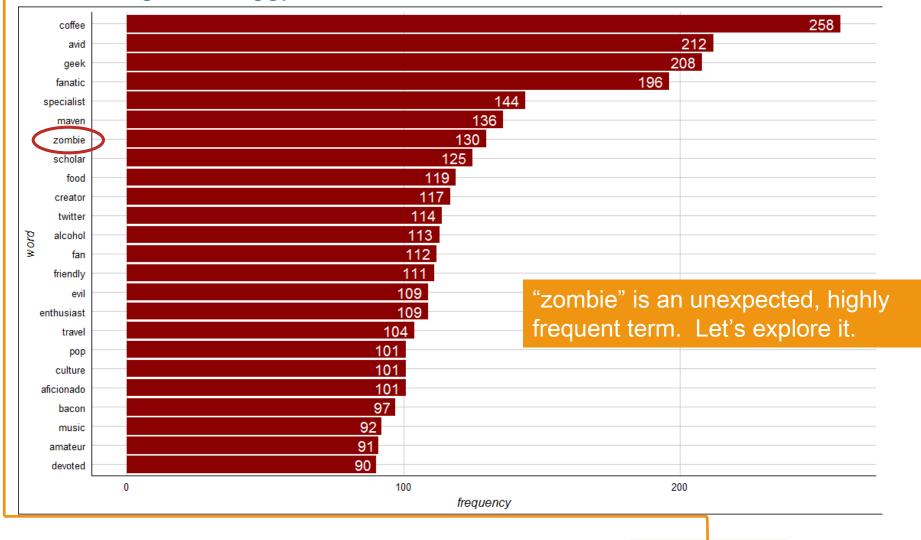
beerFreq<-data.frame(word=names(beerFreq),frequency=beerFreq)</pre>

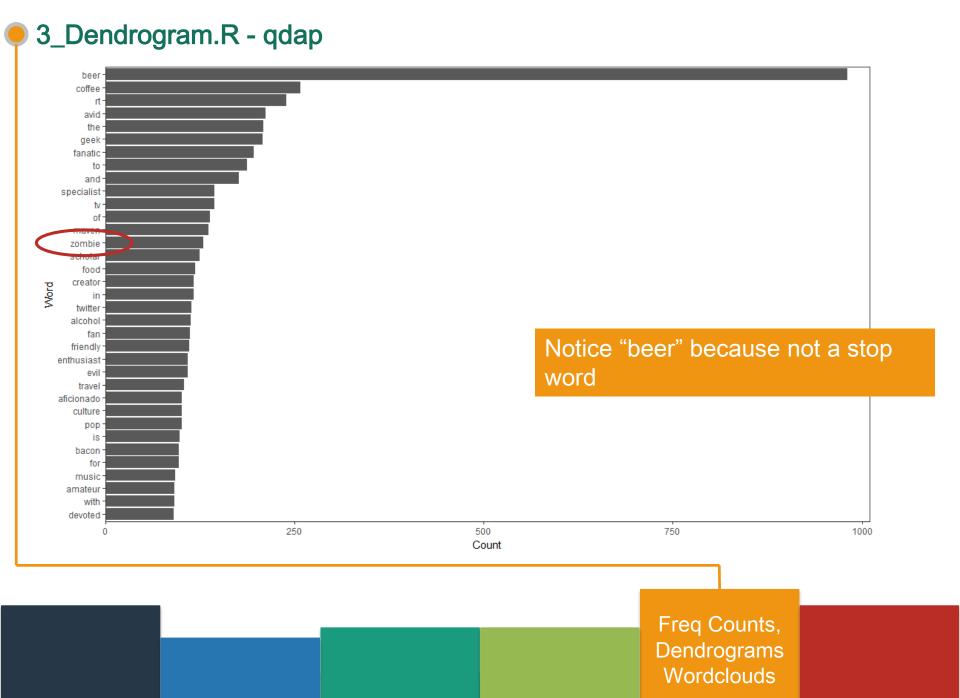
Term Document Matrix

	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

3_Dendrogram.R - ggplot



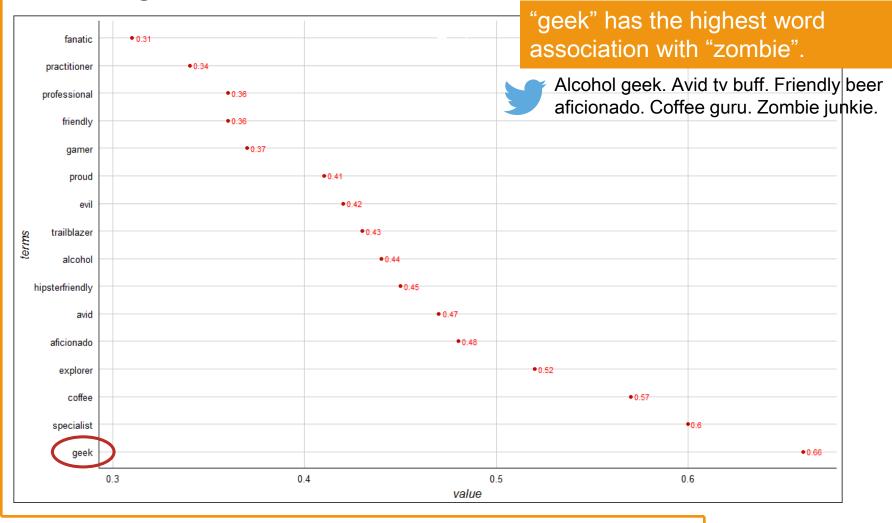


3_Dendrogram.R

Next let's explore word associations, similar to correlation

- 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.

3_Dendrogram.R



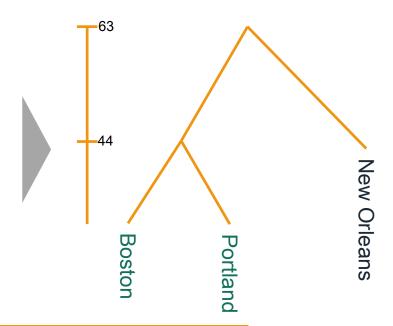
Extracting Meaning using dendrograms

Dendrograms visualize hierarchical clusters based on frequencies distances.

- 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

```
# Reduce TDM
beerTDM2 <- removeSparseTerms(beerTDM, sparse=0.97) #shoot for ~50 terms
beerTDM2<-as.data.frame(as.matrix(beerTDM2))

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

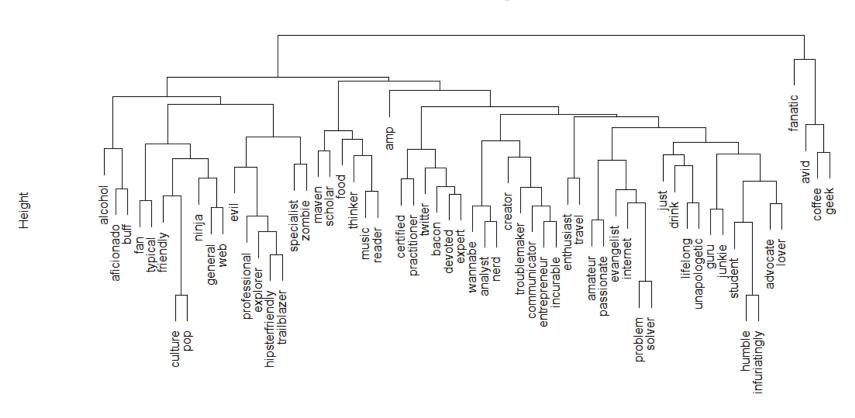
New Text Mining Concept

<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.



Base Plot of a Dendrogram

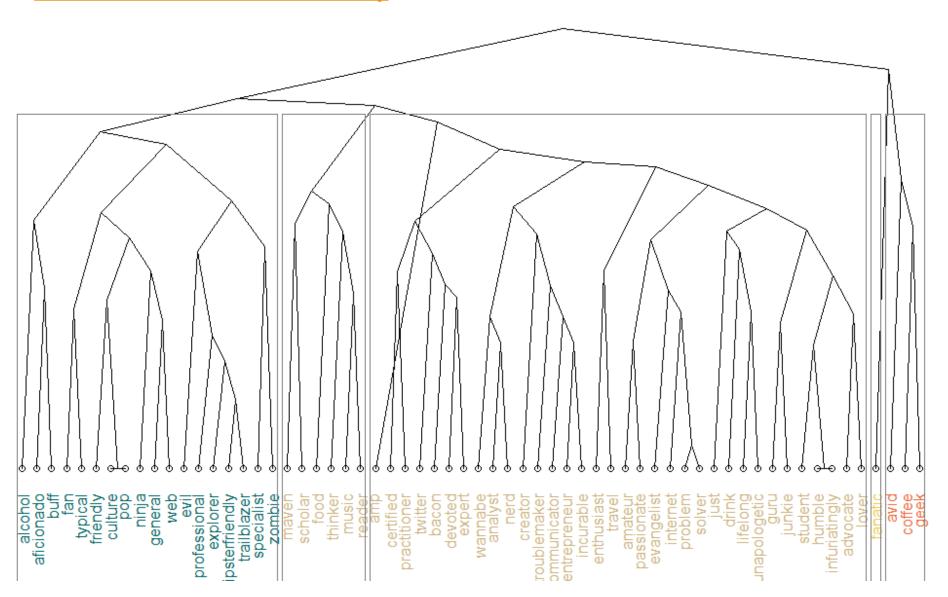
Cluster Dendrogram



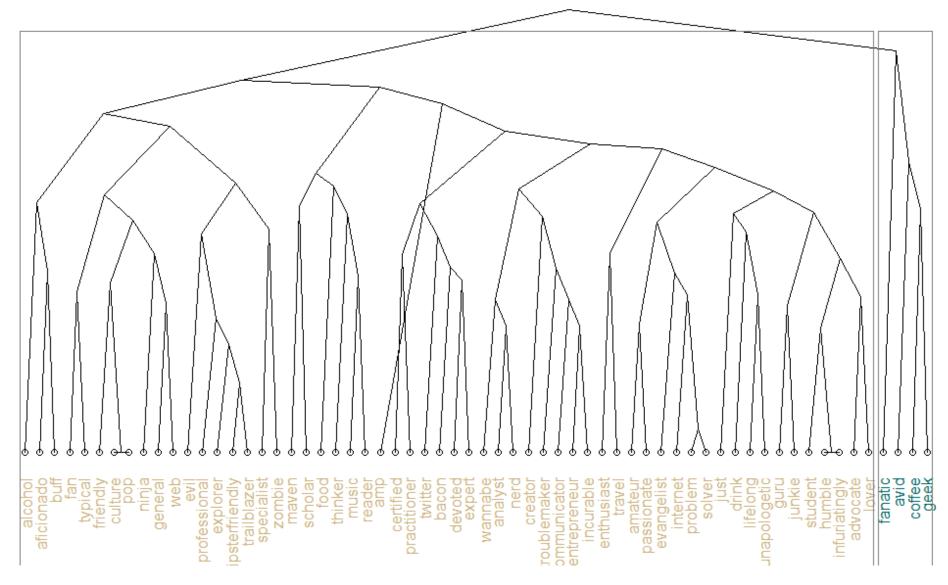
- Less visually appealing
- Clusters can be hard to read given the different heights

dist(beerTDM2) hclust (*, "complete")

Dendextend offers more flexibility



Dendextend offers more flexibility



Intro to TM w/R Deep Dive

- •What is Text Mining?
- String Manipulation & Cleaning
- •Frequency Counts, Dendrograms & Wordclouds
- Sentiment Analysis
- Topic Modeling
- OpenNLP/Named Entity Recognition

Goals:

- Learn the basics of text mining
- Apply methods to real data and uses
- Follow a business workflow for stakeholders needing text mining insights

```
#bigram token maker
bigramTokens <-function(x)
unlist(lapply(NLP::ngrams(words(x), 2), paste, collapse = " "), use.names = FALSE)</pre>
```

wineTDM<-TermDocumentMatrix(txtCorpus, control=list(tokenize=bigramTokens))</pre>

Text Mining is so fun. So do Text Mining!

Unigram

Docs Terms 1 fun. 1 mining 2 text 2

*with common stopwords

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	

New Text Mining Concept

<u>Tokenization</u>- 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.

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

```
# Get Row Sums
wineTDMv <- sort(rowSums(wineTDMm),decreasing=TRUE)
wineDF <- data.frame(word = names(wineTDMv),freq=wineTDMv)</pre>
```

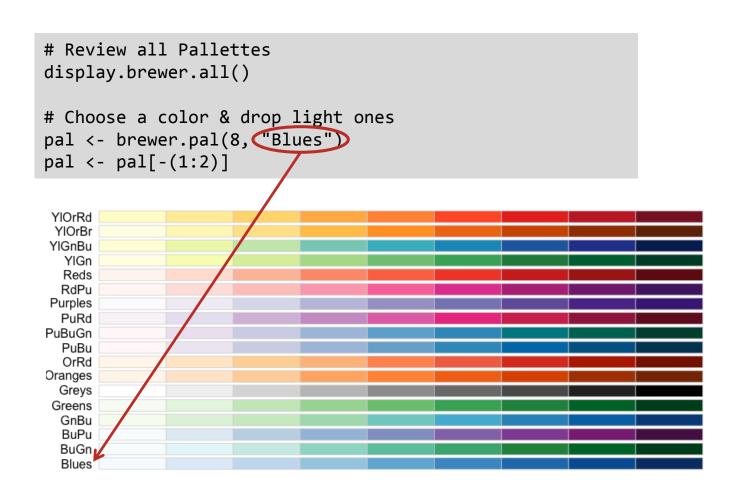
Term Document Matrix

	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

Next we need to select the colors for the wordcloud.



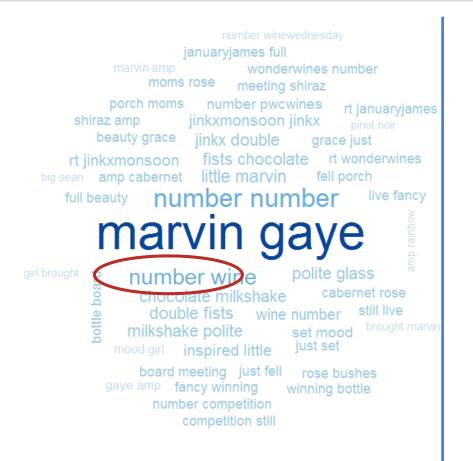
set.seed(1234)
wordcloud(wineDF\$word, wineDF\$freq, max.words=50, random.order=FALSE, colors=pal)

januaryjames full wonderwines number moms rose meeting shiraz number pwcwines porch moms rt januaryjames jinkxmonsoon jinkx shiraz amp beauty grace jinkx double grace just fists chocolate rt jinkxmonsoon rt wonderwines big sean amp cabernet fell porch live fancy number number full beauty marvin c polite glass number wine cabernet rose chocolate milkshake double fists wine number milkshake polite set mood iust set mood girl inspired little board meeting just fell rose bushes gaye amp fancy winning winning bottle number competition

competition still

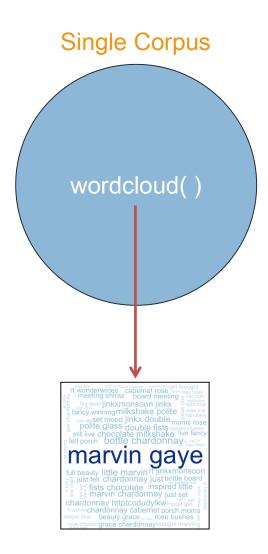
- Bigram Tokenization has captured "marvin gaye"
- A word cloud is a frequency visualization. The larger the term (or bigram here) the more frequent the term.
- You may get warnings if certain tokens are to large to be plotted in the graphics device.
- In cleanCopurs() the function ... changes numeric with the generic string "number" so be careful with your preprocessing steps!

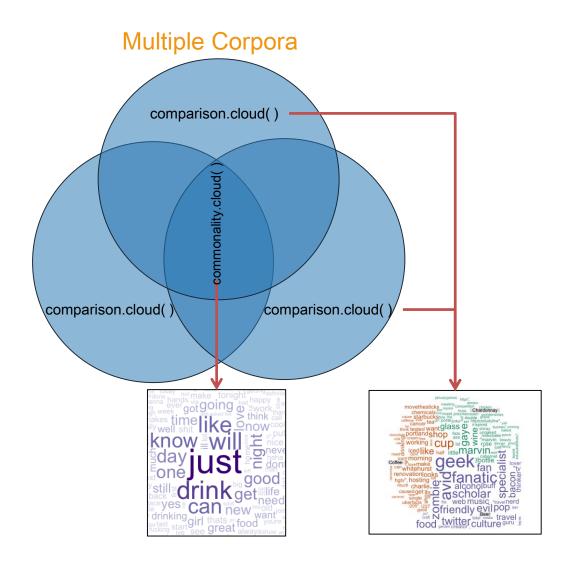
In cleanCopurs() the function replace_symbol() changes numeric with the generic string "number" so be careful with your preprocessing steps!





Types of Wordclouds





5_Other_Wordclouds.R

The challenge is working with multiple corpora efficiently. Many ways to do it...

```
# Read in multiple files as individuals
txtFiles<-c('chardonnay.csv','coffee.csv','beer.csv') #use list.files() for a
lot
for (i in 1:length(txtFiles)){
   assign(txtFiles[i], read.csv(txtFiles[i]))
   cat(paste('read completed:',txtFiles[i],'\n'))
}
# Read them into a single list with individual elements
all<-pblapply(txtFiles,read.csv)</pre>
```

Two ways to import csv files:

- 1. Each file is read in and an object created for each.
- 2. A list called "all" is created. Each list element represents a single document. Using data.table::rbindlist() one can create a single document from all files in the folder.

5_Other_Wordclouds.R

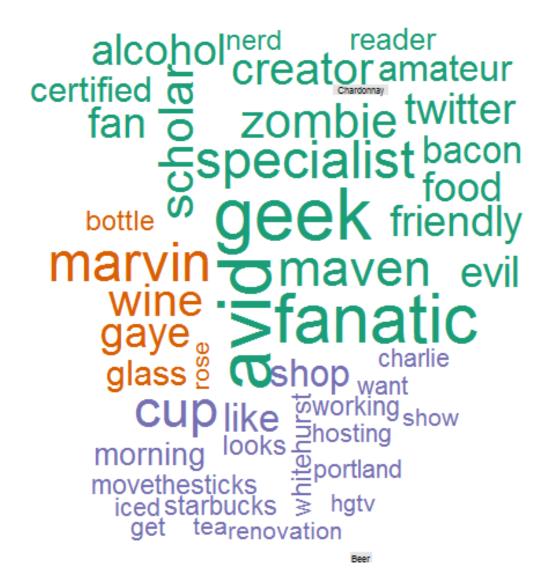
Commonality Cloud

- 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.



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.



Make comparison cloud
comparison.cloud(drinkTDMm, max.words=75, random.order=FALSE,
title.size=0.5,colors=brewer.pal(ncol(drinkTDMm),"Dark2"))

Intro to TM w/R Deep Dive

- •What is Text Mining?
- String Manipulation & Cleaning
- •Frequency Counts, Dendrograms & Wordclouds
- Sentiment Analysis
- Topic Modeling
- OpenNLP/Named Entity Recognition

Goals:

- Learn the basics of text mining
- Apply methods to real data and uses
- Follow a business workflow for stakeholders needing text mining insights

Simple Sentiment Polarity

Scoring

Surprise is a sentiment.

Hit by a bus! – Negative Polarity

Won the lottery!- Positive Polarity

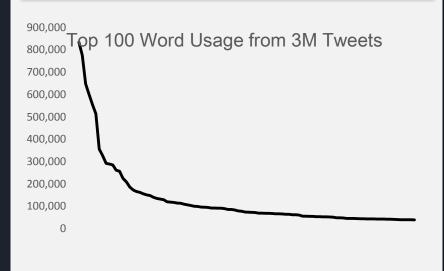
- I loathe BestBuy Service -1
- I <u>love</u> BestBuy Service. They are the <u>best</u>. +2
- I <u>like</u> shopping at BestBuy but <u>hate</u> traffic. 0

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)

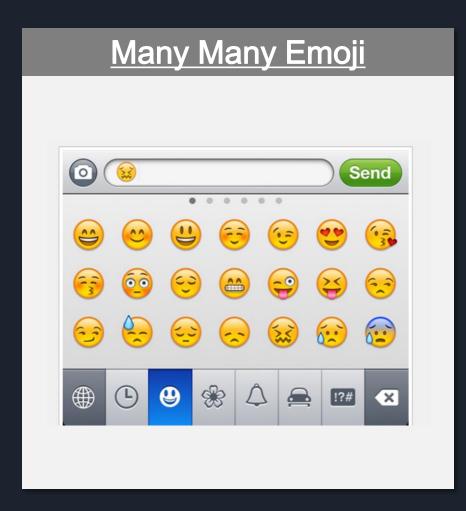
text1<-'i love St Peters University'
text2<-'this lecture is good'
text3<-'this lecture is very good'
text4<-'data science is hard I like it a little'
text5<-'data science is hard'

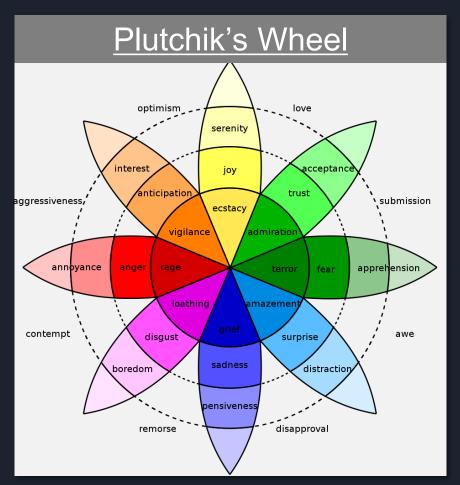
polarity(text1)
polarity(text2)
polarity(text3)
polarity(text4)
polarity(text5)</pre>
```

- <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
- Text 5: "hard" is -1/sqrt(4)=-.50

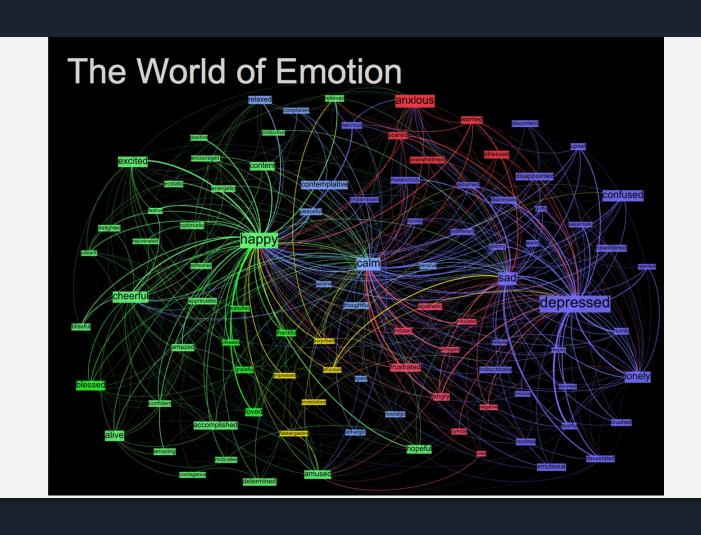
First it looks for the polarized word. Then identifies valence shifters (default 4 words before and two words after) Amplifiers are assigned +.8 and de-amplifiers weight is constrained to -1. Lastly the sum is divided by the square root of the total number of words in the passage.

In reality sentiment is more complex.





Kanjoya's Experience Corpus



Sentiment the Tidy Way!

```
##Tidy Sentiment Analysis
data(sentiments)
sentiments

#Stopwords
data(stop_words)
stop_words

#Add stopwords
custom.stopwords<-data.frame(word=c('amp','beer'),
lexicon='custom')

stop_words<-rbind(stop_words,custom.stopwords)</pre>
```

```
> sentiments
# A tibble: 23,165 \times 4
          word sentiment lexicon score
         <chr>>
                    <chr>>
                             <chr> <int>
        abacus
1
                    trust
                               nrc
       abandon
                     fear
                                       NA
                               nrc
       abandon negative
                               nrc
                                       NA
       abandon
                  sadness
                               nrc
                                       NA
     abandoned
                    anger
                                       NA
                               nrc
     abandoned
                     fear
                               nrc
                                       NA
     abandoned
                 negative
                               nrc
                                       NA
     abandoned
                  sadness
                               nrc
                                       NA
   abandonment
                    anger
                               nrc
                                       NA
10 abandonment
                     fear
                               nrc
                                       NA
  ... with 23,155 more rows
```

```
> stop_words
# A tibble: 1,151 × 2
          word lexicon
                  <chr>>
         <chr>>
1
                  SMART
           a's
                  SMART
          able
                  SMART
         about
                  SMART
         above
                  SMART
     according
                  SMART
  accordingly
                  SMART
        across
                  SMART
      actually
                  SMART
10
         after
                  SMART
# ... with 1,141 more rows
```

```
> library(tidytext)
> data("sentiments")
> sentiments
# A tibble: 27,314 x 4
          word sentiment lexicon score
         <chr>
                    <chr>
                            <chr> <int>
        abacus
 1
                    trust
                              nrc
                                      NA
                    fear
 2
       abandon
                              nrc
                                      NA
       abandon negative
 3
                              nrc
                                      NA
                 sadness
 4
       abandon
                              nrc
                                      NA
 5
     abandoned
                    anger
                              nrc
                                      NA
 6
     abandoned
                     fear
                              nrc
                                      NA
     abandoned negative
                              nrc
                                      NA
                 sadness
     abandoned
                              nrc
                                      NA
 9 abandonment
                    anger
                              nrc
                                      NA
10 abandonment
                     fear
                              nrc
                                      NA
# ... with 27,304 more rows
```

```
> stop_words
# A tibble: 1,149 x 2
          word lexicon
         <chr>>
                  <chr>>
                  SMART
            a's
                  SMART
          ab1e
 3
                  SMART
         about
                  SMART
 5
         above
                  SMART
     according
                  SMART
 7 accordingly
                  SMART
        across
                  SMART
      actually
                  SMART
10
         after
                  SMART
# ... with 1,139 more rows
```

Tidytext is part of the tidy universe including ggplot and dplyr. Code is structured so it is more easily read using the %>%. The data format is a tibble and is in "tidy" format (long form).

```
mtcars %>% group_by(cyl) %>% mutate(rank = min_rank(desc(mpg)))
```

This reads as "Using the mtcars object then group by the cyl vector then mutate a new variable called rank.



Bagaille

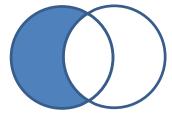
The pipe operator

%>%

Forwards an object so the code is easy to understand & concise.

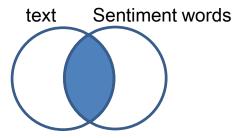
all.tidy <- all.tidy %>%
 anti_join(stop_words)

Tweet words Stop words



anti_join()

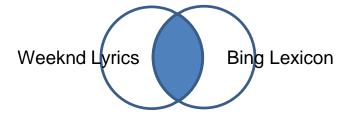
all.sentiment <- all.tidy %>%
 inner_join(nrc.lexicon) %>%
 count(tweet,sentiment) %>%
 spread(tweet, n, fill = 0)



inner_join()

```
# DTM
                                                        DTM is from the "tm" library
txtDTM<-DocumentTermMatrix(txtCorpus)</pre>
txtDTM
dim(txtDTM)
# Tidy
                                                        Easy way to make it into a tibble.
tidyCorp<-tidy(txtDTM)</pre>
tidyCorp
dim(tidyCorp)
# Get bing lexicon
# "afinn", "bing", "nrc", "loughran"
bing<-get_sentiments(lexicon = c("bing"))</pre>
head(bing)
# Perform Inner Join
                                                     Weeknd Lyrics
                                                                            Bing Lexicon
bingSent<-inner_join(tidyCorp,bing,</pre>
by=c('term'='word'))
```

The polarity function from qdap and the inner_join show similar negative results.



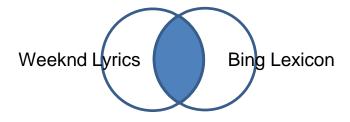
The sentiments tibble has multiple lexicons. Each can be used in an inner join to get different ways of assessing sentiment.

AFINN- Dutch researcher Words scored -5 to 5

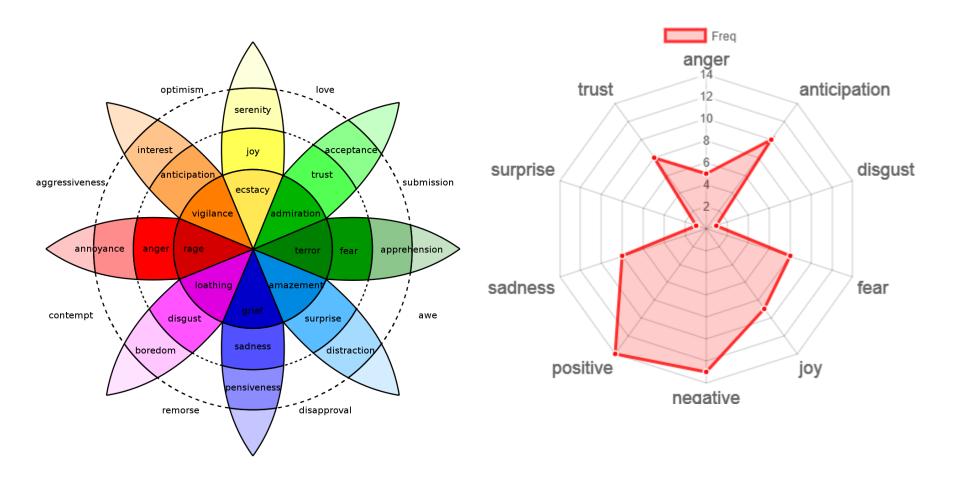
Bing- Uofl-Chi Researcher Words scored Pos/Neg

NRC – mTurk Words classified into 8 primary & pos/neg

```
> head(nrc)
# A tibble: 6 x 2
       word sentiment
      <chr>>
                 <chr>>
     abacus
                 trust
    abandon
                  fear
    abandon negative
    abandon
              sadness
5 abandoned
                 anger
6 abandoned
                 fear
```



Remember Plutchik's Wheel of emotion? Let's mimic it!



You could remove Pos/Neg before making the radarchart to more closely fit the wheel.

Intro to TM w/R Deep Dive

- •What is Text Mining?
- String Manipulation & Cleaning
- •Frequency Counts, Dendrograms & Wordclouds
- Sentiment Analysis
- Topic Modeling
- OpenNLP/Named Entity Recognition

Goals:

- Learn the basics of text mining
- Apply methods to real data and uses
- Follow a business workflow for stakeholders needing text mining insights

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

<u>on</u>

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

7_TopicModeling_Sentiment.R

LDA Topic Modeling

library(topicmodels) is not actively maintained. So, library(lda) is used here.

txtLex<-lexicalize(txt)</pre>

Reads raw text, and creates a corpus of n "articles", and complete word lexicon.

```
txtWordCount<-word.counts(txtLex$documents, txtLex$vocab)
txtDocLength<-document.lengths(txtLex$documents)</pre>
```

Word counts for each article Number of words in each article

Model Fit

7_TopicModeling_Sentiment.R

LDA Topic Modeling

One of the cool things about Ida is that you can find the prototypical document

```
# Prototypical Document
top.topic.documents(fit$document_sums,1)
```

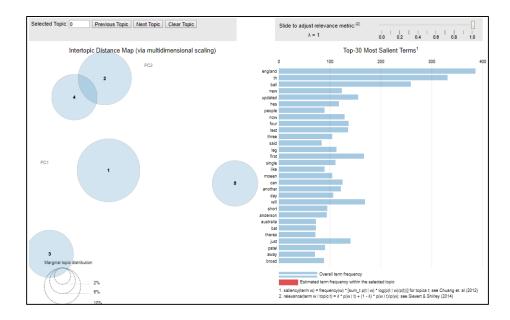
[1] 21 23 2 24 3

Since K= 5, the 21 article best represents the first topic and so on.

```
> text$webTitle[21]
[1] "A little shared knowledge goes a long way in helping to scale social enterprise"
```

7_TopicModeling_Sentiment.R

Better yet, let's explore in a dynamic gui w/ IdaVis



7_TopicModeling_Sentiment.R

TREEMAP: multi dimensional representation of the corpus attributes.





-1.5 -1.0 -0.5 0.0 0.5 1.0 1.5 polarity

- Each article is a small square
- The area of the square is related to the number of terms <u>document</u> <u>length</u>
- Color will represent quick, simple polarity from qdap
- The larger grouping is based on <u>LDA topic assignment</u>

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

7_TopicModeling_Sentiment.R

Intro to TM w/R Deep Dive

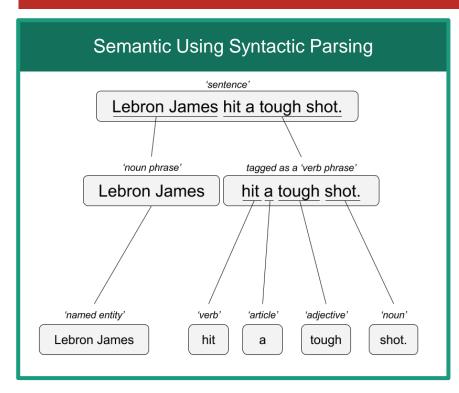
- •What is Text Mining?
- String Manipulation & Cleaning
- •Frequency Counts, Dendrograms & Wordclouds
- Sentiment Analysis
- -Topic Modeling
- OpenNLP/Named Entity Recognition

Goals:

- Learn the basics of text mining
- Apply methods to real data and uses
- Follow a business workflow for stakeholders needing text mining insights

Remember this? Text Mining Approaches

"Lebron James hit a tough shot."





Named Entity Recognition using OpenNLP



Welcome to Apache OpenNLP

The Apache OpenNLP library is a machine learning based toolkit for the processing of natural language text.





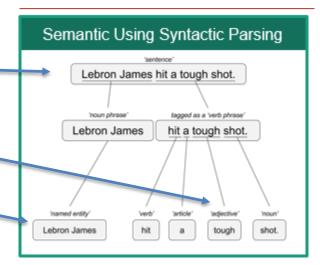
R's OpenNLP package wraps the Apache OpenNLP project. However, documentation & examples can be hard to come by.

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!

<u>Annotations</u>

- 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>
```

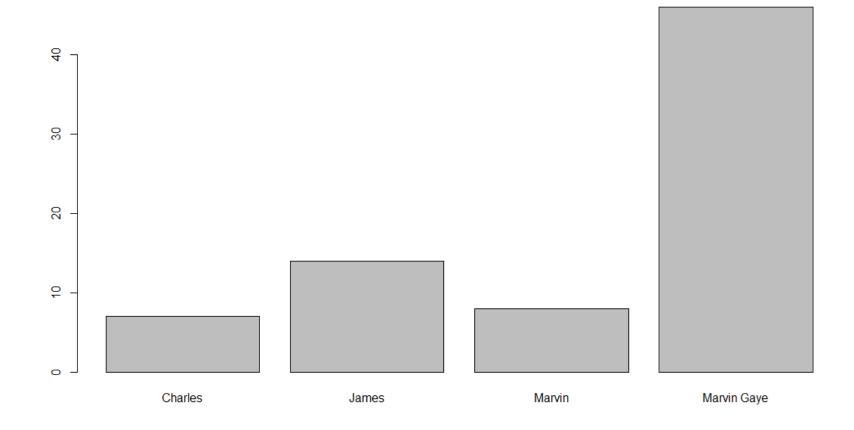
Annotate is used in ggplot2 SO you must either change the above to NLP::annotate(...) or ensure ggplot2 is not loaded in your environment!

```
#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"

```
> head(people)
[1] "Marvin Gaye" "Marvin" "LeedsVsSydney" "Marvin Gaye" "Marvin Gaye"
[6] "Marvin"
> peopleTable<-table(as.factor(people))
> topPeeps<-peopleTable[peopleTable>5]
> barplot(topPeeps)
```



Workshop Agenda

9-10:45 Intro to TM w/R 10:45-11 Break 11 – 11:45 Gather Text 11:45 – 12 Break 12-12:40 HTML Widgets & Flexdashboard

12:40-12:50 Optional Break

12:50-1 Scheduling

Goals:

- Learn the basics of text mining
- Apply methods to real data
- Follow a business workflow for stakeholders needing text mining insights

Review: read.csv()

To read a single "comma separated value" file (.csv):

```
text<-read.csv('coffee.csv', header=TRUE)</pre>
```

Part of utils package installed with base R. Used for tabled data with a "delimited" to create a data.table.

Example: 1_Keyword_Scanning.R

Some of the alternatives include:

Package::Function	Description
<pre>data.table::fread()</pre>	Fast And Friendly File Finagler
readr::read_csv()	Read A Delimited File (Including Csv & Tsv) Into A Tibble
ff::read.csv.ffdf()	"Imports csv files into ff data frames (memory- efficient storage of large data on disk and fast access functions)
<pre>sqldf::read.csv.sql()</pre>	Read File Filtered By SQL statement
<pre>bigmemory::read.big.matrix()</pre>	Package for working with "massive matrices" File Interface For A ``Big.Matrix' class

Review: read.csv() with Multiple Files*

```
Start with a vector of .csv files OR you can use txtFiles<-c('chardonnay.csv','coffee.csv','beer.csv')

OR if you have a folder with many txtFiles<-list.files(pattern = "*.csv")
```

Example: 5_Other_Wordclouds.R

A For loop will read each file individually and add it to the environment.

```
for (i in 1:length(txtFiles)){
  assign(txtFiles[i], read.csv(txtFiles[i]))
  cat(paste('read completed:',txtFiles[i],'\n'))
}
```

Using lapply or pblapply (personal preference) will put all individual files into a single list all<-pblapply(txtFiles,read.csv)

rbindlist() from data.table is an efficient way to collapse the list of documents if needed.

Plain Text files: readLines()



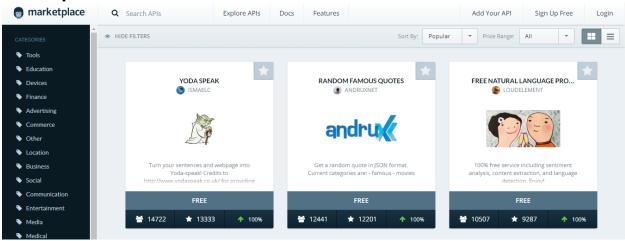
The Weeknd - Starboy (official) ft. Daft Punk

```
# Data
text<-readLines('Weeknd.txt')</pre>
```

Example: 6_TidyText_Sentiment.R

Now onto APIs!!

Mashape.com



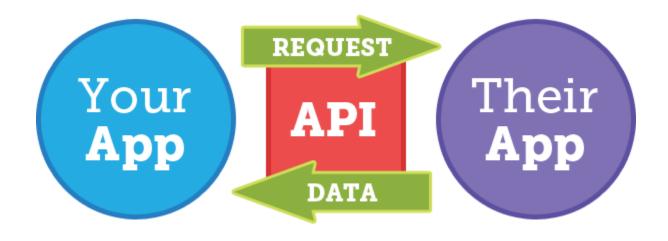
ProgrammableWeb.com



Plus literally millions of others waiting to be found!

Getting Text from an API: What is an API?

Application Program Interface (API)
Clearly defined methods of communication between various software components.

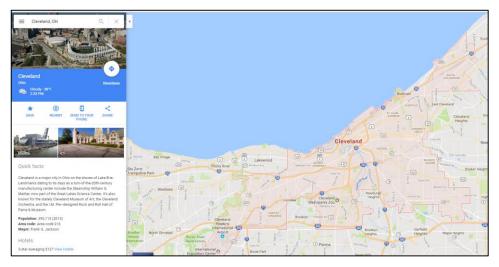


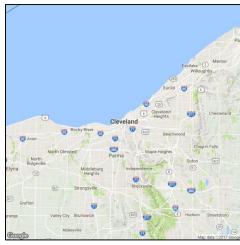
APIs are behind many of the sites you use today.



Google Maps API Services include

- Geocoding (get lat/lon from name & vice versa)
- Base map "tiles"
- Basic Geographic Info

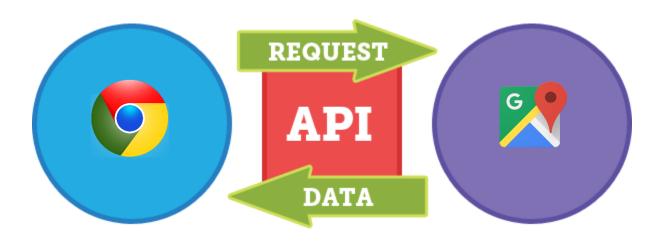




www.google.com/maps/place/Cleveland,+OH/@41.4951143,-81.8462865,11z maps.googleapis.com/maps/api/staticmap?center=cleveland,+oh&zoom=10&size=640x640&scale=2&maptype=terrain maps.googleapis.com/maps/api/geocode/json?address=cleveland&sensor=false

In our diagram, the browser/maps interface...

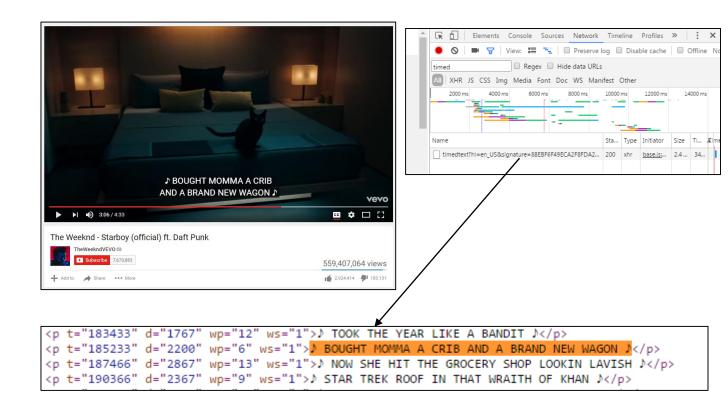
"Your Chrome browser interacts with the maps API services to render the information.



APIs can help unlock information...



Need song lyrics? Access the Closed Caption API file called "timedtext.



Two popular formats.

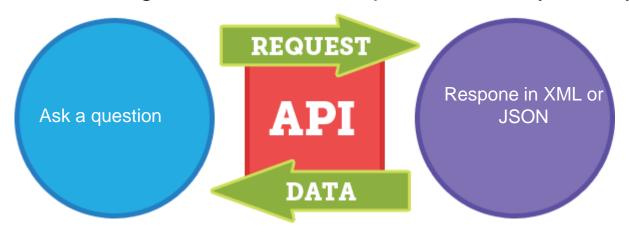
XML – Extensible Markup Language

Extensible Markup Language is a markup language that defines a set of rules for encoding documents in a format that is both human-readable and machine-readable.

Ever wonder why Excel files went from xls to xlsx? The data is stored as XML.

JSON- Javascript Object Notation

is an open-standard file format that uses human-readable text to transmit data objects consisting of attribute—value pairs and array data types.



Comparing the two formats

This XML file does not appear to have any style information associated with it.

```
▼ <GeocodeResponse>
   <status>OK</status>
 ▼<result>
    <type>locality</type>
    <type>political</type>
    <formatted address>Boston, MA, USA</formatted address>
   ▼<address component>
      <long name>Boston</long name>
      <short name>Boston</short name>
      <type>locality</type>
      <type>political</type>
    </address component>
   ▼<address component>
      <long name>Suffolk County</long name>
      <short name>Suffolk County</short name>
      <type>administrative area level 2</type>
      <type>political</type>
    </address component>
   ▼<address component>
      <long_name>Massachusetts</long_name>
      <short name>MA</short name>
      <type>administrative area level 1</type>
      <type>political</type>
    </address component>
   ▼<address component>
```

http://maps.googleapis.com/maps/api/geocode/**xml**?addre ss=boston&sensor=false

```
"results" : [
      "address_components" : [
            "long_name" : "Boston",
            "short name" : "Boston",
            "types" : [ "locality", "political" ]
            "long name" : "Suffolk County",
            "short_name" : "Suffolk County",
            "types" : [ "administrative_area_level_2", "political" ]
            "long_name" : "Massachusetts",
            "short name" : "MA",
            "types" : [ "administrative area level 1", "political" ]
            "long_name" : "United States",
            "short name" : "US",
            "types" : [ "country", "political" ]
      "formatted address" : "Boston, MA, USA",
      "geometry" : {
         "bounds" : {
            "northeast" : {
              "lat": 42.40081989999999,
               "lng" : -70.749455
            "southwest" : {
               "lat": 42.22788,
               "lng": -71.191113
         "location" : {
            "lat": 42.3600825,
            "lng" : -71.0588801
         "location_type" : "APPROXIMATE",
         "viewport" : {
            "northeast" : {
              "lat" : 42.3988669,
               "lng": -70.9232011
```

http://maps.googleapis.com/maps/api/geocode/**json**?address=boston&sensor=false

Types of API requests

HTTP GET - Use GET requests to retrieve resource representation/information only e.g. your browsers GETs google map information

HTTP POST/PUT - Use POST APIs to create new subordinate resources e.g. you post a file to kaggle, it is scored and the response is a change in your ranking

HTTP DELETE - DELETE APIs are used to delete resources

HTTP PATCH – PATCH requests are to make partial update on a resource e.g. updating an existing customer record

Generally, a response in the 200s, is a success and anything in the 400s is considered a failure for various reasons.

Let's try an API

Request Helpers:

www.Hurl.it

PostMan Chrome Add-in or software

Simple API interactions

- https://newsapi.org/
- https://pipl.com/dev/demo
- http://maps.googleapis.com/maps/api/staticmap?center= Boston+MA&zoom=10&size=640x640

4/25/2018

Files to Open

XML Example

https://github.com/kwartler/ODSC/blob/master/p2_Get_Text/GetXML.R

GetXML.R

JSON Example

https://github.com/kwartler/ODSC/blob/master/p2_Get_Text/instruction_newsAPI.R

instruction_newsAPI.R

Workshop Agenda

9-10:45 Intro to TM w/R 10:45-11 Break 11 – 11:45 Gather Text 11:45 – 12 Break 12-12:40 HTML Widgets & Flexdashboard 12:40-12:50 Optional Break

12:50-1 Scheduling

Goals:

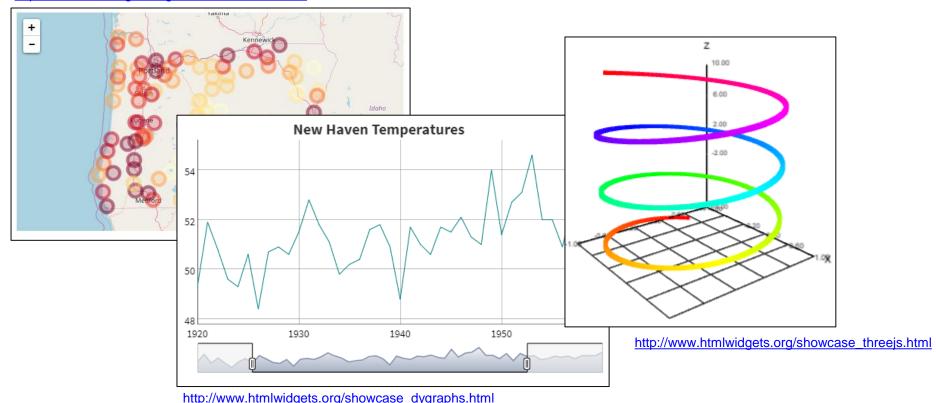
- Learn the basics of text mining
- Apply methods to real data
- Follow a business workflow for stakeholders needing text mining insights

HTML Widgets www.HTMLWidgets.org

Let's an R programmer use JavaScript libraries!

- Dynamic meaning a user can interact with the plot
- Easy to do from within the R console without JS knowledge
- "ships" with all the underlying data, so basic exploration/plotting doesn't need a Shiny or OCPU backend.

http://www.htmlwidgets.org/showcase_leaflet.html



instruction_htmlWidgets.R

FlexDashboard

http://rmarkdown.rstudio.com/flexdashboard/

Let's an R programmer use HTML with Javascript!

- Embed the htmlwidgets, or other plots
- Easy to do from within the R console without HTML knowledge
- Easily change the layout of the dashboard.
- Authoring is done in Rmarkdown (.Rmd) files by "knitting" the components
 & output is HTML for easy viewing







Open these files:

blank_dashboard.Rmd top_focal.Rmd

Workshop Agenda

9-10:45 Intro to TM w/R 10:45-11 Break 11 – 11:45 Gather Text 11:45 – 12 Break 12-12:40 HTML Widgets & Flexdashboard 12:40-12:50 Optional Break

12:50-1 Scheduling

Goals:

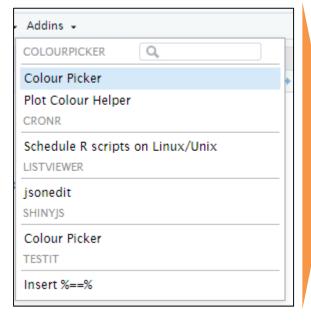
- Learn the basics of text mining
- Apply methods to real data
- Follow a business workflow for stakeholders needing text mining insights

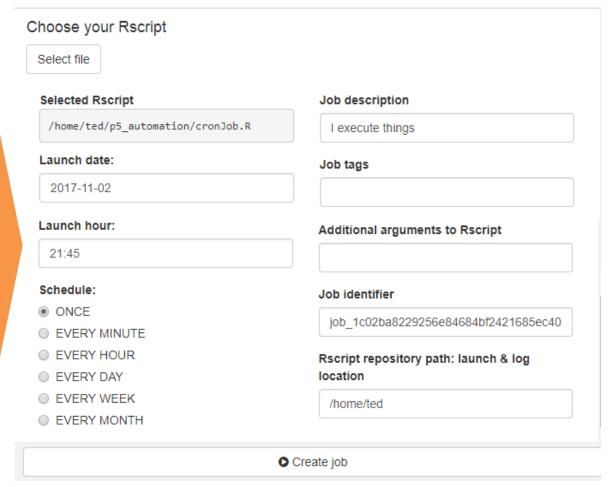
Open these files:

cronDash.Rmd cronJob.R

Use a scheduler to automate the construction & dissemination of your work!

R Studio cronR Add-In





Questions?

https://github.com/kwartler/ODSC



www.linkedin.com/in/edwardkwartler