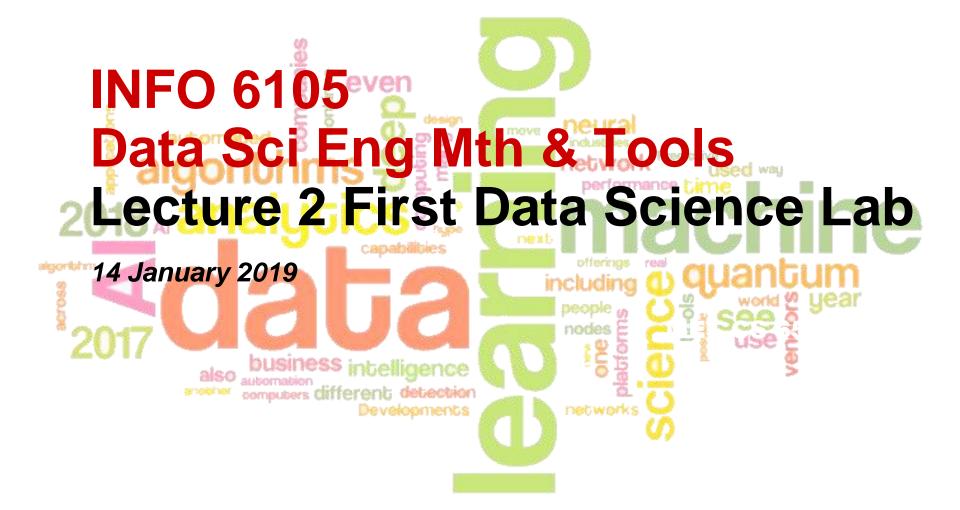


Northeastern University







What is Data Science



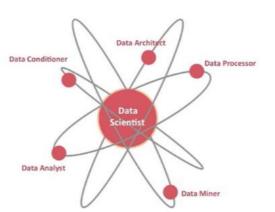
- Using automated methods to analyze sometimes massive amounts of data for the purpose of extracting knowledge learned from derived results
- Extracting resulting insights used to understand and improve:
 - Business
 - Investment and finance
 - Research
 - Health and lifestyle
 - etc

Data Science skillset



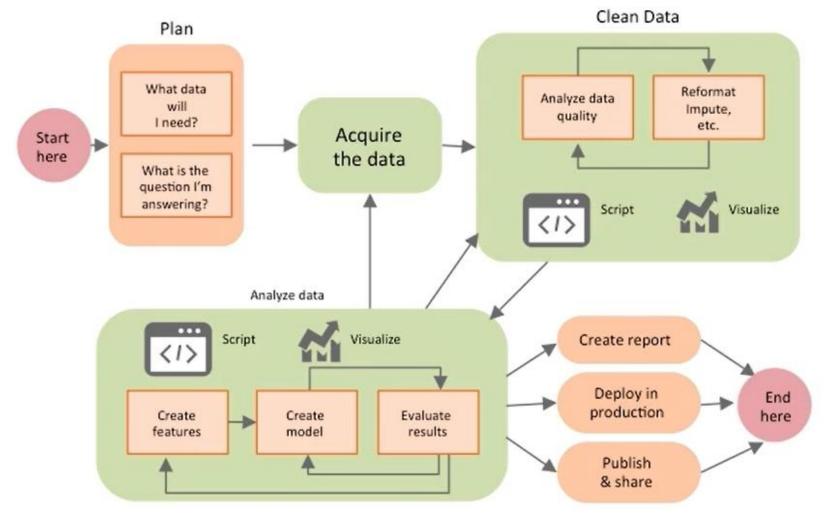
- Statistician
 - Individual who models and summarizes datasets
- Domain expertise
 - Subject matter expert
- Mathematician
 - Linear Algebra the lingua franca of Deep Learning
- Computer scientist

Writes algorithms to summarize, visualize, and store data



Data Science Process Flow

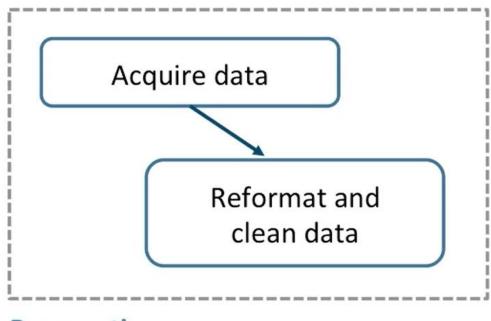




Stage 1 – Planning/Preparation



- Obtain data
- Format and clean data

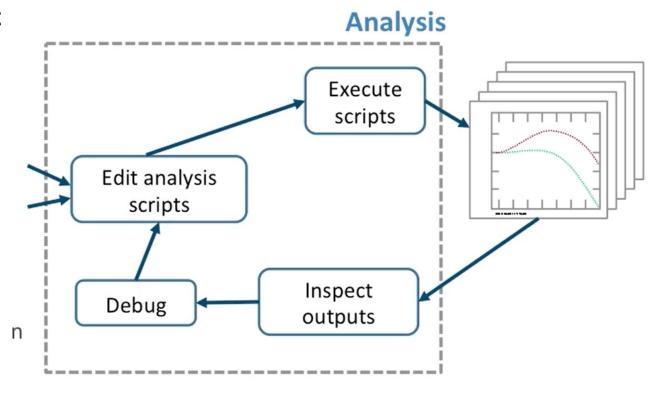


Preparation

Stage 2 - Analysis

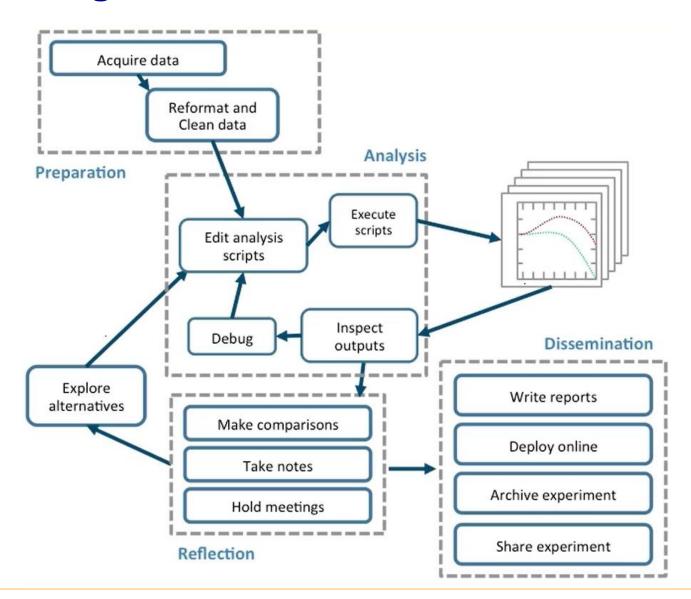


- Create scripts for analysis
- Run scripts
- Debug
- Inspect output
- Debug
- Repeat

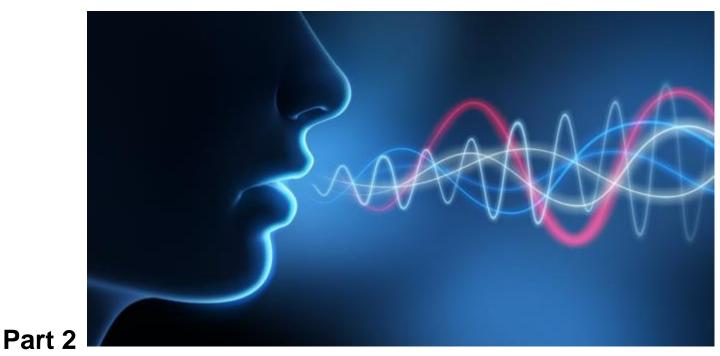


Stage 3 - Publish









NATURAL LANGUAGE PROCESSING (NLP)

What is NLP?

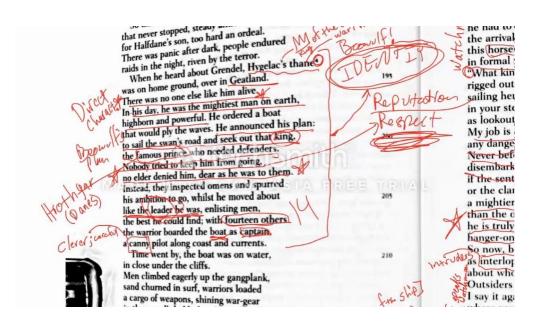


- Natural Language Processing (NLP) is a field of computer science and engineering that has developed from the study of language and computational linguistics within the field of Artificial Intelligence
- The goals of NLP are to design and build applications that facilitate human interaction with machines through the use of natural language
 - Document classification, machine translation, and speech recognition are major areas of NLP
- 2018 was a milestone year for NLP and machines, because machines achieved human parity in translation (all languages) and speech to text
 - In other words, machine perform as well as the best human experts in the field

Annotations



- It is not enough to simply provide a computer with a large amount of data and expect it to learn to speak—the data has to be prepared in such a way that the computer can more easily find patterns and inferences
- This is usually done by adding relevant metadata to a dataset. Any metadata tag used to mark up elements of the dataset is called an annotation over the input



What kinds of Annotations?



- Grammar is the name typically given to the mechanisms responsible for creating well-formed structures in language
- Most linguists view grammar as itself consisting of distinct modules or systems, either by cognitive design or for descriptive convenience These areas usually include:

Syntax

The study of how words are combined to form sentences. This
includes examining parts of speech and how they combine to
make larger constructions ("I love you", "Wo ai ni", etc..)

Semantics

 The study of meaning in language. Semantics examines the relations between words and what they are being used to represent.

Morphology

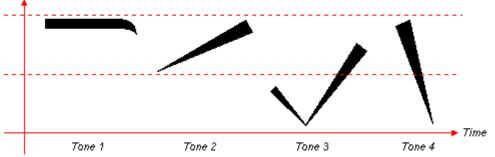
The study of units of meaning in a language. A morpheme is the smallest unit of language that has meaning or function, a definition that includes words, prefixes, affixes, and other word structures that impart meaning. In Chinese, it's an ideogram (fire)

Grammar areas (continued)



Phonology

- The study of the sound patterns of a particular language
- Aspects of study include determining which patterns are significant and have meaning (i.e., the phonemes); how syllables are structured and combined; and what features are needed to describe the discrete units (segments) in the language, and how they are interpreted
- In Chinese, it includes the 4 intonations:



Phonetics

 The study of the sounds of human speech, and how they are made and perceived. A phoneme is the term for an individual sound, and is essentially the smallest unit of human speech.

Grammar areas (continued)



Lexicon

 The study of the words and phrases used in a language, that is, a language's vocabulary.

Discourse analysis

 The study of exchanges of information, usually in the form of conversations, and particularly the flow of information across sentence boundaries.

Pragmatics

 The study of how the context of text affects the meaning of an expression, and what information is necessary to infer a hidden or presupposed meaning.

Text structure analysis

 The study of how narratives and other textual styles are constructed to make larger textual compositions

Reference: adapted from Natural Language Annotation for Machine Learning

https://www.safaribooksonline.com/library/view/natural-language-annotation/9781449332693/

Parts of Speech



- In traditional grammar, a part of speech (abbreviated as PoS or POS) is a category of words that have similar grammatical properties
- Words that are assigned to the same POS generally display similar behavior in terms of syntax—they play similar roles within the grammatical structure of sentences—and sometimes in terms of morphology, in that they undergo inflection for similar properties

Commonly listed <u>English</u> parts of speech are <u>noun</u>, <u>verb</u>, <u>adjective</u>, <u>adverb</u>, <u>pronoun</u>, <u>preposition</u>, <u>conjunction</u>, and <u>interjection</u>



Noun Adjective Verb Conjunction

Preposition Article Interjection Adverb

Parts of Speech



PARTS OF SPEECH IN ENGLISH (WORD CLASSES)

Parts of Speech



NOUN

Name of a thing, a person, an animal, a place, or an idea.

Examples: Daniel, London, table, hope - Mary uses a blue pen for her letters.

ADJECTIVE

Describes, modifies or gives more information about a noun or pronoun. Examples: cold, happy, young, two, fun - The little girl has a pink hat.

ADVFRB

Modifies a verb, an adjective or another adverb. It tells how (often), where, when. Examples: slowly, very, always, well, too-Yesterday, I ate my lunch quickly.

CONJUNCTION

Joins two words, ideas, phrases together and shows how they are connected.

Examples: and, or, but, because, until, if
- I was hot and tired but I still finished it.

PRONOUN

A pronoun is used in place of a noun or noun phrase to avoid repetition.

Examples: I, you, it, we, us, them, those
- I want her to dance with me.

VERB

Shows an action or a state of being. It can show what someone is doing or did.

Examples: go, speaking, lived, been, is

- I listen to the word and then repeat it.

PREPOSITION

Shows the relationship of a noun, noun phrase or pronoun to another word.

Examples: at, on, in, from, with, about - I left my keys on the table for you.

INTERJECTION

A word or phrase that expresses a strong emotion. It is a short exclamation.

Examples: Ouch! Hey! Wow! Oh! Ugh! - Wow! I passed my English exam.

vw.grammar.cl www.woodwardenglish.com www.vocabularv.cl

PARTS OF SPEECH

In traditional English grammar, a part of speech is a category of words that have similar grammatical properties. Parts of speech tell us how a word is used in a sentence.

Sometimes parts of speech

are called word classes.

PARTS OF SPEECH NOUN Naming word PRONOUN Replaces a noun ADJECTIVE Describes something VERB Is an action or state ADVERB Describes a werb PREPOSITION Shows a relationship CONJUNCTION Joins words or clauses INTERJECTION Expresses emotions

PARTS OF SPEECH

ADJECTIVE

Describes, modifies or gives more information about a noun or pronoun.

Examples: cold, happy, young

- The little girl has a pink hat.

girl little * hat pink*

PARTS OF SPEECH

PRONOUN
A pronoun is used in place of a noun or noun phrase to avoid repetition.
Examples: mine, yours, hers

- This bike is my blke.
- This bike is mine. Repetition

PARTS OF SPEECH

ARTICLES

A / AN: to talk about a noun in general.

THE: to talk about a specific noun.

Examples: a, an, the

- I need a dictionary.

- The dictionary needs to be in English.

The Penn Treebank



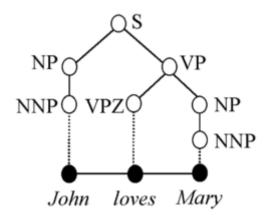
- In linguistics, a treebank is a parsed text corpus that annotates syntactic or semantic sentence structure
 - Syntax (from ancient greek: σύνταξις "coordination") is the set of rules, principles, and processes that govern the structure of sentences in a given language, usually including word order
 - Semantics (from ancient greek: σημαντικός, "significant") is the linguistic and philosophical study of meaning, in language. It is concerned with the relationship between signifiers—like words, phrases, signs, and symbols—and what they stand for, their denotation
- The most famous treebank for the English language is the Penn Treebank, which in its eight years of operation (1989-1996), produced approximately 7 million words of part-ofspeech tagged text, 3 million words of skeletally parsed text, and over 2 million words of text parsed for structure
 - Material annotated includes such wide-ranging genres as IBM computer manuals, nursing notes, Wall Street Journal articles, and transcribed telephone conversations
- https://catalog.ldc.upenn.edu/ldc99t42

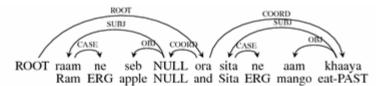
Languages have their own treebanks



There are treebanks for each language of the world

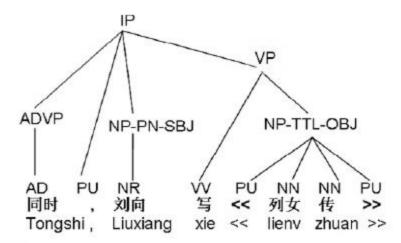
– What is yours?

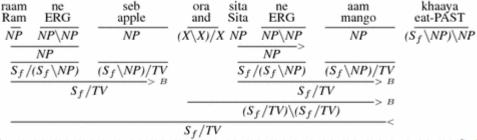




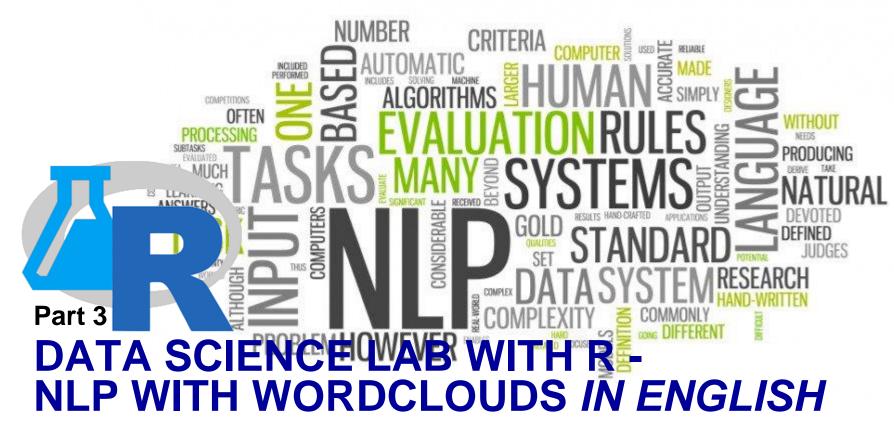
bracketed sentence:

(IP (ADVP (AD 同时)) (PU ,) (NP-PN-SBJ (NR 刘向)) (VP (VV 写) (NP-TTL-OBJ (PU 《) (NN 列女) (NN 传) (PU 》))))
English: At the same time, Liuxiang wrote <<The Biography of Strong-minded Women>>









Document-Term and Term-Document Matrices



- A document-term matrix or term-document matrix is a mathematical matrix that describes the frequency of terms that occur in a collection of documents
- In a document-term matrix, rows correspond to documents in the collection and columns correspond to terms
- In a term-document matrix, rows correspond to terms, and columns to documents
- There are various schemes for determining the value that each cell in the matrix should take
 - In our lab and for now, we are just going to count occurrences

tf-idf



Stands for Term Frequency – Inverse Document Frequency

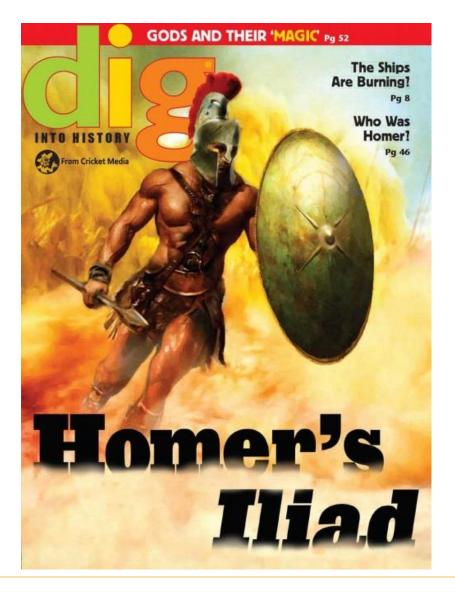
 A numerical statistic that is intended to reflect how important a word is to a document in a collection of documents (also called a *corpus*)

The tf-idf value:

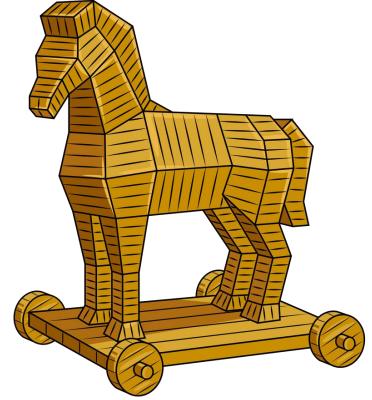
- Increases proportionally to the number of times a word appears in the document
- Offset by the frequency of the word in the corpus, which helps to adjust for the fact that some words appear more frequently in general
- So, if "Troy" appears very often in the document "Iliad", then "Troy" is a good descriptor for the document
- ...But so is the word "the"...
- However, the word "the" appears in all documents a lot
- And so its tf-idf should be low, while that of "Troy" high...
- Variations of the tf-idf weighting scheme are often used by <u>search engines</u> as a central tool in scoring and ranking a document's relevance given a user query

Homer's iliad



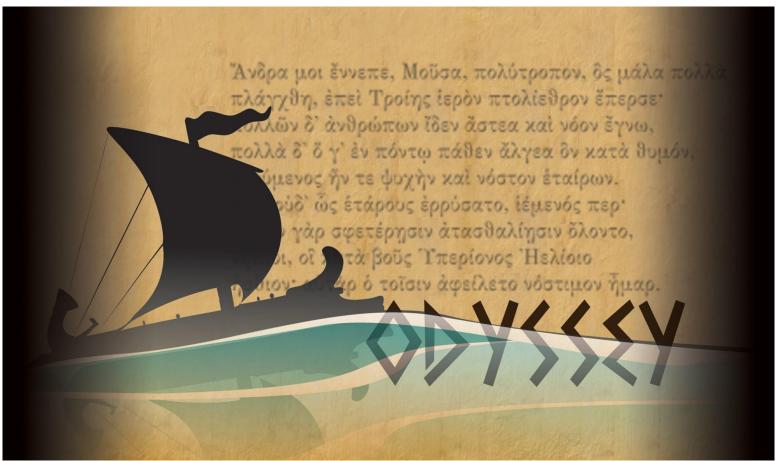


- A War in antiquity
 - Athens vs Troy
 - And lessons in humanity



Odyssey





Ulysses' return trip to Ithaca

R packages for text mining



- Our main text mining package:
 - tm (text mining)
 - https://cran.r-project.org/web/packages/tm/index.html
- Other packages we are going to use:
 - SnowballCC, RColorBrewer, ggplot2, wordcloud, biclust, cluster, igraph, fpc
- Package references on CRAN:
 - https://cran.rproject.org/web/views/NaturalLanguageProcessing.html

Loading texts



- Start by saving your text files in a folder titled texts This will be the "corpus" (body) of texts you will be mining
 - You are going to work with Homer's Iliad and the Odyssey, some of the oldest literature about Wars
 - You are going to also to work in your own language, so get started in English for now, but plan ahead..
- Download corpus from blackboard
- Mac:
 - cname <- file.path("~", "Desktop", "texts")</pre>
 - cname
- Windows:
 - cname <- file.path("C:", "texts")</pre>
 - cname
- Read your documents in the R terminal with inspect (docs)
- if you prefer to look at only one of the documents you loaded, then you can specify inspect (docs[1])

Preprocessing



- Remove numbers, capitalization, common words, punctuation, and prepare your texts for analysis
 - Remove punctuation, numbers, stopwords, common word endings, strip whitespace, convert to lower case, and combine words that should stay together
- Be sure to use the following script once you have completed preprocessing

This tells R to treat your preprocessed documents as text documents

- library(tm)
- docs <- Corpus(DirSource(cname))</pre>
- docs <- tm_map(docs, PlainTextDocument)</pre>

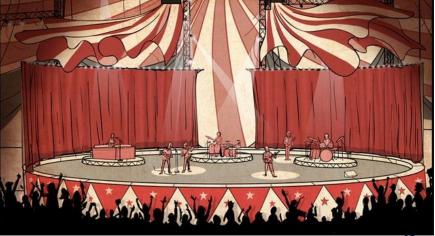


Data staging



Create a document term matrix and a term document matrix

```
> dtm <- DocumentTermMatrix(docs)</pre>
> dtm
<<DocumentTermMatrix (documents: 2, terms: 11050)>>
Non-/sparse entries: 14535/7565
Sparsity
Maximal term length: 35
Weighting
                    : term frequency (tf)
> tdm <- TermDocumentMatrix(docs)</pre>
> tdm
<<TermDocumentMatrix (terms: 11050, documents: 2)>>
Non-/sparse entries: 14535/7565
                    : 34%
Sparsity
Maximal term length: 35
Weighting
                    : term frequency (tf)
```

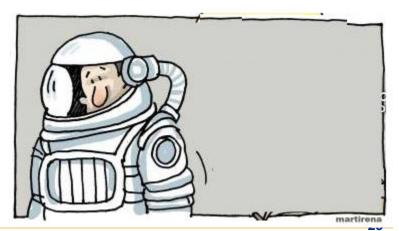


Data Exploration



- Organize terms by their frequency
 - freq <- colSums(as.matrix(dtm))</pre>
 - length(freq)
 - ord <- order(freq)</pre>
- Remove sparse terms
 - dtms <- removeSparseTerms (dtm, 0.1) # This makes a matrix that is 10% empty space, maximum
 - inspect(dtms)

 Docs god now one said shall son thi thus ulyss will content 626 573 348 211 511 420 943 621 123 188 meta 267 303 573 483 217 336 3 132 646 745
- Least frequently occurring words
 - freq[head(ord)]
- Most frequently occurring words
 - freq[head(ord)]



Frequency of frequencies



- Output is two rows of numbers, one per document
 - The top number is the frequency with which words appear
 - The bottom number reflects how many words appear that frequently
 - How many terms appear only once?
- Lowest word frequences:
 - head(table(freq), 20)
 - considering only the 20 *lowest* word frequencies, we can see that ??? terms appear only once
- Highest word frequencies:
 - tail(table(freq), 20)
 - Considering only the 20 greatest frequencies, there's usually a big disparity in how frequently some terms appear

Coarse-grained term frequencies



- For a less fine-grained look at term frequency, view a table of sparse terms:
 - freq <- colSums(as.matrix(dtms))</pre>
 - freq
 - > freq <- colSums(as.matrix(dtms))</pre> > freq abhor abandon abound abrida abroad absenc absent abus accident accommod accompani accomplish accept access accord achill

Alternatively:

- freq <- sort(colSums(as.matrix(dtm)), decreasing=TRUE)</pre>
- head(freq, 40)

head(freq, 40)													
thi	will	one	god	now	ulyss	son	thus	shall	said	great	hand	arm	jove
946	933	921	893	876	769	756	753	728	694	686	653	625	617
man	heaven	day	may	men	let	see	come	hector	can	ship	like	war	achill
595	567	534	513	499	496	489	469	461	453	450	447	440	436
hous	troy	yet	chief	first	father	greek	eye	upon	trojan	oer	make		
429	425	420	417	414	408	395	391	391	368	364	363		

Coarse-grained term frequencies



- All terms that appear frequently (250 or more times):
- findFreqTerms(dtm, lowfreq=250)

```
> findFreqTerms(dtm,lowfreq=250)
                                                                "can"
     "achill"
                    "arm"
                                   "back"
                                                 "came"
                                                                               "care"
                                                                                              "chief"
 [1]
 Г81
     "come"
                    "day"
                                   "dead"
                                                  "death"
                                                                "even"
                                                                               "everi"
                                                                                              "eve"
                    "fate"
                                                 "field"
                                                                "fiaht"
                                                                               "fire"
     "fall"
                                   "father"
                                                                                              "first"
                                                                "god"
                                                                               "good"
                                                                                             "great"
     "fli"
                    "forc"
                                   "friend"
                                                 "give"
                                   "hand"
                                                                "heart"
                                                                               "heaven"
                                                                                             "hector"
      "areek"
                    "around"
                                                 "head"
                    "high"
                                   "home"
                                                                "jove"
                                                                               "kina"
     "hero"
                                                 "hous"
                                                                                              "know"
                                                                "made"
[43]
     "let"
                    "lie"
                                   "like"
                                                 "lona"
                                                                               "make"
                                                                                             "man"
                    "men"
                                   "much"
                                                                                             "old"
F501
     "may"
                                                  "must"
                                                                "now"
                                                                               "oer"
                    "place"
                                   "plain"
                                                 "power"
                                                                "rage"
                                                                               "round"
                                                                                             "said"
      "one"
     "say"
                    "sea"
                                   "see"
                                                 "sha11"
                                                                               "son"
                                                                "ship"
                                                                                              "stand"
                                                                "tell"
                                                                               "thi"
                                                                                             "thou"
     "still"
                    "suitor"
                                   "take"
                                                 "telemachus"
     "though"
                    "thus"
                                   "till"
                                                  "time"
                                                                "trojan"
                                                                               "troy"
                                                                                              "two"
[85] "ulyss"
                    "upon"
                                   "vain"
                                                  "war"
                                                                "way"
                                                                               "went"
                                                                                              "whose"
[92] "will"
                    "word"
                                   "work"
                                                 "yet"
```

Alternatively:

wf <- data.frame(word=names(freq), freq=freq)</pre>

```
> head(wf)
- head(wf)
                           word freq
                   thi
                            thi
                                 946
                   will
                           will
                                 933
                                 921
                   one
                            one
                                 893
                            god
                   god
                                 876
                            now
                   now
                   ulyss ulyss
                                 769
```

32

Plotting with ggplot2



library(ggplot2)

```
p <- ggplot(subset(wf, freq>50), aes(word, freq))
p <- p + geom_bar(stat="identity")
p <- p + theme(axis.text.x=element_text(angle=45, hjust =1))
p</pre>
```

Fixing potential ggplot2 error

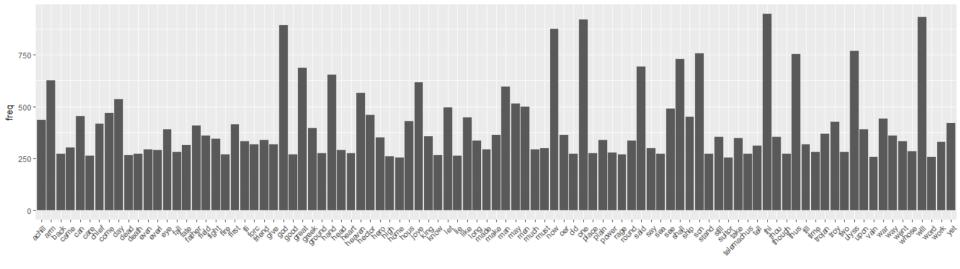


- □ library(ggplot2)
 - Error in loadNamespace(i, c(lib.loc, .libPaths()), versionCheck = vl[[i]]): namespace 'scales' 0.4.0 is being loaded, but >= 0.4.1 is required
 Error: package or namespace load failed for 'ggplot2'
- Then look at the complete error message. It lists a couple of packages which are missing or have errors. Uninstall and reinstall them. So for me:

```
remove.packages("ggplot2")
install.packages('ggplot2', dependencies = TRUE)
remove.packages("scales")
install.packages('scales', dependencies = TRUE)
library(ggplot2)
```

Words with frequencies > 250 in iliad + odyssey





Relationships between terms



- If you have a term in mind that you have found to be particularly meaningful to your analysis, then you may find it helpful to identify the words that most highly correlate with that term
 - If words always appear together, then correlation=1
 - library(tm)
 - findAssocs(dtm, "achilles", corlimit=0.98) # specifying a correlation limit of 0.98
- Can also find associations for word collections:
 - findAssocs(dtm, c("achilles", "hector"), corlimit=0.98)

Word Clouds



- Humans are generally strong at visual analytics
 - We are going to construct word clouds on our document corpus
- Load the package that makes word clouds in R
 - library (wordcloud)
- Plot words that occur at least 250 times
 - set.seed(142) #The "set.seed() function just makes the configuration of the layout of the clouds consistent each time you plot them. You can omit that part if you are not concerned with preserving a particular layout back son heaven place

Jack Sonheaven place ship heart jove god ship much day be tillground everi heart yet fight hand ship went tellold skinow fire achill hous the still be still be sea word father Ulyss stand forc though hector high per round said mayon

With Color



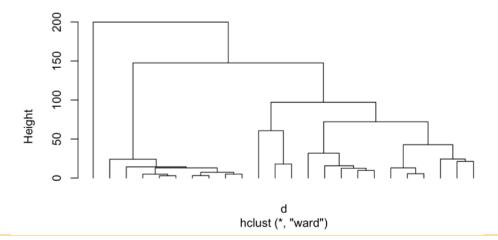
set.seed(142)
 dark2 <- brewer.pal(6, "Dark2")
 wordcloud(names(freq), freq, max.words=250, rot.per=0.2, c
 olors=dark2)



Clustering by term similarity



- First calculate distance between words & then cluster them according to similarity:
 - library(cluster)
 - dtmss <- removeSparseTerms(dtm, 0.15) # This makes a matrix that is only 15% empty space, maximum
 - inspect(dtmss)
 - d <- dist(t(dtmss), method="euclidian")</pre>
 - fit <- hclust(d=d, method="ward.D2")</pre>
 - plot(fit, hang=-1)



Ask R to help identify the clusters



- plot.new()
- □ plot(fit, hang=-1)
- groups <- cutree(fit, k=5) # "k=" defines the nu mber of clusters we are using
- rect.hclust(fit, k=5, border="red") # draw dendogr am with red borders around the 5 clusters

K-Means clustering



- Cluster words into a specified number of groups (in this case 2), such that the sum of squared distances between individual words and one of the group centers is minimized
 - library(fpc)
 - d <- dist(t(dtmss), method="euclidian")</pre>
 - kfit <- kmeans(d, 2)</pre>
 - clusplot(as.matrix(d), kfit\$cluster, color=T, shade=T, l
 abels=2, lines=0)



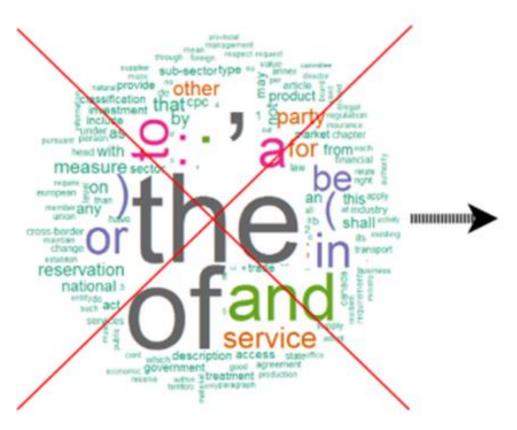


INTRODUCTION TO ADVANCED NLP: WORDCLOUDS IN YOUR LANGUAGE & FOR YOUR DATASET, WITH UDPIPE

Part 4

Better word clouds with better data science





face court road crash road toll to n koreawall stworld cup deposit guarantee child porn bus crash mental health media call st george Car crash credit crisis attack man police searchact election crisis talks missing man police probe house blaze water restrictions

Versions



- This lab may require the *latest* versions!
 - Version 3.5.0 of R runtime and 1.1.447 of RStudio
- This means you may not be able to run it with Anaconda Navigator's RStudio
 - You may have to download the latest R runtime and RStudio by hand..

Not sure, but keep that in mind.



Goal



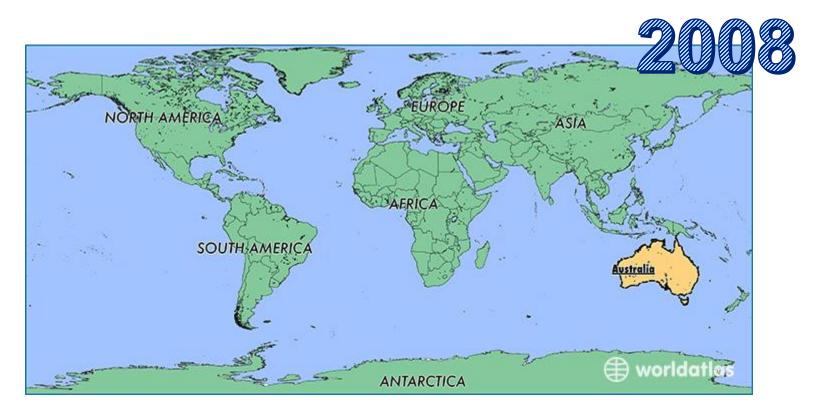
- You are going to do NLP...
 - in your language
 - on a topic of your choice
 - on a dataset you find (not the one here..)



Ok, so what are we doing here?



A friend told you he dreamt about the years 2008 and Australia, and he knows you are a Data scientist, so he asks you to tell him what happened in 2008 that was noteworthy, in Australia..



So you go look for datasets...



- You locate some news headlines published over a period of 15 years...
 - From Australian news source ABC (Australian Broadcasting Corp.)
 - At: http://www.abc.net.au/
- The Dataset is also available on Kaggle:
 - https://www.kaggle.com/therohk/million-headlines

Better use the version I put up on blackboard, instead (more data)

Package udpipe



- Then you go looking for data science libraries
 - You're still learning Python in class, so...
 - You go looking for R libraries on CRAN...

You find.. udpipe provides pretrained language models for spoken languages (not programming languages) and you can download the required model using udpipe download model()



Cleaning dataset



Read in the data

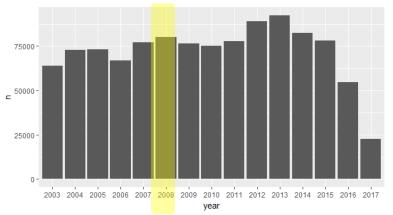
- news <- read.csv('abcnews-date-text.csv', header =
T, stringsAsFactors = F)</pre>

Do a quick exploration:

- news %>% group_by(publish_date) %>% count() %>%
 arrange(desc(n))
- What is %>%? Go to page 37 for an explanation

□ To subset data only for 2008:

- news_more_2008 <- news_more %>% filter(year
2008 & month == 10)

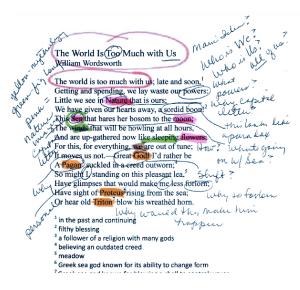




Staging with pretrained language models



- library(udpipe)
- model <- udpipe_download_model(language =
 "english")</pre>
- udmodel_english <- udpipe_load_model(file =
 'english-ud-2.0-170801.udpipe')</pre>
- s <- udpipe_annotate(udmodel_english, news_more_2008\$headline_text)
- x <- data.frame(s)</pre>



Package lattice



- Powerful and elegant high-level data visualization system inspired by Trellis graphics, with an emphasis on multivariate data
 - Sufficient for typical graphics needs, and flexible enough to handle most nonstandard requirements
 - Nice barcharts ©
 - barchart(key ~ freq, data = head(stats, 20), col =
 "magenta", main = "Keywords-simple noun phrases",
 xlab = "Frequency")
 - https://cran.r-project.org/web/packages/lattice/index.html

Factors in R

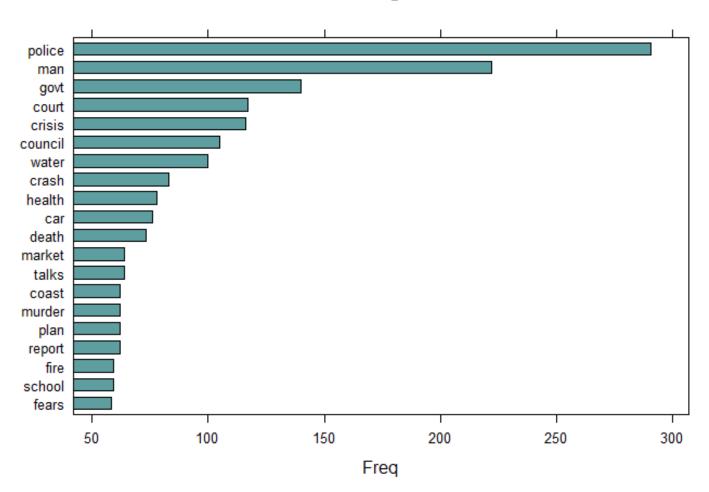


- Variables that take on a *limited* number of different values
 - Often referred to as categorical variables
- Very important in statistical modeling
 - Categorical variables enter into statistical models differently than continuous variables..
- Factors are also a very efficient way to store character values
 - Because each unique character value is stored only once, and the data itself is stored as a vector of integers
 - We'll use factors in artificial neural networks
- fdata = factor(data)
- fdata
- [1] 1 2 2 3 1 2 3 3 1 2 3 3 1
- □ Levels: 1 2 3

Most occurring nouns for 2008



Most occurring nouns



RAKE/RKEA

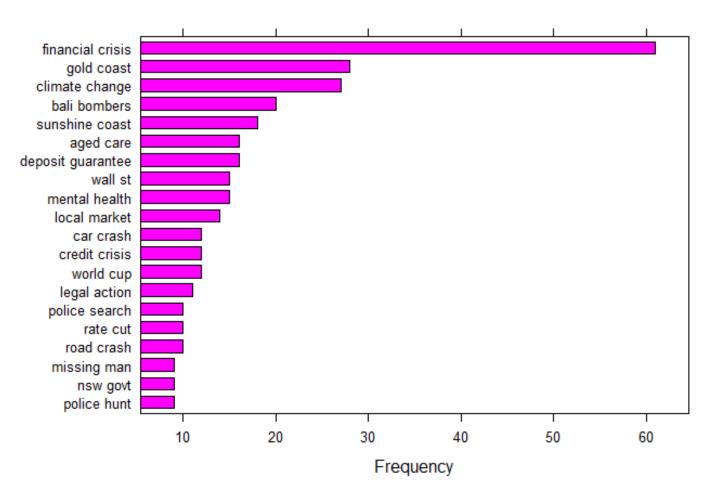


- One you get good at this, you realize scientists wrote a package that already does what you're doing..
- Rapid Keyword Extraction Algorithm (RKEA) or Rapid Algorithm for Keyword Extraction (RAKE)
 - Never figured out which one it is...

Simple noun phrases



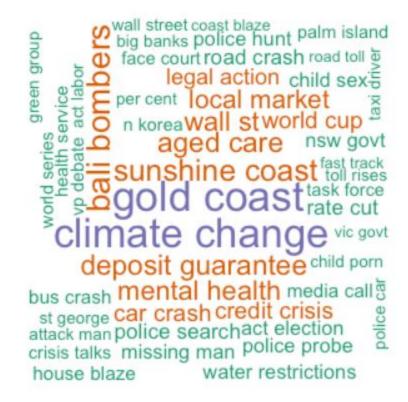
Keywords - simple noun phrases



And you give your friend his word cloud...



- library(wordcloud)
- wordcloud(words = stats\$key, freq = stats\$freq, min.freq = 3, max.words = 100, random.order = FALSE, colors = brewer.pal(6, "Dark2"))



Labs due...



- First version: Next week Monday, teams of 2 students each
- Second version: Next week Wednesday
- Code is up on Blackboard! At most 3.R files are due
 - 4-textmining-cloud.R
 - In English
 - 5-nlp.R
 - In English
 - Adapted to your language (e.g. Mandarin, Hindi, Greek, etc.) and your dataset
- If you don't need to modify the .R file, that's great, just submit a picture of your wordclouds. If you need to modify the R file, then submit the R file as well. For your language text, you will have to submit your R file, your dataset (in your language), and plots of your wordclouds
- Your dataset does not have to be as big as the abcnews dataset, but it should not be too small, either
- Best hindi & mandarin teams will present their work in class



