

# Computational text analysis: A practical overview

...

Fall 2019

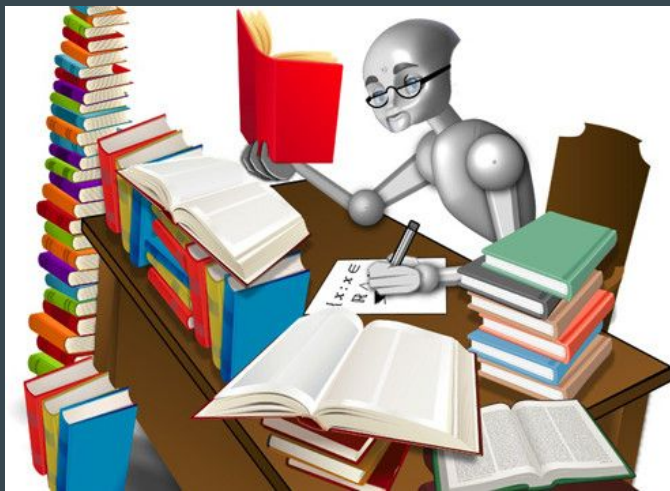
Anustubh Agnihotri

Borrowing from Geoff Bacon

With special thanks to Ben Gebre-Medhin and Laura Nelson

# Computational text analysis (CTA)

- Making computers analyze texts for us
- Closely related to natural language processing (NLP) and computational linguistics



# Types of languages

- Natural languages

*Time flies like an arrow. Fruit flies like a banana.*

- Artificial languages

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \dots + \beta_n X_n + \epsilon$$

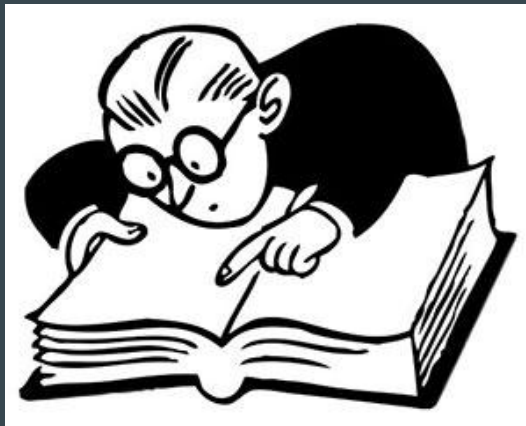
```
import scipy  
from scipy import sparse
```

```
n = 200000  
matrix = scipy.sparse.rand(n, n, density=.001)  
print(matrix)
```

# How do humans analyze texts?

*We need to steer clear of this poverty of ambition, where people want to drive fancy cars and wear nice clothes and live in nice apartments but don't want to work hard to accomplish these things. Everyone should try to realize their full potential.*

- Barack Obama

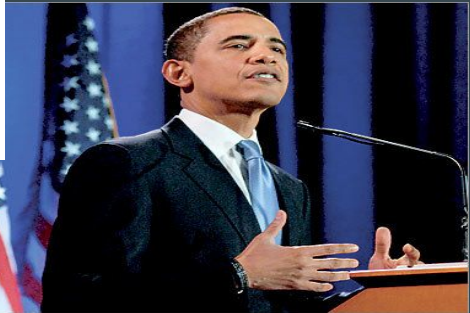
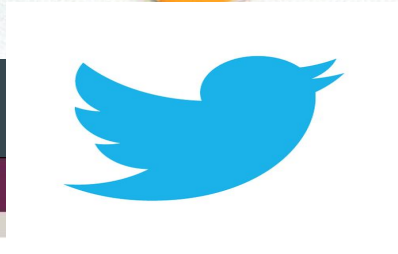
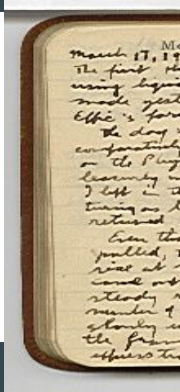


Close reading

# The promise of distant reading

- Scale/speed
- Reproducibility
- Does not have human biases
- Has other (at times unknown) biases
- Consistent

# Text as data



## Exhibit Feedback

### 1. Please explain below:

What did you think overall?

What would you improve?

# What can you do with CTA?

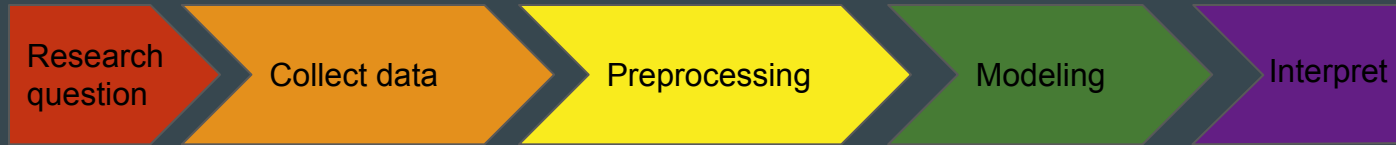
- Estimate political ideology from Twitter
- Uncover government censorship
- Relate the stock market to sentiment in the media
- Study why particular papers get cited
- Detect impending disease epidemics
- Determine who actually wrote something
- ...

# CTA lifecycle

...



# CTA lifecycle



- Cooking analogy

# Research question

- Domain specific
- Possible answer is encoded in text
- E.g. *What are the early warning symptoms of depression?*
- *How did different European nations react to the election of Trump?*
- *Do Twitter users react differently to mass shootings based on the ethnicity of the perpetrator?*
- *What distinguishes different styles of hip-hop?*
- *Have any of these essays been plagiarized?*

# What do we want?

- Machine readable
- Plain text = just the readable characters (i.e. letters, punctuation) plus new line.
- Not Word or pdf
- Sometimes useful to have some structure: csv, xml, html, json
- If many files, we'd like some structure (subdirectories or useful file names)

# How do we get it?

- You (your collaborator or a stranger on the Internet) already have it\*
  - Web scraping
  - API
  - OCR
- 
- \*Still important to know how the data was collected

# What is preprocessing?

- Tokenization = separating running text into words
- Sentence segmentation/tokenization = separating words into sentences
- Text normalization = dealing with upper/lower case, spelling mistakes, removing special characters, replacing URLs, numbers, etc.
- Remove “stop words”
- Stemming/lemmatization = removing morphological affixes
- POS tagging = assigning a part-of-speech category to each word
- Syntactic parsing = assigning a (normally graph or tree) structure to a sentence
- Chunking = shallow version of syntactic parsing
- Named entity recognition = identifying the proper nouns in a text
- ...

# Why do we do preprocessing?

- Because later methods require preprocessed data as input
  - Counting words requires having already identified the words of a text
  - Knowing the POS of a word might help us in knowing whether it is important (modal *can* vs noun *can*)
- Because we gain intuition about our data
  - It forces us to look at the data
  - Often coupled with exploratory data analysis (EDA)
  - We might find out that all the reviews are exactly 500 characters long, which suggests some have been truncated.

# Modeling

- In this session, we'll focus on topic modeling and text classification.
- Topic modeling: input = many different texts, output = what each text is about
  - Newspaper articles
  - Emails
- Classification: input = many different texts and hand-labeled categories, output = something that can take in unlabeled texts and predict the category.
  - Hand label a bunch of documents, train a computer to mimic your hand coding
  - Spam / ham
  - positive/negative reviews
- The big difference is the need for labeled data (unsupervised vs supervised)
- Text = document, could be a review, newspaper article, journal article, whole book, tweet, ...

# Clarifying some terminology

- Text = document, a single unit under study
- Corpus = a collection of documents
- Lexical  $\approx$  fancy name for word
- N-gram: sequences of adjacent words of length  $N$ 
  - *poverty*: uni-gram
  - *poverty of*: bi-gram
  - *poverty of ambition*: tri-gram
- Bag of words



# The Bag of Words Representation

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet!

15



it	6
I	5
the	4
to	3
and	3
seen	2
yet	1
would	1
whimsical	1
times	1
sweet	1
satirical	1
adventure	1
genre	1
fairy	1
humor	1
have	1
great	1
...	...

# But there's more out there

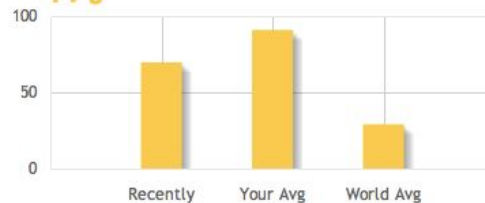
- Dictionary methods
  - Lists of positive/negative words
- Document-Term Matrix (DTM)
  - Words are rows, documents are columns, entries are number of times word appears in document.
- Term Frequency-Inverse Document Frequency (TF-IDF)
  - Modification of DTM
  - Entries are scaled by how common a word is across the whole corpus
- Clustering
  - Of DTM or TF-IDF

# Dictionary method

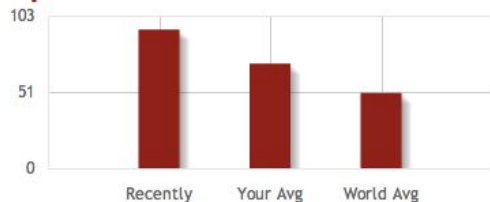
## Affectionate



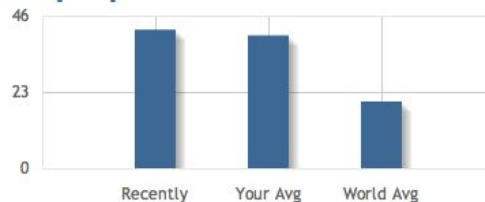
## Happy



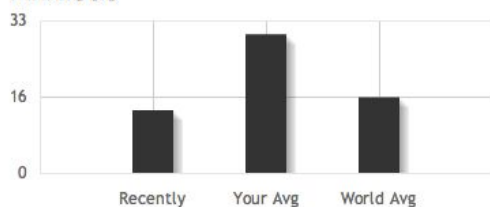
## Upset



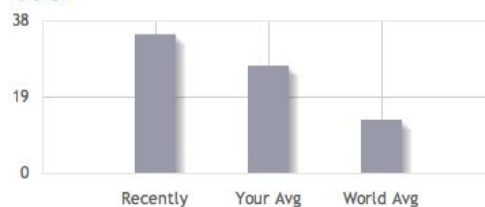
## Self-Expressive



## Anxious



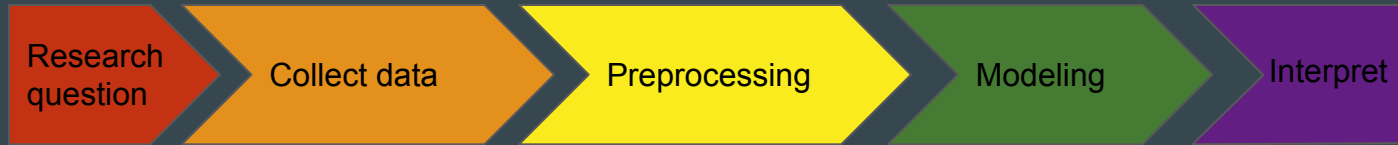
## Sad



# Now what?

- Use the insight from the preprocessing and modeling stage to understand the initial question
- In classification, this could be looking at words with high coefficients.
  - E.g. *People with depressive symptoms are more likely to talk about abstract concepts and use more negation.*
- In topic modeling, this could be qualitative inspection of the topics.
  - E.g. *In Germany the main themes centered around taxes and immigration, while in France people were more concerned about racism.*

# CTA lifecycle



# Questions?



# Yeah, but how do I do it?

- Yourself
  - This workshop
  - Other resources (in a few slides)
  - [NLTK](#), [SpaCy](#), [gensim](#), [polyglot](#), [TextBlob](#), etc.
- Using third-party APIs
  - [Google](#), [Amazon](#), [Microsoft](#) and [IBM](#)
- Work with someone
  - Fellow researchers/staff
  - D-Lab consulting

Where to go from here?

...



# Further resources

- D-Lab: <http://dlab.berkeley.edu/>
  - Regular expressions in Python
  - NLP with nltk/SpaCy
  - Consulting
- CTAWG: <http://dlabctawg.github.io/>
- [Lectures from Stanford's NLP class](#)
- Jurafsky and Martin [textbook](#)