

# Introduction to text mining (with R and tidyverse)

Peeter Tinit (University of Tartu/Tallinn University)

03.04.2019 at Rīgas Tehniskā Universitāte

# What is text mining

- (textual) data <-> questions
- Questions matter.
- Technology can be easy (if data is available)

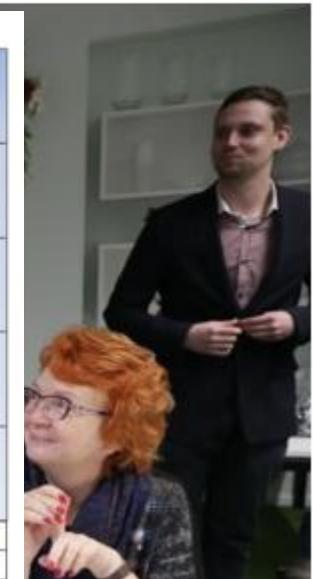
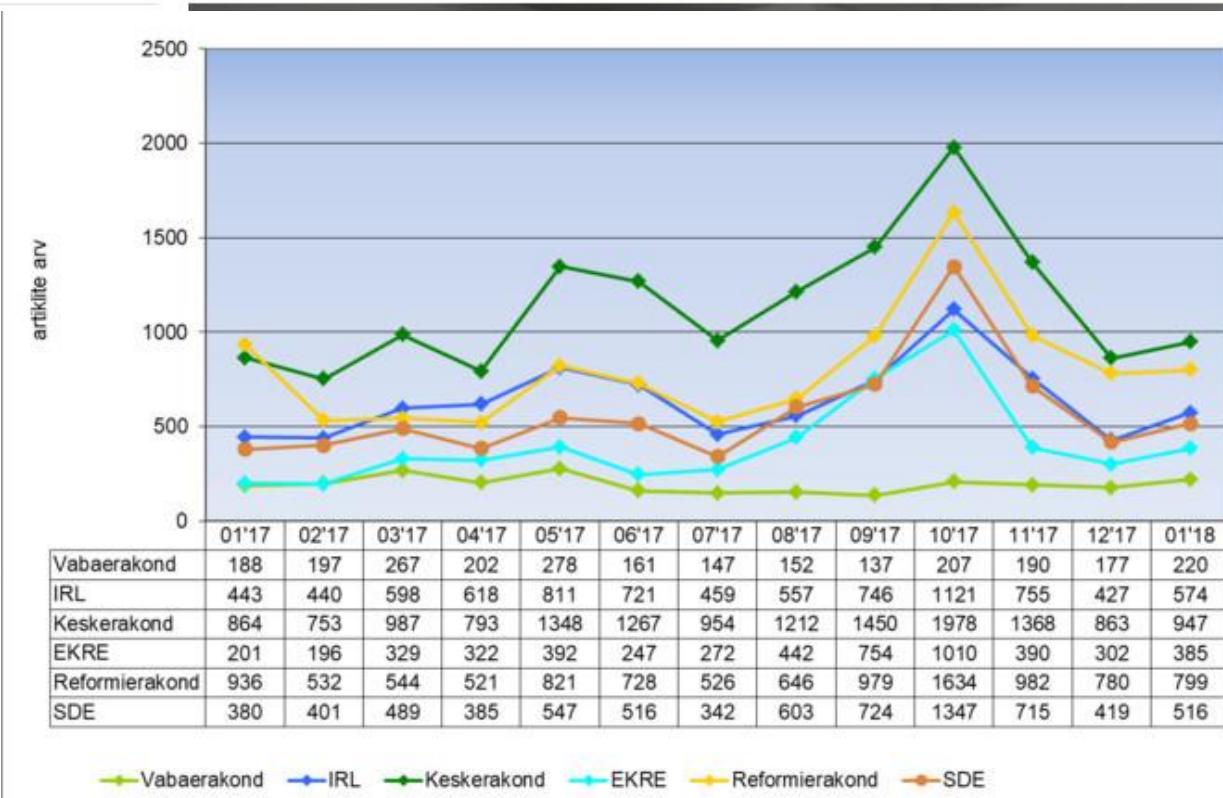
# Examples of text mining

# Media monitoring

## Jaanuaris pälvis meedias enim tähelepanu Keskerakond

EESTI

18.02.2018 16:51



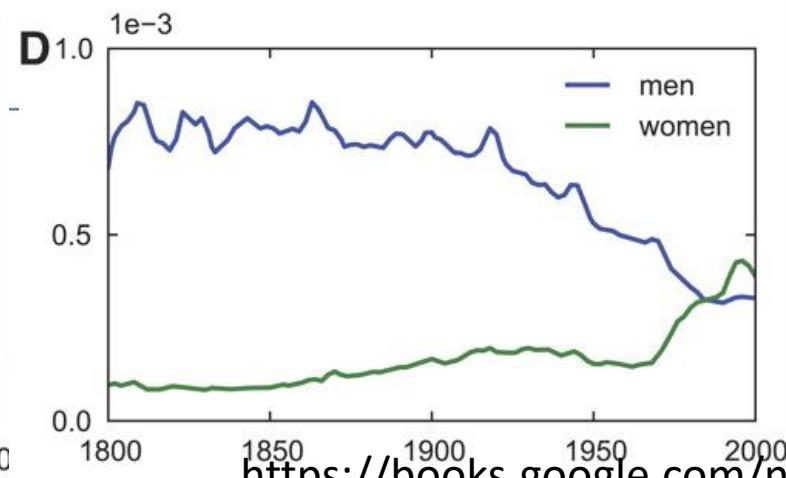
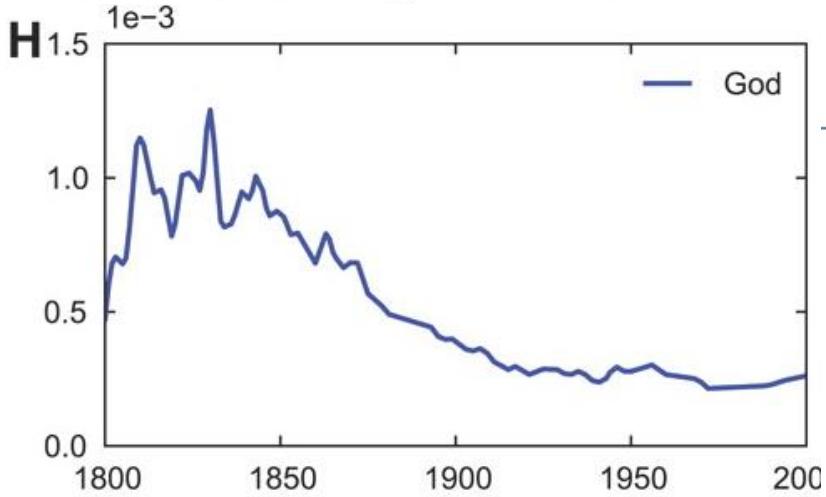
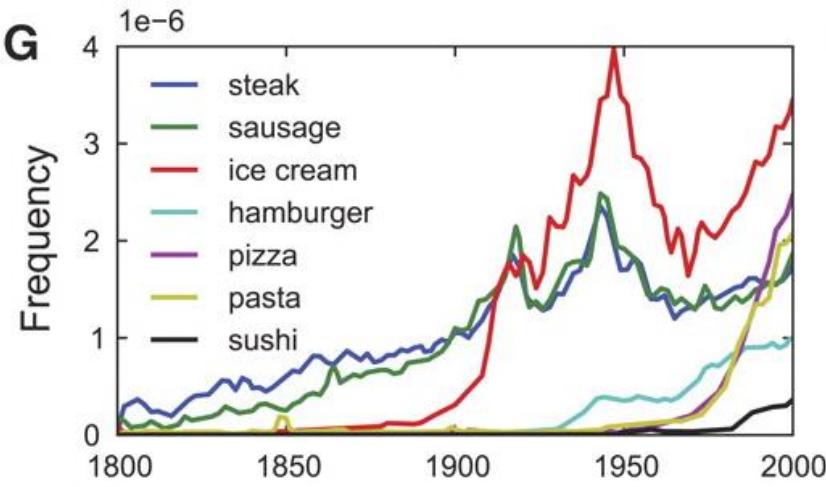
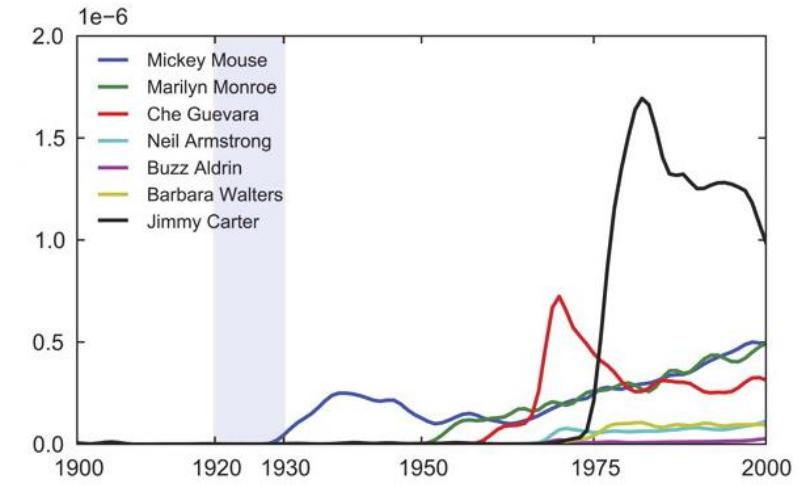
# Media monitoring

## Erakondade kajastatus meedias

Kui palju on erakondi sel kuul (eelmisel kuul) uudistes mainitud?

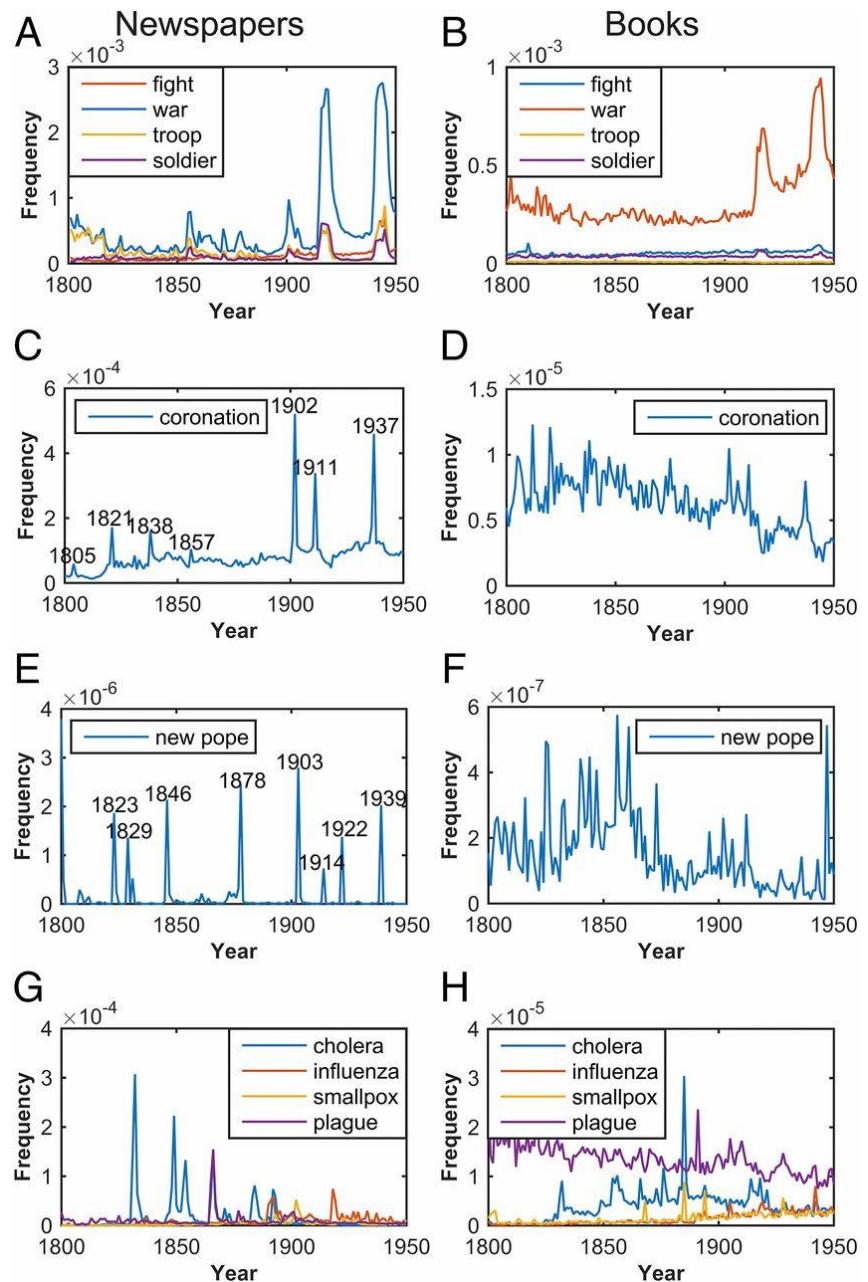
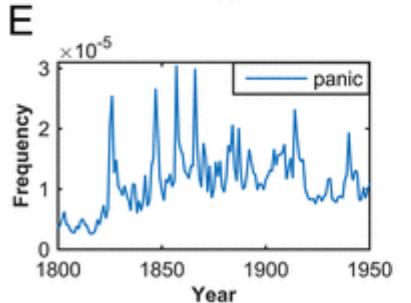
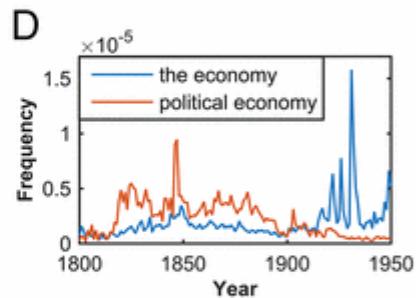
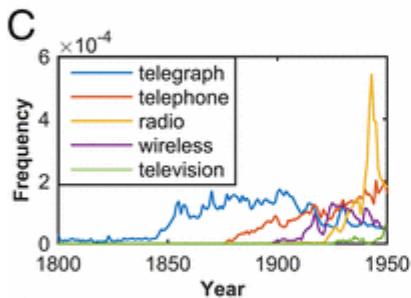
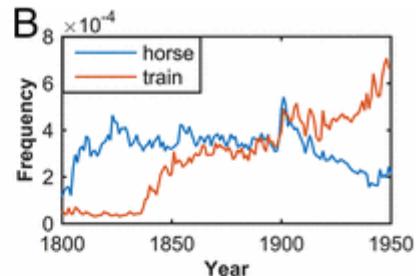
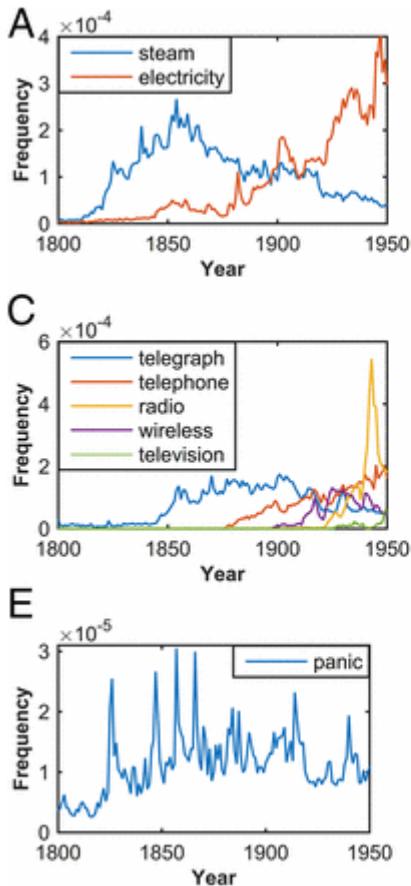
Erakond	Mainimisi kokku	Negatiivseid	Positiivseid
 Reformierakond	145 (157)	33 (28)	33 (41)
 KESKERAKOND	319 (288)	62 (50)	77 (65)
 IRL	101 (116)	24 (26)	22 (34)
 SDE	217 (185)	34 (38)	52 (41)

# „All digitized texts“



<https://books.google.com/ngrams/>

# 150 years of news



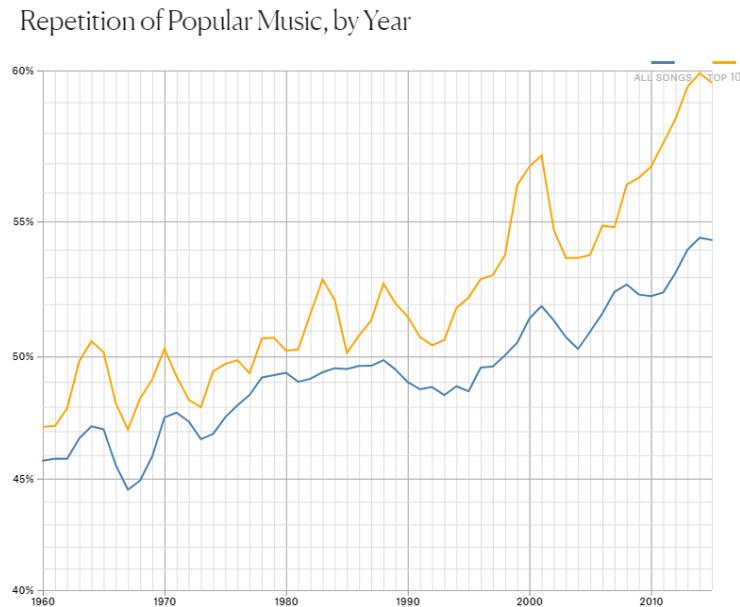
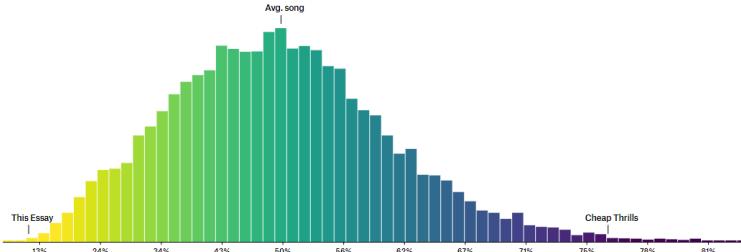
# Are Pop Lyrics Getting More Repetitive?

- 15 000 songs of top 100-s 1957-2017
- .zip the text

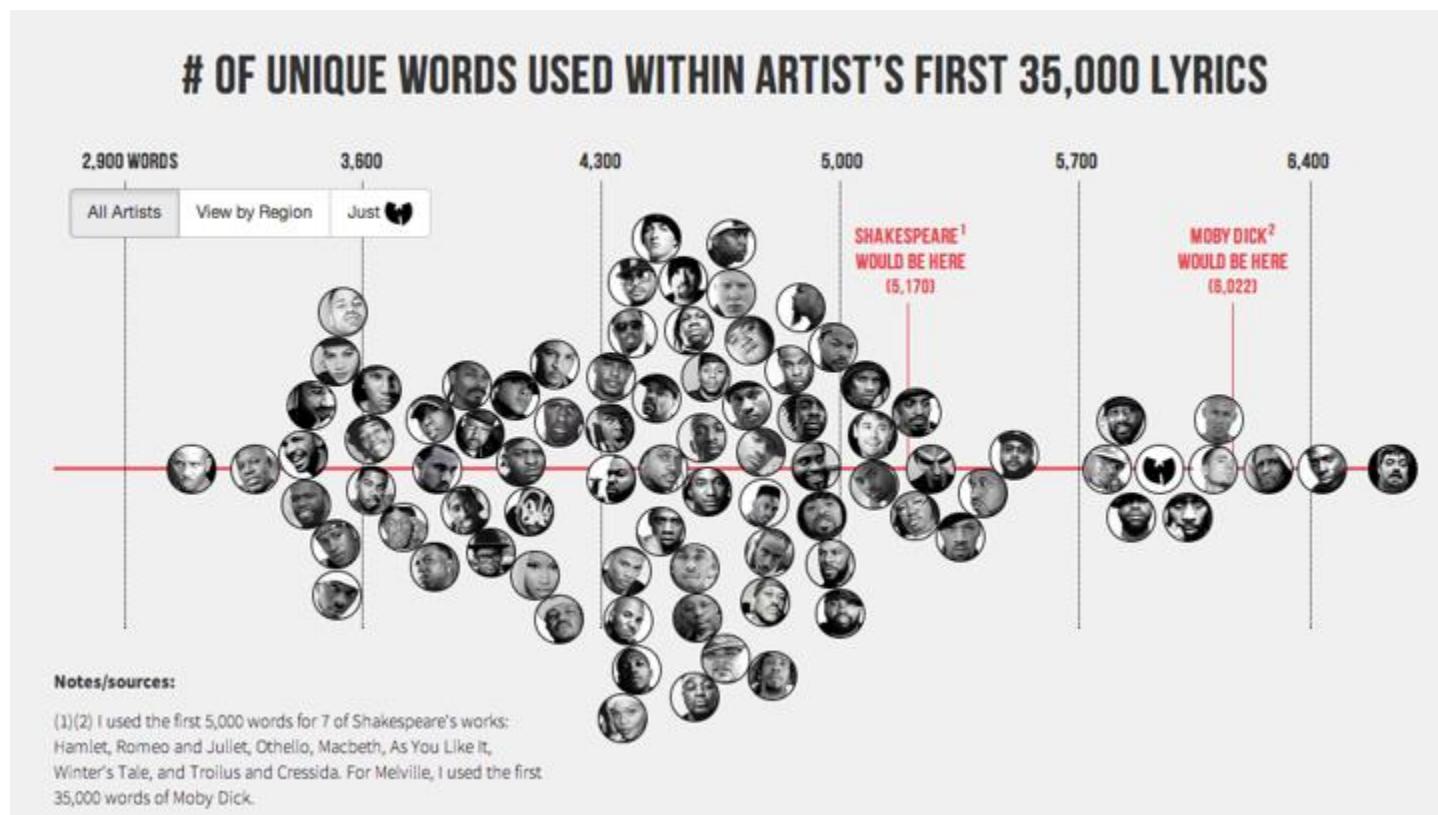
```
baby I don't need dollar bills to have fun tonight  
I love cheap thrills!  
baby I don't need dollar bills to have fun tonight  
I love cheap thrills!  
I don't need no money  
as long as I can feel the beat  
I don't need no money  
as long as I keep dancing
```

```
tonight I need dollar bills  
I don't keep fun  
cheap thrills long to feel money  
the bills don't need to be dancing baby  
fun dollar dancing thrills the baby I need  
don't have fun  
no no don't have dancing fun tonight  
beat the can as I don't feel thrills  
love the dancing money
```

29.5% Size Reduction

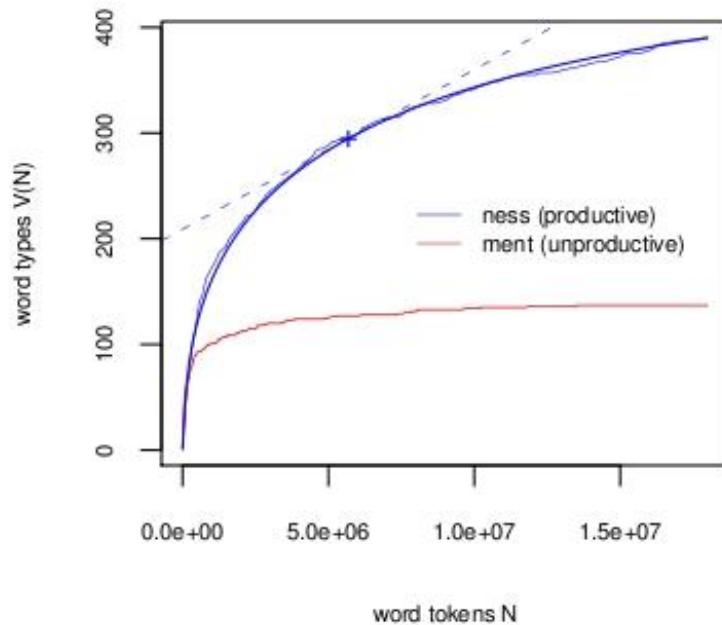


# Vocabulary of rap artists

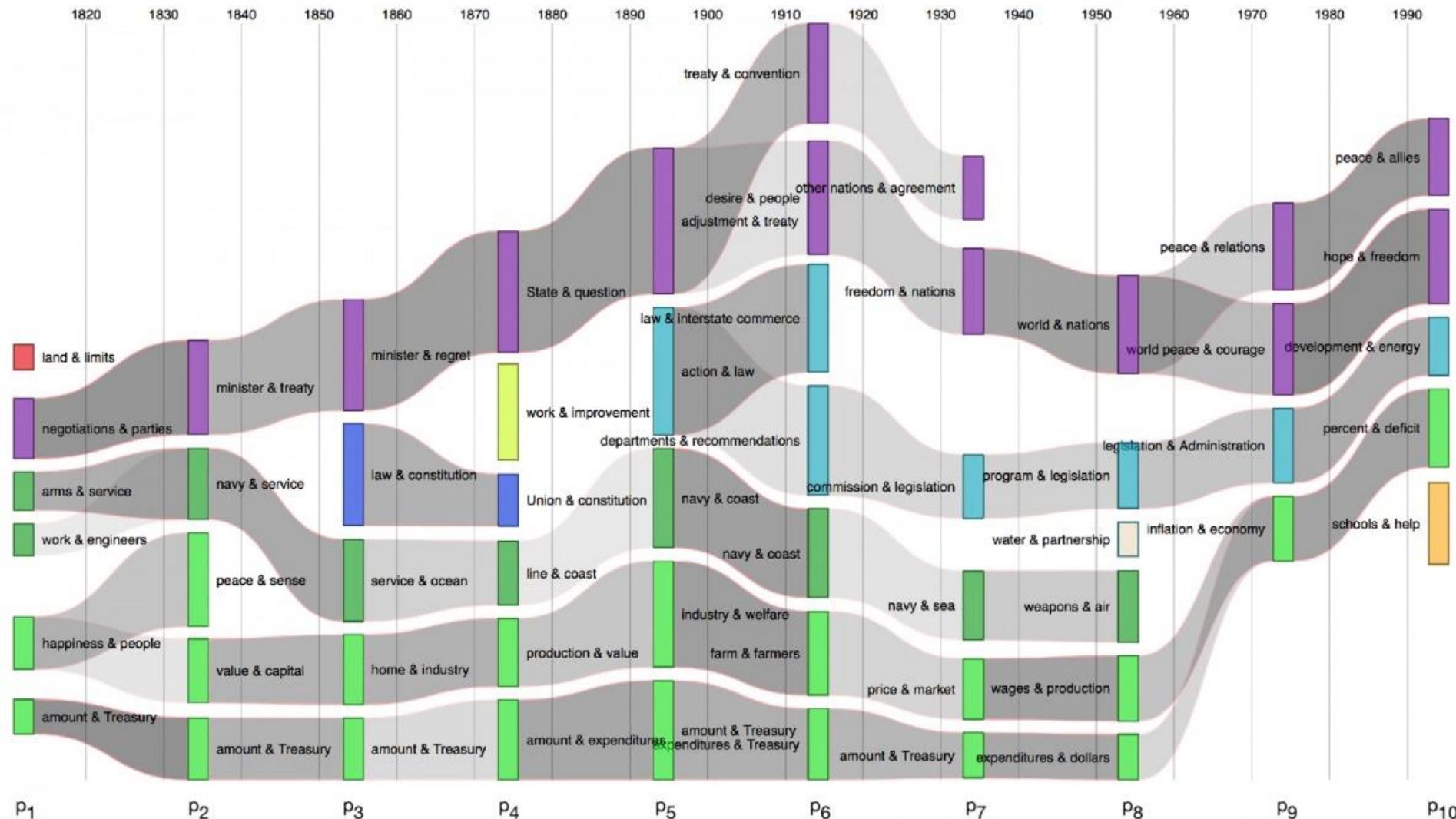


# Morphological productivity

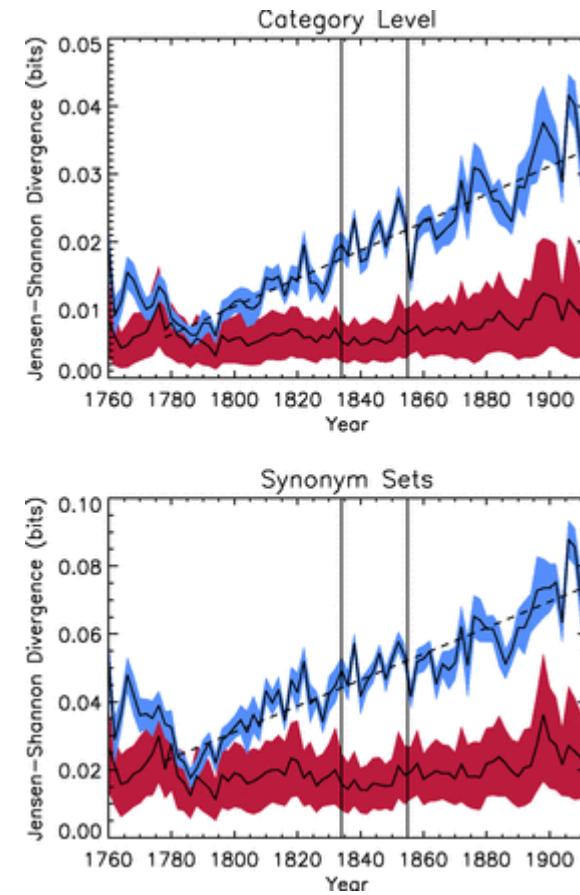
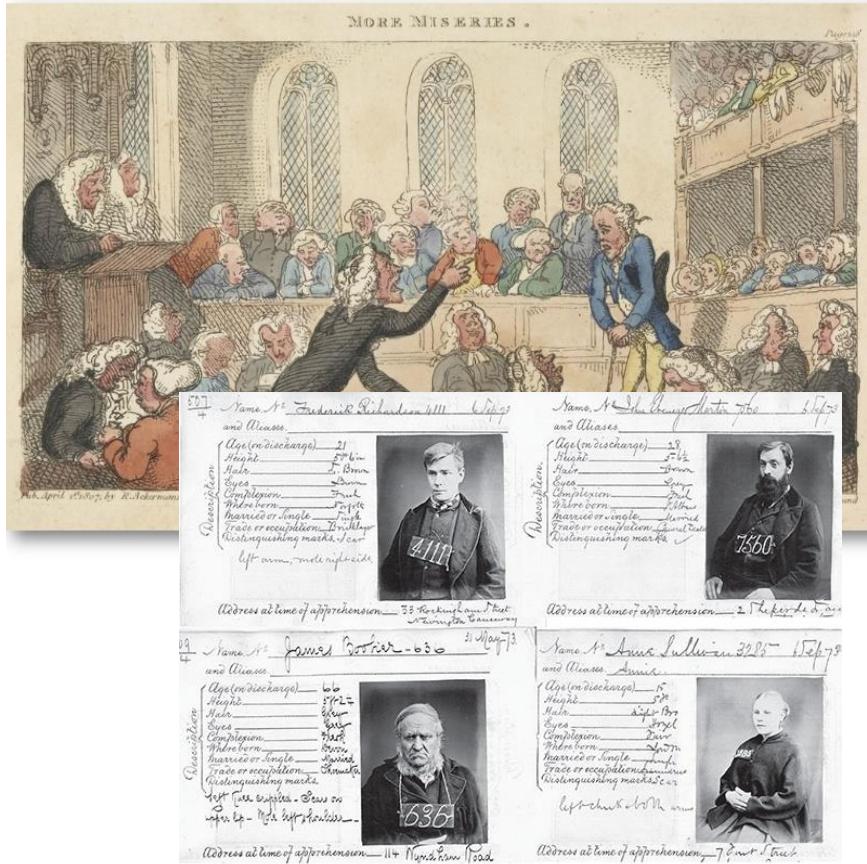
- Number of types / number of tokens
- ‘-ness’ – (*goodness*)
- ‘-ment’ – (*improvement*)



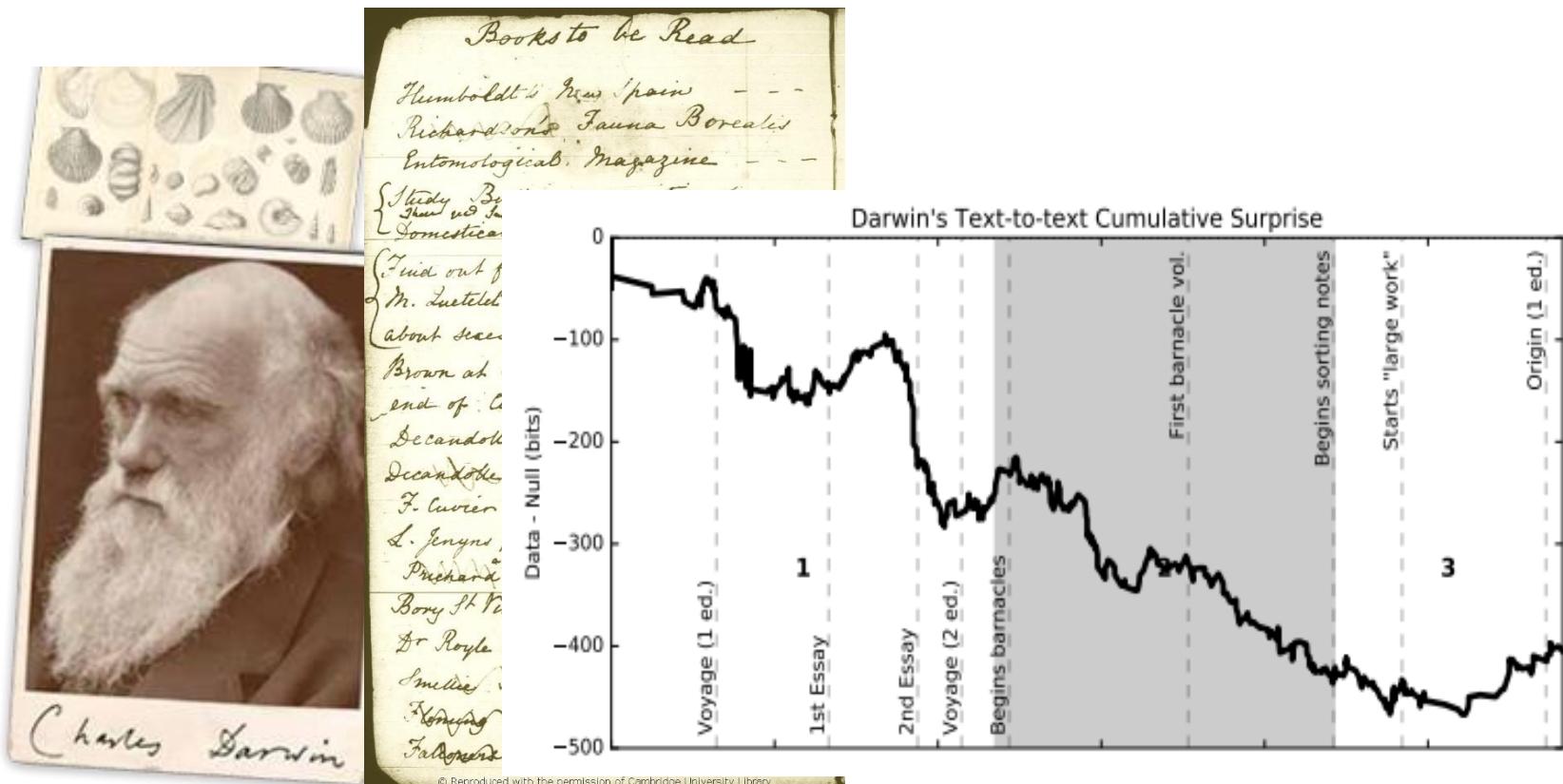
# Presidential speeches in USA



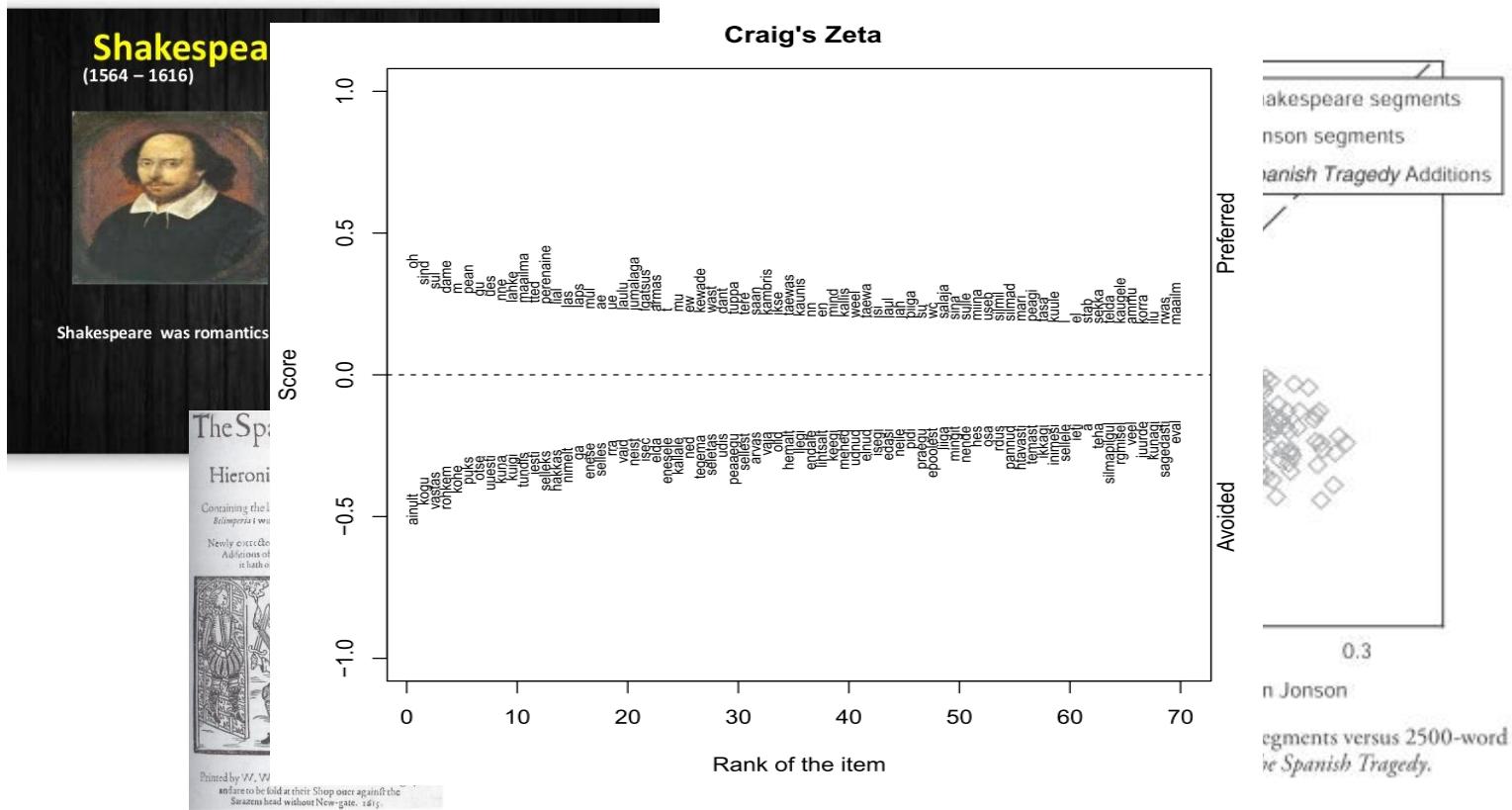
# Civilizing process



# Darwin's reading habits



# Finding the author

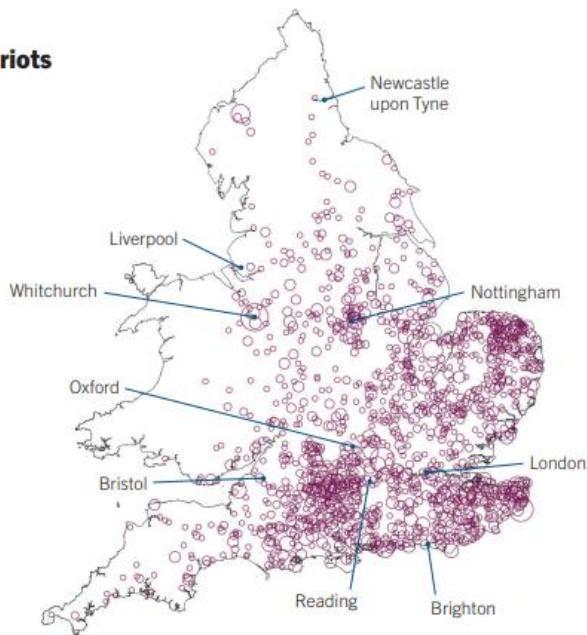


# Deep mining history

- Advertisements in newspapers -> diffusion of technology
- „Swing riots“ 1829-1830: machines, unemployment and crops.

**Figure 1:**  
**Location of Swing riots**

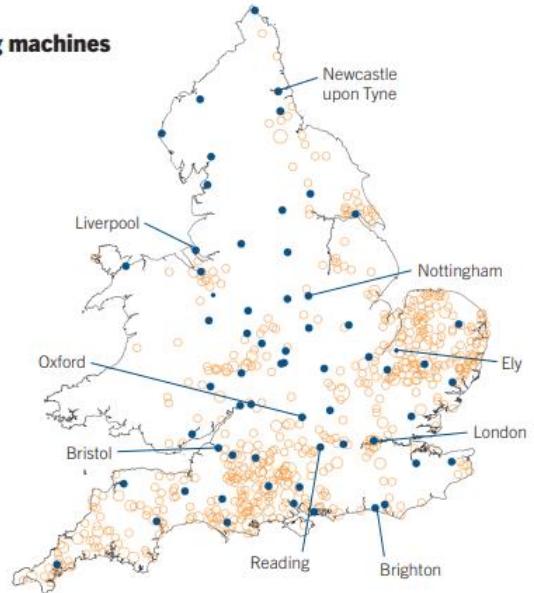
- 1 Riot
- 26 Riots  
e.g. Whitchurch (Shropshire)



Notes: Location of Swing riots. Purple circles identify parishes with Swing riots; the size of the circles is proportional to the number of episodes recorded in each parish.

**Figure 2:**  
**Location of threshing machines**

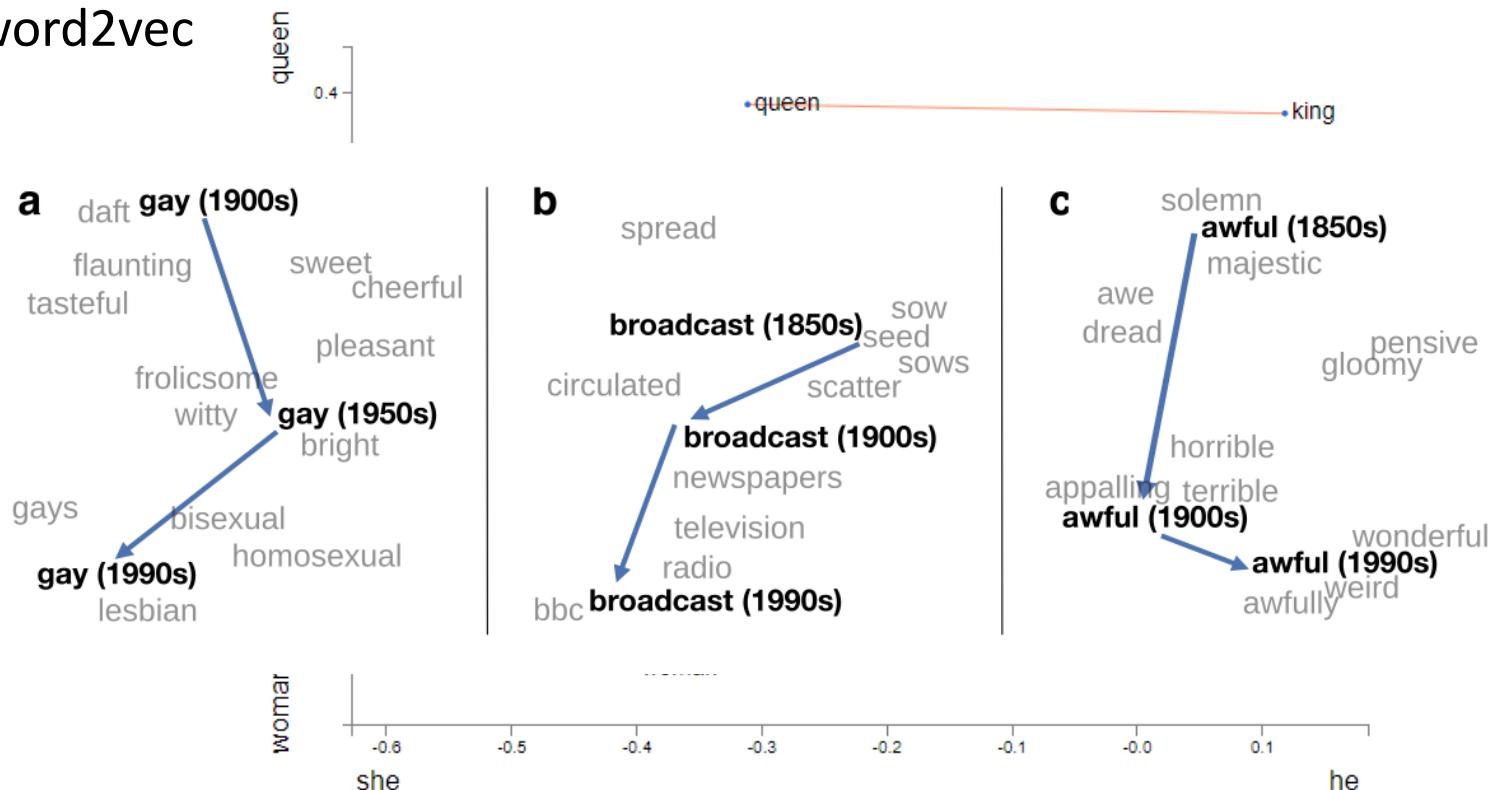
- 1 Machine
- 5 Machines  
e.g. Ely (Cambridgeshire)
- Newspaper ads



Notes: Location of threshing machines. Orange circles identify parishes with threshing machines; the size of the circles is proportional to the number of machines we found. Blue dots show cities that published newspaper advertisements at the time.

# Word semantics

- word2vec



# Plagiarism and bot detection

"In the matter of restoring Internet freedom. I'd like to recommend the commission to undo The Obama/Wheeler power grab to control Internet access. Americans, as opposed to Washington bureaucrats, deserve to enjoy the services they desire. The Obama/Wheeler power grab to control Internet access is a distortion of the open Internet. It ended a hands-off policy that worked exceptionally successfully for many years with bipartisan support.",

"Chairman Pai: With respect to Title 2 and net neutrality. I want to encourage the FCC to rescind Barack Obama's scheme to take over Internet access. Individual citizens, as opposed to Washington bureaucrats, should be able to select whichever services they desire. Barack Obama's scheme to take over Internet access is a corruption of net neutrality. It ended a free-market approach that performed remarkably smoothly for many years with bipartisan consensus.",

"FCC: My comments re: net neutrality regulations. I want to suggest the commission to overturn Obama's plan to take over the Internet. People like me, as opposed to so-called experts, should be free to buy whatever products they choose. Obama's plan to take over the Internet is a corruption of net neutrality. It broke a pro-consumer system that performed fabulously successfully for two decades with Republican and Democrat support.",

"Mr Pai: I'm very worried about restoring Internet freedom. I'd like to ask the FCC to overturn The Obama/Wheeler policy to regulate the Internet. Citizens, rather than the FCC, deserve to use whichever services we prefer. The Obama/Wheeler policy to regulate the Internet is a perversion of the open Internet. It disrupted a market-based approach that functioned very, very smoothly for decades with Republican and Democrat consensus.",

"FCC: In reference to net neutrality. I would like to suggest Chairman Pai to reverse Obama's scheme to control the web. Citizens, as opposed to Washington bureaucrats, should be empowered to buy whatever products they prefer. Obama's scheme to control the web is a betrayal of the open Internet. It undid a hands-off approach that functioned very, very successfully for decades with broad

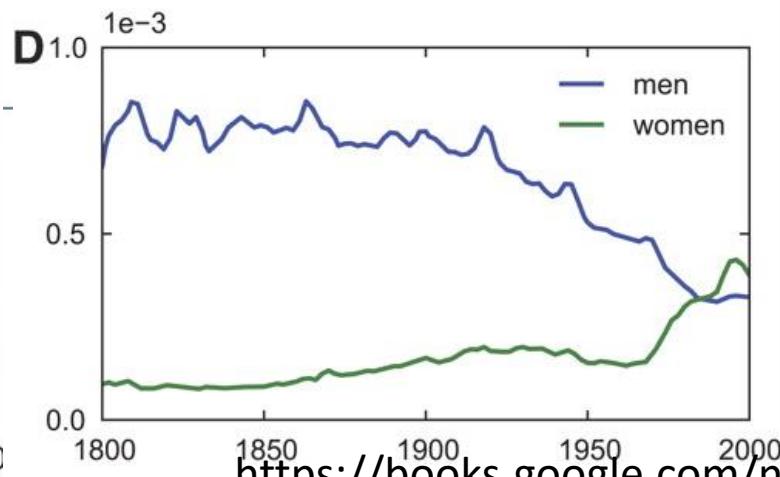
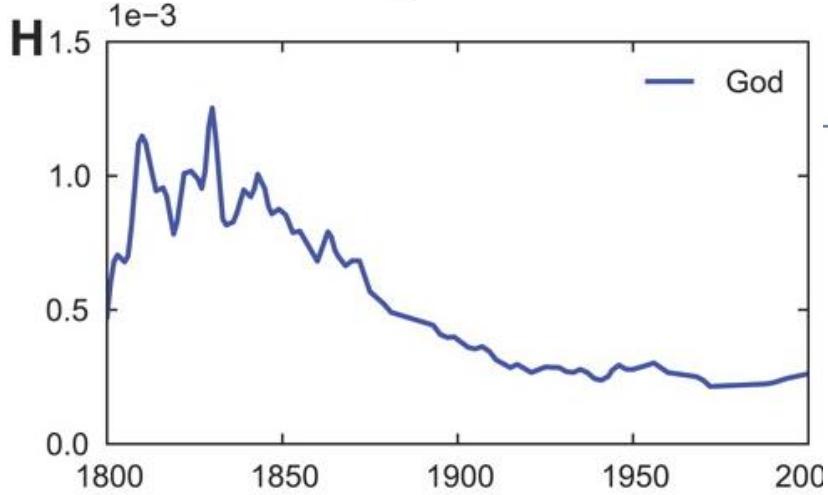
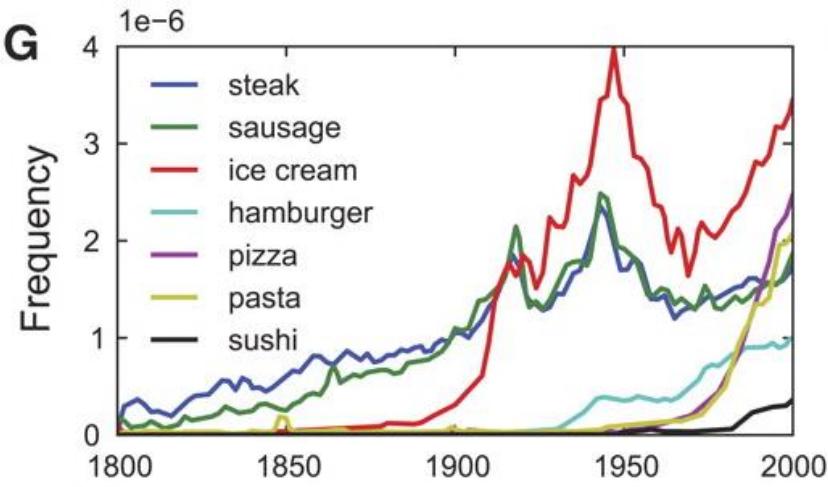
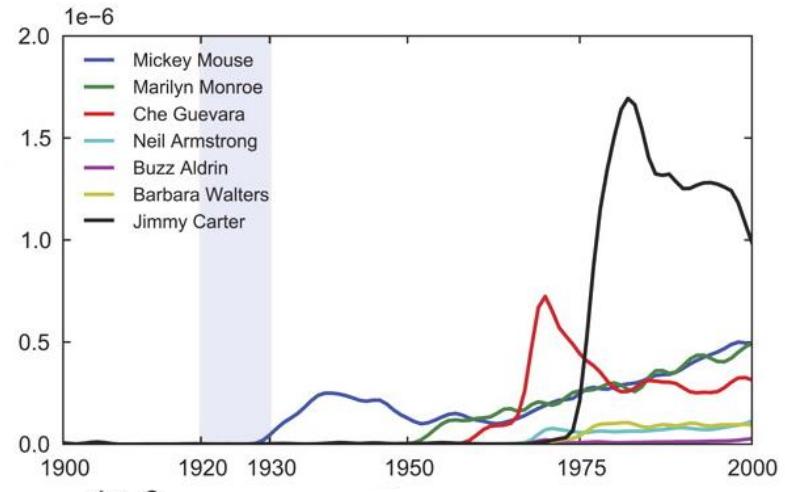
# Basic methods

# Basic methods

- 1. Frequencies
- 2. Comparing texts
- 3. Text-internal structure
- 4. Lexicon-based
- 5. Word distributions
- 6. Metainformation!

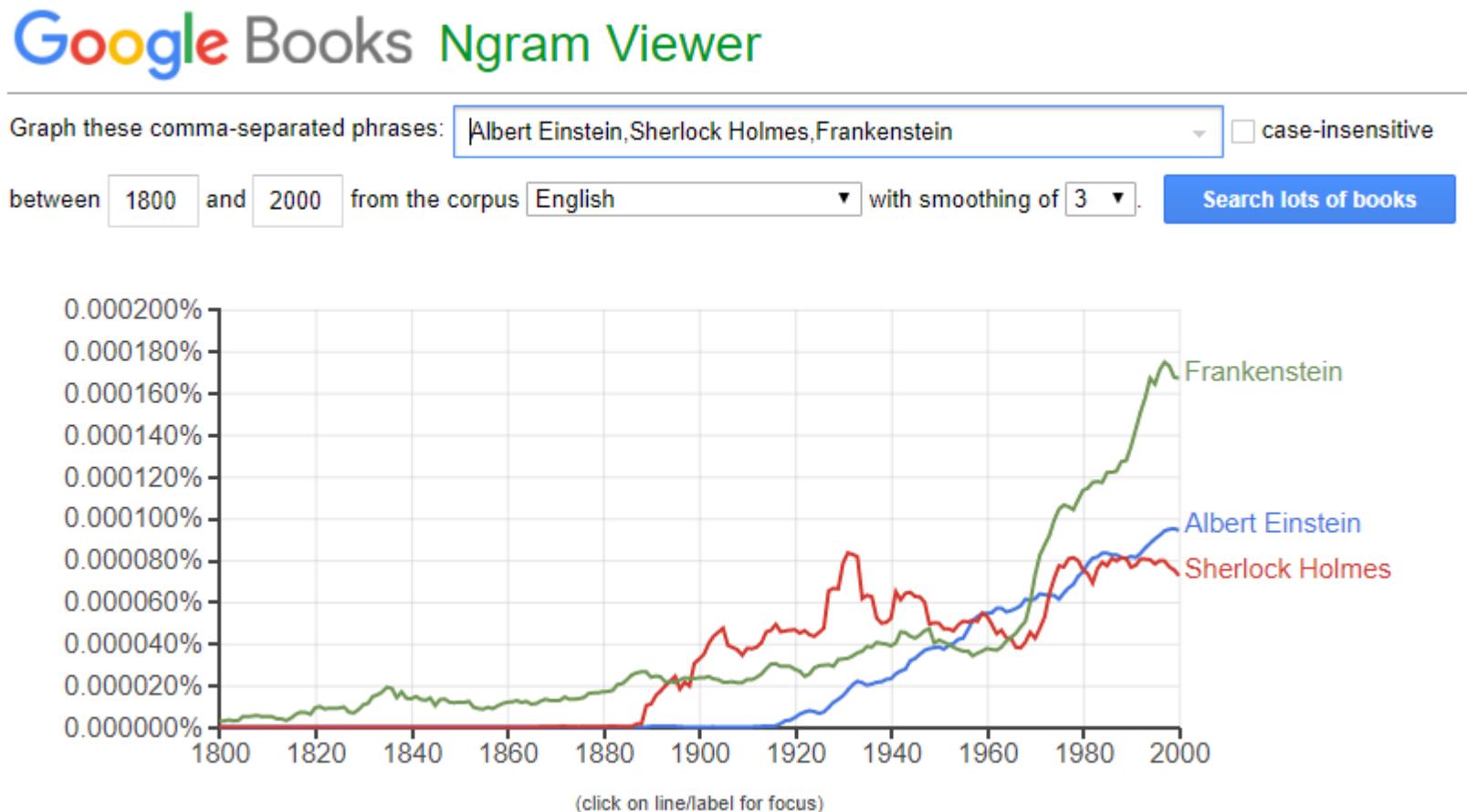
# Frequencies

# Frequencies



<https://books.google.com/ngrams/>

# Google Ngrams (1-grams, 2-grams, 3-grams...)



<https://books.google.com/ngrams>

# Counts vs frequencies

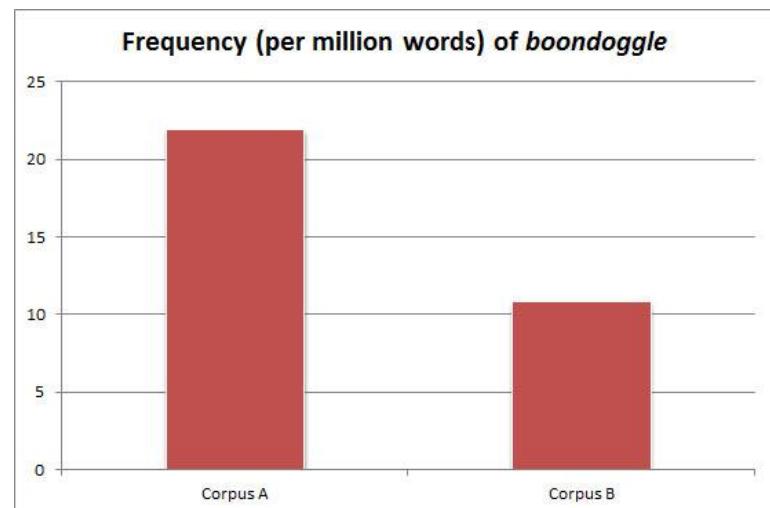
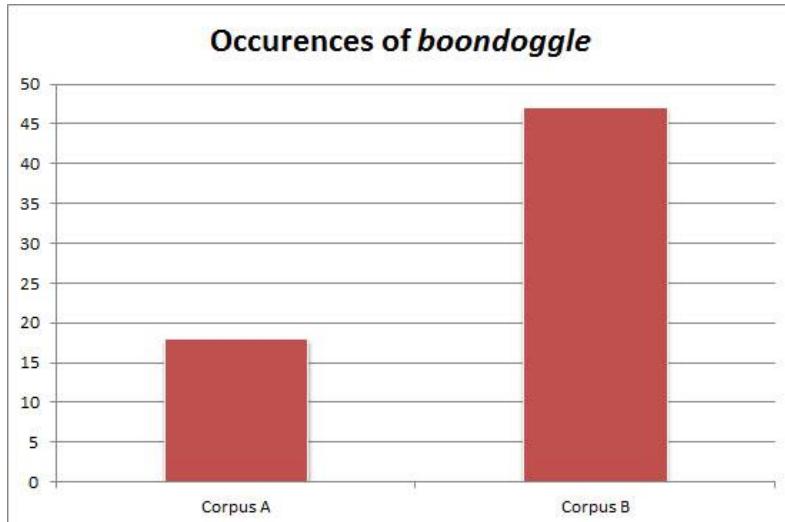
- The 1-gram ‘boondoggle’ in corpus
- (A **boondoggle** is a project that is considered a waste of both time and money, yet is often continued due to extraneous policy or political motivations.)

$$\frac{18}{821,273} = \frac{x}{1,000,000}$$

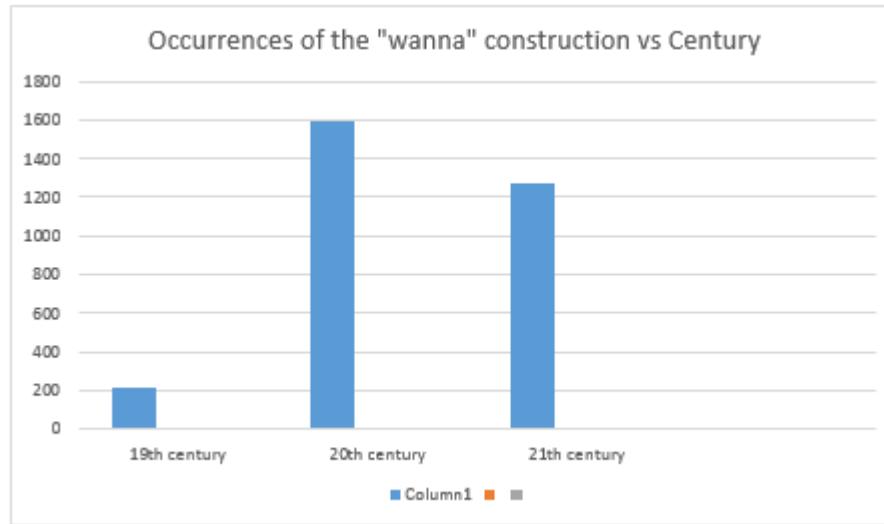
$$x (821,273) = 18 (1,000,000)$$

$$x = 18 (1,000,000) / (821,273)$$

$$x = \sim 22$$



# Quick test of attention



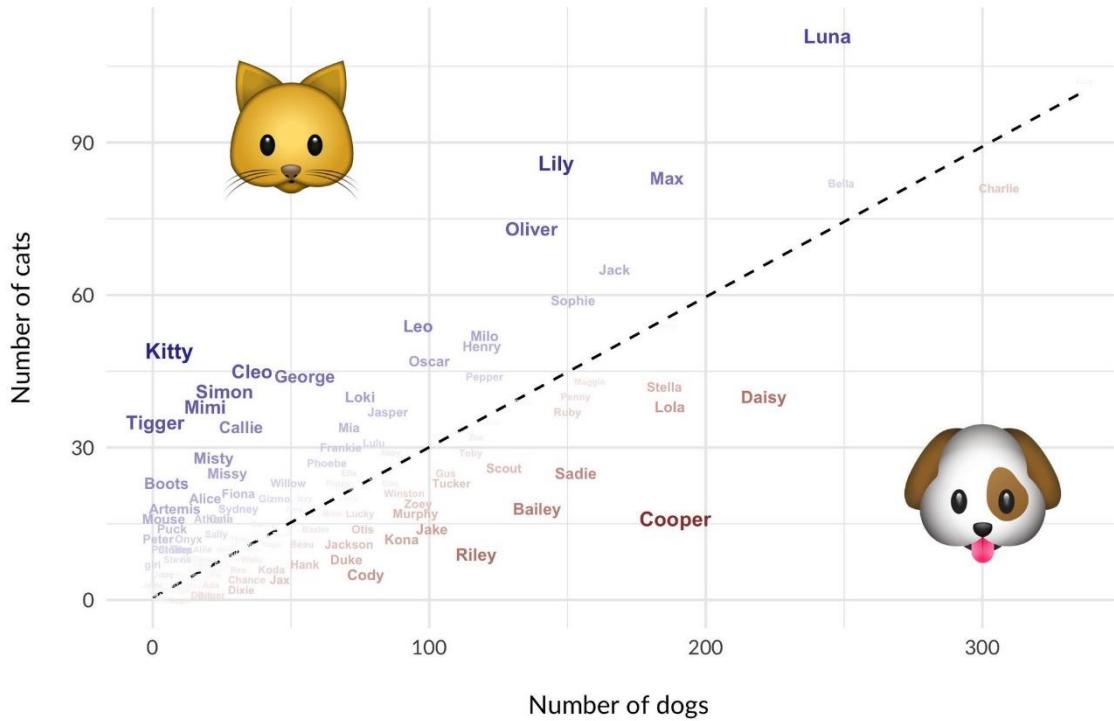
[StackExchange: How do I compare data from two corpora correctly?](#)

# Comparing texts

# Comparing frequencies

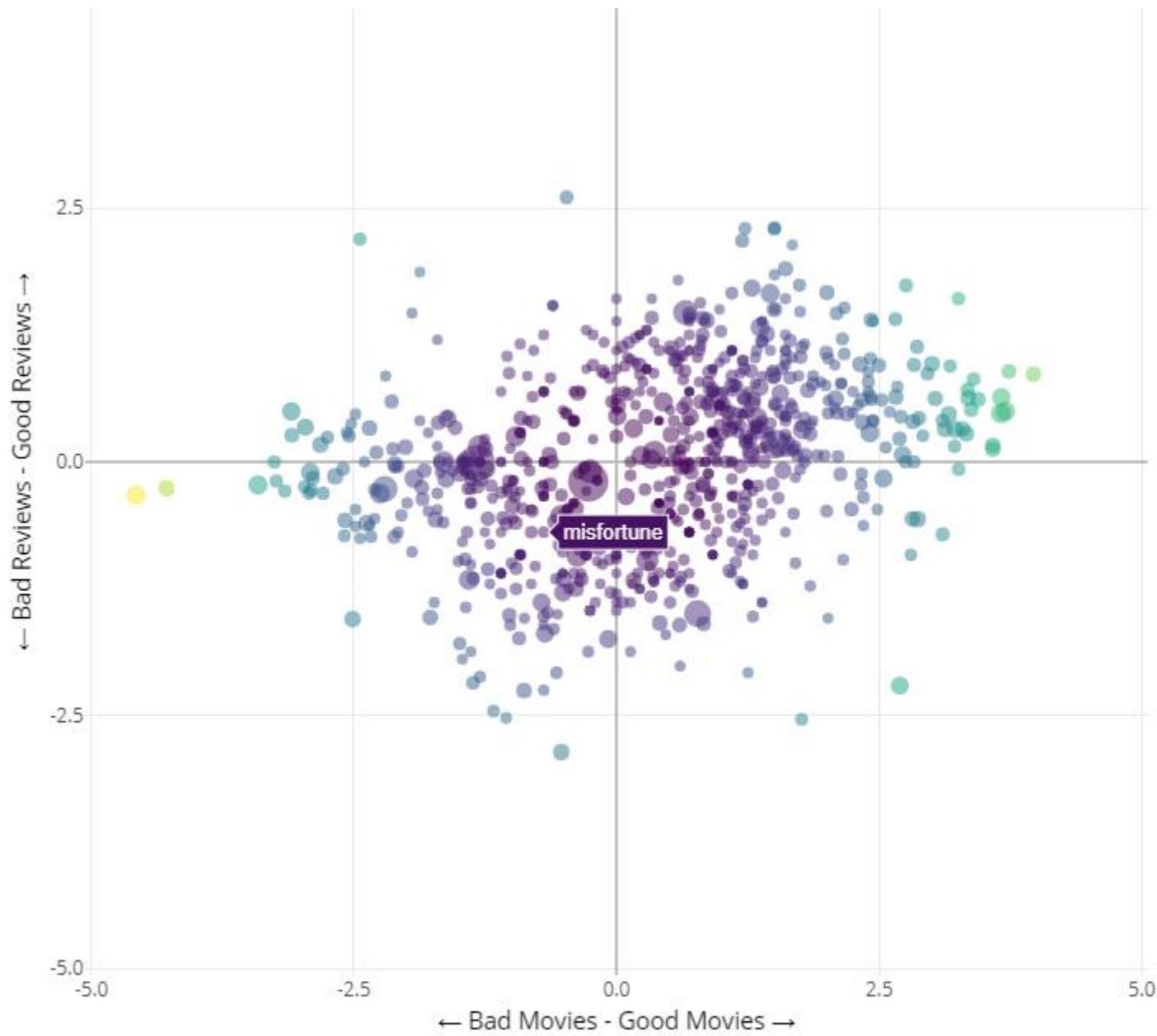
Is your pet's name more 'cat' or more 'dog'?

## Seattle's most popular pet names in 2018

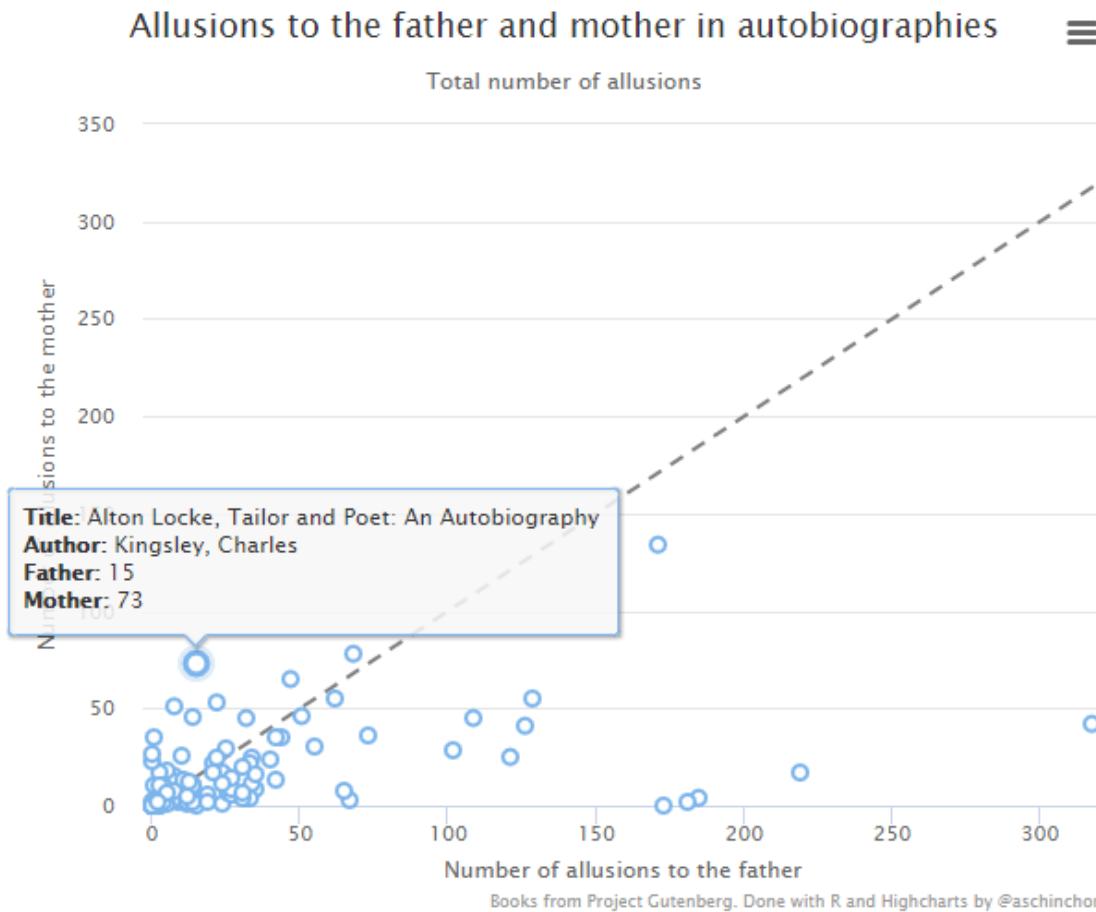


Source: #tidytuesday  
Graphic: @JoshuaFeldman

<https://twitter.com/JoshuaFeldman/status/1110604459669897218>

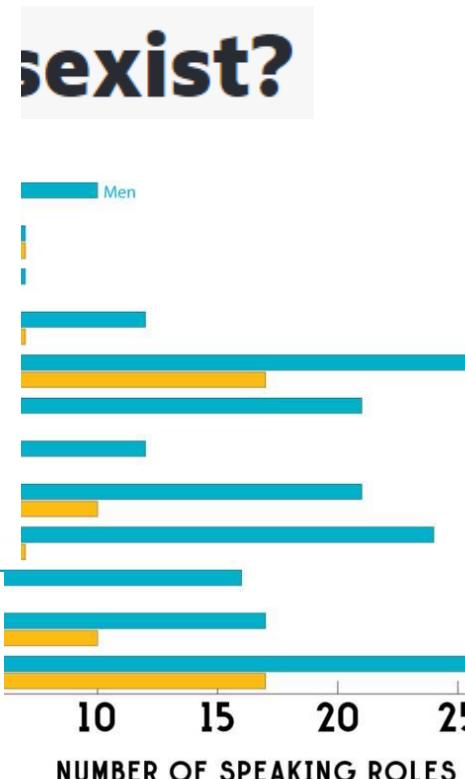
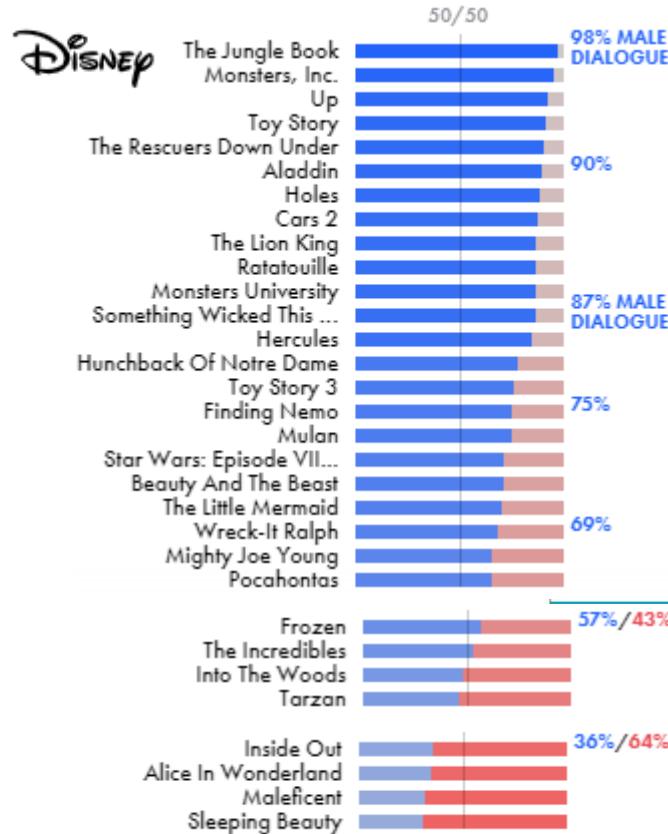


# Lexicon+comparison



# Text-internal structure

# Who speaks when



Source: Hannah Anderson, The Pudding

100% of Words  
are Male

90%

75%

60%

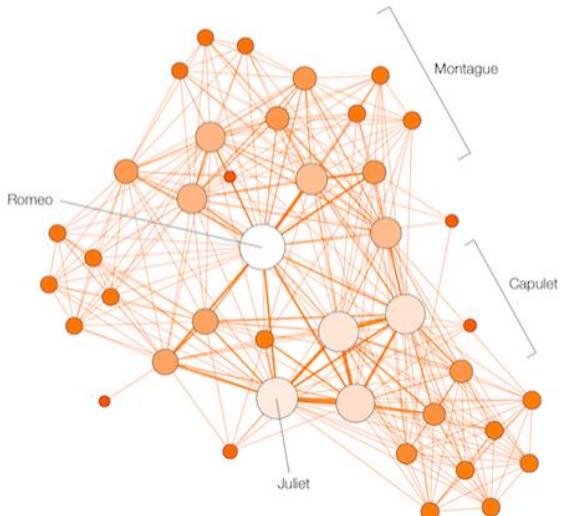
50/50  
Gender Split

100% of Words  
are Female

THE TRUMAN SHOW

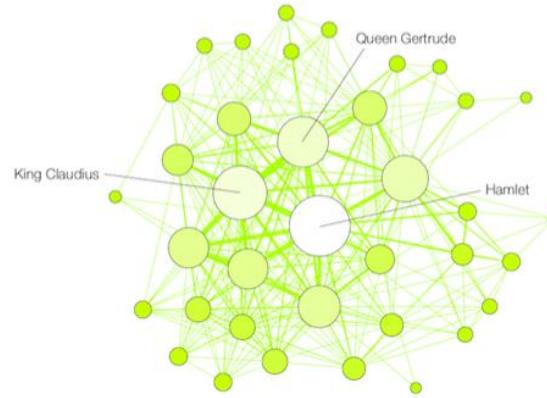
<https://pudding.cool/2017/03/film-dialogue/>

# Co-occurrences in scenes



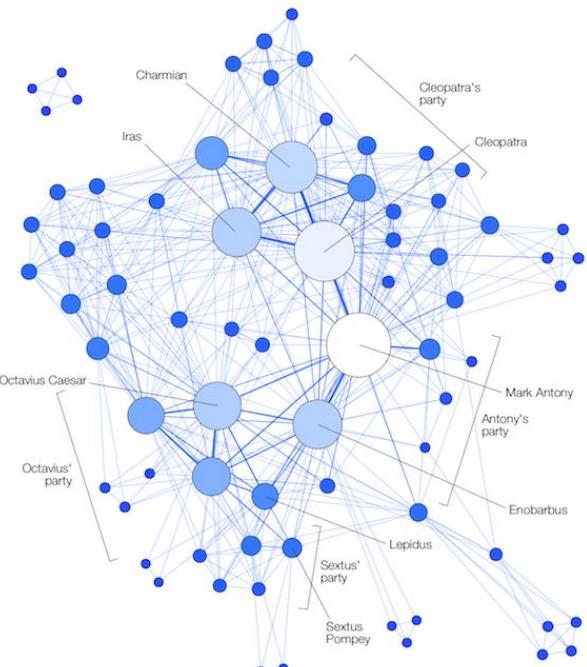
ROMEO AND JULIET

Number of characters **41** | **37%** Network density



HAMLET

Number of characters **37** | **39%** Network density

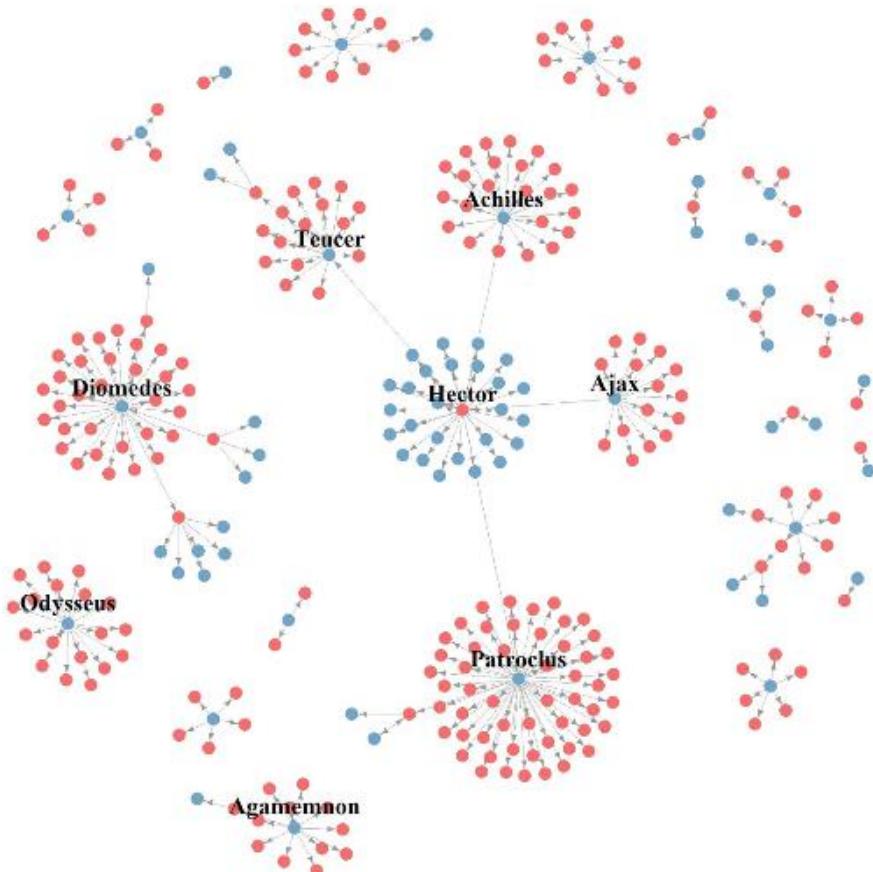
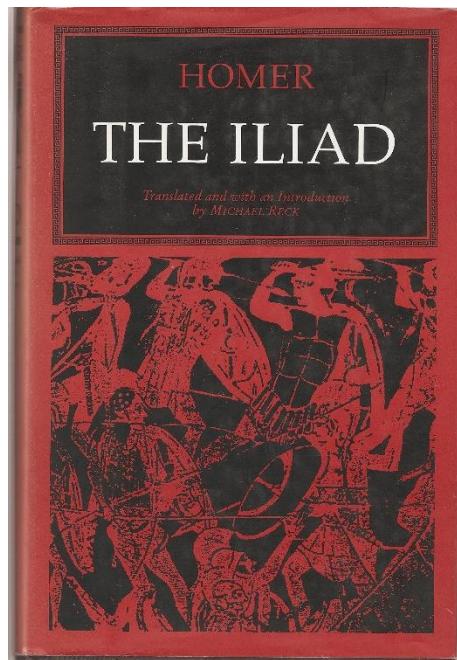


ANTONY AND CLEOPATRA

Number of characters **74** | **17%** Network density

# Cooccurrences in scenes

- Fights in Iliad



# Lexicon-based

# Lexicon based

- Sentiment analysis

API TEST TOOL

English ▾ Sentiment ▾ Graphical ▾

I **①** really enjoyed using the **①** Canon Ixus in Madrid on March 4. The **②** Panasonic Lumix **②** is a bit disappointing, but the **③** Canon **③** camera is **④** not bad **④** at all. All I want when taking photos is point it and then just press the button. For only 200 dollars, a **④** really fair **④** price, this **⑤** camera **⑤** is **⑥** perfect **⑥** for me. Besides, I have had a **⑥** good **⑥** customer **⑥** service **⑥** experience **⑥**. **⑦** John Faraday was **⑦** very nice!

LEGEND color key SENTIMENT

- Sentiment topic
- Positive sentiment text
- Negative sentiment text
- Text and topic link

ANALYZE TEXT ► RESET ⏪

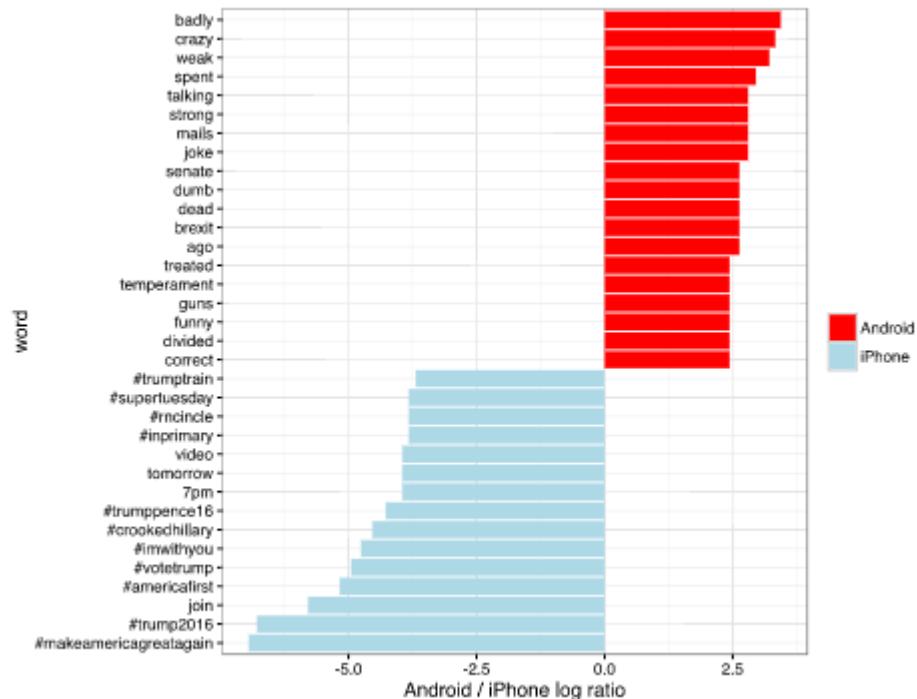


Bitext tool

# Political tweets

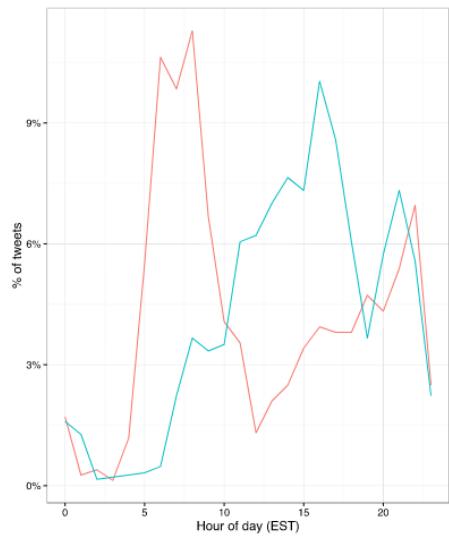
## Text analysis of Trump's tweets confirms he writes only the (angrier) Android half

I don't normally post about politics (I'm not particularly savvy about polling, which is where data science has had the largest impact on politics). But this weekend I saw a hypothesis about Donald Trump's twitter account that simply begged to be investigated with data:

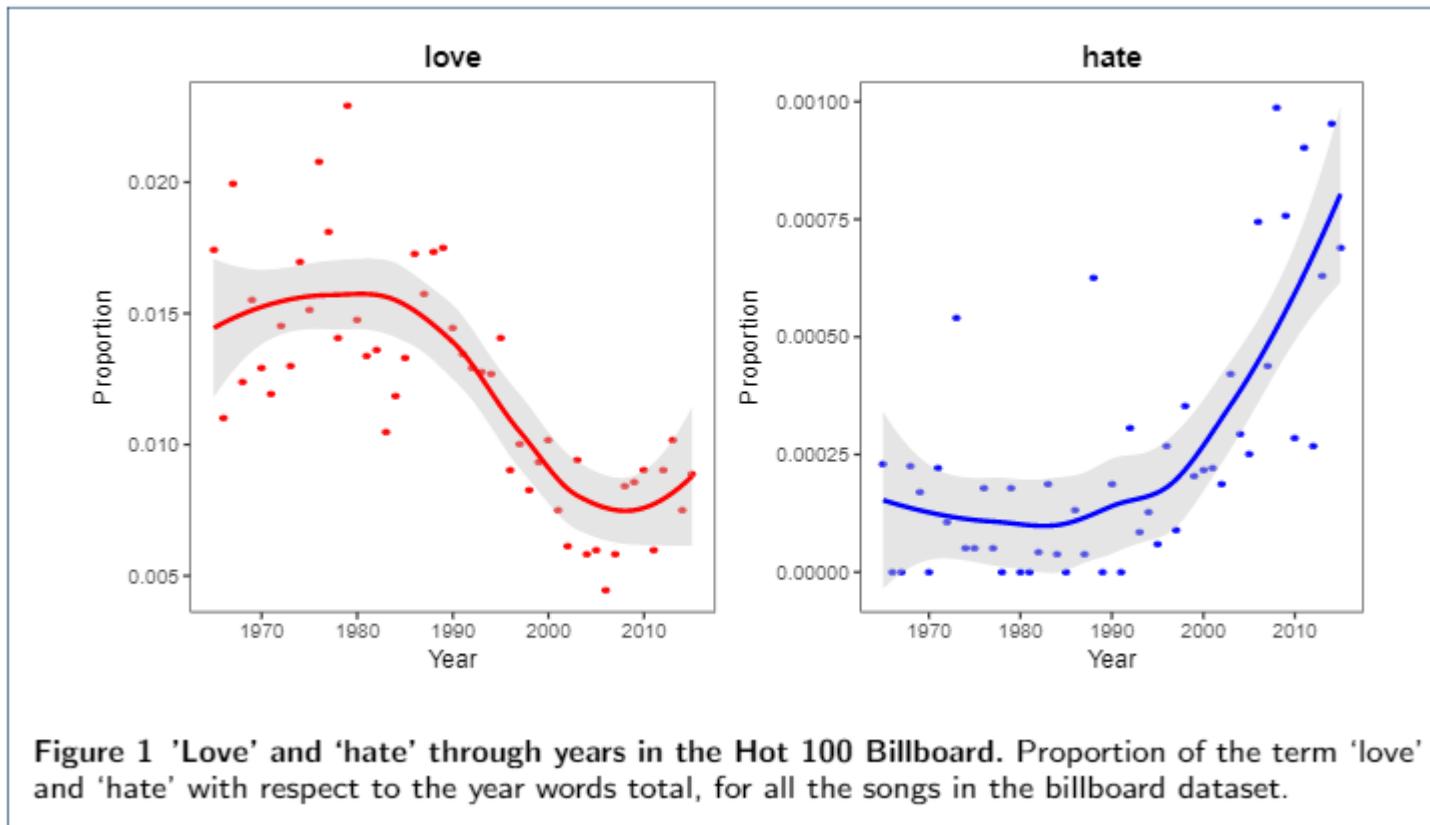


A screenshot of Donald J. Trump's Twitter profile. It shows a photo of him at a campaign rally, a bio reading "THANK YOU WINDHAM, NEW HAMPSHIRE JULY 6, 2016 #TRUMPTRAIN", his handle @realDonaldTrump, and a tweet he posted at 5:19 AM on August 7, 2016, from Windham, NH. The tweet reads: "Thank you Windham, New Hampshire! #TrumpPence16 #MAGA". It has 2,338 replies, 6,044 retweets, and 19,778 likes.

A screenshot of Donald J. Trump's Twitter profile showing another tweet. The bio remains the same. The tweet reads: "The media is going crazy. They totally distort so many things on purpose. Crimea. nuclear. "the baby" and so much more. Very".



# Sentiment across texts



Morin & Acerbi 2016 Birth of the cool

Brand et al. 2018 Cultural evolution of emotional expression in 50 years of song lyrics

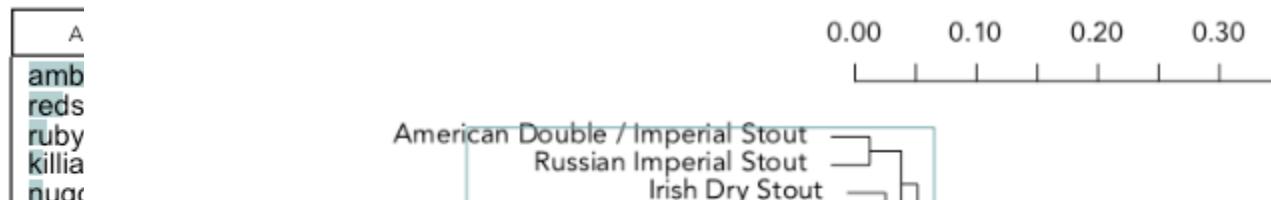
# Geo-locations

- Which locations are mentioned in text

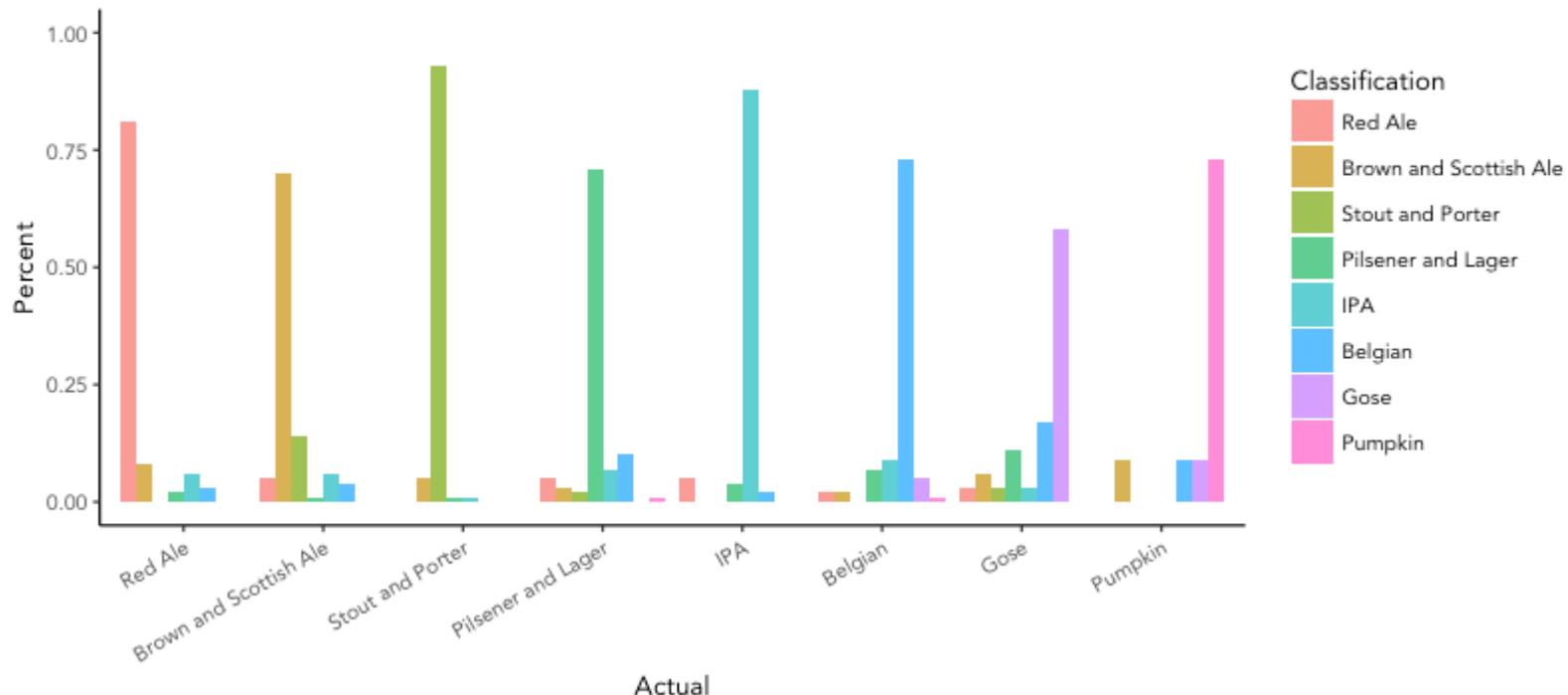


# Word distributions

## Top TF-IDF Terms for Selected Beer Styles

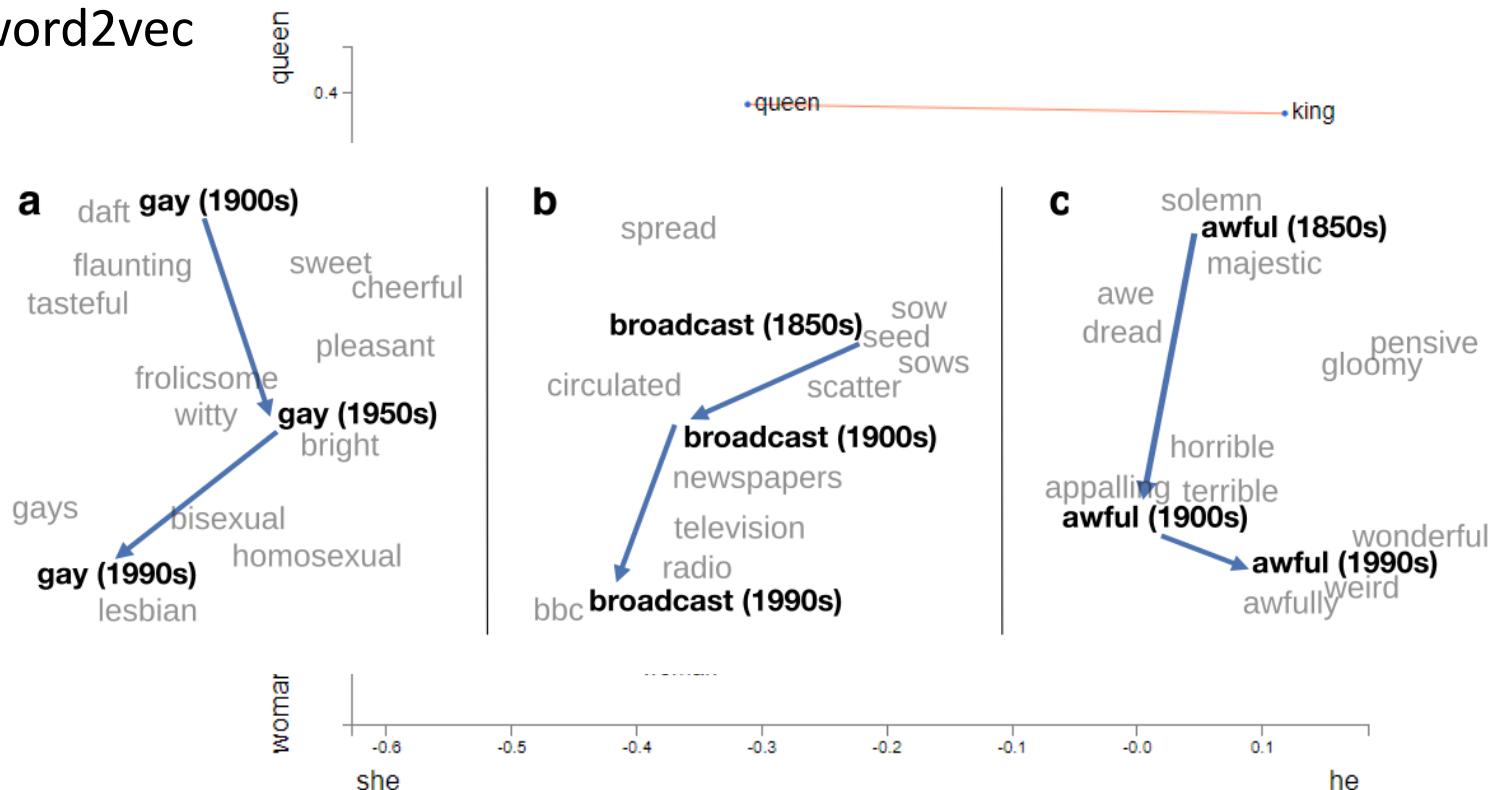


## Accuracy of KNN Classification into HClusters

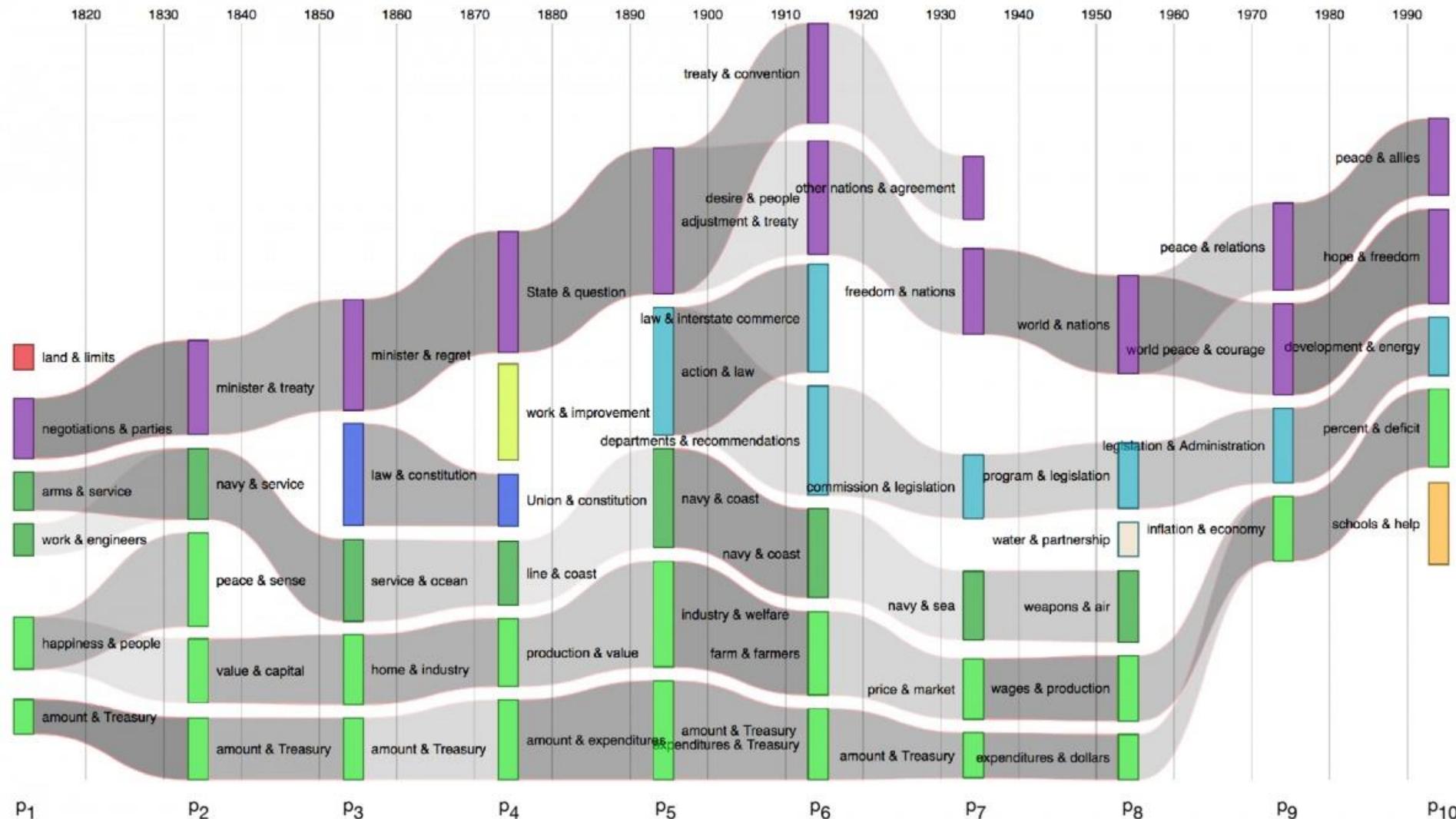


# Word semantics

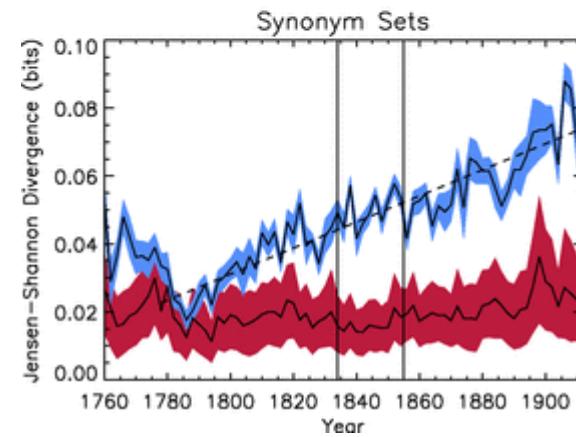
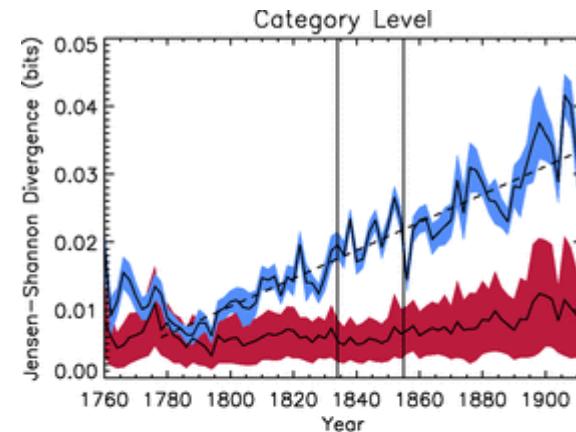
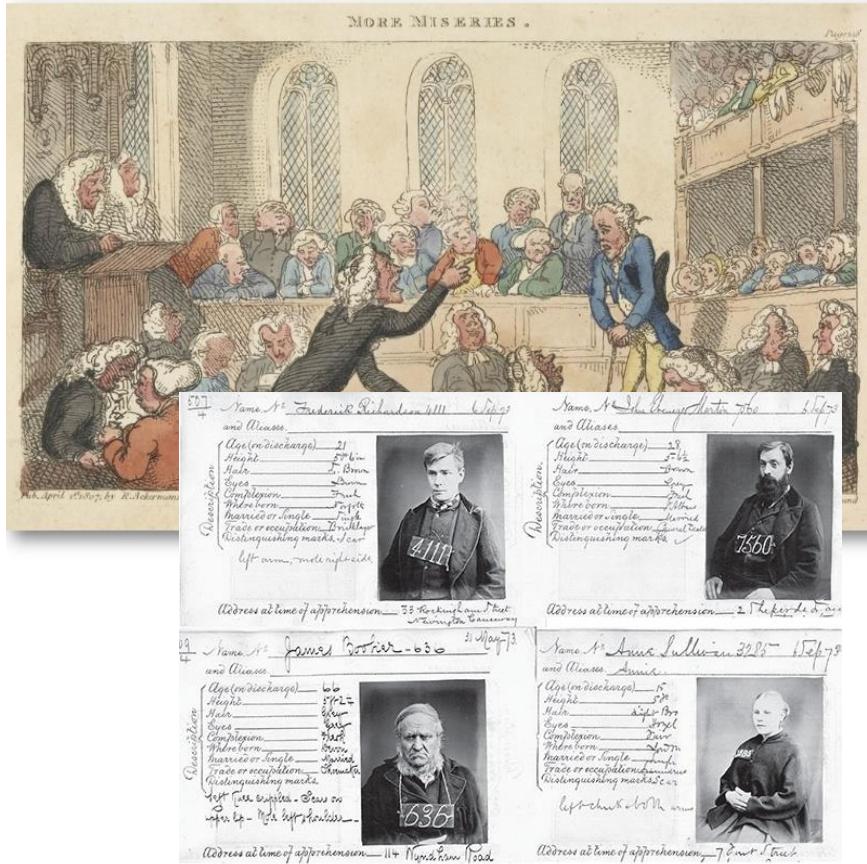
- word2vec



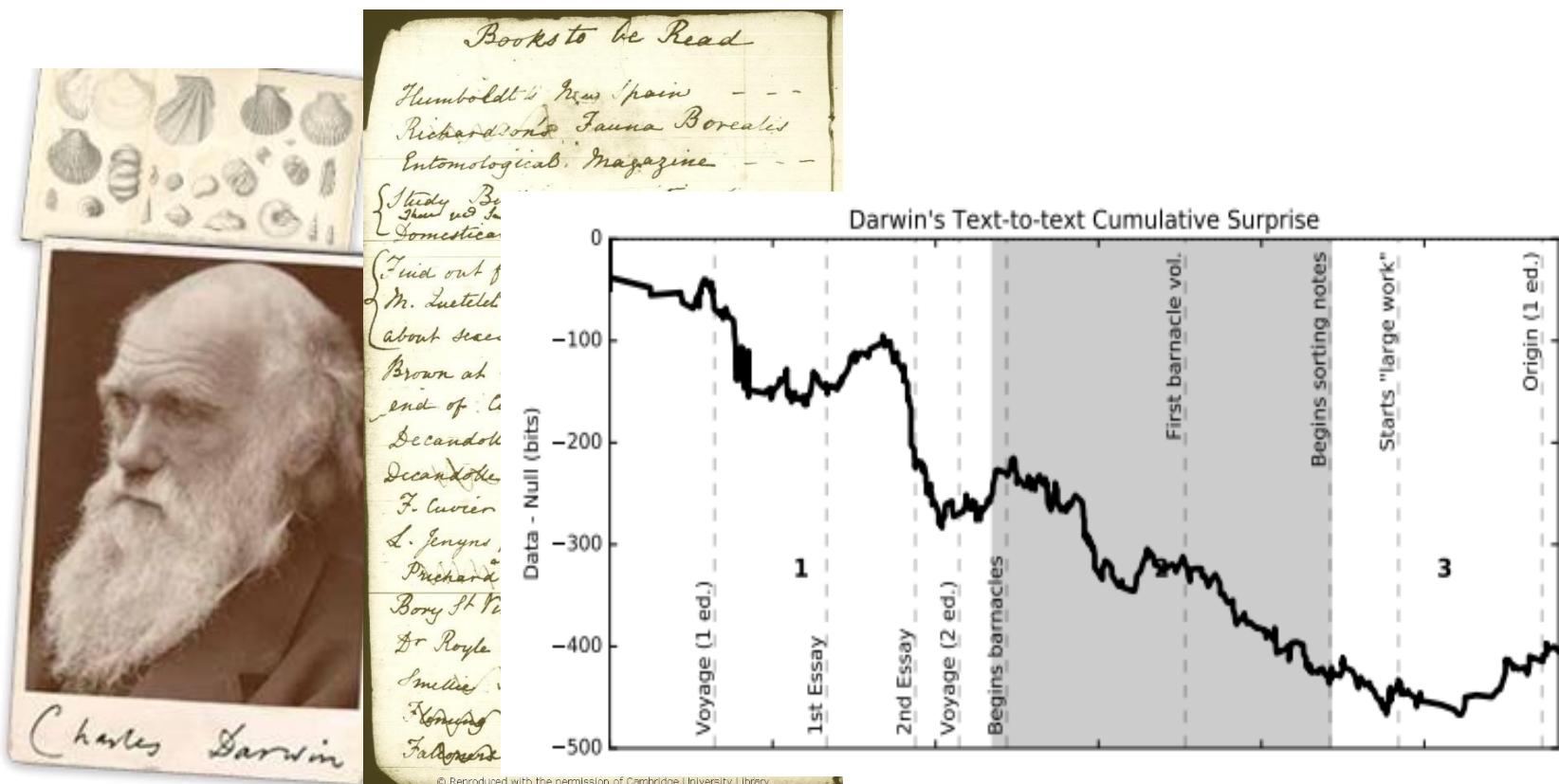
# Presidential speeches in USA



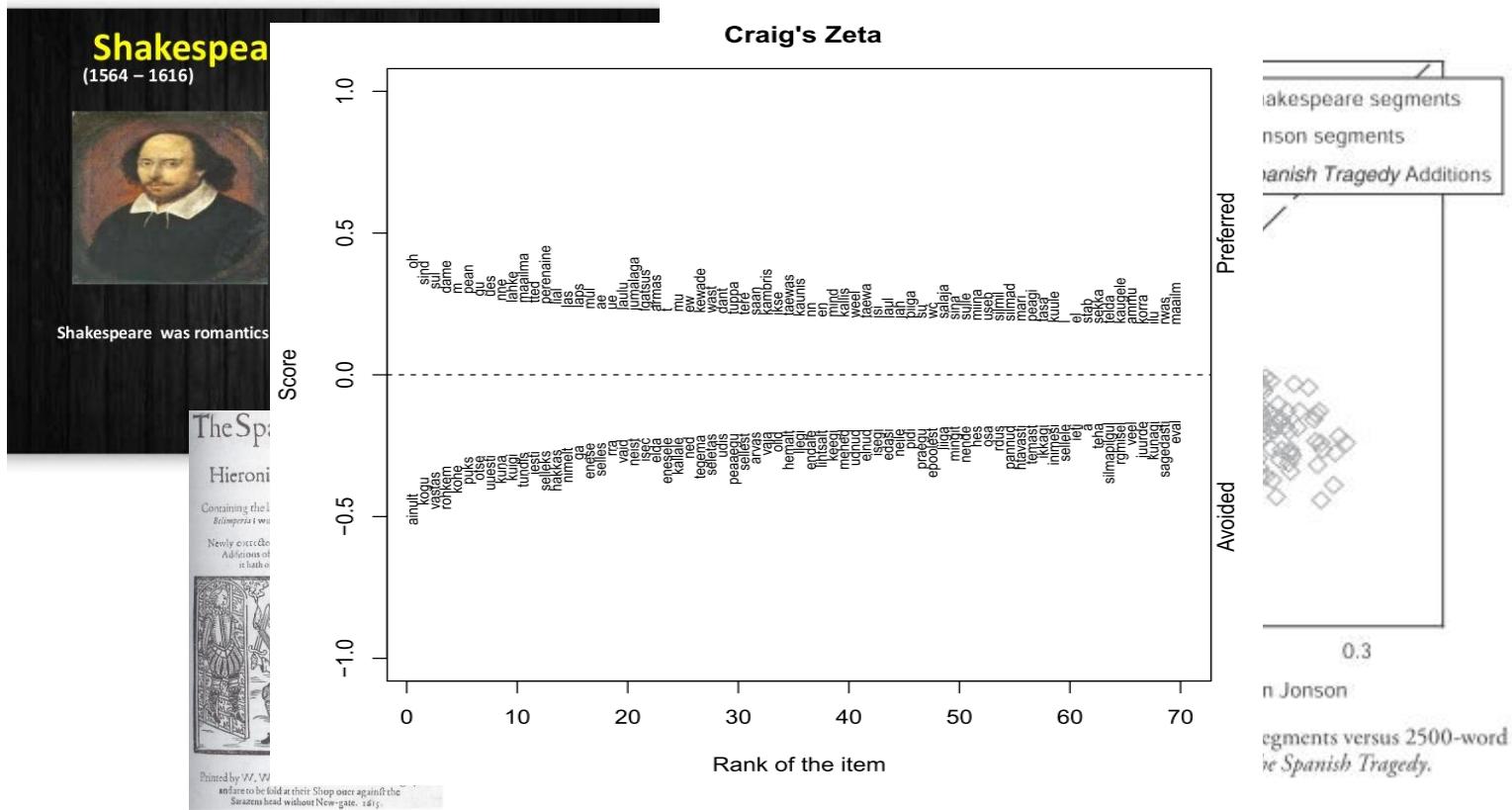
# Civilizing process



# Darwin's reading habits



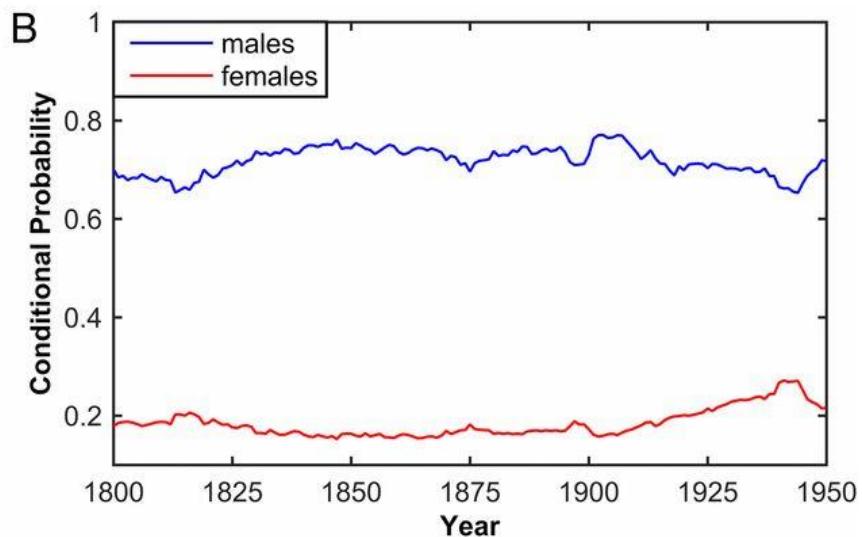
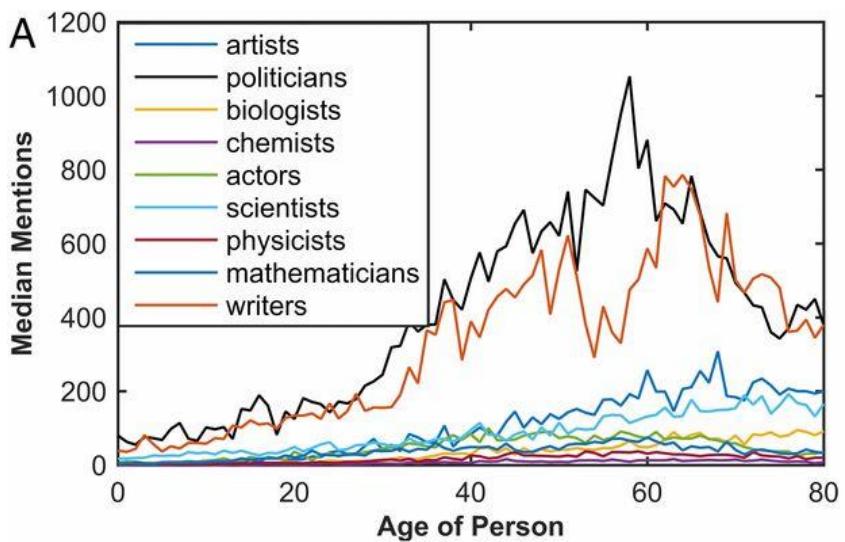
# Finding the author



Metainformation!

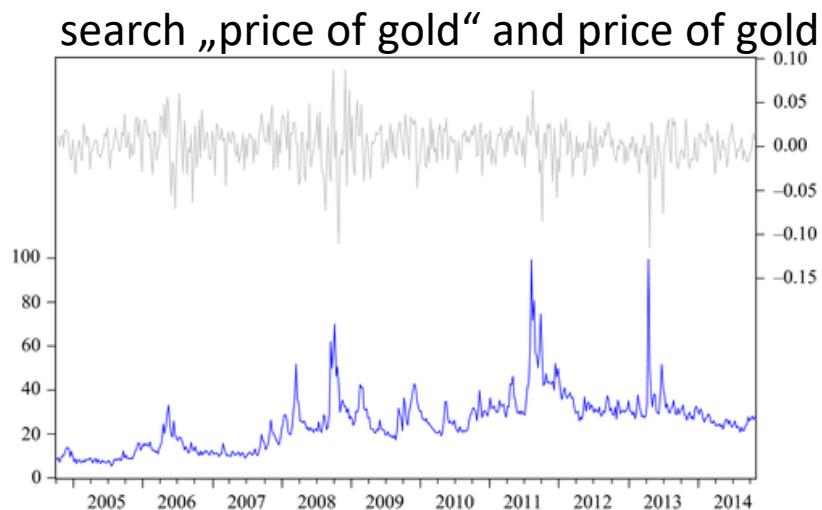
# 150 years of newspapers

- Linked to biographical data

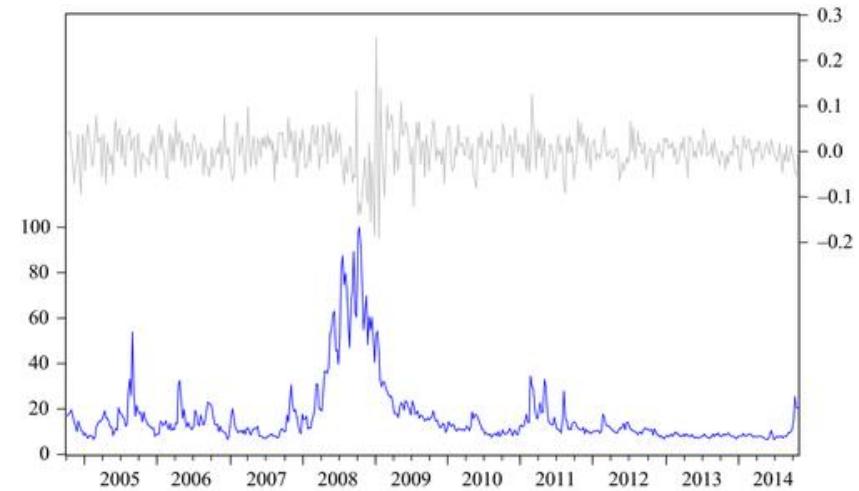


# Google searches

- Linked to exact time and economy

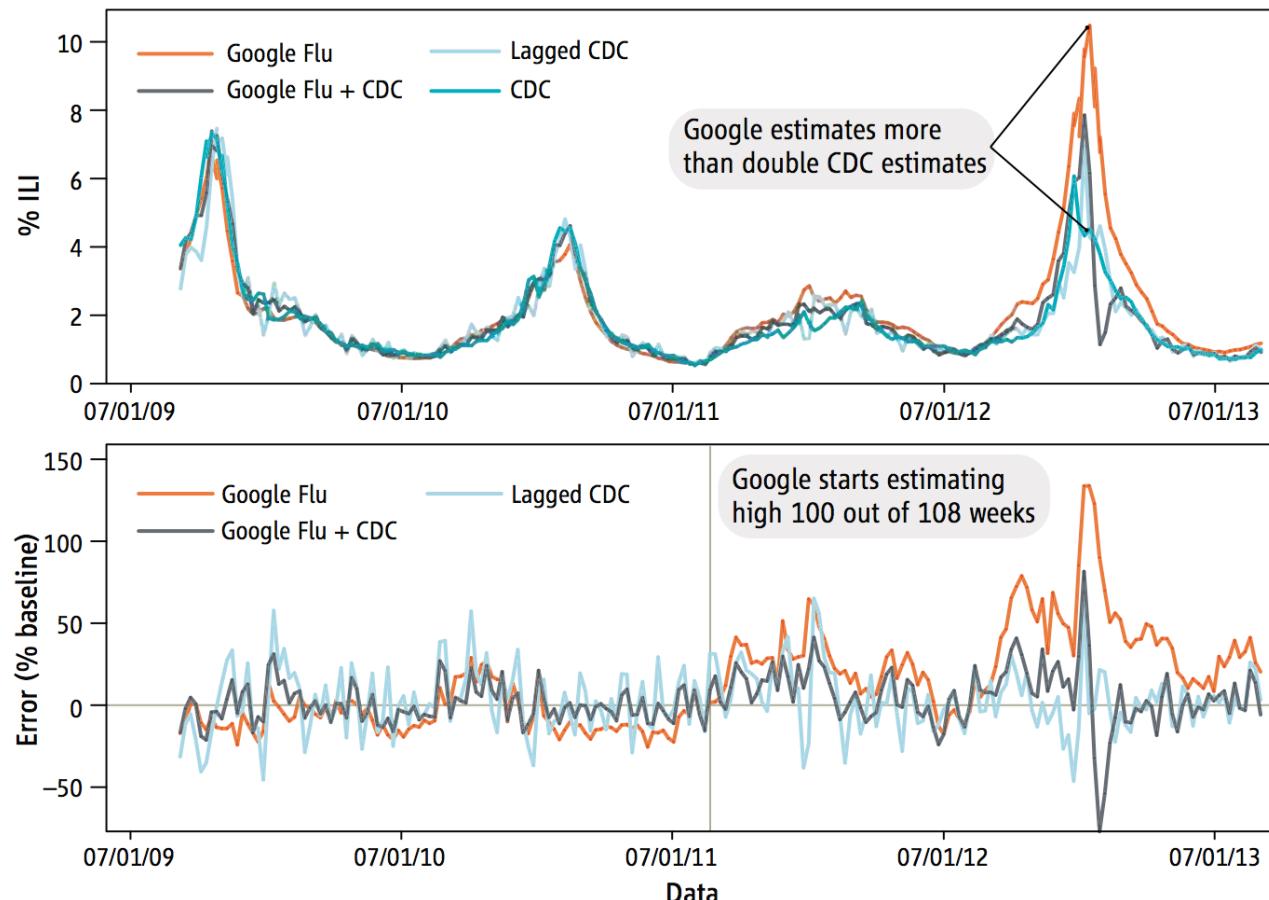


search „price of oil“ and price of oil



# Google Flu trends

- Search for „flu“ linked with medical estimates

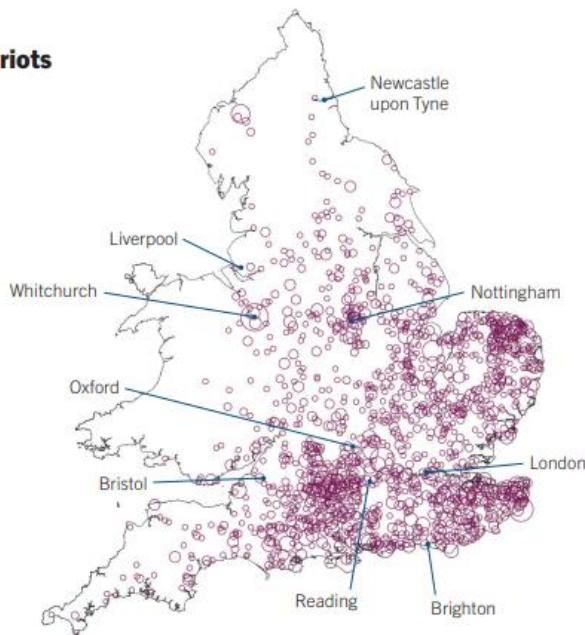


# Deep mining history

- Advertisements in newspapers -> diffusion of technology
- „Swing riots“ 1829-1830: machines, unemployment and crops.

**Figure 1:**  
**Location of Swing riots**

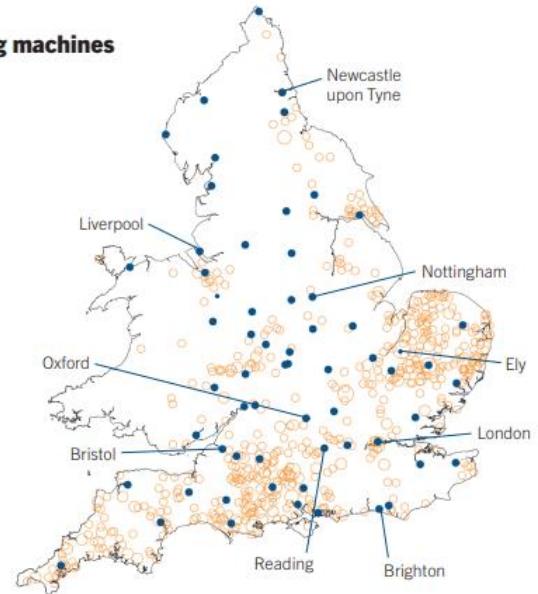
- 1 Riot
- 26 Riots  
e.g. Whitchurch (Shropshire)



Notes: Location of Swing riots. Purple circles identify parishes with Swing riots; the size of the circles is proportional to the number of episodes recorded in each parish.

**Figure 2:**  
**Location of threshing machines**

- 1 Machine
- 5 Machines  
e.g. Ely (Cambridgeshire)
- Newspaper ads



Notes: Location of threshing machines. Orange circles identify parishes with threshing machines; the size of the circles is proportional to the number of machines we found. Blue dots show cities that published newspaper advertisements at the time.