



HOW TO TURN WORDS INTO STATISTICS?

Sandeep Soni

09/03/2024

Today's slides draw heavily from David Bamman's slides on the same topic

AGENDA

- What is a word?
- What is a vocabulary?
- How to segment text into words?
- What are the issues that we should care about?
- How to compare different text collections?



Dog



WORDS REPRESENT CONCEPTS

- Words are symbolic representations of a concept



WORDS REPRESENT CONCEPTS

- Words are symbolic representations of a concept
 - Concise



WORDS REPRESENT CONCEPTS

- Words are symbolic representations of a concept
 - Concise
 - Shared



WORDS REPRESENT CONCEPTS

- Words are symbolic representations of a concept
 - Concise
 - Shared
 - Memorable



WORDS REPRESENT CONCEPTS

- Words are symbolic representations of a concept
 - Concise
 - Shared
 - Memorable
- Words are a form of dimensionality reduction



WHAT ARE WORDS?

- May the Force be with you
- I'm going to make him an offer he can't refuse.
- And I says, “What kinda action?”
- I'da rather die-died then do that

Word boundaries are not clear-cut in other languages.

WHAT ABOUT THESE?

WHAT ABOUT THESE?

- #MarchMadness

WHAT ABOUT THESE?

- #MarchMadness
- @KamalaHarris

WHAT ABOUT THESE?

- #MarchMadness
- @KamalaHarris
- 9423164156

WHAT ABOUT THESE?

- #MarchMadness
- @KamalaHarris
- 9423164156
- <https://www.google.com>

TYPES AND TOKENS

TYPES AND TOKENS

- Types: Symbolic representation of real-world concepts
- Tokens: Instances of types
- Tokenization: process to segment text (or speech) into tokens

TYPES AND TOKENS

- Types: Symbolic representation of real-world concepts
- Tokens: Instances of types
- Tokenization: process to segment text (or speech) into tokens

“to be or not to be”

TYPES AND TOKENS

- Types: Symbolic representation of real-world concepts
- Tokens: Instances of types
- Tokenization: process to segment text (or speech) into tokens

“to be or not to be”

Types	to, be, or, not	4
Tokens	to, be, or, not, to, be	6

VOCABULARY

- Types are the unique tokens in a text
- Types are what we typically use to define a lexicon or a vocabulary or a dictionary.



HOW CAN WE SEGMENT TEXT INTO TOKENS?

JUST SEARCH?

JUST SEARCH?

Start with a long vocabulary and then lookup

JUST SEARCH?

JUST SEARCH?

- Fast: Linear time to tokenize any text
- Easy: Just lookup in a dictionary

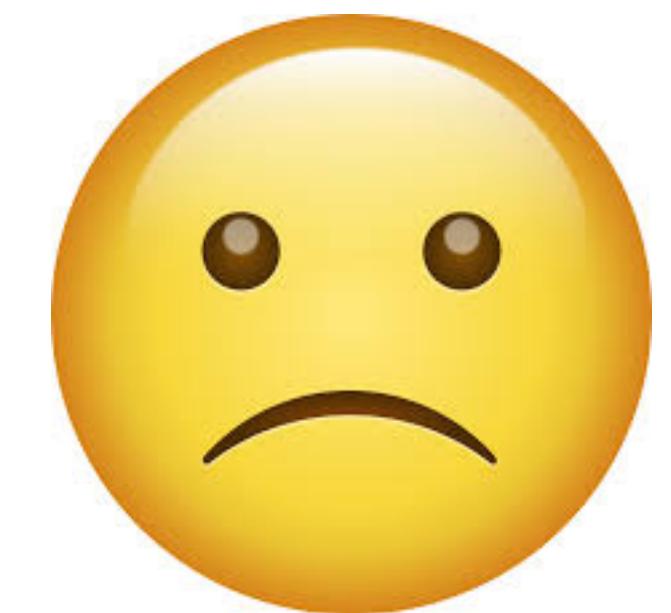


JUST SEARCH?

- Fast: Linear time to tokenize any text
- Easy: Just lookup in a dictionary



- Not robust: the vs them
- Brittle: supercalifragilisticexpialidocious
- Unmanageable: move, moves, moving



WHITESPACE

```
text.split(" ")
```

It is a truth universally acknowledged, that a single man in possession of a good fortune, must be in want of a wife.

WHITESPACE

```
text.split(" ")
```

It is a truth universally acknowledged, that a single man in possession of a good fortune, must be in want of a wife.

- Delimiting on whitespace often misses to separate punctuations from words
- grep -E -o 'wife.?' 853 wife
439 wife,
229 wife.
93 wife'
61 wife
43 wife;
36 wife?
28 wife
27 wife-
17 wife!
7 wifey
7 wifel
7 wife'
4 wife_
2 wifer
2 wife-
1 wifeh
1 wife

PUNCTUATION

- Naively removing punctuation can lead to loss of information about:
 - sentence boundaries
 - pauses
 - parentheticals or asides

WHITESPACE

WHITESPACE

- Fast
- Easy
- Less brittle



WHITESPACE

- Fast
- Easy
- Less brittle



- Lossy or noisy segmentation



REGULAR EXPRESSIONS

- Instead of separating tokens by whitespace, define a pattern for a token
- A pattern like 'wife.' is called a regular expression

REGULAR EXPRESSIONS

- More formally, regular expressions are a compact language to specify search patterns in text

/single/

It is a truth universally acknowledged, that a **single** man in possession of a good fortune, must be in want of a wife.

REGULAR EXPRESSIONS

- More formally, regular expressions are a compact language to specify search patterns in text

/ ?n /

It is a truth **un**iversally acknowledged, that a **sing**le **man** in possession of a good **fortu**ne, must be in **want** of a wife.

REGULAR EXPRESSIONS

- More formally, regular expressions are a compact language to specify search patterns in text

/ [Ii]t/

It is a truth universally acknowledged, that a single man in possession of a good fortune, must be in want of a wife.

REGULAR EXPRESSIONS

regex	matches	doesn't match
/the/	the, isothermally	The
/ [Tt]he/	the, isothermally, The	
/\b[Tt]he\b/	the, The	—The

Slide credit to David Bamman

REGULAR EXPRESSIONS

- Use brackets to specify a range or specific alternations

Regex	Range	Example matches
[Ii]t	It or it	It, it
[a-z]	{a,b,c,...,z}	for, she,hhh
[0-9]	{0,1,2,3,...,9}	123, 0123
[A-Z]	{A,B,C,...,Z}	FOR, SHE, HHH
[a-zA-Z0-9]	{a,b,...,z,A,B,...,Z,0,1,...,9}	1aSBv2

REGULAR EXPRESSIONS

Term	Meaning	Sample regex	Matches
+	one or more	he+y	hey, heeeeeey
?	optional	colou?r	color, colour
*	zero or more	toys*	toy, toys, toysss

Slide credit to David Bamman

SYMBOLS

Symbol	Function
\b	Word boundary (zero width)
\d	Any decimal digit (equivalent to [0-9])
\D	Any non-digit character (equivalent to [^0-9])
\s	Any whitespace character (equivalent to [\t\n\r\f\v])
\S	Any non-whitespace character (equivalent to [^\t\n\r\f\v])
\w	Any alphanumeric character (equivalent to [a-zA-Z0-9_])
\W	Any non-alphanumeric character (equivalent to [^a-zA-Z0-9_])
\t	The tab character
\n	The newline character

COMPOSITION OF REGEXES

- You can create complex regular expressions from simple ones by using disjunction

/single| [Ii]t

It is a truth universally acknowledged, that a **single** man in possession of a good fortune, must be in want of a wife.

PYTHON

- Use the `re` module for regex searches
- `re.findall(regex, text)` will find non-overlapping matches
- `re.findall("[It]", "If it is meant to be it is meant to be")`
- The function returns a list of all the matches: `[It, it]`
- Also see `re.match`, `re.compile`
- <https://www.dataquest.io/wp-content/uploads/2019/03/python-regular-expressions-cheat-sheet.pdf>

BACK TO TOKENIZATION

<https://spacy.io/usage/spacy-101#annotations-token>

SENTENCE SEGMENTATION

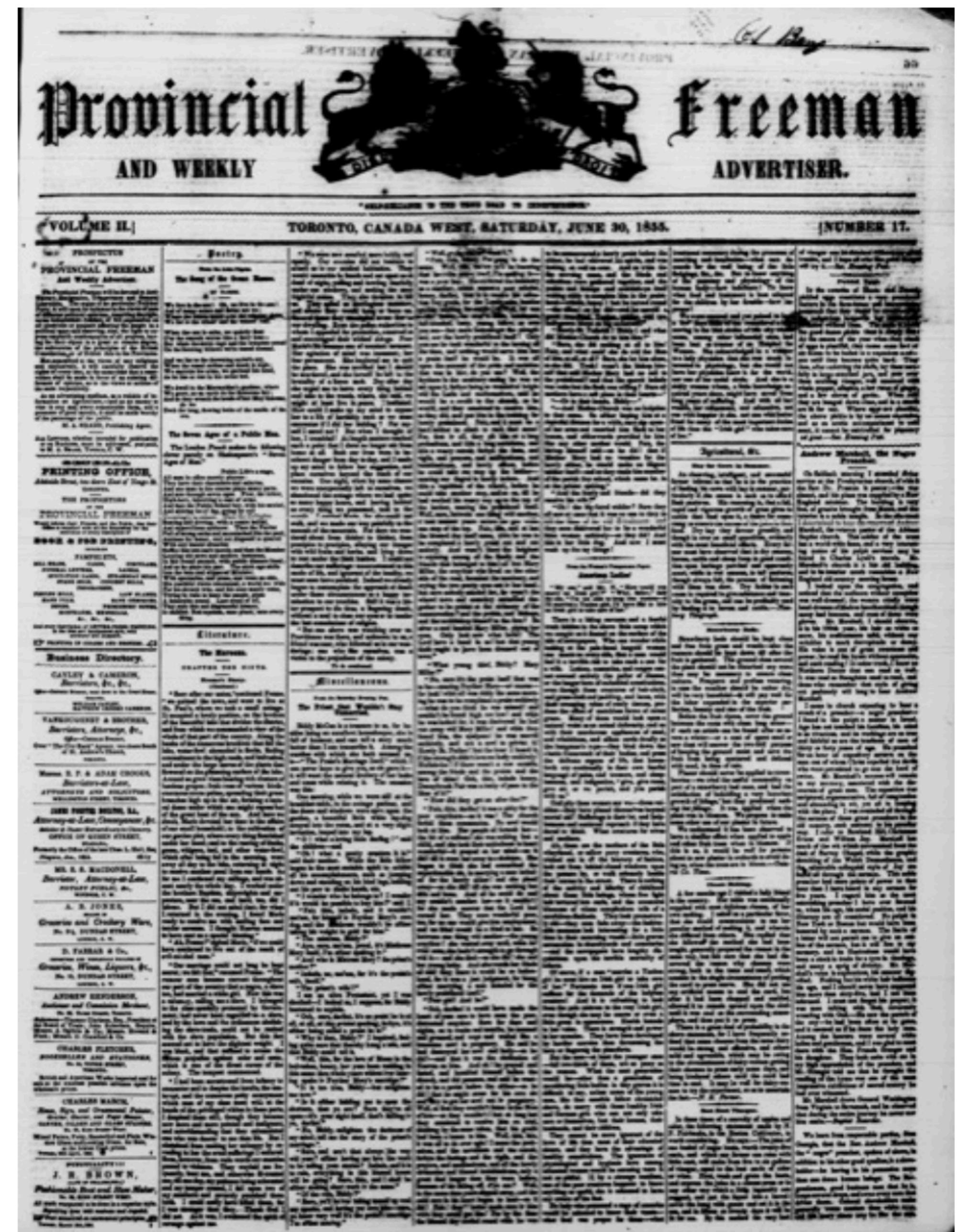
- This is usually not done using regular expressions
- He lives in the U.S. John, however, lives in Canada.
- Typically done using a combination of dependency parsing and unsupervised learning

STEMMING AND LEMMATIZATION

- Languages have rich inflectional and derivational morphology
- Stemming is a heuristic to clip the suffixes
 - *arguing, argues, argue, argued* → *argu*
- Lemmatization is more principled as it collapses inflectional forms to a typical dictionary entry
 - *arguing, argues, argue, argued* → *argue*

POST PROCESSING

- What counts as tokens can have important consequences
- In historical texts, tokenization errors need to be fixed as a post-processing step
- senatoradmits -> senator admits

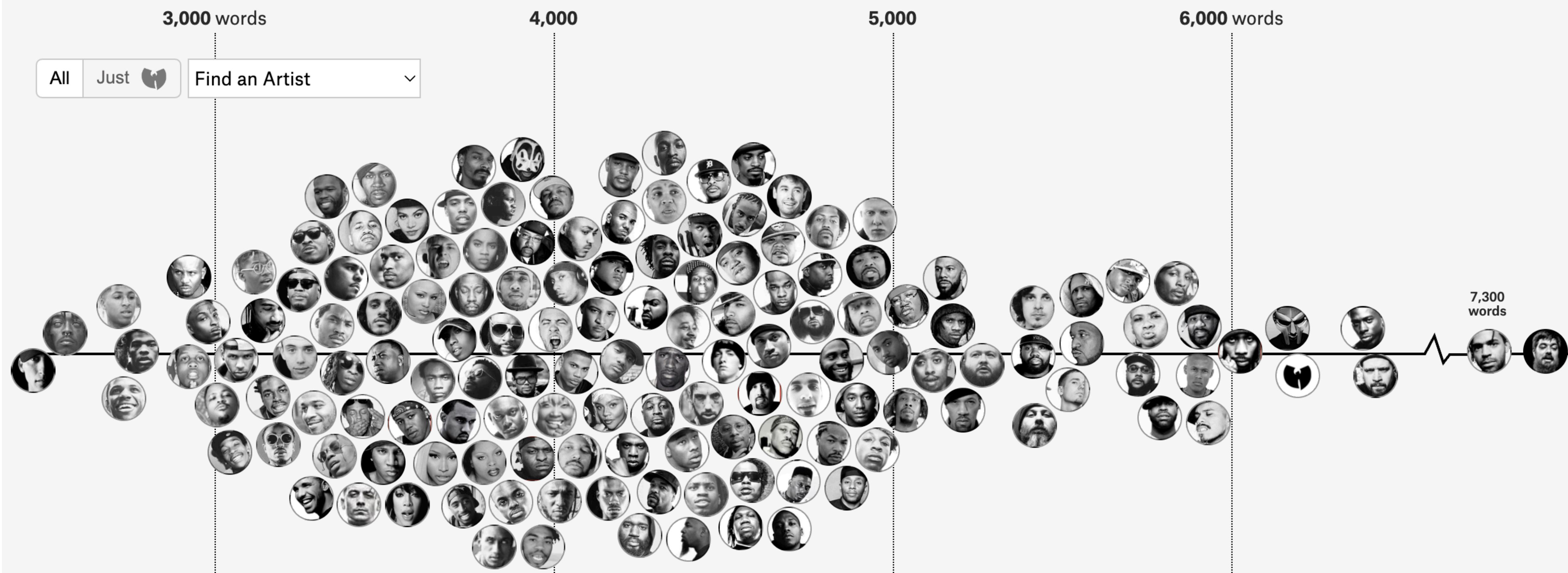


WORD ANALYSIS

- Vocabulary size and diversity
- Subword analysis
- Word frequency

VOCABULARY SIZE

of Unique Words Used Within Artist's First 35,000 Lyrics

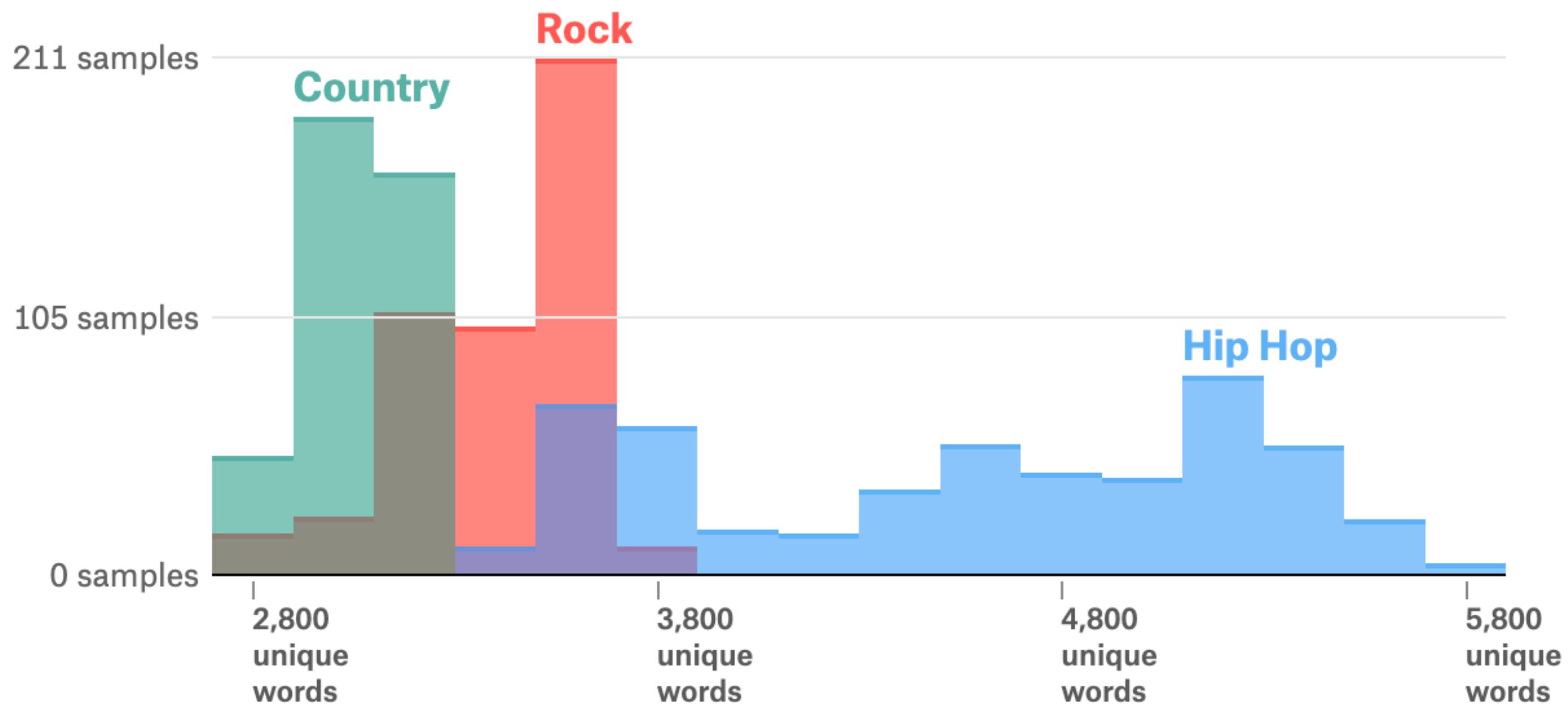


of Unique Words Used Within Artist's First 35,000 lyrics

BY ERA¹

1980s | 1990s | 2000s | 2010s

Run-D.M.C.											
2Pac	Biz Markie										
Big L	Ice T										
Insane Clown...	Rakim										
MC Lyte	Brand Nubian										
Foxy Brown	Geto Boys										
Juvenile	Three 6 Mafia	Beastie Boys									
Master P	UGK	Jay-Z	Big Daddy Kane								
Salt-n-Pepa	Dizzee Rascal	Mobb Deep	LL Cool J								
Snoop Dogg	Jadakiss	Outkast	Busta Rhymes								
Eve	Kano	Public Enemy	Cypress Hill								
Gucci Mane	Lil' Kim	Cam'ron	De La Soul								
Kanye West	Nelly	Eminem	Fat Joe								
Lil Wayne	Rick Ross	The Game	Gang Starr								
Bone Thugs-n...	Missy Elliot	T.I.	Joe Budden	KRS-One							
50 Cent	Trick Daddy	2 Chainz	Kevin Gates	Method Man							
Juicy J	Trina	A\$AP Ferg	Royce da 5'9	A Tribe Call...							
Drake	Young Jeezy	Big KRIT	Tech n9ne	Atmosphere							
Future	Big Sean	Brockhampton	Twista	Ludacris	Common						
DMX	Kid Cudi	BoB	Cupcakke	Ab-Soul	Das EFX	Del the Funk...					
21 Savage	Kid Ink	Childish Gam...	Hopsin	A\$AP Rocky	Mos Def	The Roots					
A Boogie wit...	Kodak Black	G-Eazy	Jay Rock	Danny Brown	E-40	Blackalicious					
Lil Baby	Lil Yachty	J Cole	Kendrick Lamar	Death Grips	Goodie Mob	Canibus					
Lil Durk	Logic	Machine Gun ...	Mac Miller	Denzel Curry	Kool G Rap	Ghostface Ki...					
Wiz Khalifa	Migos	Meek Mill	ScHoolboy Q	\$uicideboy\$	Nas	Immortal Tec... GZA					
Lil Uzi Vert	YG	Nicki Minaj	Tyga	Flatbush Zom...	Talib Kweli	Rzaekwon					
NF	YoungBoy Nev...	Russ	Vince Staples	Tyler the Cr...	Brother Ali	Immortal Tec... GZA					
<2,675 unique words	2,675-3,050 unique words	3,050-3,425 unique words	3,425-3,800 unique words	3,800-4,175 unique words	4,175-4,550 unique words	4,550-4,925 unique words	4,925-5,300 unique words	5,300-5,675 unique words	5,675-6,050 unique words	6,050-6,425 unique words	6,425+ unique words



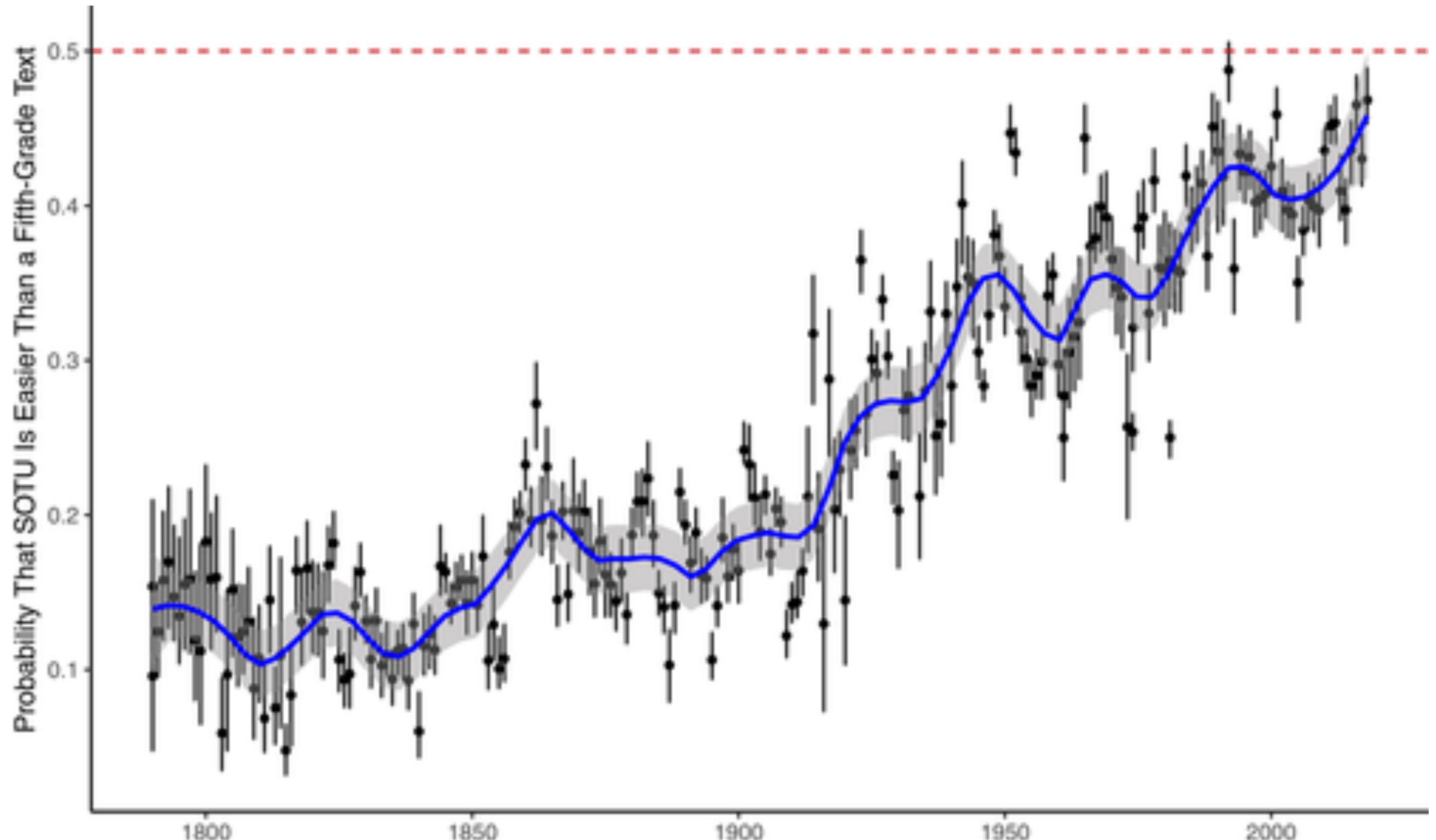
The contraction in vocabulary could be because of the style and structure of some genres

I dumbed down for my audience to double my dollars
 They criticized me for it, yet they all yell "holla"
 If skills sold, truth be told, I'd probably be
 Lyrically Talib Kweli
 Truthfully I wanna rhyme like Common Sense
 But I did 5 mil - I ain't been rhyming like Common since

Or it could be that artists make a genuine choice to make songs simpler and shorter

WORD CHARACTERISTICS

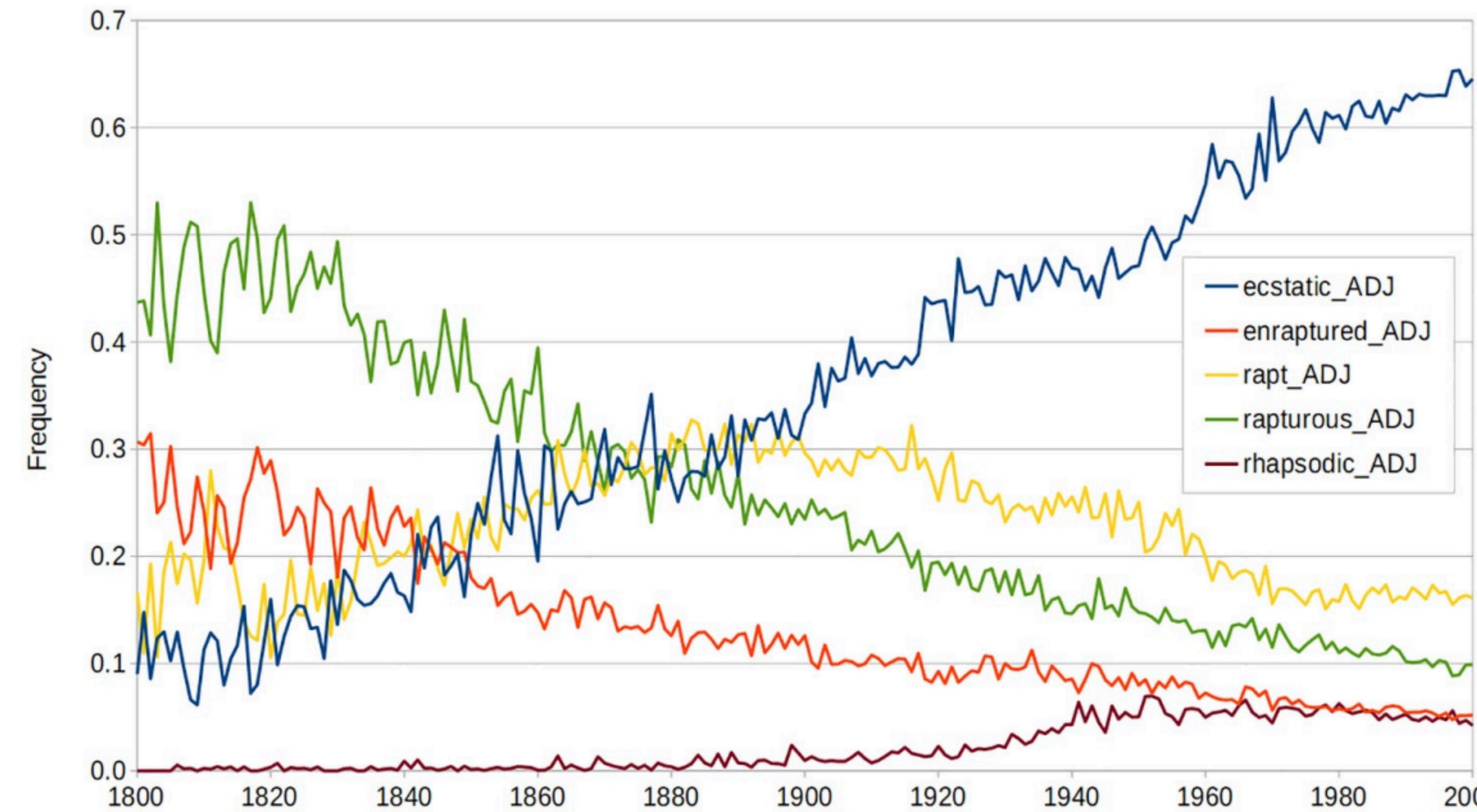
Source of Complexity
Long Words
Mean characters per word
Words with at least 7 characters
Words with at least 6 characters
Mean syllables per word
Words with at least 3 syllables
Words with fewer than 3 syllables
Words with 2 syllables
Words with 1 syllable
Rare Words
Google Books baseline usage
Brown corpus baseline usage



Benoit et. al. (2019) Measuring and Explaining Political Sophistication through Textual Complexity

WORD FREQUENCY

- We can create frequency profiles of words and compare them based on their frequencies



HOW TO COMPARE GROUPS USING WORDS?

PROBLEM SETUP

- Assume we have text from two “groups”
 - Twitter in 2014 Vs Twitter in 2024
 - Fiction Vs Non-fiction
 - Donald Trump Speeches Vs Speeches from all other presidents
- Find what differentiates the groups using words and their statistics

DIFFERENCE IN PROPORTIONS

w	Indexes a word from our lexicon
k_i	Indexes a category
c_{w,k_i}	Count of w in category k_i
$f_{w,k_i} = \frac{c_{w,k_i}}{\sum_w c_{w,k_i}}$	Normalized count or proportion of w in k_i

$f_{w,k_1} - f_{w,k_2}$ is the difference in proportions across two groups

DIFFERENCE IN PROPORTIONS

DIFFERENCE IN PROPORTIONS

w="students"

law of large numbers is
a rule in math that
students learn in school
It is not a law though

k₁="maths"

The law school teaches
law to **students** who
want to learn the law
and practice law in their
future careers

k₂="legal"

DIFFERENCE IN PROPORTIONS

w="students"

law of large numbers is
a rule in math that
students learn in school
It is not a law though

$k_1 = \text{"maths"}$

$c_{\text{students, maths}} = 1$

The law school teaches
law to **students** who
want to learn the law
and practice law in their
future careers

$k_2 = \text{"legal"}$

$c_{\text{students, legal}} = 1$

DIFFERENCE IN PROPORTIONS

w="students"

law of large numbers is
a rule in math that
students learn in school
It is not a law though

k₁="maths"

$$c_{\text{students, maths}} = 1$$

$$f_{\text{students, maths}} = \frac{1}{20}$$

The law school teaches
law to **students** who
want to learn the law
and practice law in their
future careers

k₂="legal"

$$c_{\text{students, legal}} = 1$$

$$f_{\text{students, legal}} = \frac{1}{20}$$

DIFFERENCE IN PROPORTIONS

w="students"

law of large numbers is
a rule in math that
students learn in school
It is not a law though

k₁="maths"

^cstudents, maths = 1

*f*students, maths = $\frac{1}{20}$

$$f_{\text{students, maths}} - f_{\text{students, legal}} = \frac{1}{20} - \frac{1}{20} = 0$$

The law school teaches
law to **students** who
want to learn the law
and practice law in their
future careers

k₂="legal"

^cstudents, legal = 1

*f*students, legal = $\frac{1}{20}$

DIFFERENCE IN PROPORTIONS

w=“law”

law of large numbers is
a rule in math that
students learn in school
It is not a law though

k_1 =“maths”

$c_{\text{law, maths}} = ?$

$f_{\text{law, maths}} = ?$

$f_{\text{law, maths}} - f_{\text{law, legal}} = ?$

The law school teaches
law to students who
want to learn the law
and practice law in their
future careers

k_2 =“legal”

$c_{\text{law, legal}} = ?$

$f_{\text{law, legal}} = ?$

DIFFERENCE IN PROPORTIONS

w=“law”

law of large numbers is
a rule in math that
students learn in school
It is not a law though

k₁=“maths”

$${}^c\text{law, maths} = 2$$

$$f_{\text{law, maths}} = \frac{2}{20}$$

$$f_{\text{law, maths}} - f_{\text{law, legal}} = \frac{2}{20} - \frac{4}{20} = -0.1$$

The law school teaches
law to students who
want to learn the law
and practice law in their
future careers

k₂=“legal”

$${}^c\text{law, legal} = 4$$

$$f_{\text{law, legal}} = \frac{4}{20}$$

DIFFERENCE IN PROPORTIONS

- Simple and easy to measure and interpret
- Overemphasizes common words; for common words, there differences are also large
- No correction for chance or determination of statistical significance

$$\chi^2$$

$$\chi^2$$

Does the word “robot” occur **significantly** more frequently in science fiction?

$$\chi^2$$

Does the word “robot” occur **significantly** more frequently in science fiction?

	robot	\neg robot	
sci-fi	104	1004	= 10.3%
\neg sci-fi	2	13402	= 0.015%

Slide credit: David Bamman's Info 256 class

$$\chi^2$$

We can calculate the following statistic, which is the sum of squared difference between the observed value in each cell and the expected value assuming independence

$$\chi^2 = \sum_{i,j} \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

χ^2

	robot	\neg robot	sum	frequency
sci-fi	104	1004	1108	0.076
\neg sci-fi	2	13402	13404	0.924
sum	106	14406		
frequency	0.007	0.993		

Assuming independence:

$$\begin{aligned} P(\text{robot, scifi}) &= P(\text{robot}) \times P(\text{scifi}) \\ &= 0.007 \times 0.076 = 0.00053 \end{aligned}$$

Among 14512 words, we would expect to see 7.69 occurrences of *robot* in sci-fi texts.

	robot	\neg robot	$P(\text{scifi})$	$P(\neg\text{scifi})$
sci-fi	7.69	1095.2	0.076	
\neg sci-fi	93.9	13315.2		0.924

$P(\text{robot})$ $P(\neg\text{robot})$

0.007	0.993
-------	-------

$$\chi^2$$

	robot	\neg robot
sci-fi	104	1004
\neg sci-fi	2	13402

	robot	\neg robot
sci-fi	7.69	1095.2
\neg sci-fi	93.9	13315.2

Left is observed counts; right is expected counts assuming complete independence

$$\chi^2 = \sum_{i,j} \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

$$\chi^2$$

How different are these two tables?

	robot	\neg robot
sci-fi	104	1004
\neg sci-fi	2	13402

	robot	\neg robot
sci-fi	7.69	1095.2
\neg sci-fi	93.9	13315.2

Left is observed counts; right is expected counts assuming complete independence

$$\chi^2 = \sum_{i,j} \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

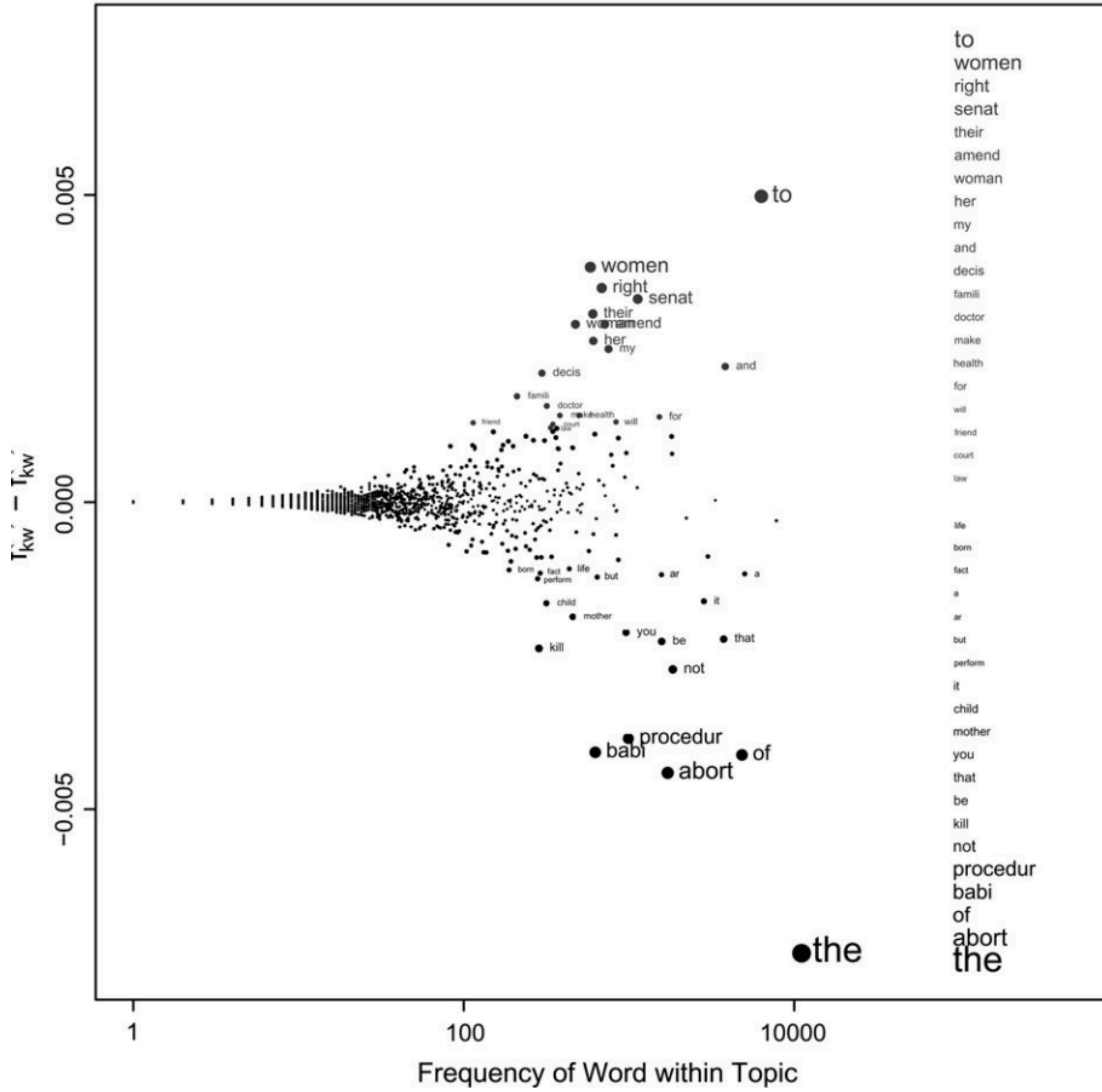
$$\chi^2$$

- Useful statistic to find differentiating markers
- We have a way to test for statistical significance
- Assumes each word is independent from the others

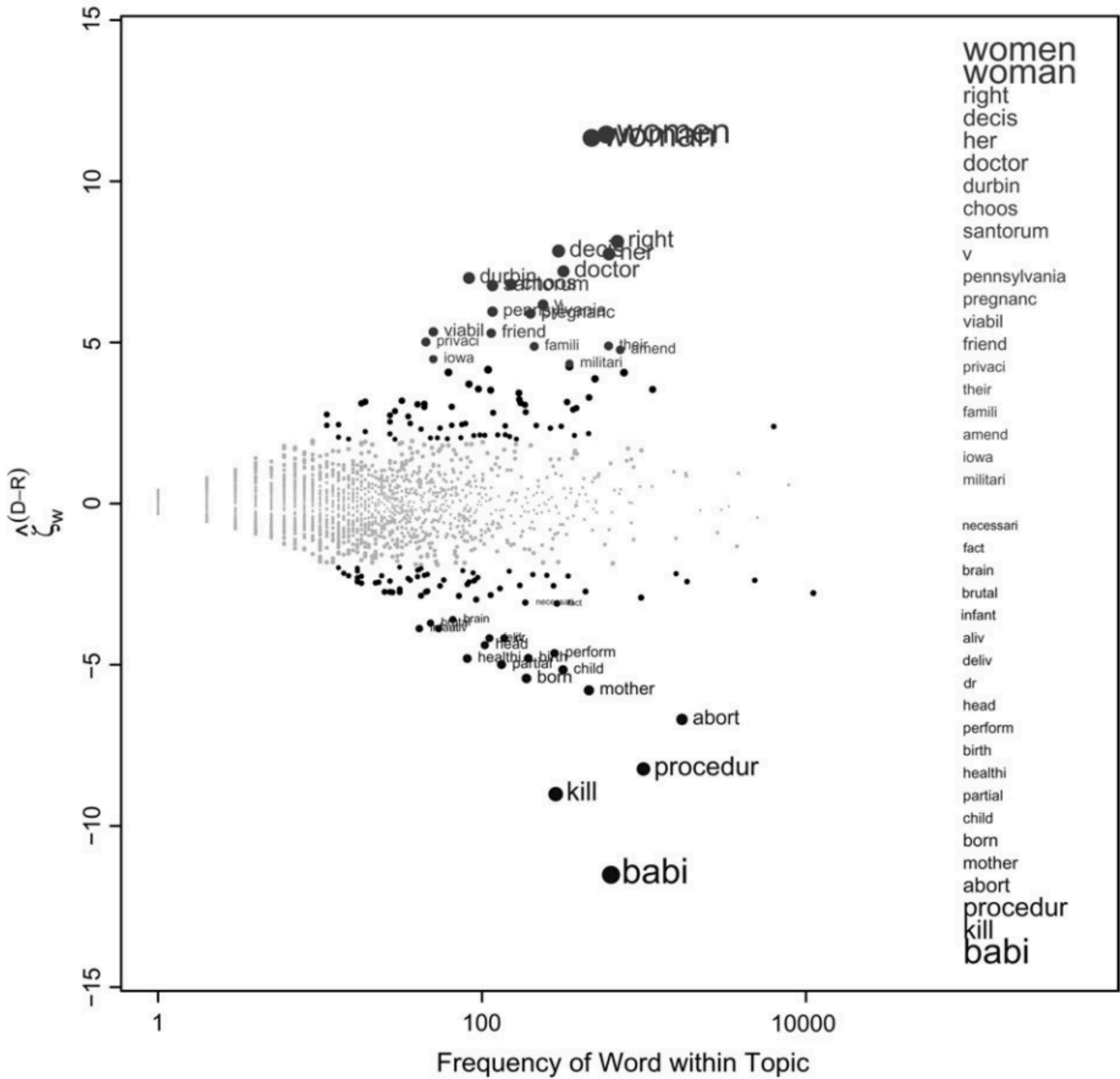
OTHER METHODS

- Many other methods to characterize differences
- Model based methods that assume parametric distributions and Bayesian priors (e.g., Monroe et. al. 2009)
- Unsupervised and Bayesian (e.g., SAGE; Eisenstein et. al. 2011)
- Supervised learning to learn precise features that are informative in separating categories (e.g., Underwood et. al. 2018)

**Partisan Words, 106th Congress, Abortion
(Difference of Proportions)**



**Partisan Words, 106th Congress, Abortion
(Weighted Log-Odds-Ratio, Informative Dirichlet Prior)**



JOURNAL OF THE AMERICAN STATISTICAL ASSOCIATION

Number 502

JUNE, 1963

Volume 58

INFERENCE IN AN AUTHORSHIP PROBLEM^{1,2}

- Who wrote the disputed 12 federalist papers?
- Compare the frequency of some basic stylistic words (e.g., upon)
- Answer: Madison

A comparative study of discrimination methods applied to the authorship of the disputed *Federalist* papers

FREDERICK MOSTELLER

Harvard University

and

Center for Advanced Study in the Behavioral Sciences

AND

DAVID L. WALLACE

University of Chicago

This study has four purposes: to provide a comparison of discrimination methods; to explore the problems presented by techniques based strongly on Bayes' theorem when they are used in a data analysis of large scale; to solve the authorship question of *The Federalist* papers; and to propose routine methods for solving other authorship problems.

Word counts are the main tool used for discrimination. Since the topic written about heavily influences the rate with which a word is used, care in selection of words is necessary. The filler words of the language such as *an*, *of*, and *upon*, and, more generally, articles, prepositions, and conjunctions provide fairly stable rates, whereas more meaningful words like *war*, *executive*, and *legislature* do not.

After an investigation of the distribution of these counts, the authors execute an analysis employing the usual discriminant function and an analysis based on Bayesian methods. The conclusions about the authorship problem are that Madison rather than Hamilton wrote all 12 of the disputed papers.

The findings about methods are presented in the closing section on conclusions.

This report, summarizing and abbreviating a forthcoming monograph [8], gives some of the results but very little of their empirical and theoretical foundation. It treats two of the four main studies presented in the monograph, and none of the side studies.

¹ This work has been facilitated by grants from The Ford Foundation, the Rockefeller Foundation, and from the National Science Foundation NSF G-13040 and G-10368, contracts with the Office of Naval Research Nonr 1806(37) and 2121(09), and the Laboratory of Social Relations, Harvard University. The work was done in part at the Massachusetts Institute of Technology Computation Center, Cambridge, Massachusetts, and at the Center for Advanced Study in the Behavioral Sciences, Stanford, California. Permission is granted for reproduction in whole or in part for purposes of the United States Government.

² Presented at a session of Special Papers Invited by the Presidents of The American Statistical Association, The Biometric Society (ENAR), and The Institute of Mathematical Statistics at the statistical meetings in Minneapolis, Minnesota, September 9, 1962.

IN CLASS

- Tokenization demo
- Distinctive terms demo