



FINDING DISTINCTIVE TERMS

Sandeep Soni

01/25/2024

STORY SO FAR

- Words are linguistic units that reference concepts
- Type Vs Token
- 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

STORY SO FAR

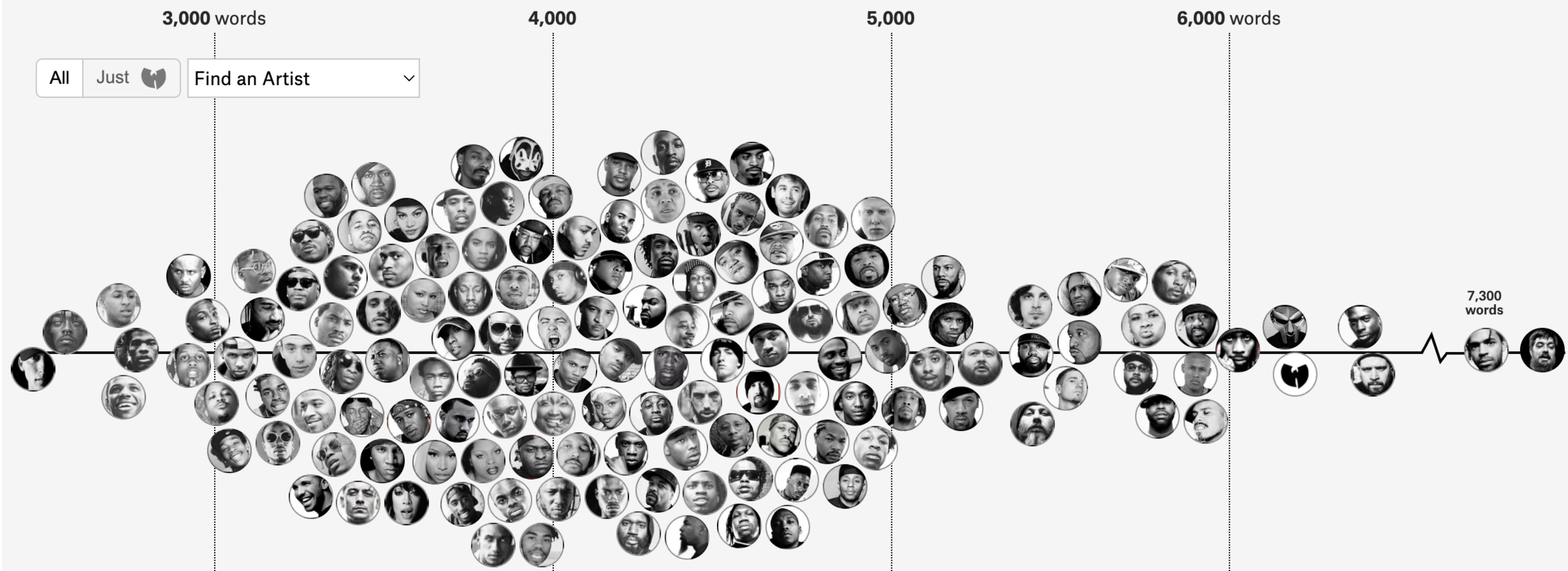
- Tokenization is a typical pre-processing step
- Tokenization is done by matching regular expressions
- Lexicon or vocabulary is a collection of word types
- Stemming/lemmatization can reduce the vocabulary size

WORD ANALYSIS

- Vocabulary size and diversity
- Subword analysis
- Word frequency

VOCABULARY SIZE

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

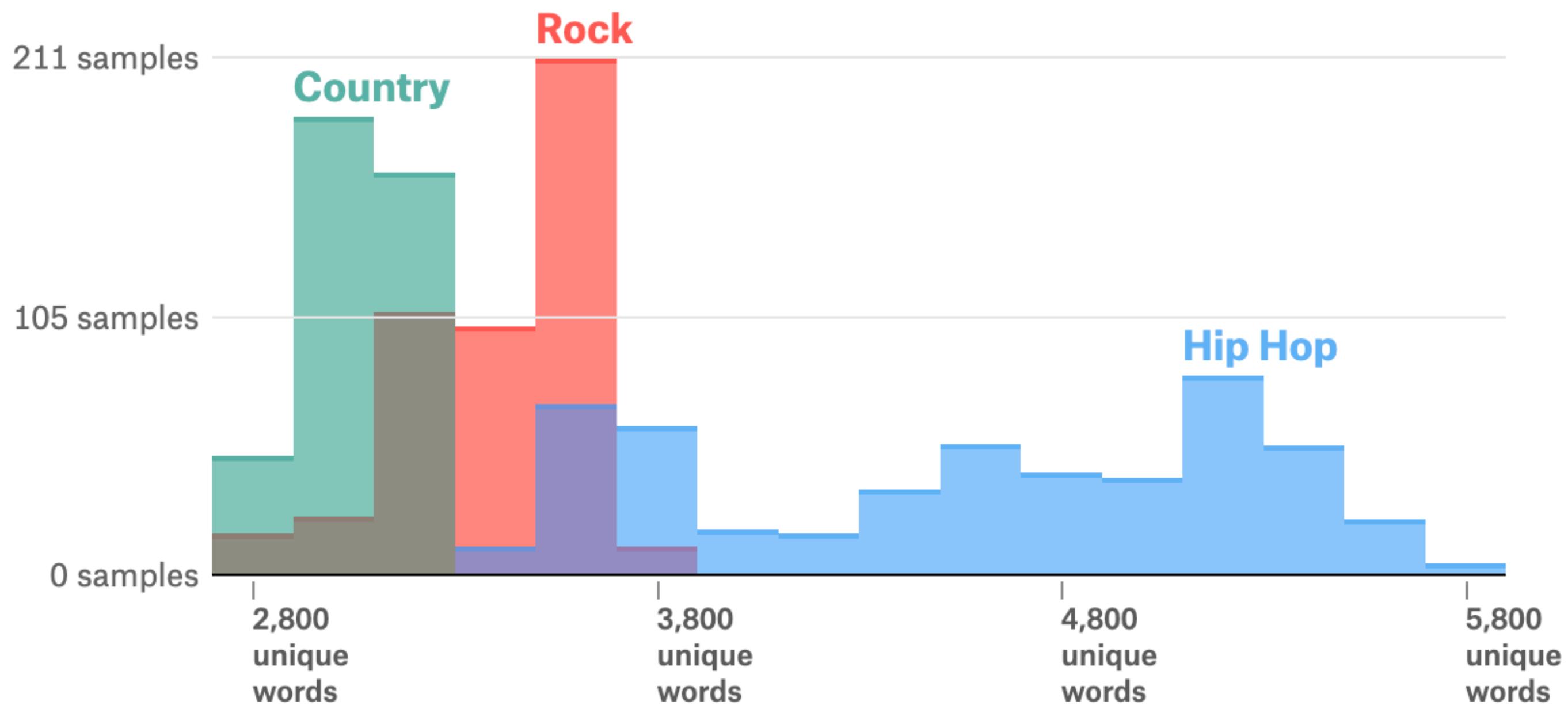


of Unique Words Used Within Artist's First 35,000 lyrics

BY ERA¹

1980s | 1990s | 2000s | 2010s

Run-D.M.C.											
2Pac			Biz Markie			Ice T			Rakim		
Big L			Brand Nubian			Geto Boys			Beastie Boys		
Insane Clown...			MC Lyte			Ice Cube			Big Daddy Kane		
Foxy Brown			Three 6 Mafia			Jay-Z			LL Cool J		
Juvenile			UGK			Dizzee Rascal			Outkast		
Master P			Jadakiss			Public Enemy			Busta Rhymes		
Salt-n-Pepa			Eve			Kano			Cypress Hill		
Snoop Dogg			Gucci Mane			Lil' Kim			Cam'ron		
Kanye West			Kanye West			Nelly			De La Soul		
Lil Wayne			Lil Wayne			Rick Ross			Eminem		
Bone Thugs-n...			Missy Elliot			T.I.			Fat Joe		
50 Cent			Trick Daddy			2 Chainz			The Game		
Juicy J			Trina			A\$AP Ferg			Gang Starr		
Drake			Young Jeezy			Big KRIT			KRS-One		
Future			Big Sean			Brockhampton			Joe Budden		
DMX			Kid Cudi			BoB			Kevin Gates		
21 Savage			Kid Ink			Childish Gam...			Royce da 5'9		
A Boogie wit...			Kodak Black			G-Eazy			Tech n9ne		
Lil Baby			Lil Yachty			J Cole			Atmosphere		
Lil Durk			Logic			Kendrick Lamar			Twista		
Wiz Khalifa			Migos			Machine Gun ...			Ludacris		
Lil Uzi Vert			YG			Meek Mill			Common		
NF			Travis Scott			Nicki Minaj			Del the Funk...		
YoungBoy Nev...			Young Thug			Russ			The Roots		
Run-D.M.C.			Brand Nubian			Geto Boys			Blackalicious		
Ice T			Public Enemy			Public Enemy			Canibus		
Rakim			Kanye West			Eminem			Ghostface Ki...		
Brand Nubian			Fat Joe			The Game			Immortal Tec...		
Geto Boys			Gang Starr			KRS-One			GZA		
Public Enemy			KRS-One			Joe Budden			Common		
Public Enemy			KRS-One			Kevin Gates			Das EFX		
Public Enemy			Method Man			Royce da 5'9			E-40		
Public Enemy			A Tribe Call...			Tech n9ne			Mos Def		
Public Enemy			Atmosphere			Ab-Soul			Lupe Fiasco		
Public Enemy			Ludacris			A\$AP Rocky			Goodie Mob		
Public Enemy			Common			Mos Def			Kool G Rap		
Public Enemy			Das EFX			E-40			Kool Keith		
Public Enemy			E-40			Nas			Nas		
Public Enemy			Goodie Mob			Goodie Mob			Raekwon		
Public Enemy			Kool G Rap			Kool G Rap			Redman		
Public Enemy			Kool Keith			Kool Keith			Raekwon		
Public Enemy			Nas			Nas			Immortal Tec...		
Public Enemy			Raekwon			Raekwon			GZA		
Public Enemy			Immortal Tec...			Immortal Tec...			Jean Grae		
Public Enemy			GZA			GZA			Wu-Tang Clan		
Public Enemy			Jean Grae			Jean Grae			Killah Priest		
Public Enemy			Wu-Tang Clan			Killah Priest			Jedi Mind Tr...		
Public Enemy			Killah Priest			Jedi Mind Tr...			Aesop Rock		
Public Enemy			MF DOOM			MF DOOM			Busdriver		
<2,675 unique words			2,675-3,050 unique words			3,050-3,425 unique words			3,425-3,800 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											



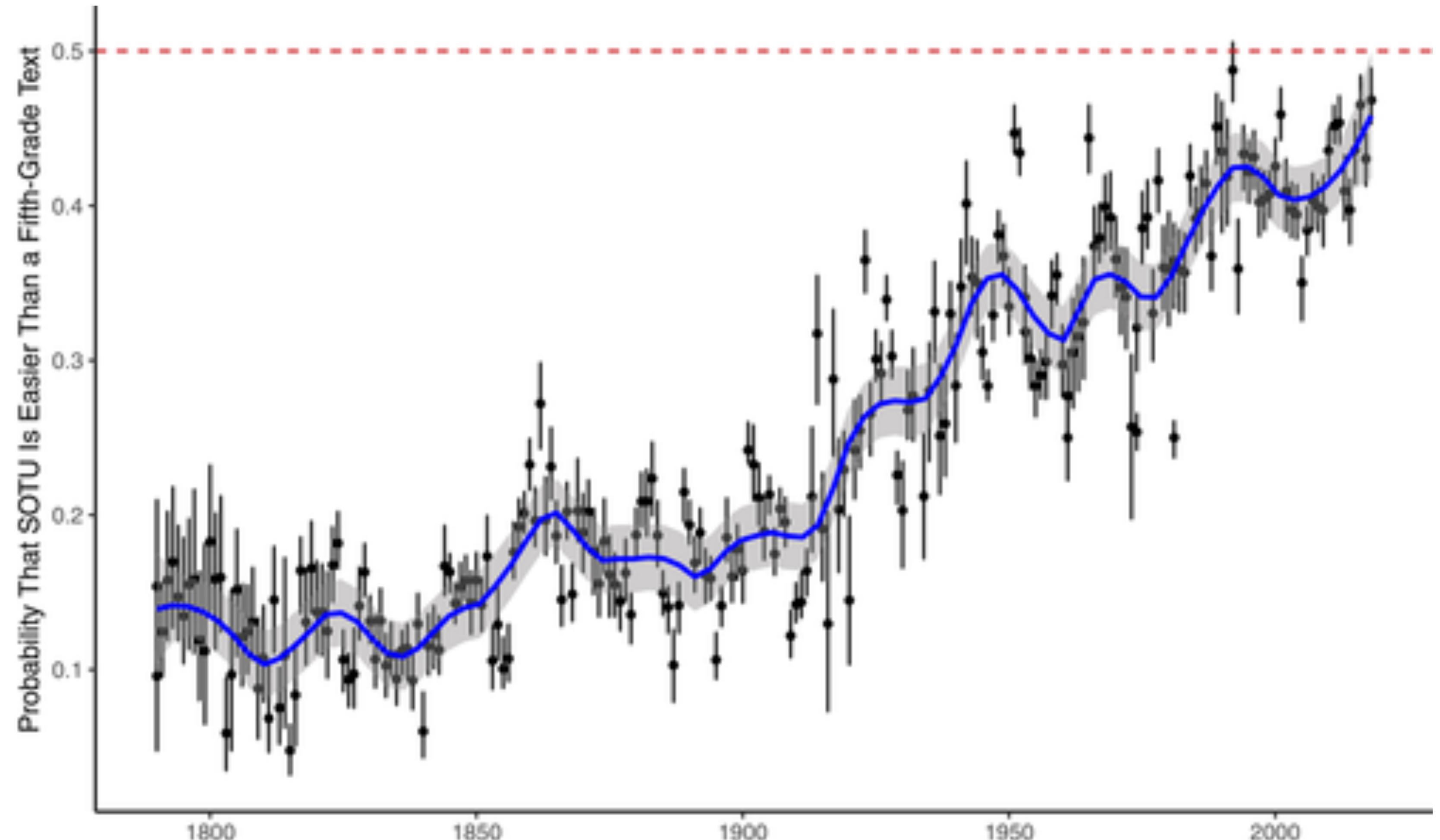
The contraction in vocabulary could be because of the style and structure of some genres

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

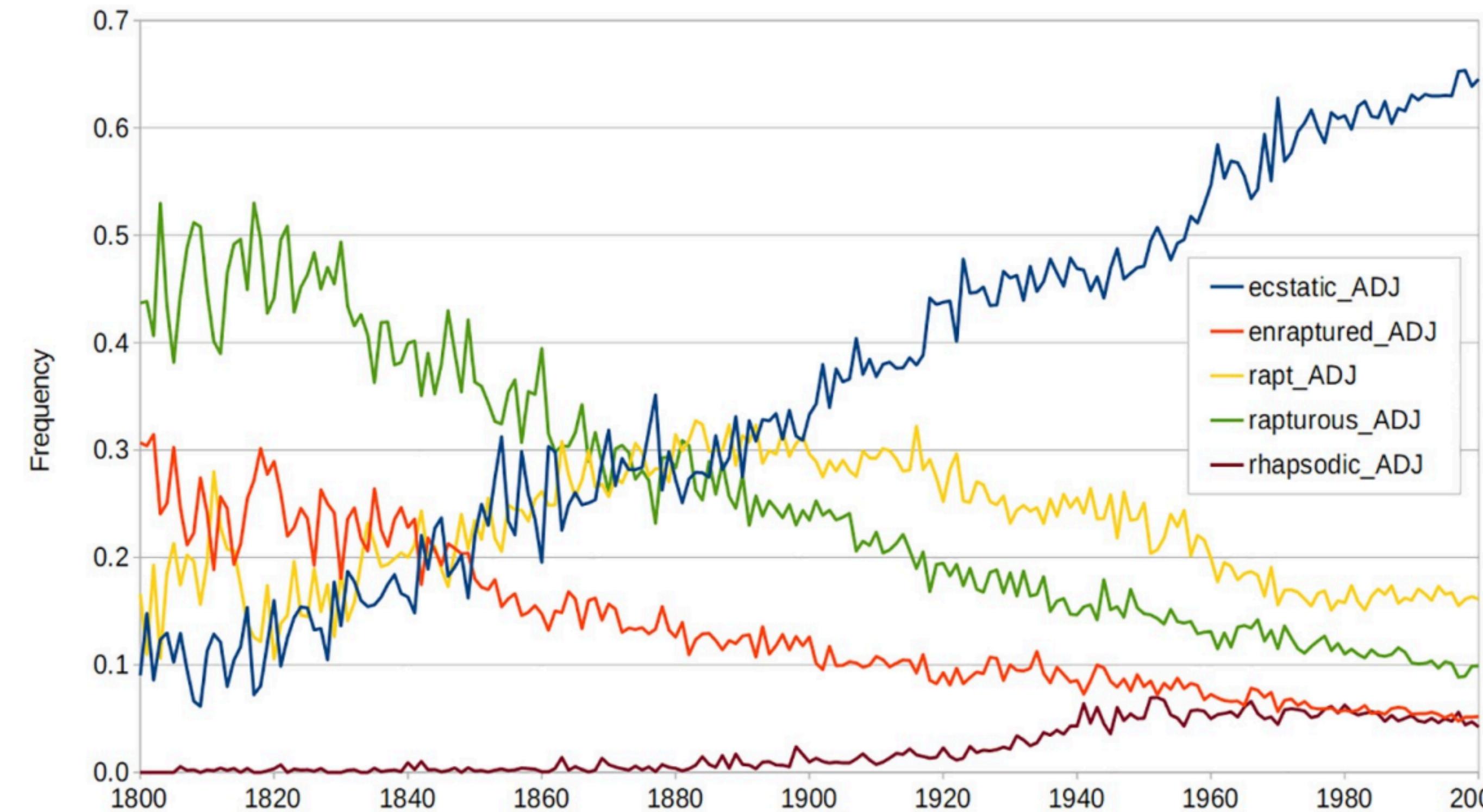
WORD CHARACTERISTICS

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



WORD FREQUENCY

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



QUESTION FOR THE DAY

“How to compare groups using words?”

PROBLEM SETUP

- Given two text collections corresponding to two “groups”
 - SotU speeches from 1800s Vs 1900s
 - Fiction Vs Non-fiction
 - r/politics Vs everything else
- Find if the two groups are different and what words differentiate the two groups

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"

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

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

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

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

k₂="legal"

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

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

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
-------	-------

χ^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

$$\chi^2$$

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

$$\chi^2$$

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

POINTWISE MUTUAL INFORMATION

- An information theoretic approach
- Generalizes to more than 2 categories
- Statistical significance can be calculated using non-parametric methods

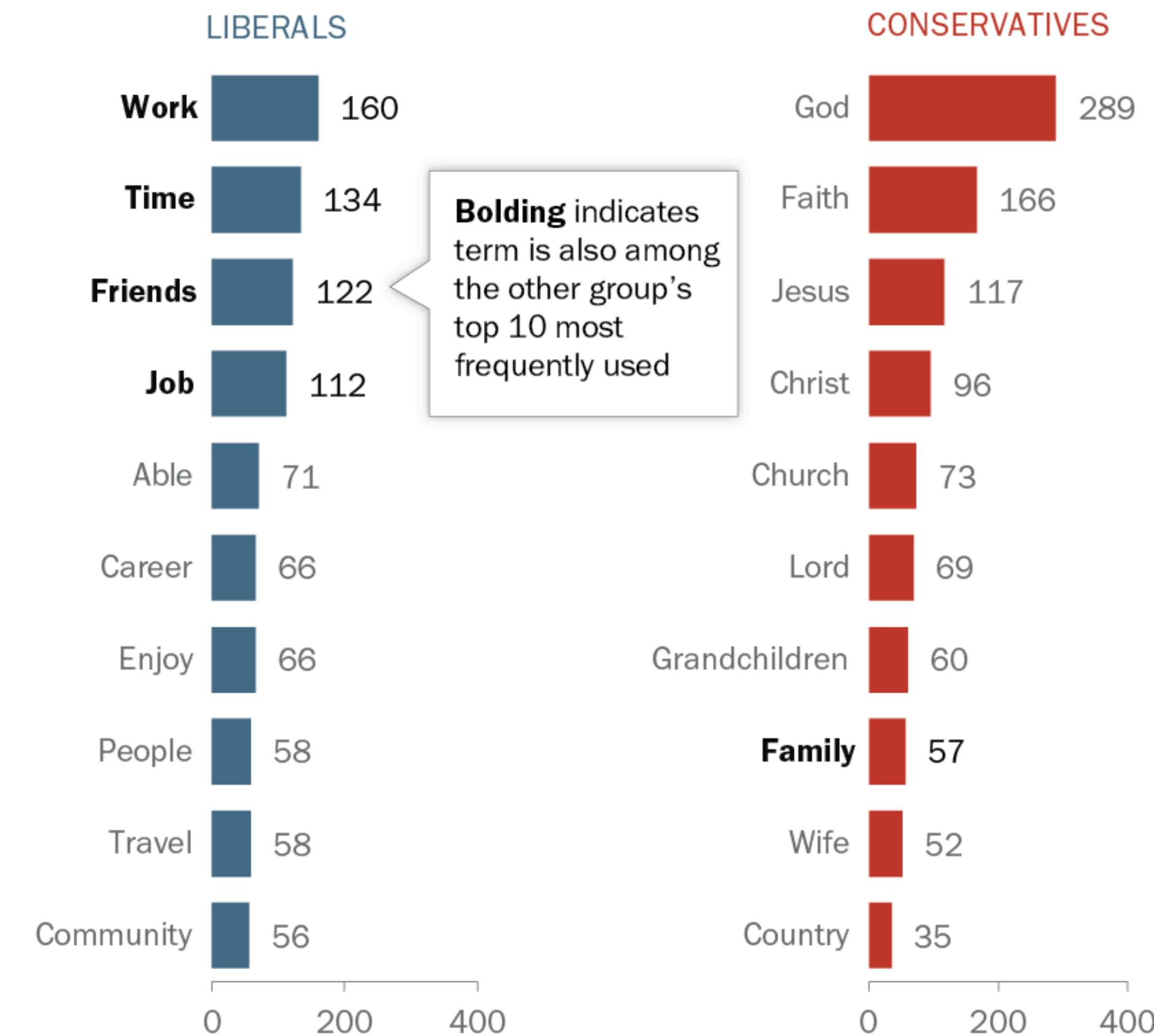
$$\text{PMI}(w, c) = \log \frac{\text{Probability of } w \text{ and } c \text{ co-occurring}}{\text{p}(w) \text{ p}(c)}$$

The diagram illustrates the components of the Pointwise Mutual Information (PMI) formula. At the top, a bracket labeled "Probability of w and c co-occurring" points down to the numerator. Below the numerator, a blue arrow labeled "A word w " points to the first argument of the PMI function. A blue bracket labeled "Probability of w occurring" points up to the term $\text{p}(w)$. To the right of the numerator, a grey box contains the function $p(w, c)$. Below this box, an orange arrow labeled "A category c " points to the second argument of the PMI function. An orange bracket labeled "Probability of c occurring" points up to the term $\text{p}(c)$.

Source: Bestvater and Shah; Pew research article

Identifying terms used more often by one group than another doesn't always indicate distinctiveness

Top 10 terms used **more** frequently by ___, by difference in word count



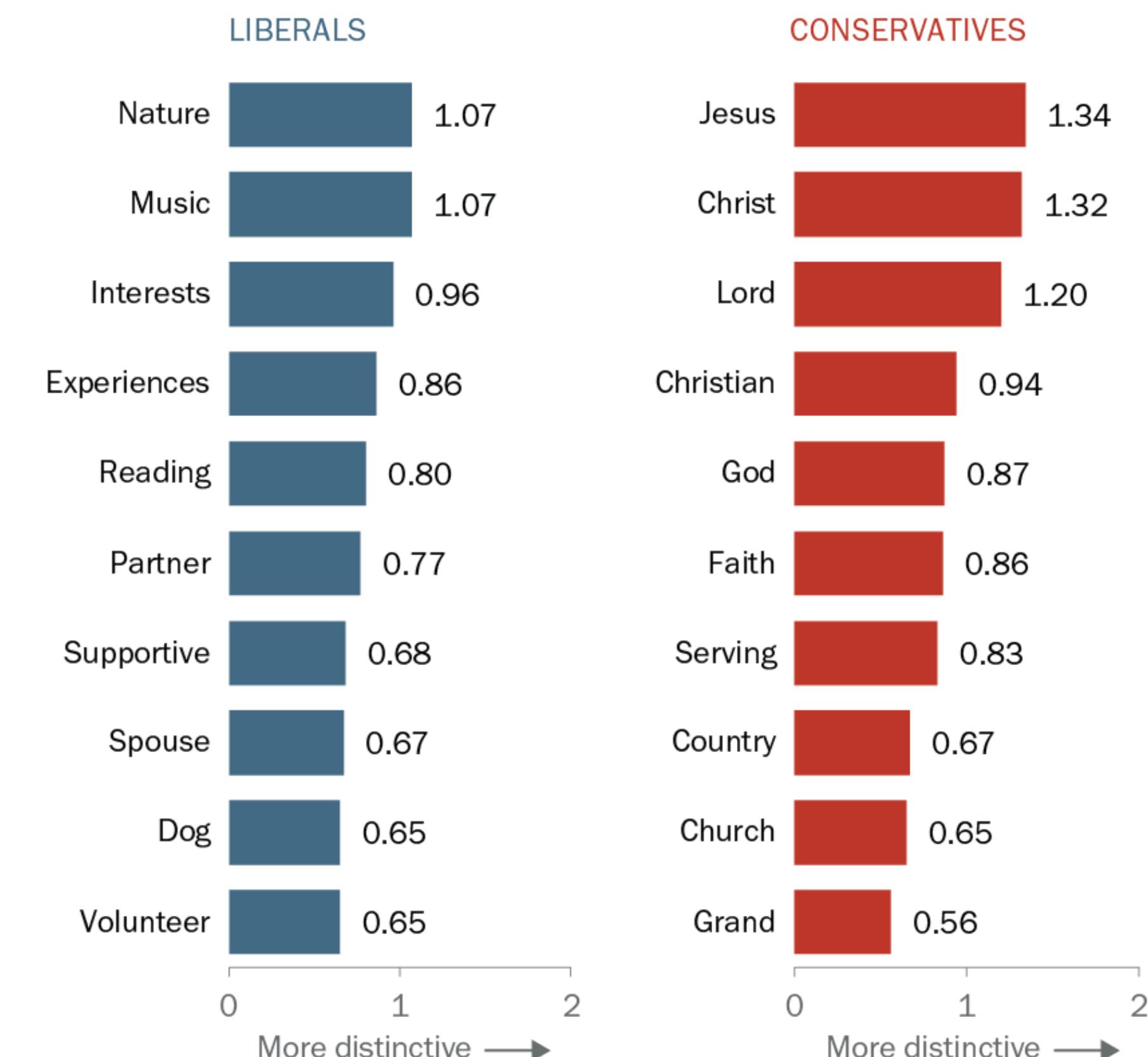
Notes: Terms used in open-ended responses to the question, "What about your life do you currently find meaningful, fulfilling, or satisfying? What keeps you going, and why?" Reported numbers are unweighted.

Source: Survey conducted Sept. 14-28, 2017, among U.S. adults.

PEW RESEARCH CENTER

Pointwise mutual information helps identify words that are used distinctively by one group, not just used more

Most distinctive terms used by each group, by PMI score



Notes: Terms used in open-ended responses to the question, "What about your life do you currently find meaningful, fulfilling, or satisfying? What keeps you going, and why?" Reported numbers are unweighted.

Source: Survey conducted Sept. 14-28, 2017, among U.S. adults.

PEW RESEARCH CENTER

Pointwise mutual information emphasizes different words than simple term frequencies or frequency differences

Top 10 terms used by conservatives, as measured by ...

Rank	Frequency	Difference in frequency	PMI
1	Family	God	Jesus
2	God	Faith	Christ
3	Friends	Jesus	Lord
4	Children	Christ	Christian
5	Health	Church	God
6	Love	Lord	Faith
7	Time	Grandchildren	Serving
8	Work	Family	Country
9	Faith	Wife	Church
10	Job	Country	Grand
.....			
60		Family	

Notes: Terms used in open-ended responses to the question, “What about your life do you currently find meaningful, fulfilling, or satisfying? What keeps you going, and why?” Reported numbers are unweighted.

Source: Survey conducted Sept. 14-28, 2017, among U.S. adults.

Pointwise mutual information emphasizes different words than simple term frequencies or frequency differences

Top 10 terms used by conservatives, as measured by ...

Rank	Frequency	Difference in frequency	PMI
1	Family	God	Jesus
2	God	Faith	Christ
3	Friends	Jesus	Lord
4	Children	Christ	Christian
5	Health	Church	God
6	Love	Lord	Faith
7	Time	Grandchildren	Serving
8	Work	Family	Country
9	Faith	Wife	Church
10	Job	Country	Grand
.....			
60		Family	

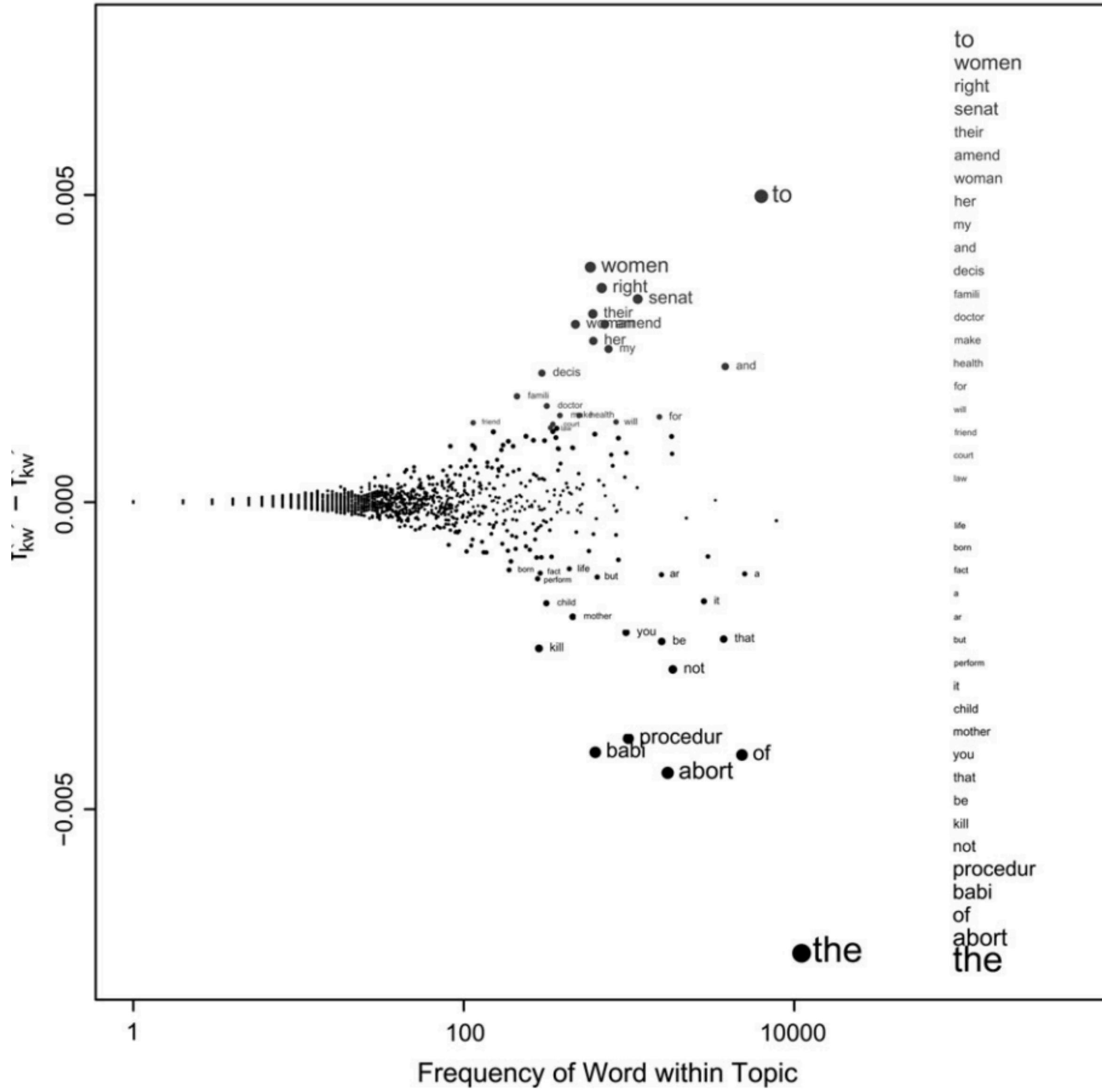
Notes: Terms used in open-ended responses to the question, “What about your life do you currently find meaningful, fulfilling, or satisfying? What keeps you going, and why?” Reported numbers are unweighted.

Source: Survey conducted Sept. 14-28, 2017, among U.S. adults.

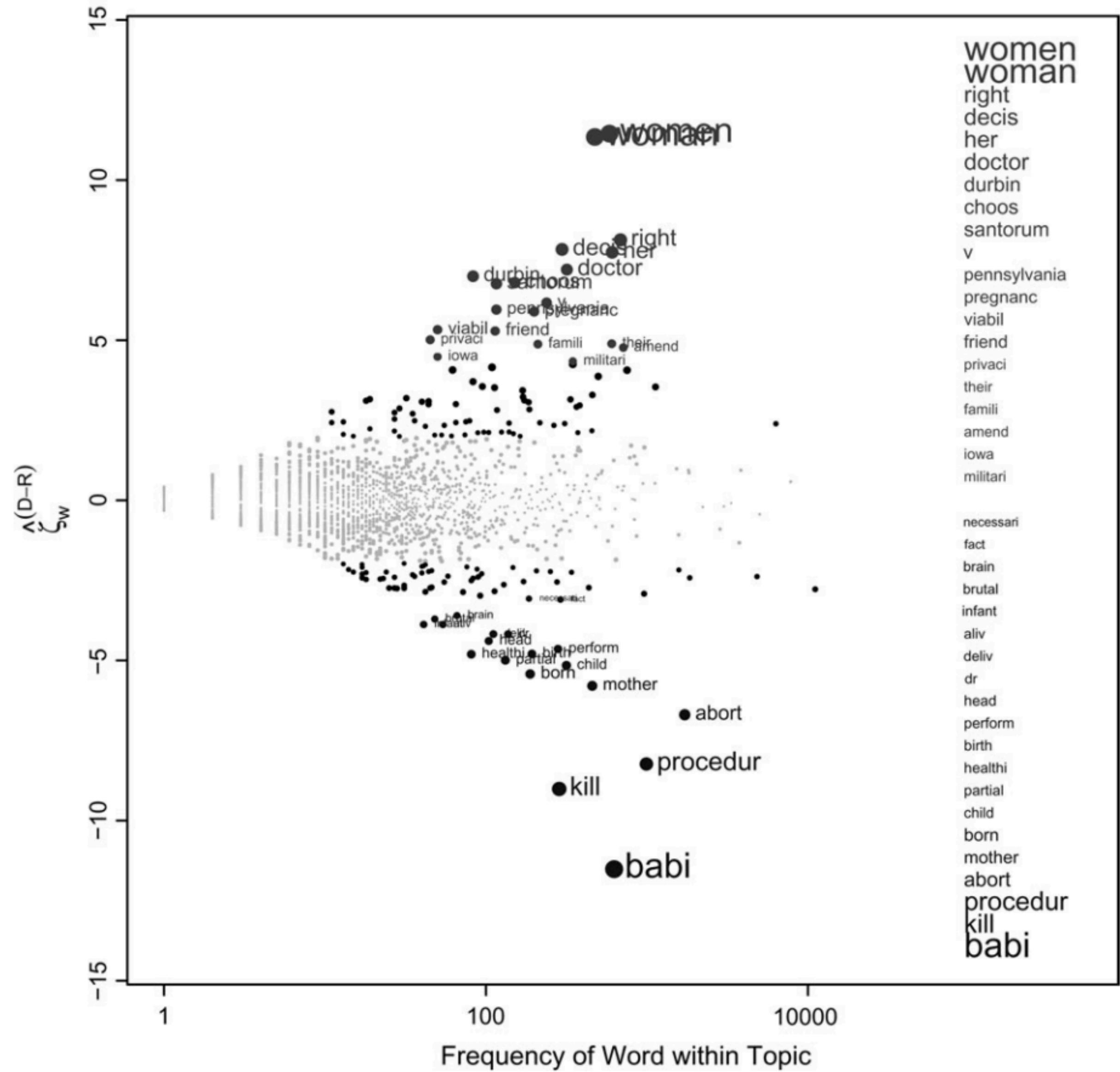
OTHER METHODS

- 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

- Distinctive terms demo