



# FINDING DISTINCTIVE TERMS

Sandeep Soni

---

09/13/2023

# CLASS LOGISTICS

---

- Hw4 released

# FROM LAST CLASS

---

- Bag of words demo
- I've added code to create a counts term-document matrix

# QUESTIONS FOR THE DAY

---

“How to calculate the similarity between documents?”

“How to compare groups using words?”

# AGENDA

---

- How to measure the “similarity” or “differences” between documents?
- How to measure differences between groups?
- How to find words that differentiate between two groups?

# DISTANCE METRICS

---

- Many ways to compare documents. A few examples:
  - Levenshtein distance
  - Jaccard index
  - Cosine similarity

# LEVENSHTEIN DISTANCE

- Objective: Transform one document to another by editing it
- Assign cost  $c$  for an edit operation
  - Insertion
  - Deletion
  - Replacement
- $\text{Dist}(\text{Doc1}, \text{Doc2}) = \min (\sum_i c_i)$

Sherlock Holmes      Doc 1

She locked Homes      Doc 2

Sherlock Holmes      Doc1

Sherlock Holmes      delete "l"

Sherlock Homes      replace r by ""

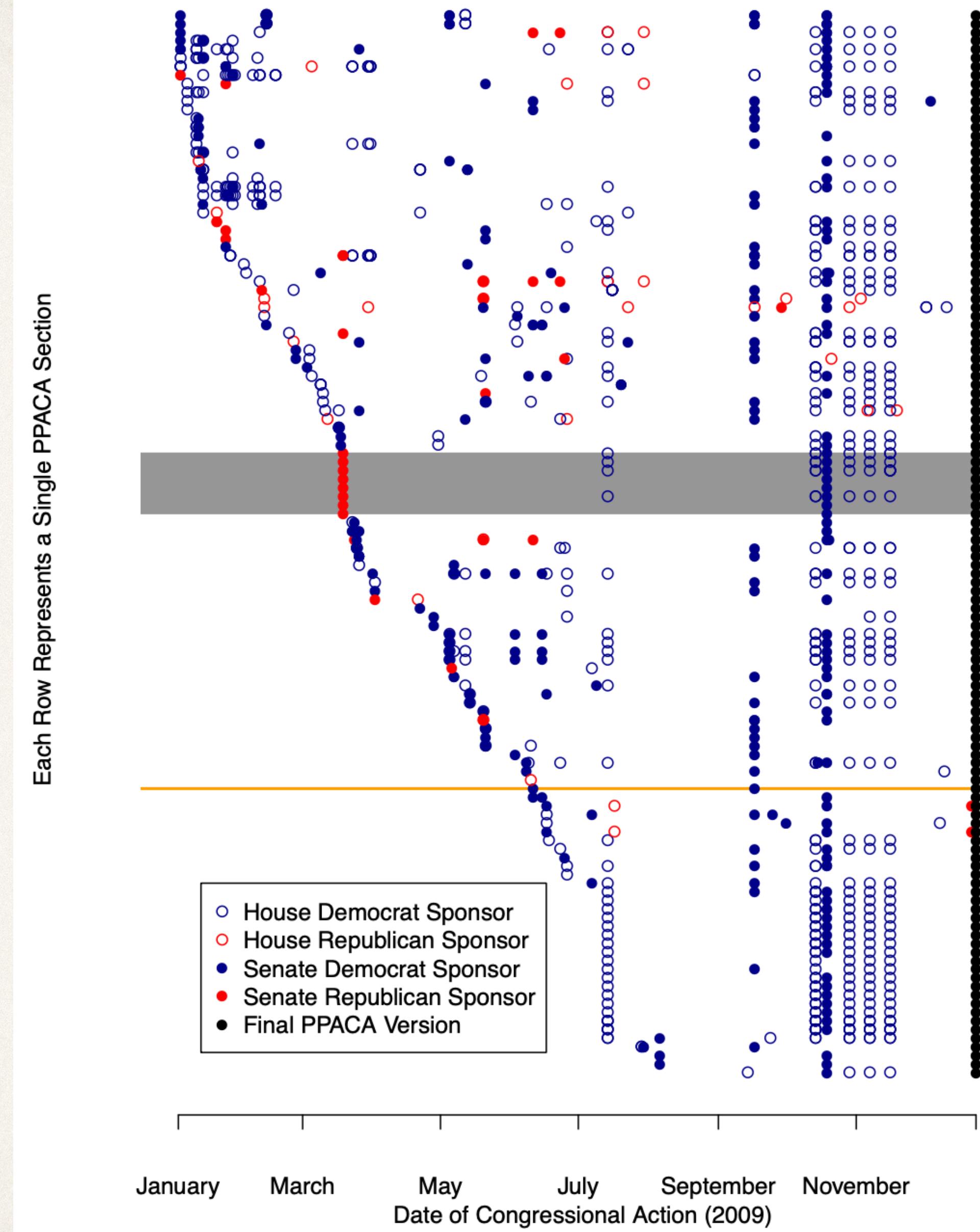
She locke Homes      insert "e"

She locked Homes      insert "d"

She locked Homes      Doc2

- Distance can be calculated using dynamic programming
- Many variations such as Needleman-Wunsch and Smith-Waterman
- Useful to quantify sequence alignment or text reuse

## Tracing Policy Ideas



# JACCARD INDEX

---

- Objective: measure commonality between documents
- Treat documents as sets of tokens
- $J(Doc1, Doc2) = \frac{|Doc1 \cap Doc2|}{|Doc1 \cup Doc2|}$

To Sherlock Holmes  
she is  
always *the* woman.

Doc 1

Sherlock Holmes is  
always right.

Doc 2

.  $J(Doc1, Doc2) = \frac{4}{9}$

- Many variations and generalizations (e.g., overlap coefficient, Tversky index)

Num	Obsecro Te 1
1	Obsecro Te domina sancta maria mater dei pietate plenissima summi
2	regis filia mater gloriosissima mater orphanorum consolatio
3	desolatorum via errantium <b>salus in te</b> sperantium virgo ante
4	partum virgo in partu <b>et</b> virgo post partum <b>Fons misericordie</b>
5	fons salutis et gratie fons pietatis et leticie fons consolationis
6	et indulgencie Per illam sanctam ineffabilem leticiam
7	qua exultavit spiritus tuus in illa hora quando tibi per gabrielem
8	annunciatus filius dei fuit
9	Et per illud divinum mysterium quod tunc operatus est spiritus sanctus

Obsecro Te 2
Obsecro Te domina sancta maria mater dei pietate plenissima summi regis filia mater gloriosissima mater orphanorum consolatio desolatorum via errantium <b>salus et spes in te</b> sperantium virgo ante partum virgo in partu virgo post partum fons salutis et gratie fons pietatis et leticie fons consolationis et indulgencie <b>Et</b> per illam sanctam inestimabilem leticiam qua exultavit spiritus tuus in illa hora quando tibi per gabrielem <b>archangelum</b> annuciatus <b>et conceptus</b> filius dei fuit Et per illud divinum mysterium quod tunc operatus est spiritus sanctus <b>in te</b>

# DOCUMENTS AS VECTORS

- Let  $d$  be a document, then  $d \in \mathbb{R}^{|V|}$
- If  $D$  is a set of documents, then  $D = [d^{(1)}; d^{(2)}; \dots; d^{(|D|)}] \in \mathbb{R}^{|D| \times |V|}$

To Sherlock Holmes  
she is  
always *the* woman. I  
have seldom heard  
him mention her  
under any other  
name. In his eyes  
she eclipses and  
predominates the  
whole of her sex.

Document



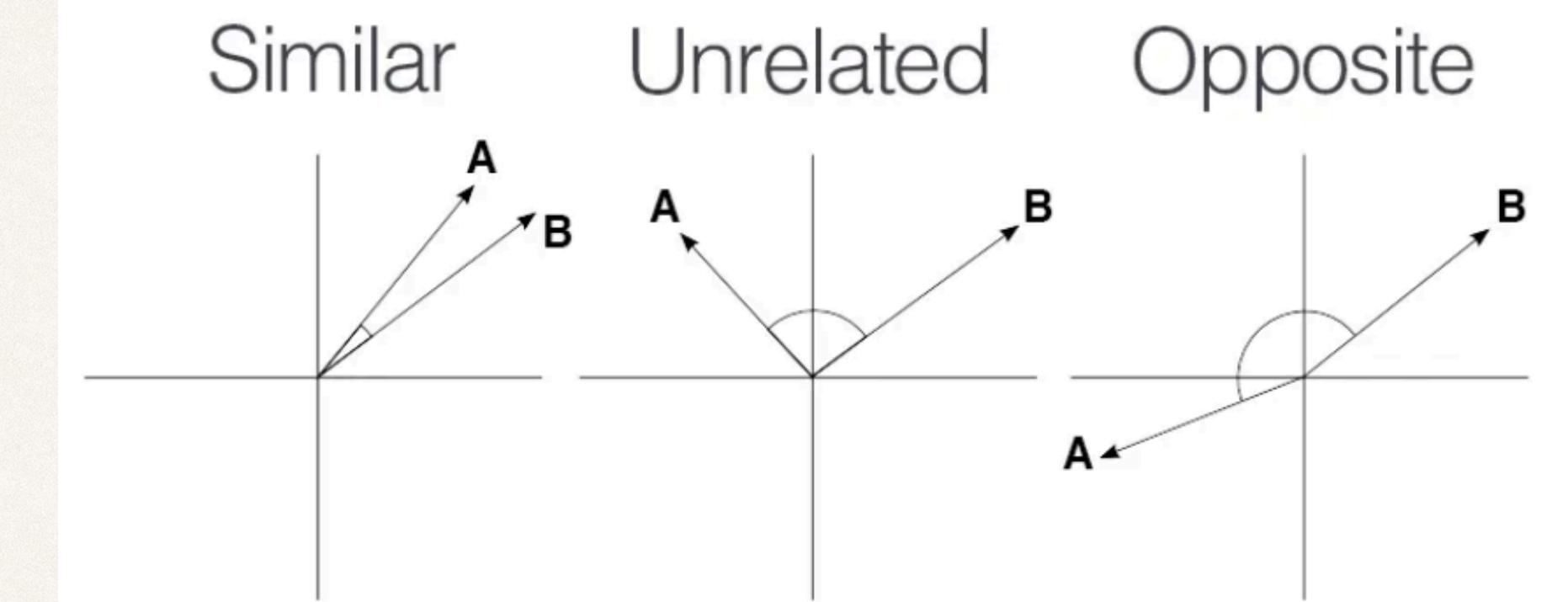
0	You
1	Sherlock
1	Holmes
0	methods
0	Watson
2	the
0	!
1	I
...	...
...	...
0	said
0	he

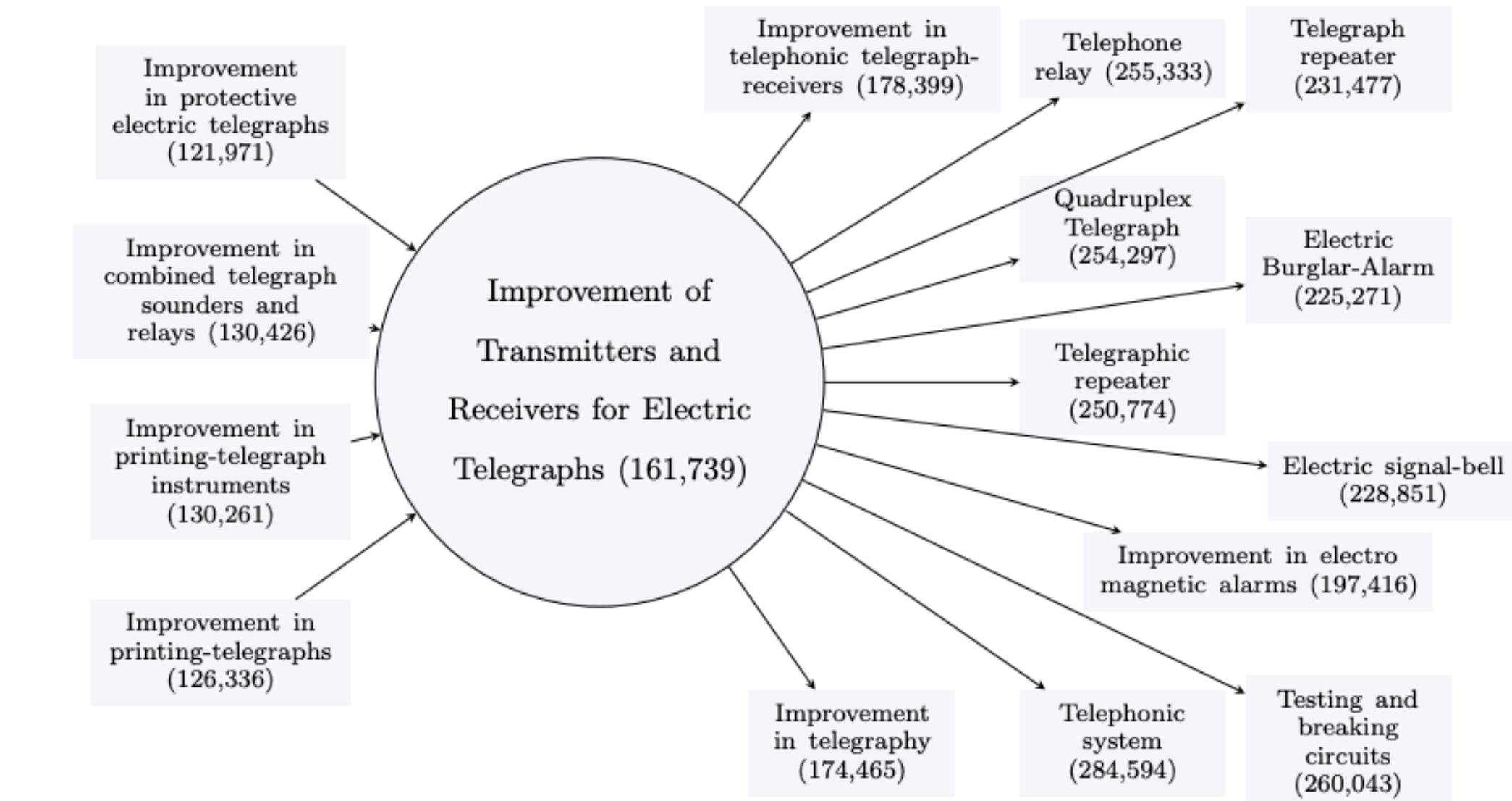
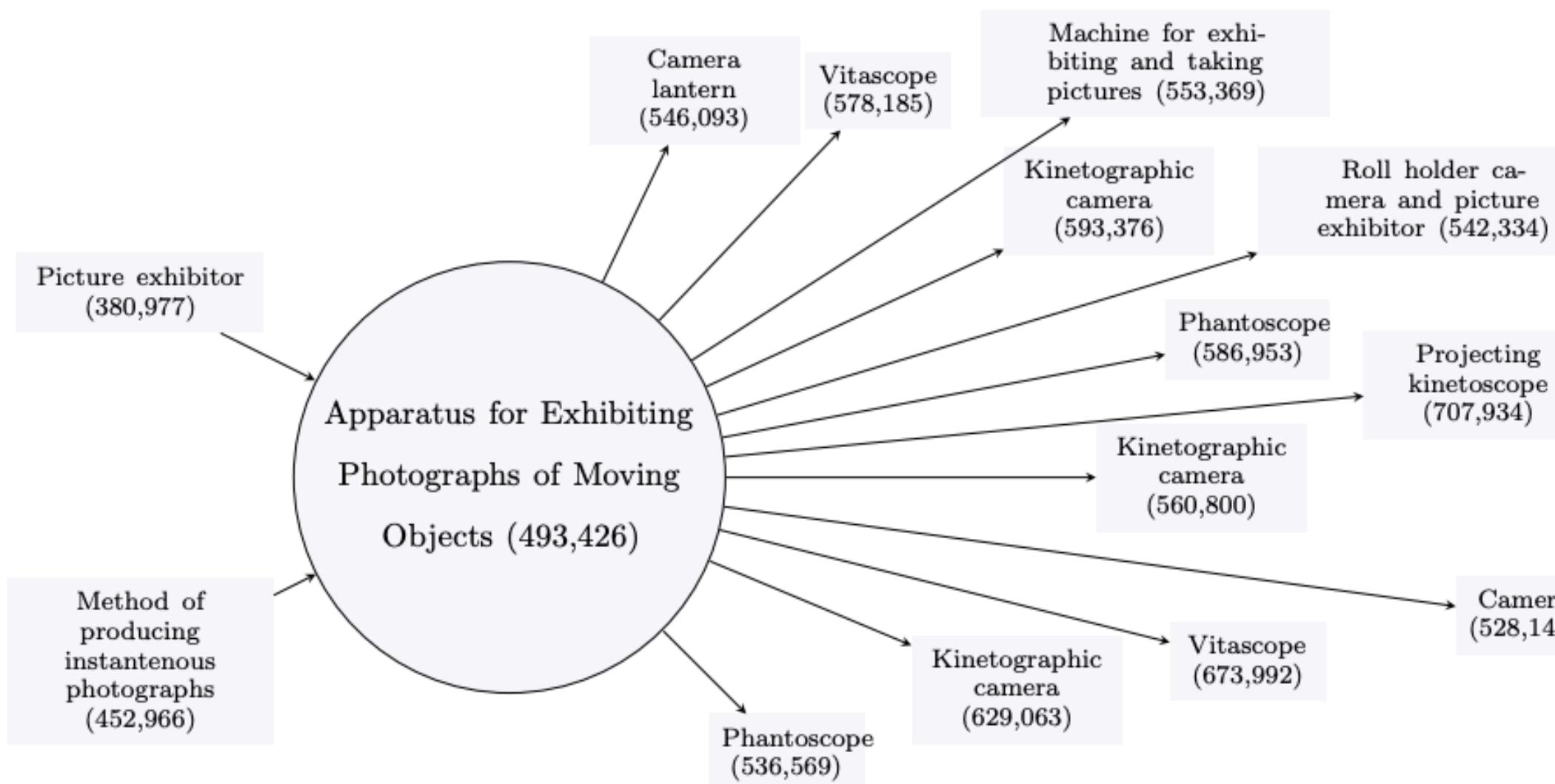
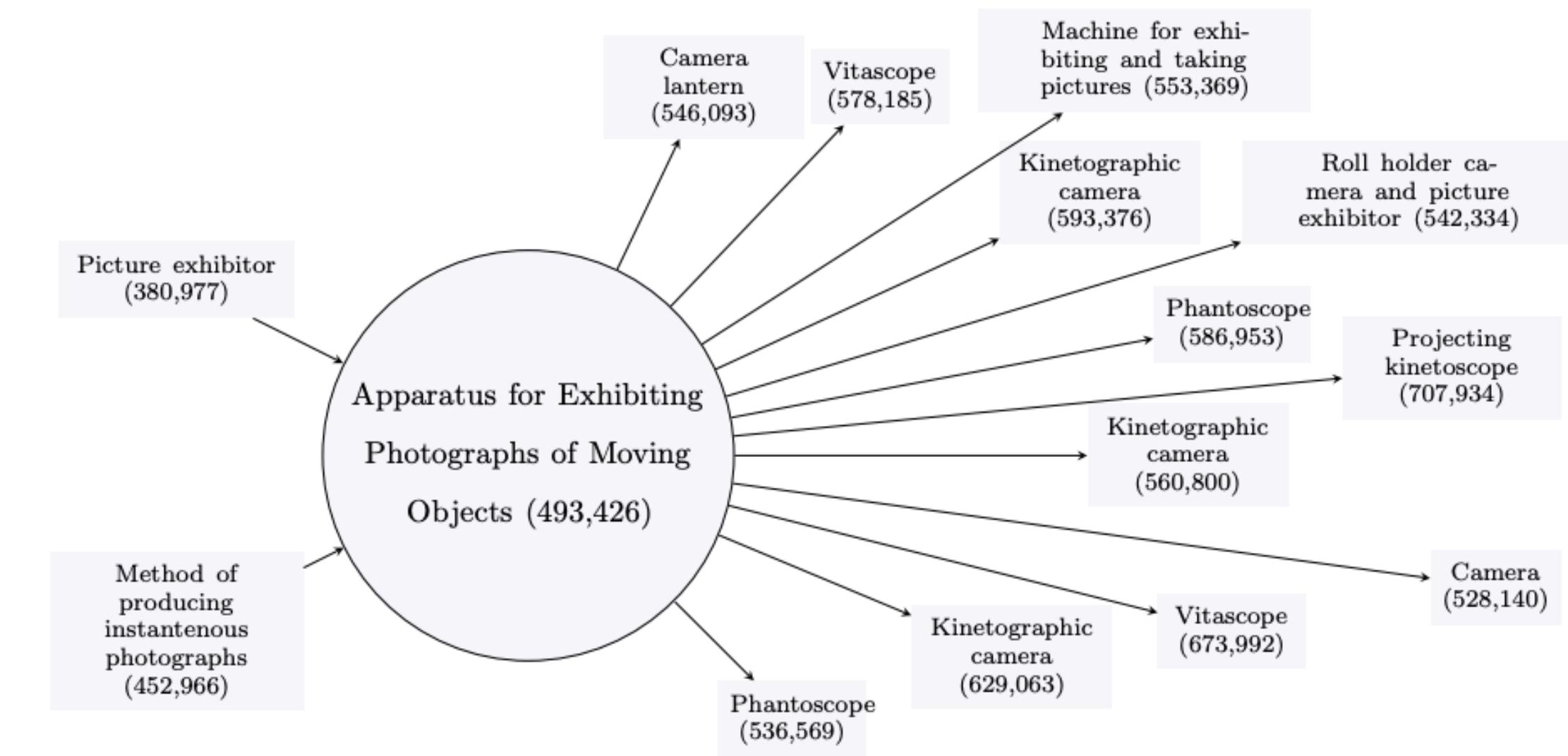
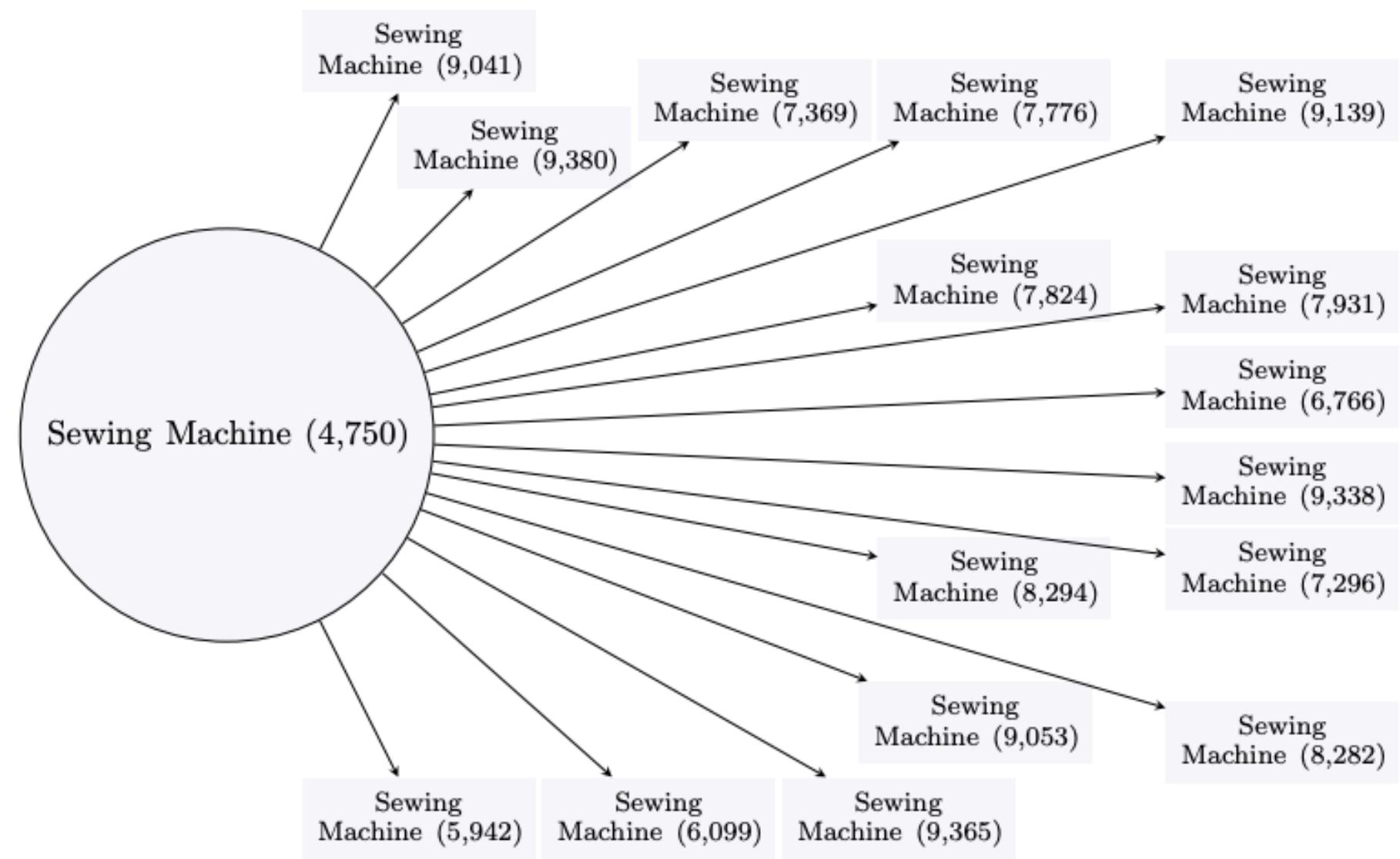
Vector

# COSINE SIMILARITY

---

- Objective: Measure similarity by treating documents as vectors
- $\text{CosSim}(\text{Doc1}, \text{Doc2}) = \frac{\mathbf{d}^{(1)} \cdot \mathbf{d}^{(2)}}{\|\mathbf{d}^{(1)}\| \|\mathbf{d}^{(2)}\|}$





# DISTANCE METRICS

---

- Many ways to compare documents. A few examples:
  - Edit based (e.g., Levenshtein distance)
  - Token based (e.g., Jaccard index)
  - Vector based (e.g., Cosine similarity)

# METADATA

---

# METADATA

---

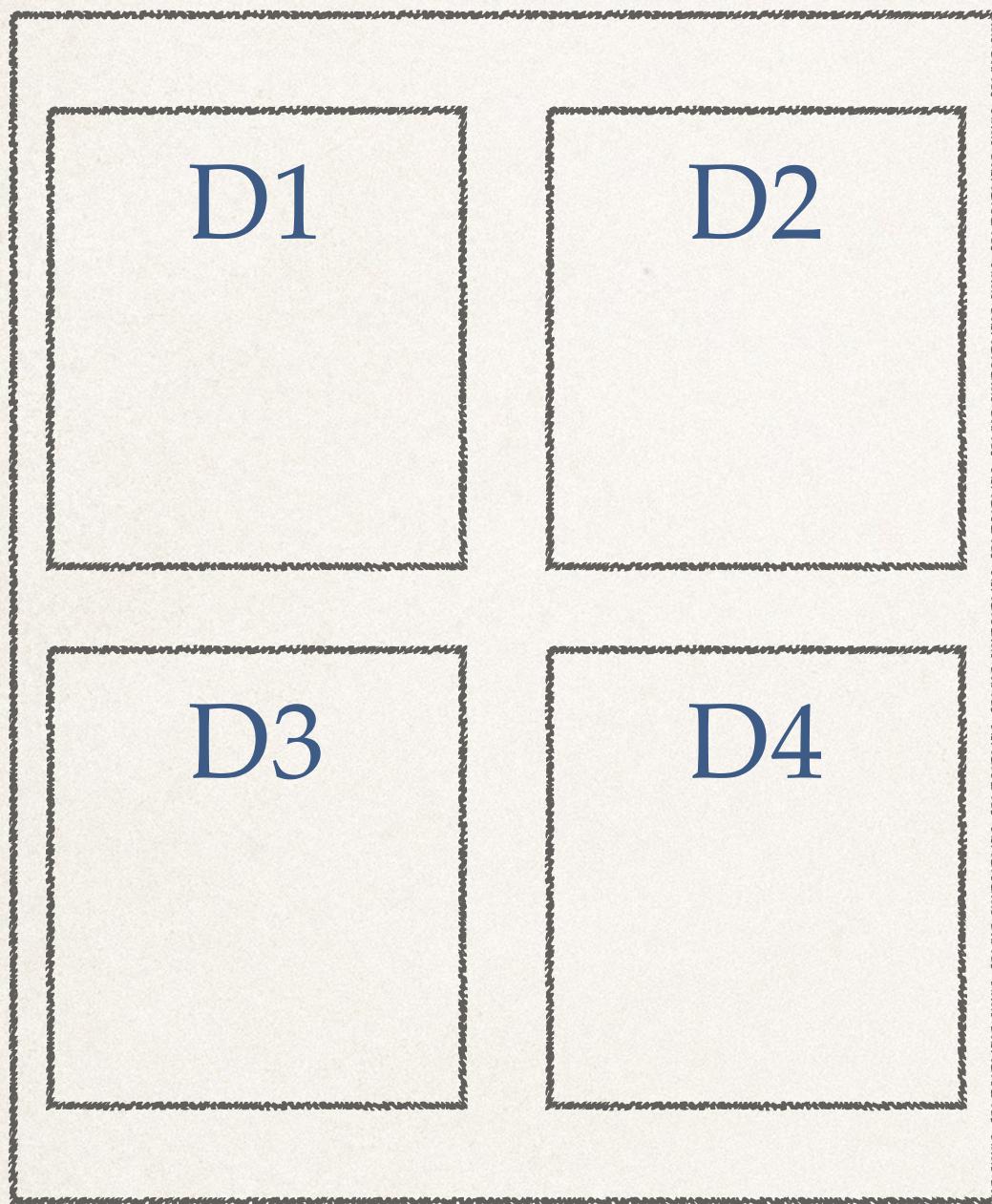
- Documents have metadata (e.g., author, date, genre)

# METADATA

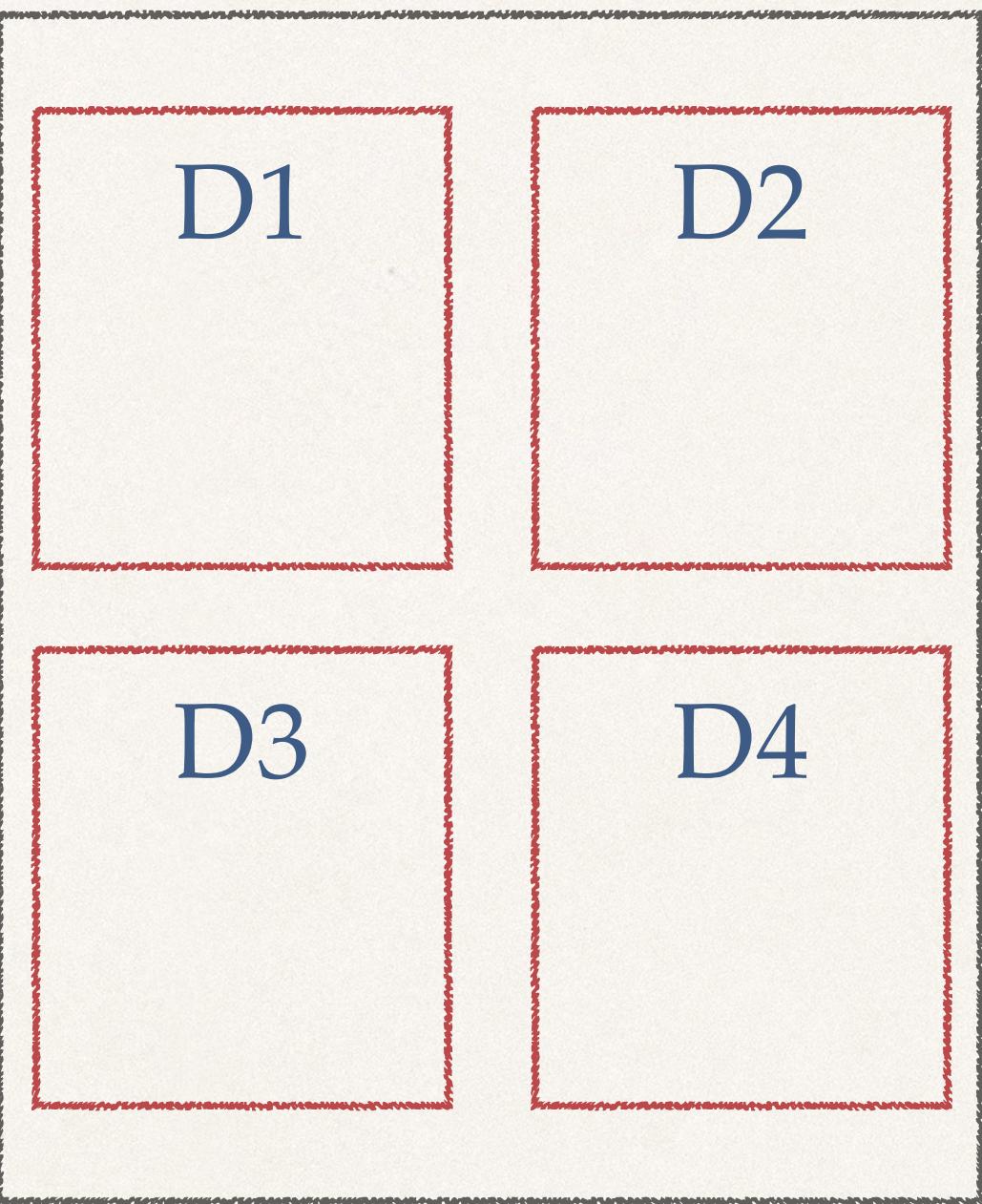
---

- Documents have metadata (e.g., author, date, genre)
- How to quantify differences between groups?

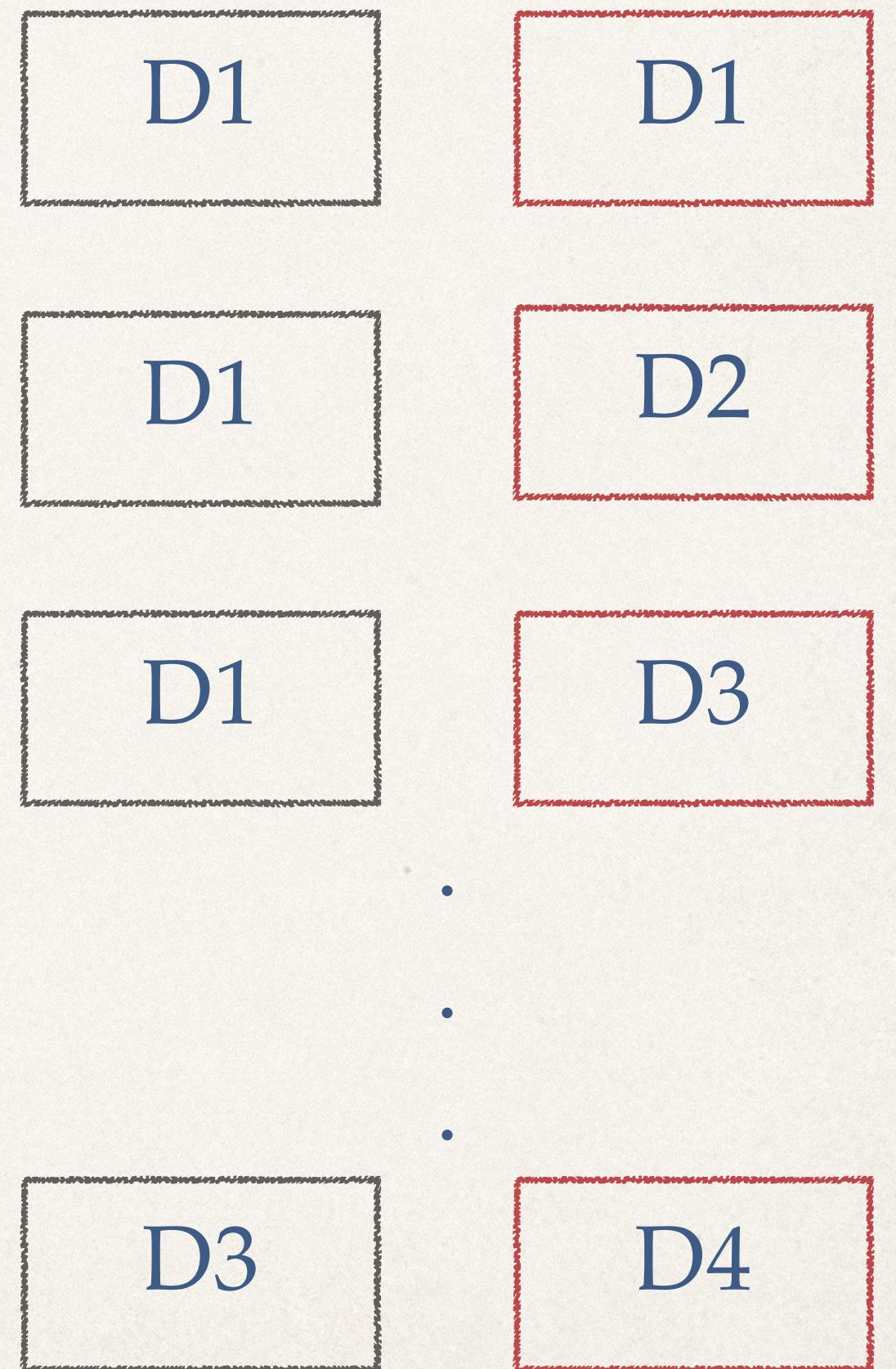
# PAIRWISE METRICS



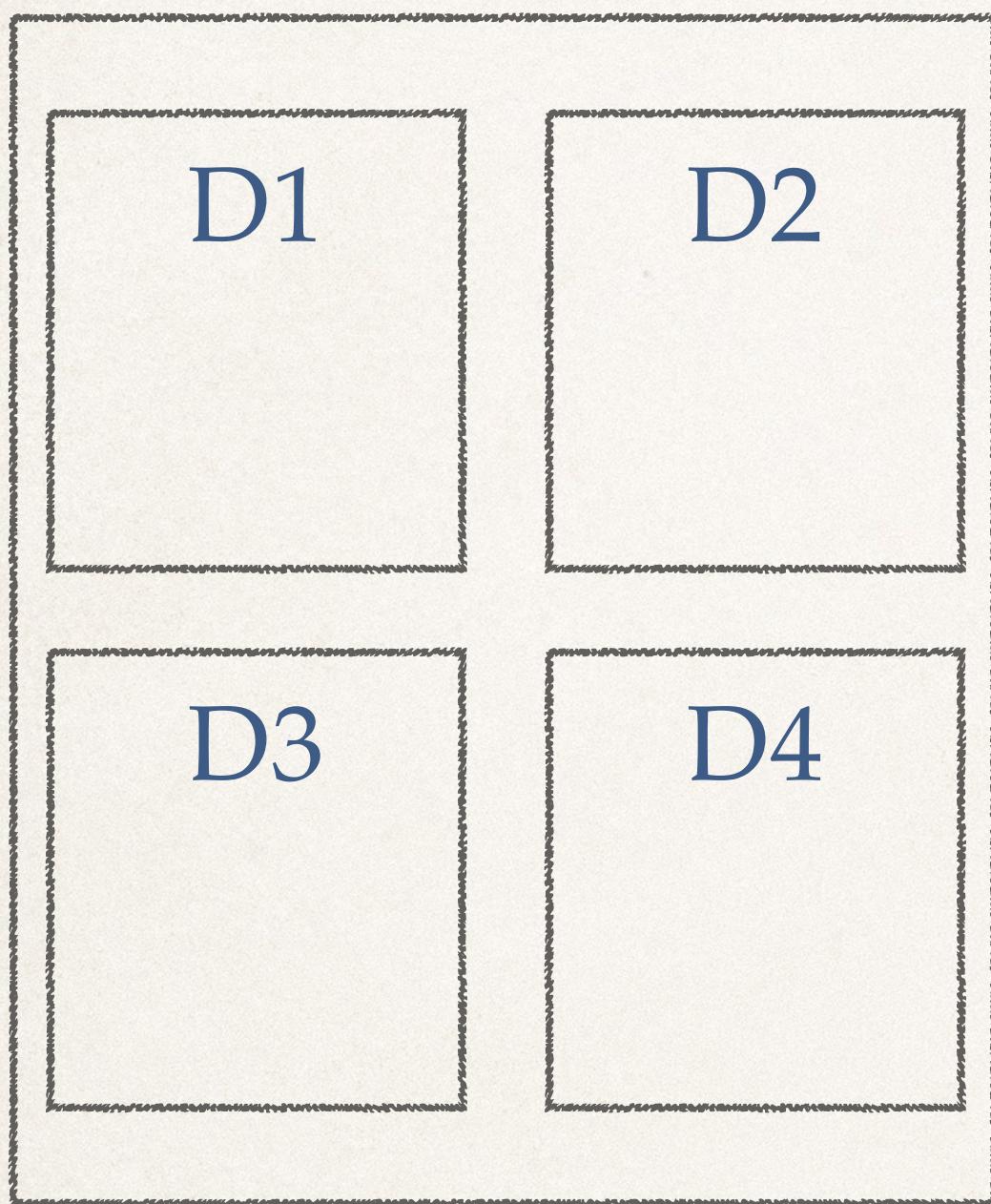
Corpus for group 1



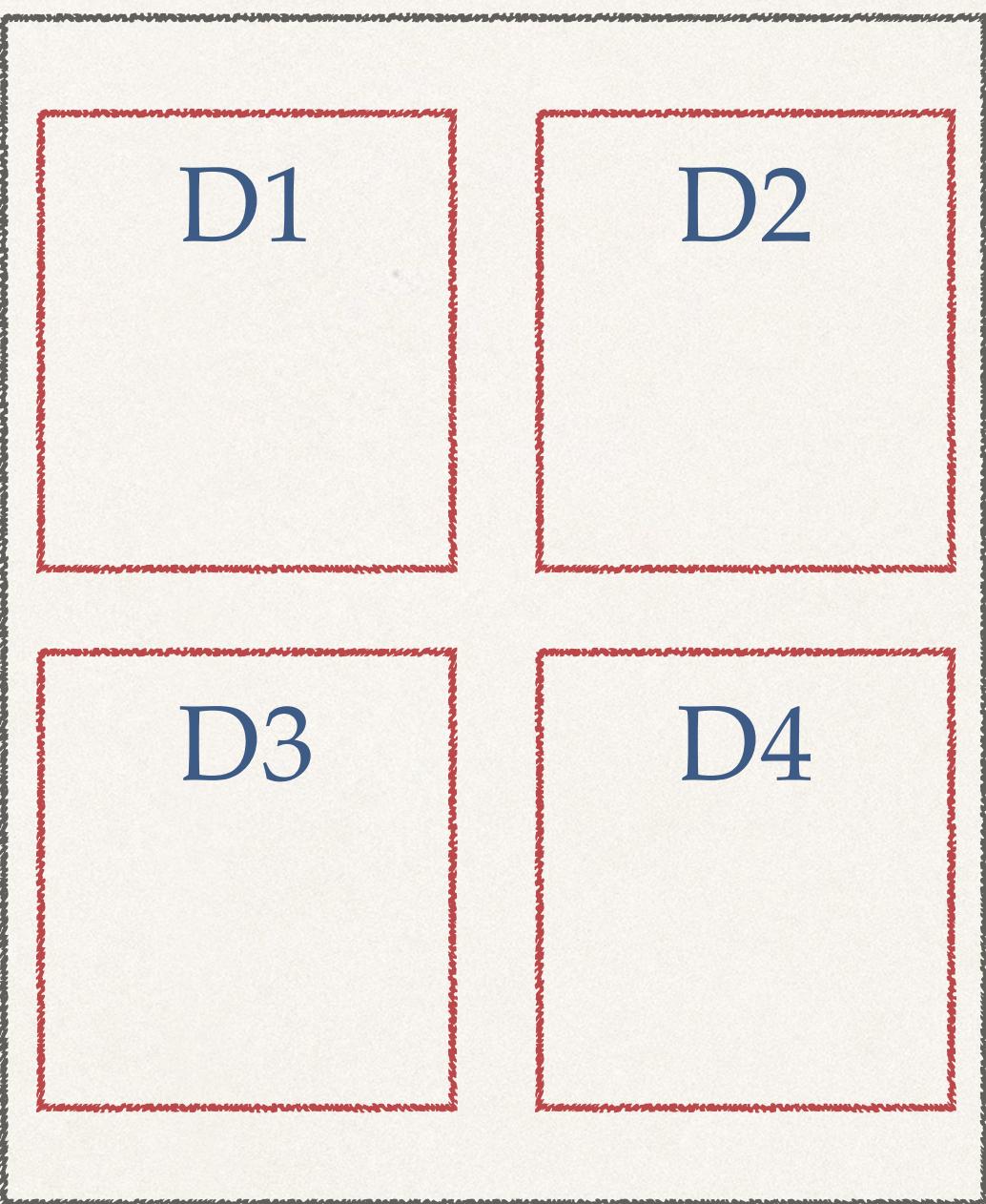
Corpus for group 2



# PAIRWISE METRICS

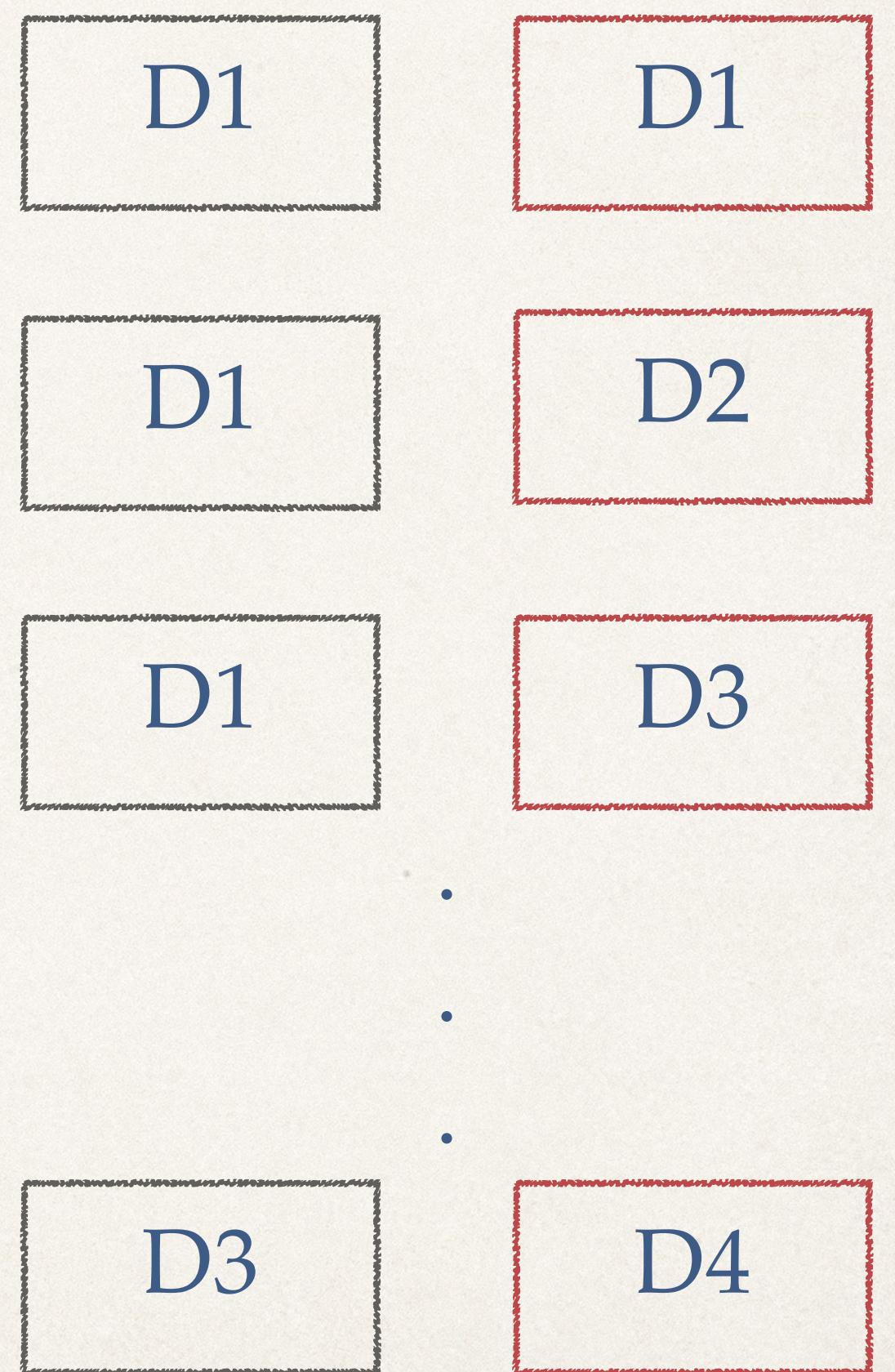


Corpus for group 1



Corpus for group 2

- Average similarity between document pairs as measure of similarity between groups

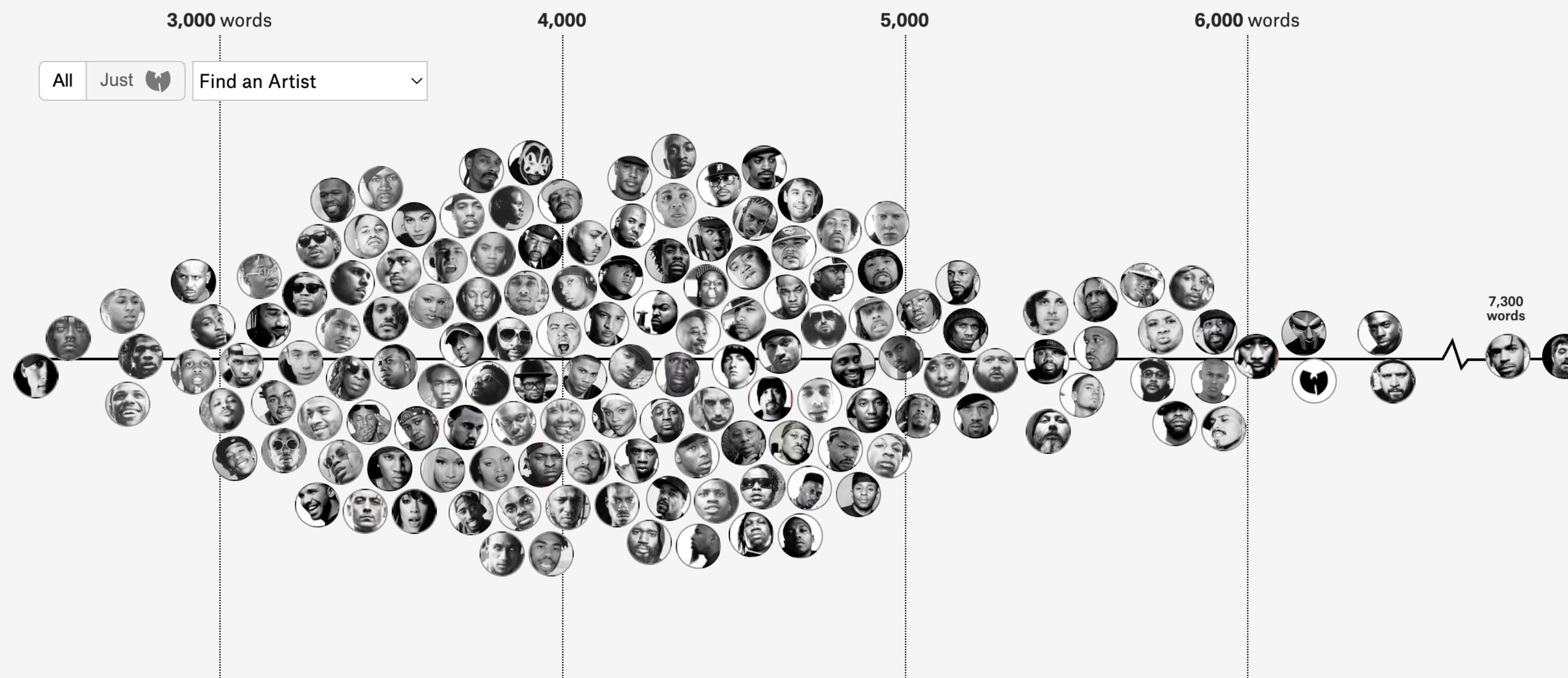


<code>metrics.pairwise.additive_chi2_kernel(X[, Y])</code>	Compute the additive chi-squared kernel between observations in X and Y.
<code>metrics.pairwise.chi2_kernel(X[, Y, gamma])</code>	Compute the exponential chi-squared kernel between X and Y.
<code>metrics.pairwise.cosine_similarity(X[, Y, ...])</code>	Compute cosine similarity between samples in X and Y.
<code>metrics.pairwise.cosine_distances(X[, Y])</code>	Compute cosine distance between samples in X and Y.
<code>metrics.pairwise.distance_metrics()</code>	Valid metrics for pairwise_distances.
<code>metrics.pairwise.euclidean_distances(X[, Y])</code>	Compute the distance matrix between each pair from a vector array X and Y. sklearn.metrics.pairwise.chi2_kernel
<code>metrics.pairwise.haversine_distances(X[, Y])</code>	Compute the Haversine distance between samples in X and Y.
<code>metrics.pairwise.kernel_metrics()</code>	Valid metrics for pairwise_kernels.
<code>metrics.pairwise.laplacian_kernel(X[, Y, gamma])</code>	Compute the laplacian kernel between X and Y.
<code>metrics.pairwise.linear_kernel(X[, Y, ...])</code>	Compute the linear kernel between X and Y.
<code>metrics.pairwise.manhattan_distances(X[, Y, ...])</code>	Compute the L1 distances between the vectors in X and Y.
<code>metrics.pairwise.nan_euclidean_distances(X)</code>	Calculate the euclidean distances in the presence of missing values.
<code>metrics.pairwise.pairwise_kernels(X[, Y, ...])</code>	Compute the kernel between arrays X and optional array Y.
<code>metrics.pairwise.polynomial_kernel(X[, Y, ...])</code>	Compute the polynomial kernel between X and Y.
<code>metrics.pairwise.rbf_kernel(X[, Y, gamma])</code>	Compute the rbf (gaussian) kernel between X and Y.
<code>metrics.pairwise.sigmoid_kernel(X[, Y, ...])</code>	Compute the sigmoid kernel between X and Y.
<code>metrics.pairwise.paired_euclidean_distances(X, Y)</code>	Compute the paired euclidean distances between X and Y.
<code>metrics.pairwise.paired_manhattan_distances(X, Y)</code>	Compute the paired L1 distances between X and Y.
<code>metrics.pairwise.paired_cosine_distances(X, Y)</code>	Compute the paired cosine distances between X and Y.
<code>metrics.pairwise.paired_distances(X, Y, *[...,])</code>	Compute the paired distances between X and Y.
<code>metrics.pairwise_distances(X[, Y, metric, ...])</code>	Compute the distance matrix from a vector array X and optional Y.
<code>metrics.pairwise_distances_argmin(X, Y, *[...,])</code>	Compute minimum distances between one point and a set of points.
<code>metrics.pairwise_distances_argmin_min(X, Y, *)</code>	Compute minimum distances between one point and a set of points.
<code>metrics.pairwise_distances_chunked(X[, Y, ...])</code>	Generate a distance matrix chunk by chunk with optional reduction.

The different pairwise metrics in sklearn

# FINDING DIFFERENTIATORS

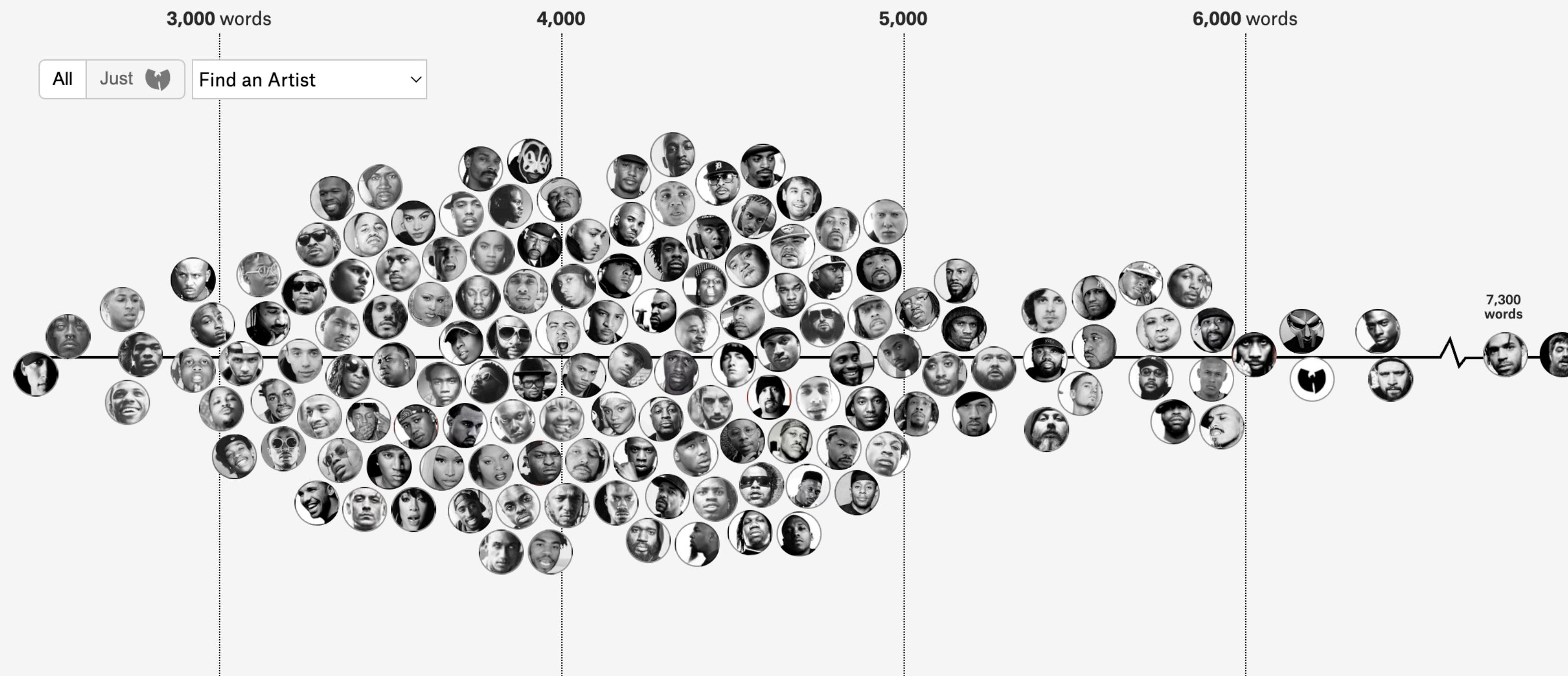
# of Unique Words Used Within Artist's First 35,000 Lyrics



Source: <https://pudding.cool/projects/vocabulary/index.html>

# FINDING DIFFERENTIATORS

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



Source: <https://pudding.cool/projects/vocabulary/index.html>

We are often interested in finding the difference markers

## DIFFERENCE IN PROPORTIONS

---

- If  $w$  is the word,  $k$  is the category and  $c_w$  is the count
- $f_{w,k} = \frac{C_{w,k}}{\sum_w C_{w,k}}$  is the normalized count or proportion of occurrence of the word  $w$  in category  $k$
- $f_{w,k_1} - f_{w,k_2}$  gives the difference in proportions across two groups

# 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

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

# $\chi^2$

	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

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

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

$$\chi^2$$

---

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

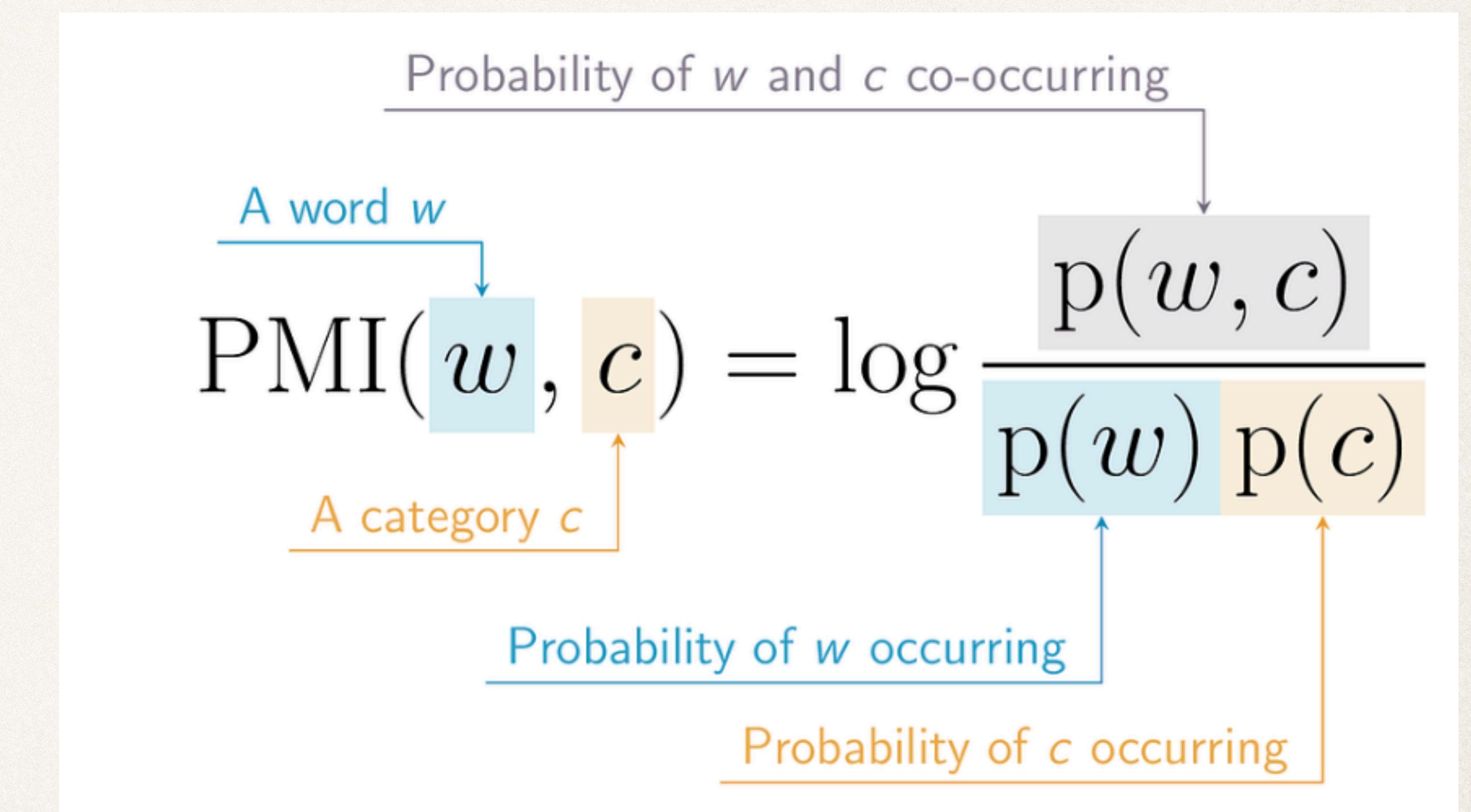
$$\chi^2$$

---

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

# PMI

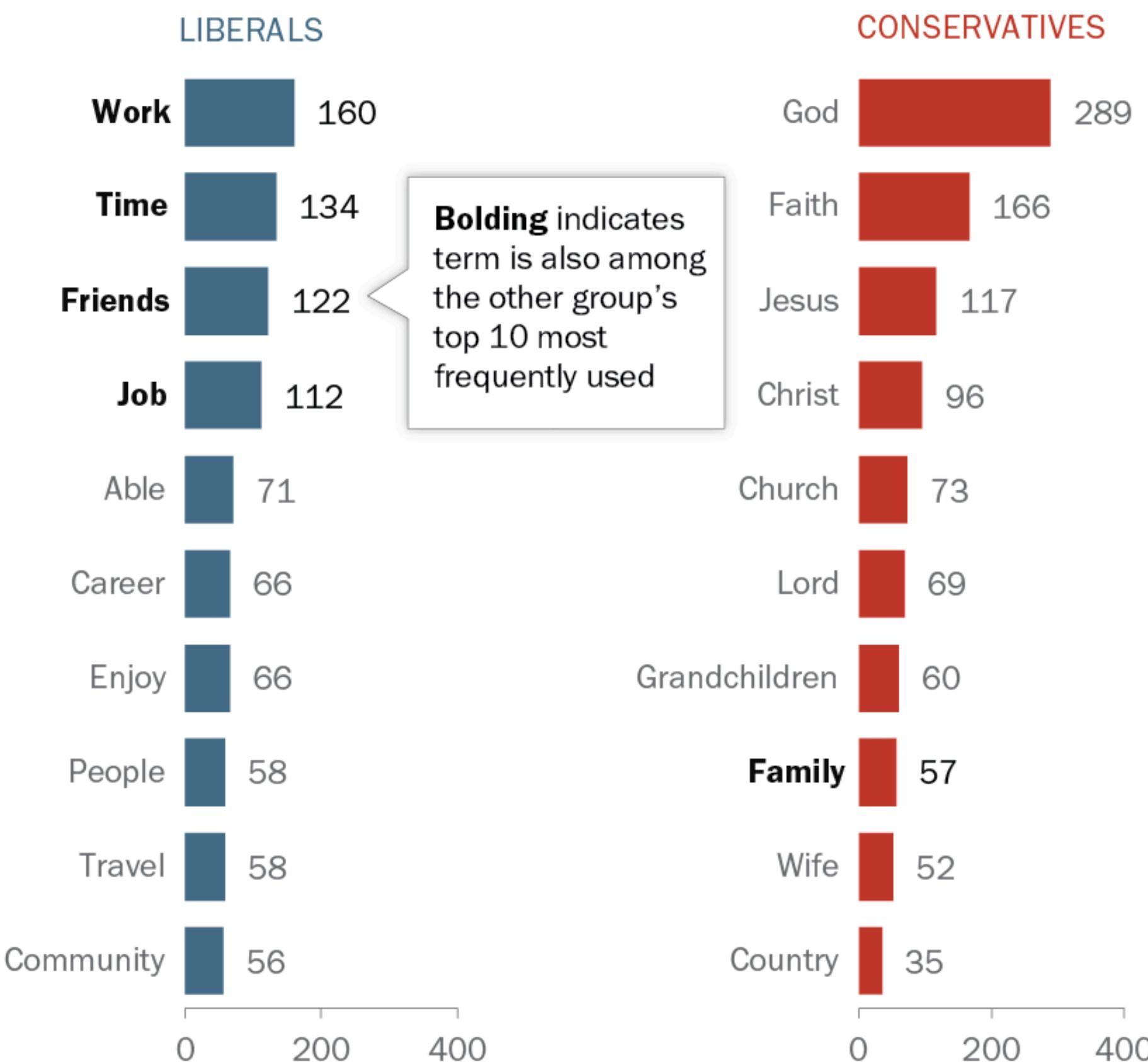
- A more information theoretic approach
- Generalization to more than 2 categories
- Statistical significance can be calculated using non-parametric methods



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



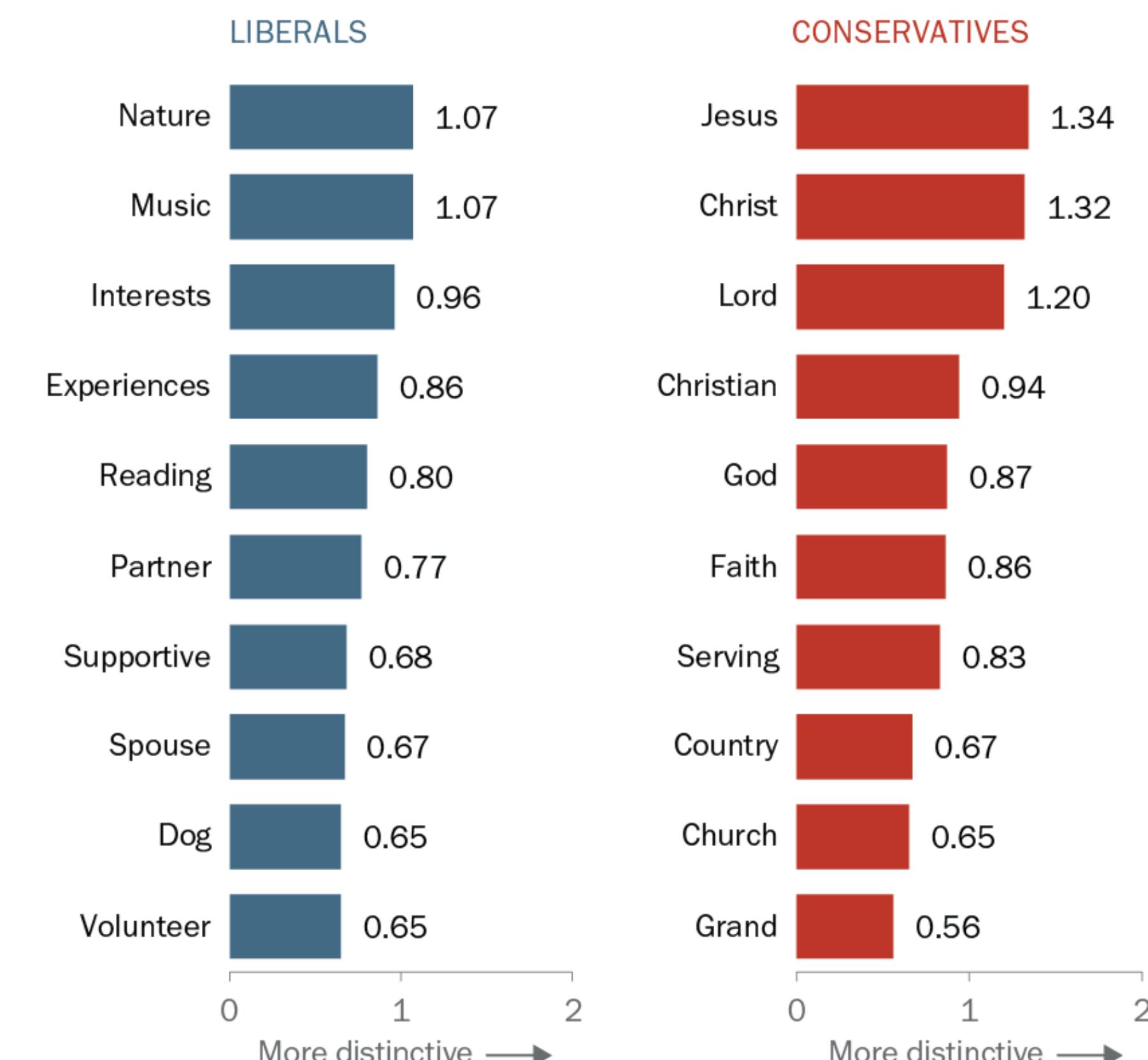
Notes: Terms used in open-ended responses to the question, "What about your life do you currently find meaningful, fulfilling, or satisfying? What keeps you going, and why?" Reported numbers are unweighted.

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	<b>Family</b>	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	<b>Family</b>	Country
9	Faith	Wife	Church
10	Job	Country	Grand
.....			
60		<b>Family</b>	

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	<b>Family</b>	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	<b>Family</b>	Country
9	Faith	Wife	Church
10	Job	Country	Grand
.....			
60		<b>Family</b>	

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)

## IN CLASS

---

- Finding distinctive terms demo