



# Natural Language Processing

Mehmet Can Yavuz, PhD

Adapted from Info 256 - David Bamman, UC Berkeley





## Panel B: Phrases Used More Often by Republicans

### *Two-Word Phrases*

stem cell	personal accounts	retirement accounts
natural gas	Saddam Hussein	government spending
death tax	pass the bill	national forest
illegal aliens	private property	minority leader
class action	border security	urge support
war on terror	President announces	cell lines
embryonic stem	human life	cord blood
tax relief	Chief Justice	action lawsuits
illegal immigration	human embryos	economic growth
date the time	increase taxes	food program

### *Three-Word Phrases*

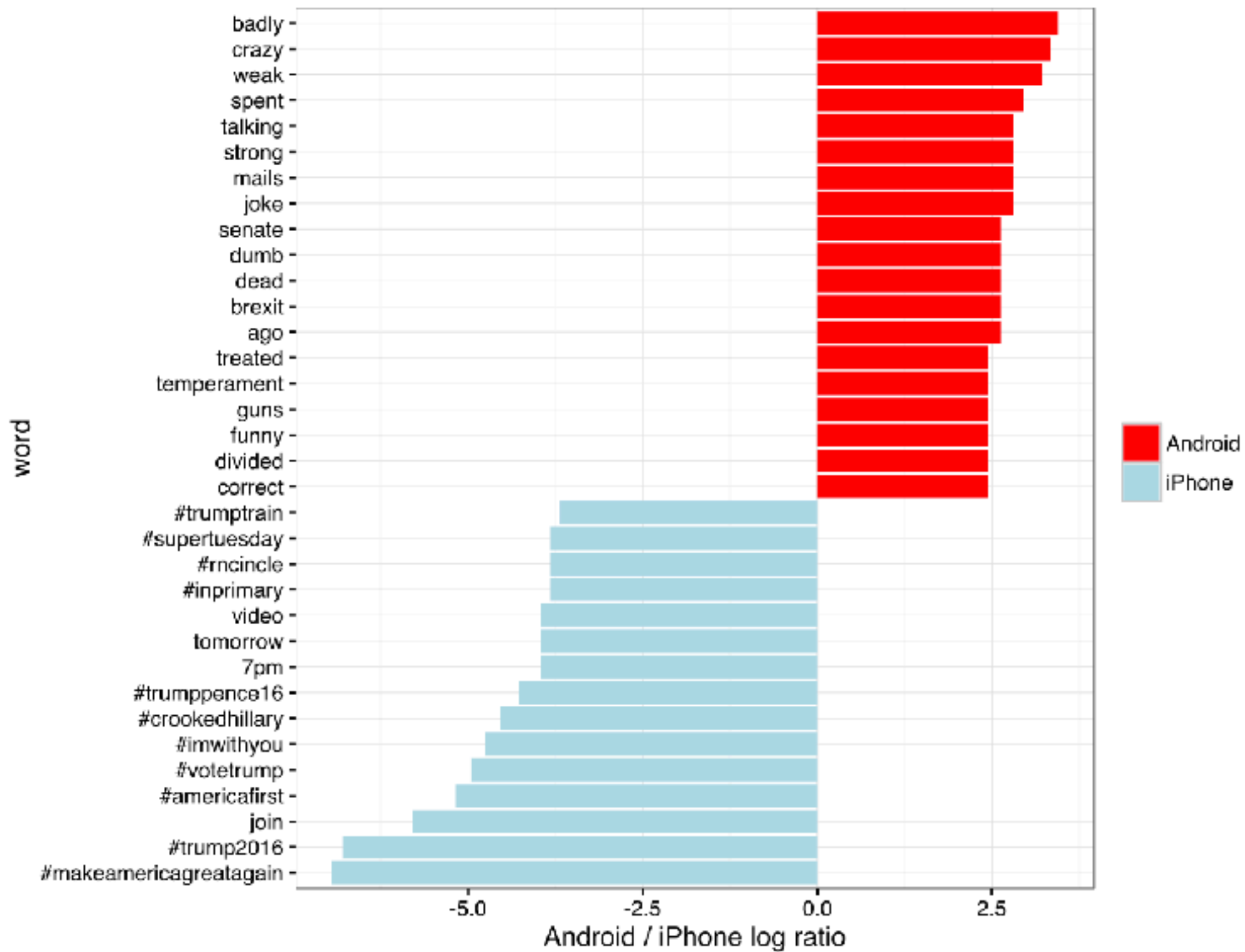
embryonic stem cell	Circuit Court of Appeals	Tongass national forest
hate crimes legislation	death tax repeal	pluripotent stem cells
adult stem cells	housing and urban affairs	Supreme Court of Texas
oil for food program	million jobs created	Justice Priscilla Owen
personal retirement accounts	national flood insurance	Justice Janice Rogers
energy and natural resources	oil for food scandal	American Bar Association
global war on terror	private property rights	growth and job creation
hate crimes law	temporary worker program	natural gas natural
change hearts and minds	class action reform	Grand Ole Opry
global war on terrorism	Chief Justice Rehnquist	reform social security

Gentzkow and Shapiro (2006), "What Drives Media Slant? Evidence from U.S. Daily Newspapers," *Econometrica*





Which are the words most likely to be from Android and most likely from iPhone?



<http://varianceexplained.org/r/trump-tweets/>

# Distinctive terms

- Finding distinctive terms is useful:
  - As a data exploration exercise to understand larger trends in individual word differences).
  - As a pre-processing step of feature selection.
- When the two datasets are  $A$  and  $\neg A$ , these terms also provide insight into what  $A$  is about.
- Many methods for finding these terms! (Developed in NLP, corpus linguistics, political science, etc.)

# Difference in proportions

For word  $w$  written by author with label  $k$  (e.g., {democrat, republican}), define the frequency to be the normalized count of that word

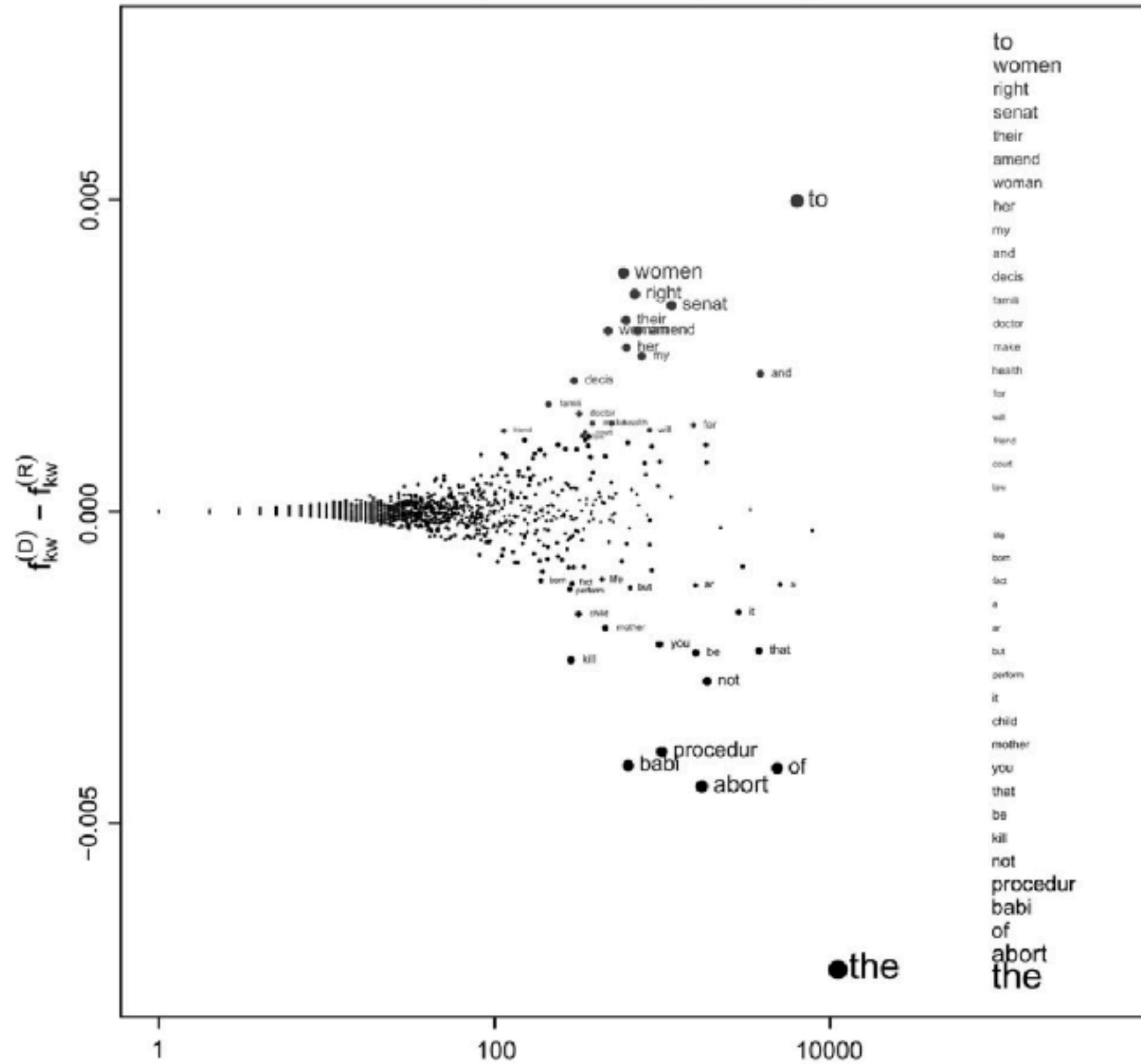
$$f_{w,k} = \frac{C(w, k)}{\sum_{w'} C(w', k)}$$

count of word  $w$  in group  $k$

count of all words in group  $k$

$$f_{w,k=\text{dem}} - f_{w,k=\text{repub}}$$

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





# Difference in proportions

- The difference in proportions is a conceptually simple measure and easily interpretable.
- Drawback: tends to emphasize words with high frequency (where even comparatively small differences in word usage between groups is amplified).
- Also, no measure whether a difference is statistically meaningful. We have **uncertainty** about the what the true proportion is for any group.

$$\chi^2$$

- $\chi^2$  (chi-square) is a statistical test of dependence—here, dependence between the two variables of word identity and corpus identity.
- For assessing the difference in two datasets, this test assumes a 2x2 contingency table:

	word	$\neg$ word
corpus 1	7	104023
corpus 2	104	251093

$$\chi^2$$

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

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



$$\chi^2$$

For each cell in contingency table, sum the squared difference between observed value in cell and the expected value assuming independence.

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

	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:

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

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

	robot	¬robot		
sci-fi	7.69	1095.2	$P(\text{scifi})$	0.076
¬sci-fi	93.9	13315.2	$P(\neg \text{scifi})$	0.924
			$P(\text{robot})$	$P(\neg \text{robot})$
			0.007	0.993



$$\chi^2$$

- What  $\chi^2$  is asking is: how different are the observed counts different from the counts we would expect given complete independence?

	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

$$\chi^2$$

- With algebraic manipulation, simpler form for 2x2 table O (cf. Manning and Schütze 1999)

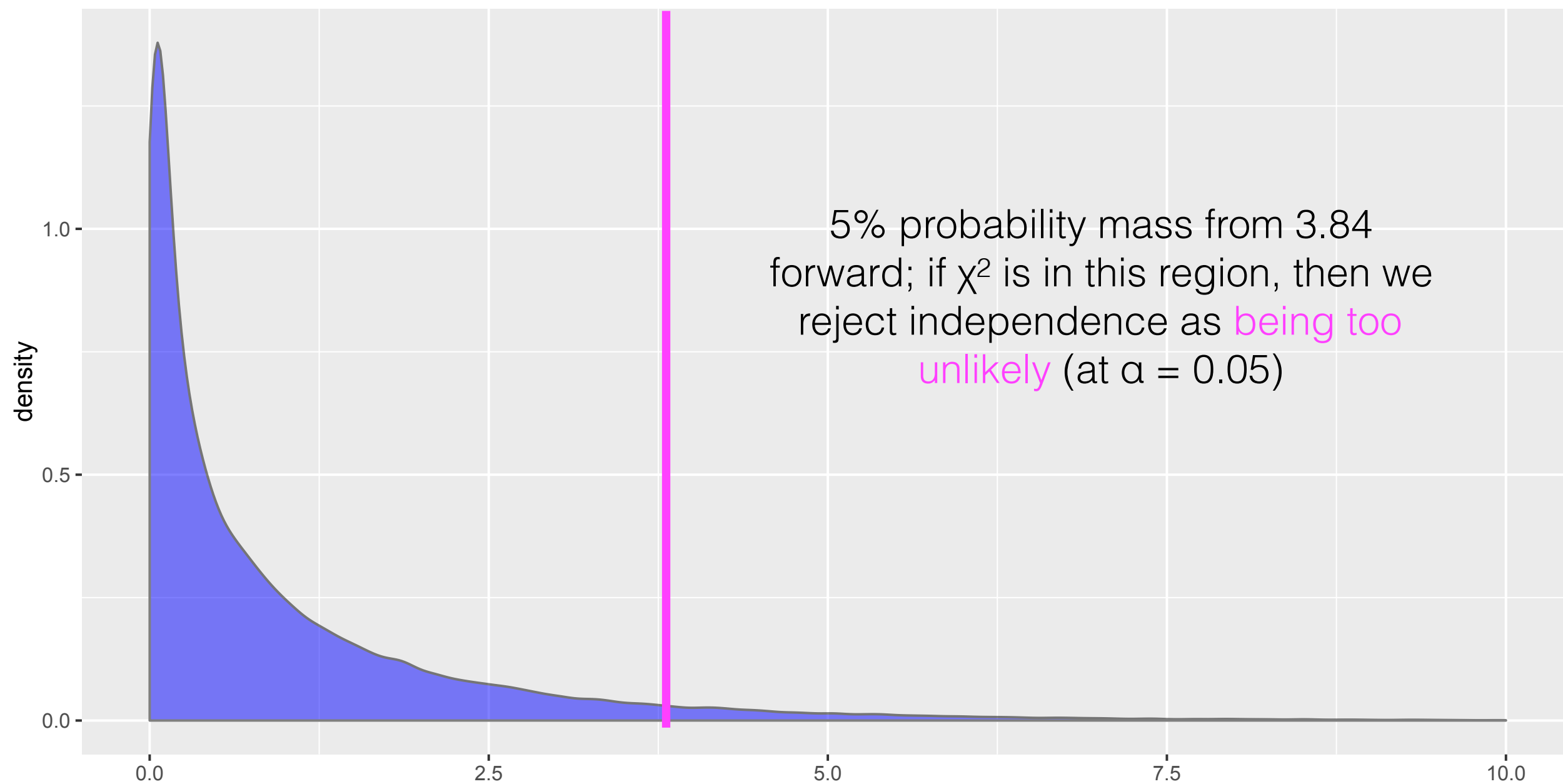
$$\chi^2 = \frac{N(O_{11}O_{22} - O_{12}O_{21})^2}{(O_{11} + O_{12})(O_{11} + O_{21})(O_{12} + O_{22})(O_{21} + O_{22})}$$

$$\chi^2$$

- The  $\chi^2$  value is a statistic of dependence with a probability governed by a  $\chi^2$  distribution; if this value has low enough probability in that measure, we can reject the null hypothesis of the independence between the two variables.



$$\chi^2$$



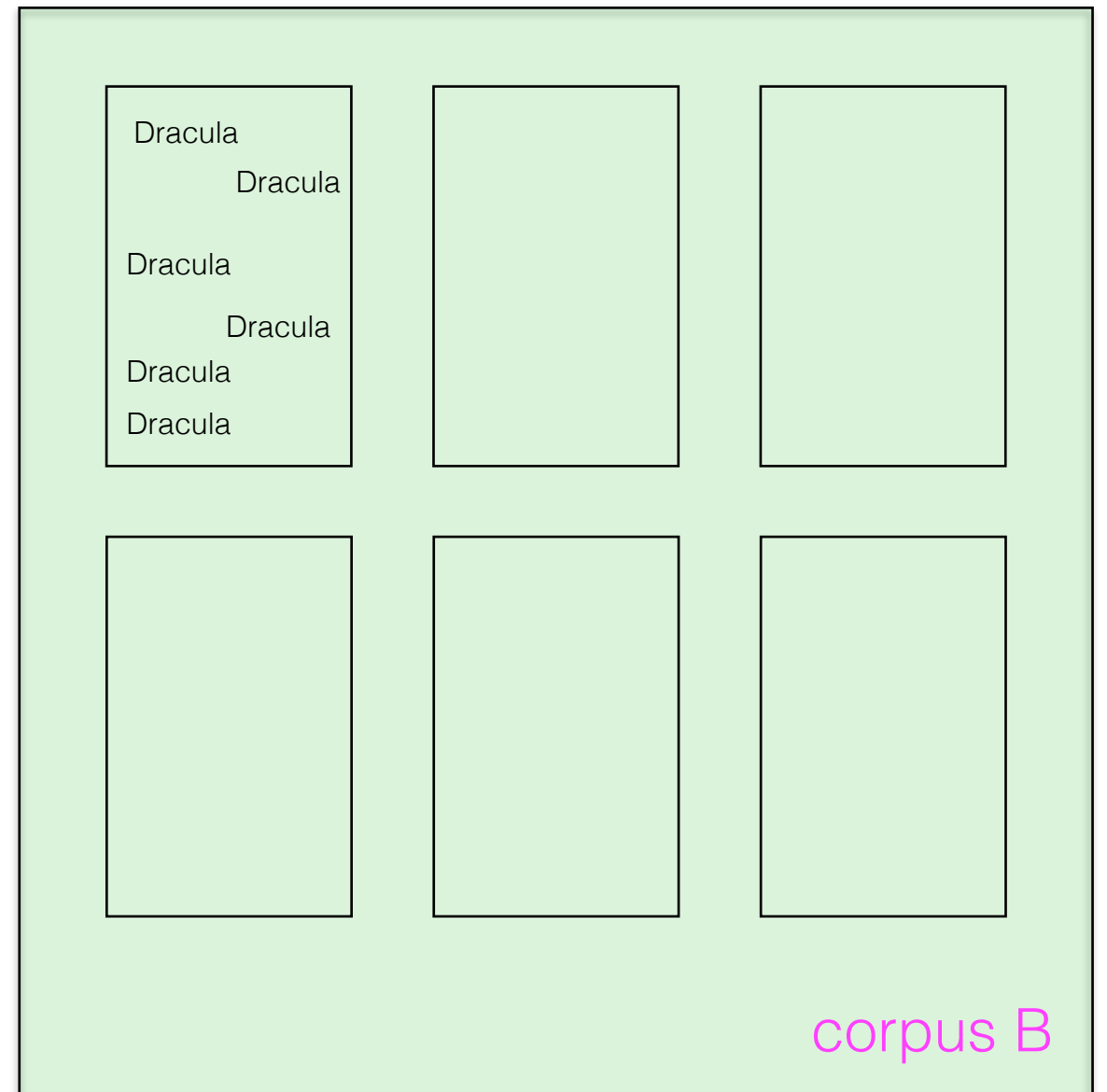
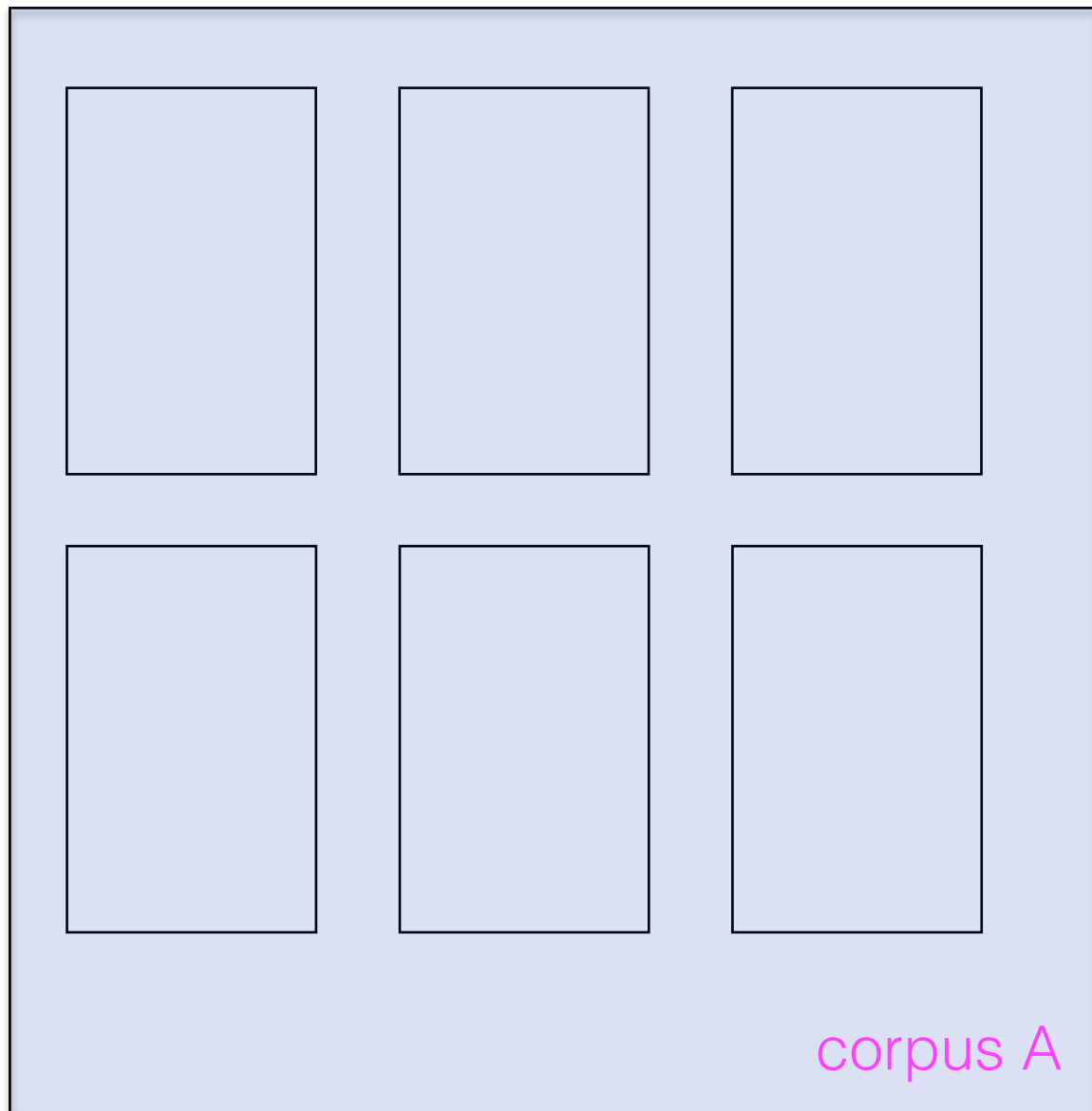
$$\chi^2$$

- Chi-square is ubiquitous in corpus linguistics (and in NLP as a measure of collocations).
- A few caveats for its use:
  - Each cell should have an *expected* count of at least 5
  - Each observation is independent

$$\chi^2$$

- A drawback, however, is due to the burstiness of language: the tendency for the same words to clump together in texts.
- Chi-square is testing for independence of two variables (word identity and corpus identity), but it **assumes** each mention of the word is independent from the others.





- Is Dracula really a word that distinguishes these two corpora?
- It distinguishes one text, but otherwise doesn't appear in the corpus at all.

# Mann-Whitney rank sums test

- Mann-Whitney is a test of the difference in some quantity of interest in two datasets. Null hypothesis: if you select a random sample from group A and another from group B, just as likely that A will be greater than B as less than B.

A	A	A	A	A	A	A	A
1	2	1	4	3	2	0	1

B	B	B	B	B	B
8	4	9	7	6	10

# Mann-Whitney

A	A	A	A	A	A	A	A
1	2	1	4	3	2	0	1

B	B	B	B	B	B
8	4	9	7	2	10

A	A	A	A	A	A	B	A	B	B	B	B	B	B
0	1	1	1	2	2	2	3	4	4	7	8	9	10
1	2	3	4	5	6	7	8	9	10	11	12	13	14

ranks

# Mann-Whitney

A	A	A	A	A	A	A	A
1	2	1	4	3	2	0	1

B	B	B	B	B	B
8	4	9	7	2	10

A	A	A	A	A	A	B	A	B	B	B	B	B	B
0	1	1	1	2	2	2	3	4	4	7	8	9	10
1	2	3	4	5	6	7	8	9	10	11	12	13	14

ranks

$$R_1 = 7+9+10+11+12+13+14 = 76$$

# Mann-Whitney

$$R_1 = 7+9+10+11+12+13+14 = 76$$

$$U_1 = R_1 - \frac{n_1(n_1 + 1)}{2}$$

- Once we have this U value, we can ask whether it's significantly different from the average value we would expect if there's no difference between the two groups at all.



A	A	A	A	A	A	A	A
1	2	1	4	3	2	0	1

B	B	B	B	B	B
8	4	9	7	6	10

- In corpus linguistics, each measurement is the count of a word in a fixed-sized **chunk** of text (e.g., 500 words).
- This lets us accommodate a more realistic assumption about the burstiness of language.

A	A	A	A	A	A	A	A
0	0	0	500	0	0	0	0

B	B	B	B	B	B
0	0	0	0	0	0

500 mentions  
of Dracula in  
one book

Not a significant  
difference in ranks

A	B	A	A	B	A	A	B	A	B	A	B	A	B
0	0	0	0	0	0	0	0	0	0	0	0	0	500
0	0	3	4	5	6	7	8	9	10	11	12	13	14

# Other methods

- There are many other methods for learning distinguishing words between two corpus; major classes:
  - Model-based methods that assume parametric forms + Bayesian priors (for smoothing) [Monroe et al. 2009]
  - Methods using classification to learn informative features that separate classes.