



# Natural Language Processing

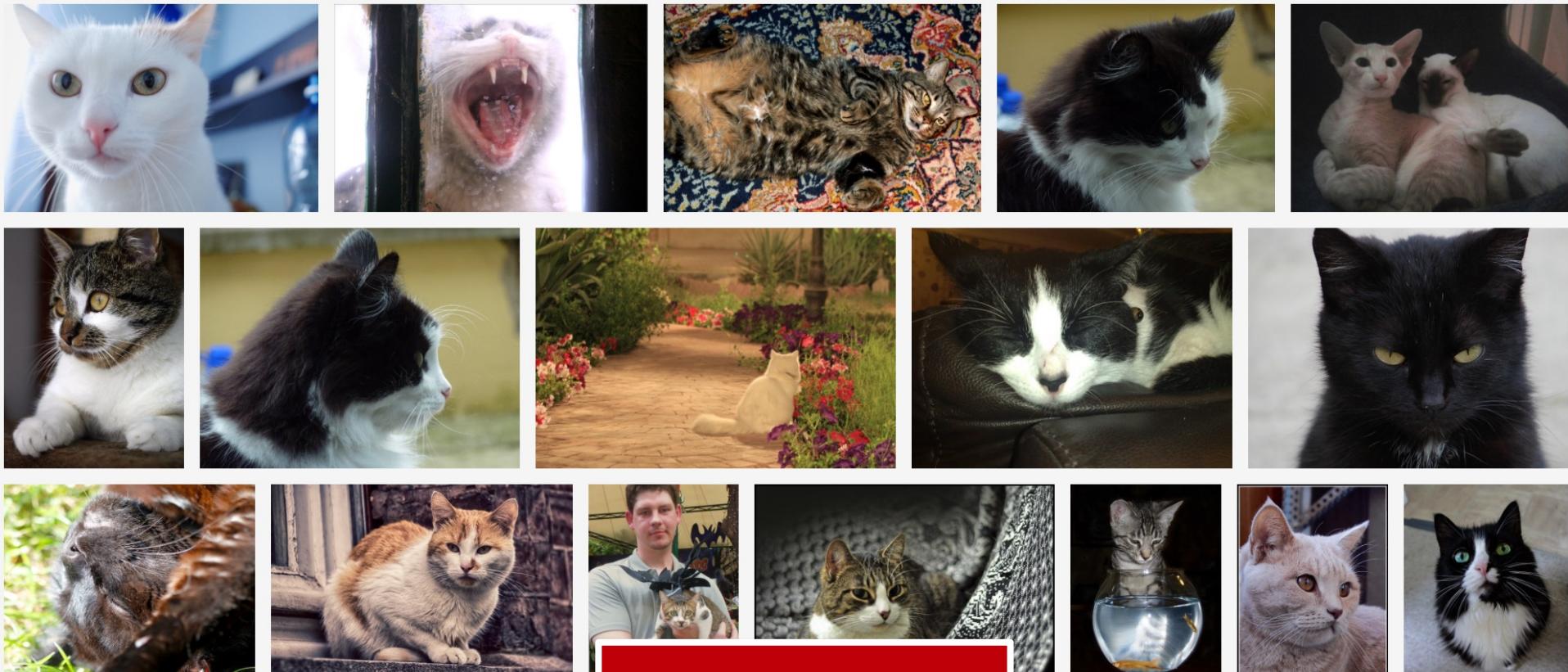
## 02: Basic Text Processing

---

Philipp Schaer, Technische Hochschule Köln, Cologne, Germany

Version: 2022-04-01

# Words as dimensionality reduction



“cats”

# Words

- One morning I shot an elephant in my pajamas
- I didn't shoot an elephant
- **Imma** let you finish but Beyonce had one of the best videos of all time
- 一天早上我穿着睡衣射了一只大象

# Words

- @phschaer have you seen this :) http://popvssoda.com

Tokenization  
before Twitter

@  
phschaer  
have  
you  
seen  
this  
:  
)  
http  
:  
//popvssoda.com

# Emoticons

Icon												Meaning	
:-( :)	:-[ :]	:3- :3	:>- >	8-) 8)	:{-} :}	:o) :o	:c) :c	:^) :^	=] =]	=) =)	Smiley or happy face. <a href="#">[4]</a> <a href="#">[5]</a> <a href="#">[6]</a>		
:D-D :D 8D	8-D xD	x-D xD	X-D XD	=D	=3	B^D						Laughing, <a href="#">[4]</a> big grin, <a href="#">[5]</a> <a href="#">[6]</a> laugh with glasses, <a href="#">[7]</a> or wide-eyed surprise <a href="#">[8]</a>	
:-))												Very happy or double chin <a href="#">[7]</a>	
:-( :( :c	:c-	:< :<	:-[ :]	:	>:[ :{	:{	:@)	>:( :					Frown, <a href="#">[4]</a> <a href="#">[5]</a> <a href="#">[6]</a> sad, <a href="#">[9]</a> angry, <a href="#">[7]</a> pouting
:-( :('												Crying <a href="#">[9]</a>	
:'-) :')												Tears of happiness <a href="#">[9]</a>	
D-': D-:	D:< D:	D: D8	D; D;	D= DX									Horror, disgust, sadness, great dismay <a href="#">[5]</a> <a href="#">[6]</a> (right to left)
:O-O :O	:o- :o	:0 8-0	>:O									Surprise, <a href="#">[3]</a> shock, <a href="#">[4]</a> <a href="#">[10]</a> yawn <a href="#">[11]</a>	
:-* :*	:x												Kiss
:)-) :)*	*-) *)	:]- :]	;^) ;^	:-, :-,	;D								Wink, <a href="#">[4]</a> <a href="#">[5]</a> <a href="#">[6]</a> smirk <a href="#">[10]</a> <a href="#">[11]</a>
:P-P :P	X-P XP	x-p xp	:p- :p	:P- :P	:p- :b	:b- :b	d: d	=p =p	>:P >:P				Tongue sticking out, cheeky/playful, <a href="#">[4]</a> blowing a raspberry

# Types and tokens

- Type = abstract descriptive concept
- Token = instantiation of a type

To be or not to be

- 6 tokens (to, be, or, not, to, be)
- 4 types (to, be, or, not)
- Types = The **vocabulary**; the unique tokens.

# Types and tokens

- Type = abstract descriptive concept
- Token = instantiation of a type

How can we use types and tokens to measure vocabulary richness?

# Whitespace

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

- As much mud in the streets as if the waters had but newly retired from the face of the earth, and it would not be wonderful to meet a Megalosaurus, forty feet long or so, waddling like an elephantine lizard up Holborn Hill.

# Whitespace

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

- As much mud in the streets as if the waters had but newly retired from the face of the **earth**, and it would not be wonderful to meet a **Megalosaurus**, forty feet long or **so**, waddling like an elephantine lizard up Holborn **Hill**.

What do we lose with  
whitespace tokenization?

368	earth
135	earth,
68	earth.
26	earth
24	earth.
18	earth."
16	earth;
14	earth,
9	earth's
5	earth!"
5	earth!
4	earth;
4	earth,"
3	earth."
3	earth?

3	earth!"
2	earth--to
2	earth--if
2	earth--and
2	earth:
2	earth,'
1	earth-worms,
1	earth-worm.
1	earth--which
1	earth--when
1	earth--something
1	earth-smeared,
1	earth-scoops,
1	earth's
1	earth--oh,

# Punctuation

- We typically don't want to just strip all punctuation, however.
  - Punctuation signals **boundaries** (sentence, clausal boundaries, parentheticals, asides)
  - Some punctuation has **expressive force**, like exclamation points (!) and question marks (?)
  - Emoticons are strong signals of e.g. sentiment
- Most tokenization algorithms (for languages typically delimited by whitespace) use **regular expressions** to segment a string into discrete tokens.

# Regular expressions

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

- Python
  - **re.findall(regex, text)** finds all non-overlapping matches for a target regex.
  - `re.findall(r'[Tt]he', "The dog barked at the cat")`
  - `[“The”, “the”]`

# NLTK for the win

```
import nltk  
tokens=nltk.word_tokenize(text)
```

Tokenizes following the conventions of the **Penn Treebank**:

- punctuation split from adjoining words double quotes ("") changes to forward/ backward quotes based on their location in word ('`the'')
- verb contractions + 's split into separate
- tokens: (did\_n't, children\_ 's)

# NLTK for the win

```
import nltk  
tokens=nltk.word_tokenize(text)
```

- Penn Treebank tokenization is important because a lot of downstream NLP is trained on annotated data that uses Treebank tokenization!

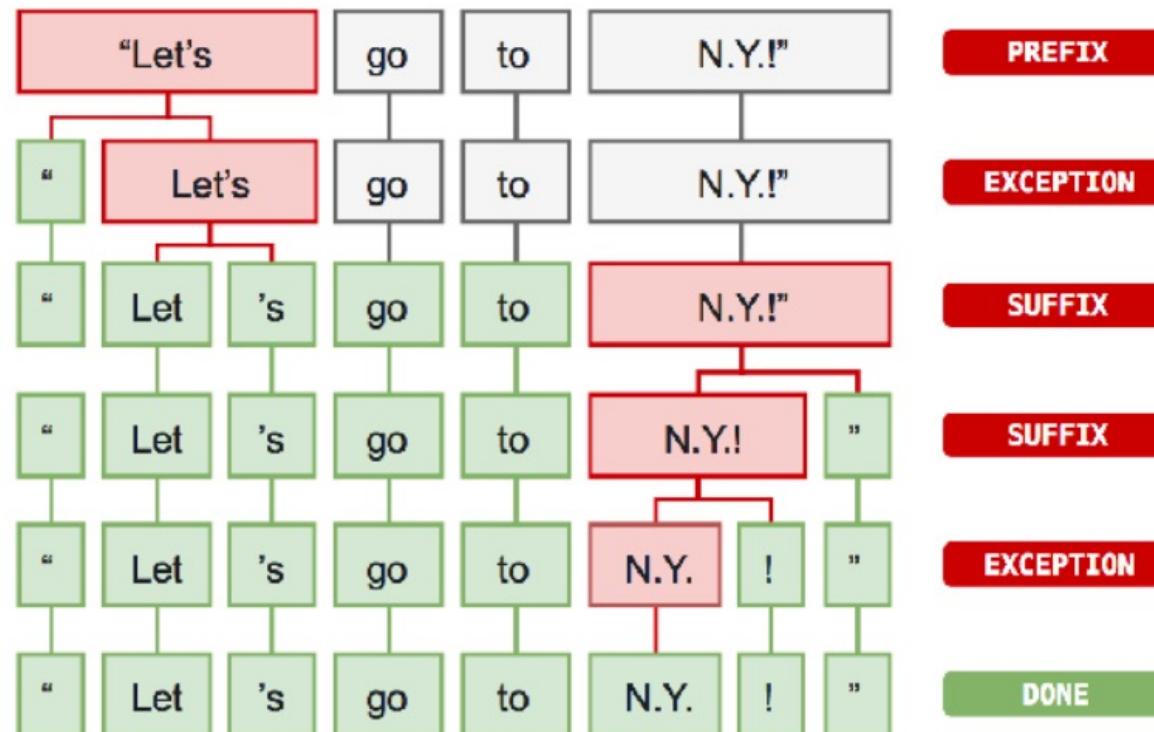
PRP VBD RB VB DT NN .  
I did n't see the parade .

PRP ??? VB DT NN .  
I didn't see the parade .

[https://www.ling.upenn.edu/courses/Fall\\_2003/ling001/penn\\_treebank\\_pos.html](https://www.ling.upenn.edu/courses/Fall_2003/ling001/penn_treebank_pos.html)

# Spacy

```
import spacy
nlp = spacy.load('en')
tokens=[token.text for token in nlp(text)]
```



# Sentence segmentation

- Word tokenization presumes a preprocessing step of sentence segmentation – Identifying the boundaries between sentences.
- Lots of NLP operates at the level of the sentence (POS tagging, **parsing**), so really important to get it right.
- Harder to write regexes to delimit these, since there are many cases where the usual delimiters (periods, question marks) serve double duty.

# Sentence segmentation

- “Do you want to go?” said Jane.
- Mr. Collins said he was going.
- He lives in the U.S. John, however, lives in Canada.

# Sentence segmentation

- **NLTK**: Punkt sentence tokenizer – unsupervised method to learn common abbreviations, collocations, sentence-initial words. Can be trained on data from new domain.
- <https://www.nltk.org/api/nltk.tokenize.html#module-nltk.tokenize.punkt>
- Kiss, Tibor and Strunk, Jan (2006): Unsupervised Multilingual Sentence Boundary Detection (*Computational Linguistics*)
  
- **spaCy**: Relies on dependency parsing to find sentence boundaries.

# Stemming and lemmatization

- Many languages have some inflectional and derivational morphology, where similar words have similar forms:  
  
organizes, organized, organizing
- Stemming and lemmatization reduce this variety to a single common **base form**.

# Stemming

- Heuristic process for chopping off the inflected suffixes words  
organizes, organized, organizing → **organ**
- Lower precision, higher recall
- Remember Information Retrieval, **Porter stemmer**: Sequence of rules for removing suffixes from words
  - EMENT → Ø
  - SSES → SS
  - IES → I
  - SS → Ø
  - S → Ø

# Lemmatization

- Using morphological analysis to return the dictionary form of a word (the entry in a dictionary you'd find all forms under)

organizes, organized, organizing → **organ**

```
import spacy  
nlp = spacy.load('en')  
lemmas=[token.lemma_ for token in nlp(text)]
```

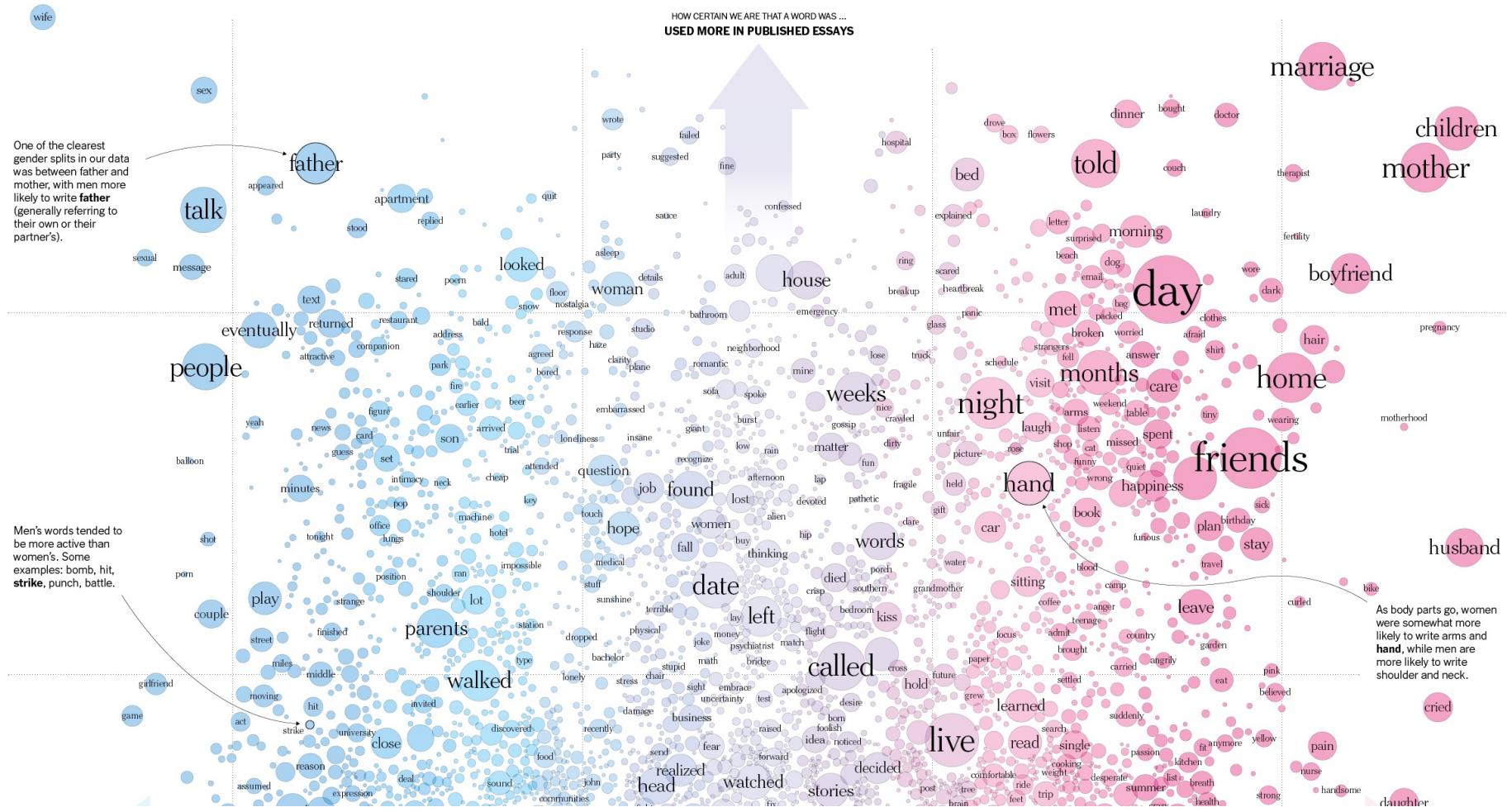
# Difficulties

- When does **punctuation** disrupt the desired boundaries of a token?

Emoticons	:) :D \o/ o_O
URLs	<a href="http://www.google.com">http://www.google.com</a>
Prices	\$19.99
Decimals	19.99
Hyphens	state-of-the-art
Usernames	@phschaer
Hashtags	#blacklivesmatter

# Finding Distinctive Terms

# The Words Men and Women Use When They Write About Love



### 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



 **Donald J. Trump** @realDonaldTrump  
Good luck #TeamUSA  
#OpeningCeremony #Rio2016  
[pic.twitter.com/mS8qsQpJPh](https://pic.twitter.com/mS8qsQpJPh)

27,391 Likes      8,392 Retweets

Aug 5, 2016 at 8:59 PM      via Twitter for iPhone



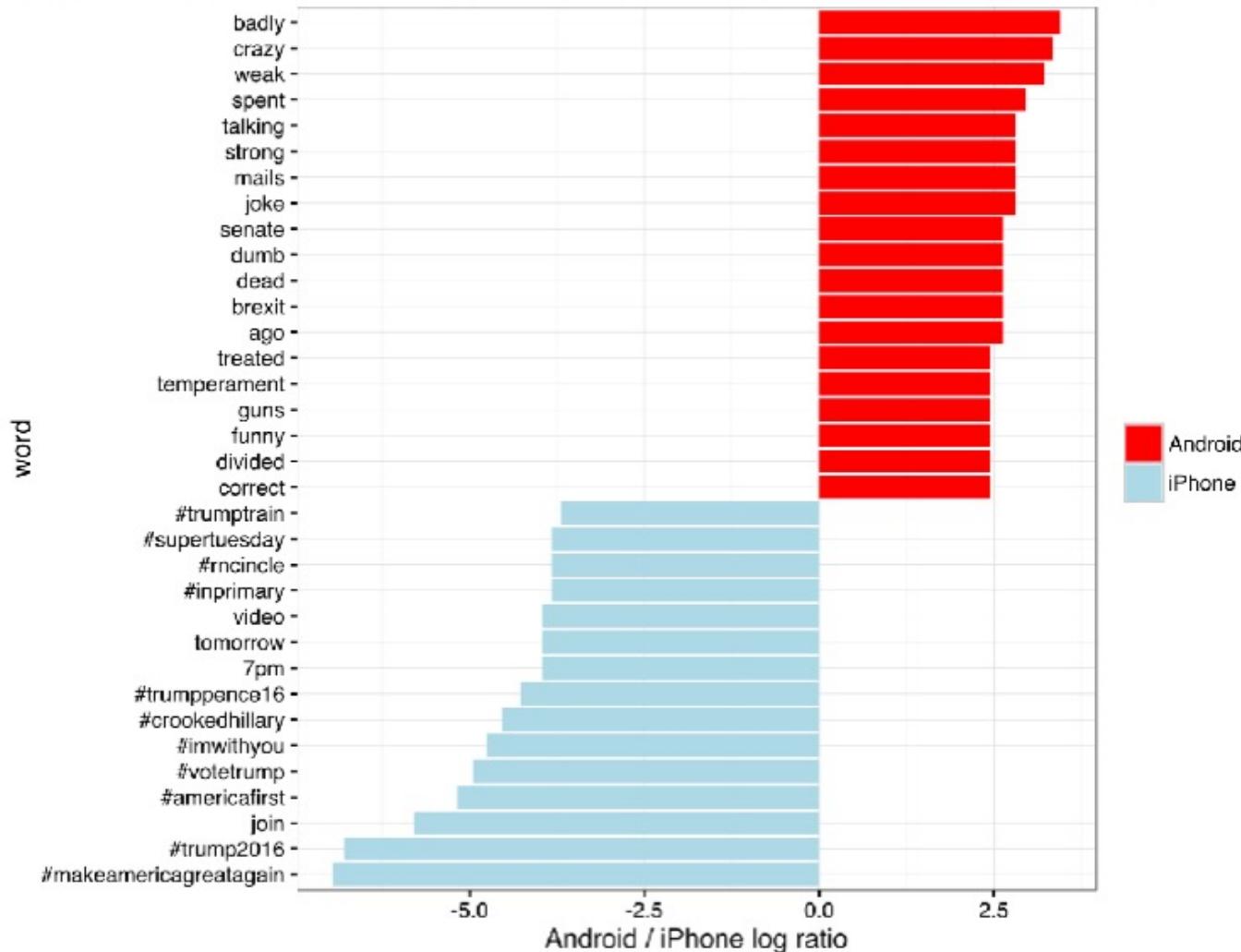
 **Donald J. Trump** @realDonaldTrump  
Heading to New Hampshire - will be  
talking about Hillary saying her brain  
SHORT CIRCUITED, and other things!

4,451 Likes      1,480 Retweets

Aug 6, 2016 at 11:11 AM      via Twitter for Android



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



# Distinctive terms

- Finding distinctive terms is useful:
  - As a pre-processing step of feature selection.
  - As a data exploration exercise to understand larger trends in individual word differences).
- When the two datasets are **A** and **¬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

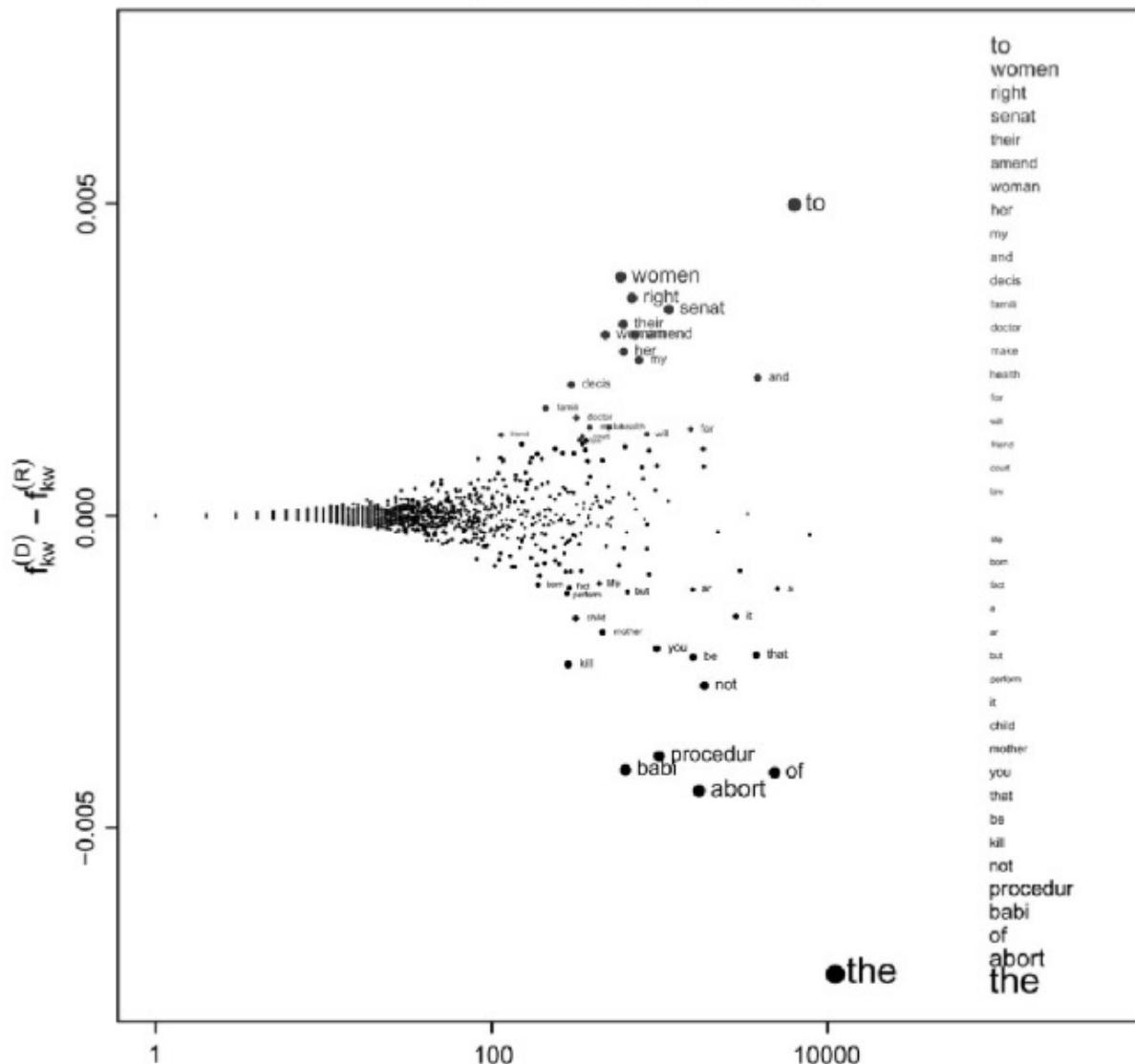
$$f_{w,k} = \frac{C(w, k)}{\sum_{w'} C(w', k)}$$

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.

# X<sup>2</sup>

- X<sup>2</sup> (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	¬word
corpus 1	7	104023
corpus 2	104	251093

# X<sup>2</sup>

- 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 ( $O_{ij}$ )** in cell and the **expected value ( $E_{ij}$ )** 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	1108/14512
$\neg$ sci-fi	2	13402	13404	0.924	
sum		106	14406	14512	
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 14,512 words, we would expect to see 7.69 occurrences of *robot* in sci-fi texts.

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

	$P(\text{robot})$	$P(\neg\text{robot})$
frequency	0.007	0.993

# X<sup>2</sup>

- What X<sup>2</sup> is asking is:  
How different are the observed counts from the counts we would expect given complete independence?

	robot	¬robot
sci-fi	104	1004
¬sci-fi	2	13402

	robot	¬robot
sci-fi	7.69	1095.2
¬sci-fi	93.9	13315.2

# $\chi^2$

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

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

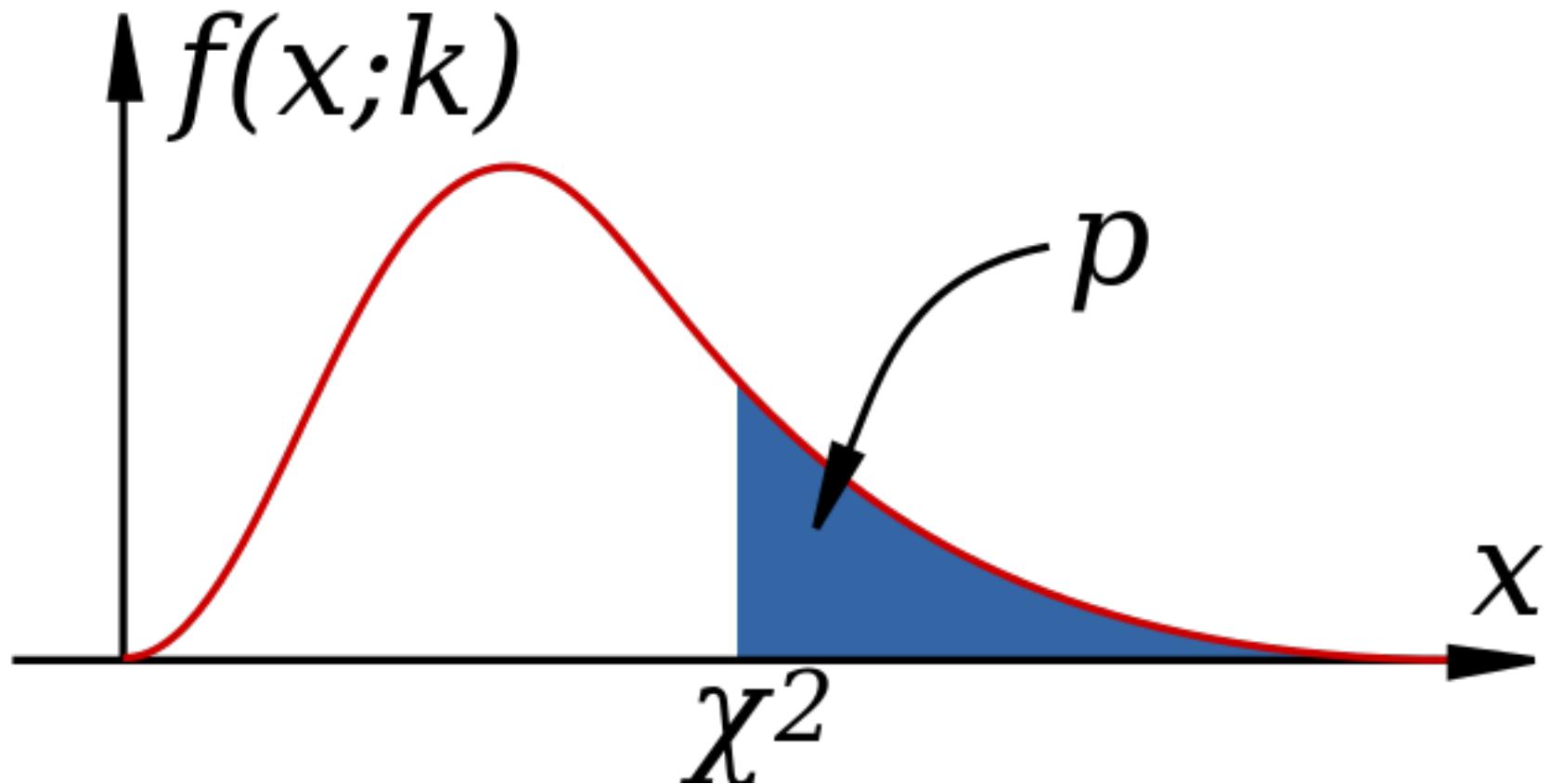
- $N = O_{11} + O_{12} + O_{21} + O_{22}$
  - The chi-square statistic  $\chi^2$  is **1239.5312**.
- |               |                 | robot           | $\neg$ robot |
|---------------|-----------------|-----------------|--------------|
| sci-fi        | O <sub>11</sub> | O <sub>12</sub> |              |
|               | O <sub>21</sub> | O <sub>22</sub> |              |
| $\neg$ sci-fi |                 |                 |              |

# $\chi^2$

- 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.
- $k$  = degrees of freedom (#rows-1)(#cols-1)

<b>k</b>	<b>0.995</b>	<b>0.99</b>	<b>0.975</b>	<b>0.95</b>	<b>0.90</b>	<b>0.10</b>	<b>0.05</b>	<b>0.025</b>	<b>0.01</b>	<b>0.005</b>
<b>1</b>	---	---	0.001	0.004	0.016	2.706	<b><u>3.841</u></b>	5.024	6.635	7.879
<b>2</b>	0.010	0.020	0.051	0.103	0.211	4.605	5.991	7.378	9.210	10.597
<b>3</b>	0.072	0.115	0.216	0.352	0.584	6.251	7.815	9.348	11.345	12.838

# $\chi^2$ distribution

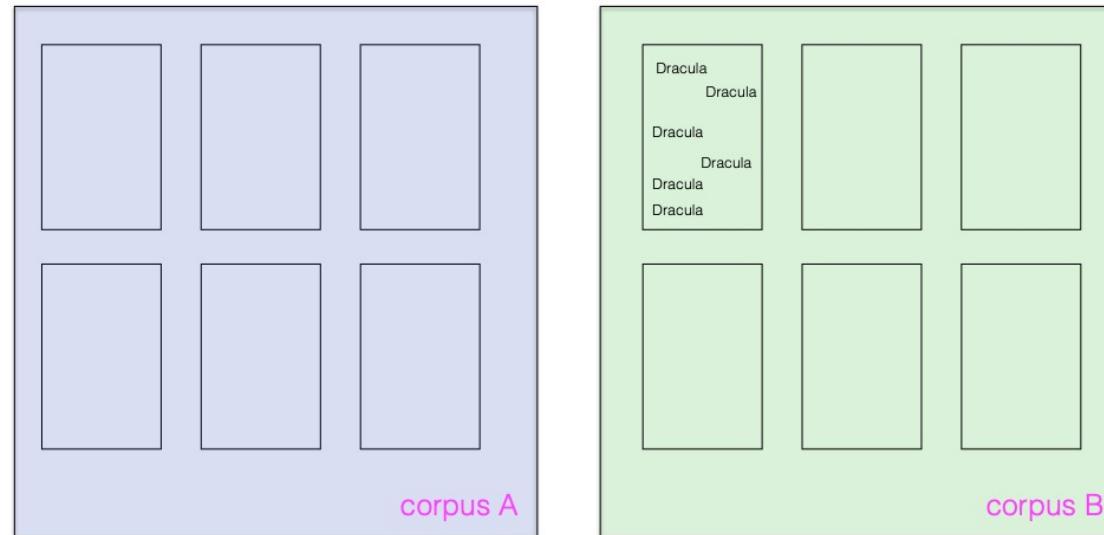


# X<sup>2</sup>

- 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
- 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.

# X<sup>2</sup>

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



- Alternatives to X<sup>2</sup>: Mann-Whitney or G-Test

# Mental homework

- Hypothesize terms that will be different between groups in the ESUPOL dataset (male vs. female, CDU vs. SPD, etc.)
- Execute chi-square to find terms that are different
- Online  $\chi^2$   
<https://www.socscistatistics.com/tests/chisquare/default2.aspx>