

# Automated Text Analysis in R



D-Lab INTENSIVE  
Instructor: Laura Nelson  
June 3-4, 2014

# Theory

- Artificial Language

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n + \varepsilon$$

→ 1.	$(x)(Q \supset Fx)$	
2.	$Q \supset Fx$	I, UI
→ 3.	$Q$	
4.	$Fx$	2, 3, M.P.
5.	$(x)Fx$	4, UG
6.	$Q \supset (x)Fx$	3 – 5, C.P.
7.	$(x)(Q \supset (x)Fx) \supset [Q \supset (x)Fx]$	1 – 6, C.P.

```
import scipy
from scipy import sparse
```

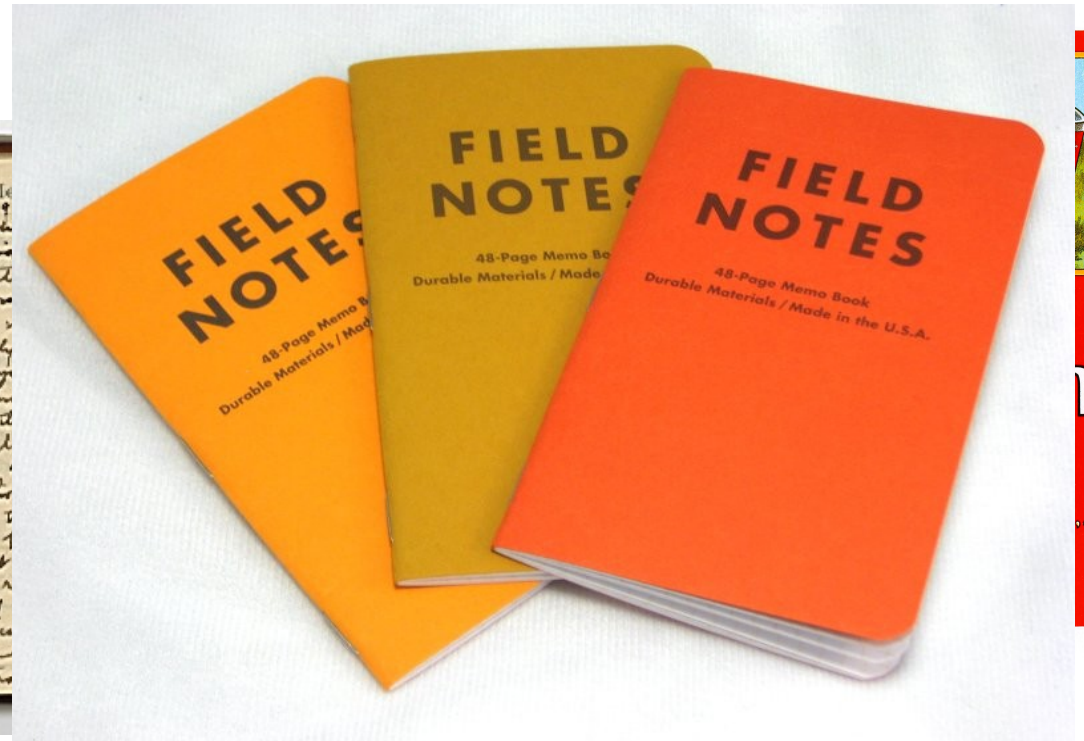
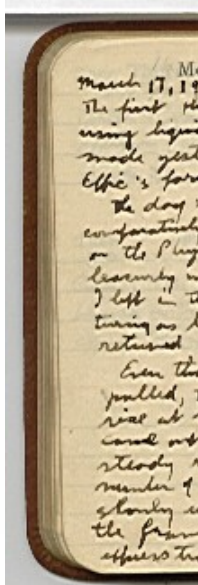
```
n = 200000
matrix = scipy.sparse.rand(n, n, density=.001)
print matrix
```

- Natural Language

- “Time **flies like** an arrow. Fruit **flies like** a banana.”

- What does “cute” mean? **Keen**, or **pretty**?

# Natural Language as Data

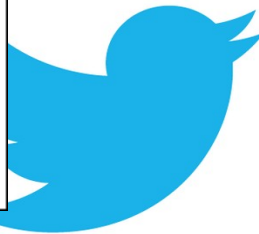


## Exhibit Feedback

### 1. Please explain below:

What did you think overall?

What would you improve?



# How Do We Analyze Text?

“We need to steer clear of this poverty of ambition, where people want to drive fancy cars and wear nice clothes and live in nice apartments but don't want to work hard to accomplish these things. Everyone should try to realize their full potential.” -Barack Obama



# How Do We Analyze Text?

“We need to steer clear of this poverty of ambition, where people want to drive fancy cars and wear nice clothes and live in nice apartments but don't want to work hard to accomplish these things. Everyone should try to realize their full potential.” -Barack Obama



# How Do We Analyze Text?

“We **need to steer** clear of this poverty of ambition, where people **want to drive** fancy cars and **wear** nice clothes and **live** in nice apartments but don't **want to work** hard **to accomplish** these things. Everyone should **try to realize** their full potential.” -Barack Obama

# Why Use Computer-Assisted Methods?

- Speed
  - Humans are slow
  - Text is becoming large
- Reliability / Reproducibility
- Validity
  - Expanded memory
  - Unburdened by bias

**Does not remove the need for interpretation!**

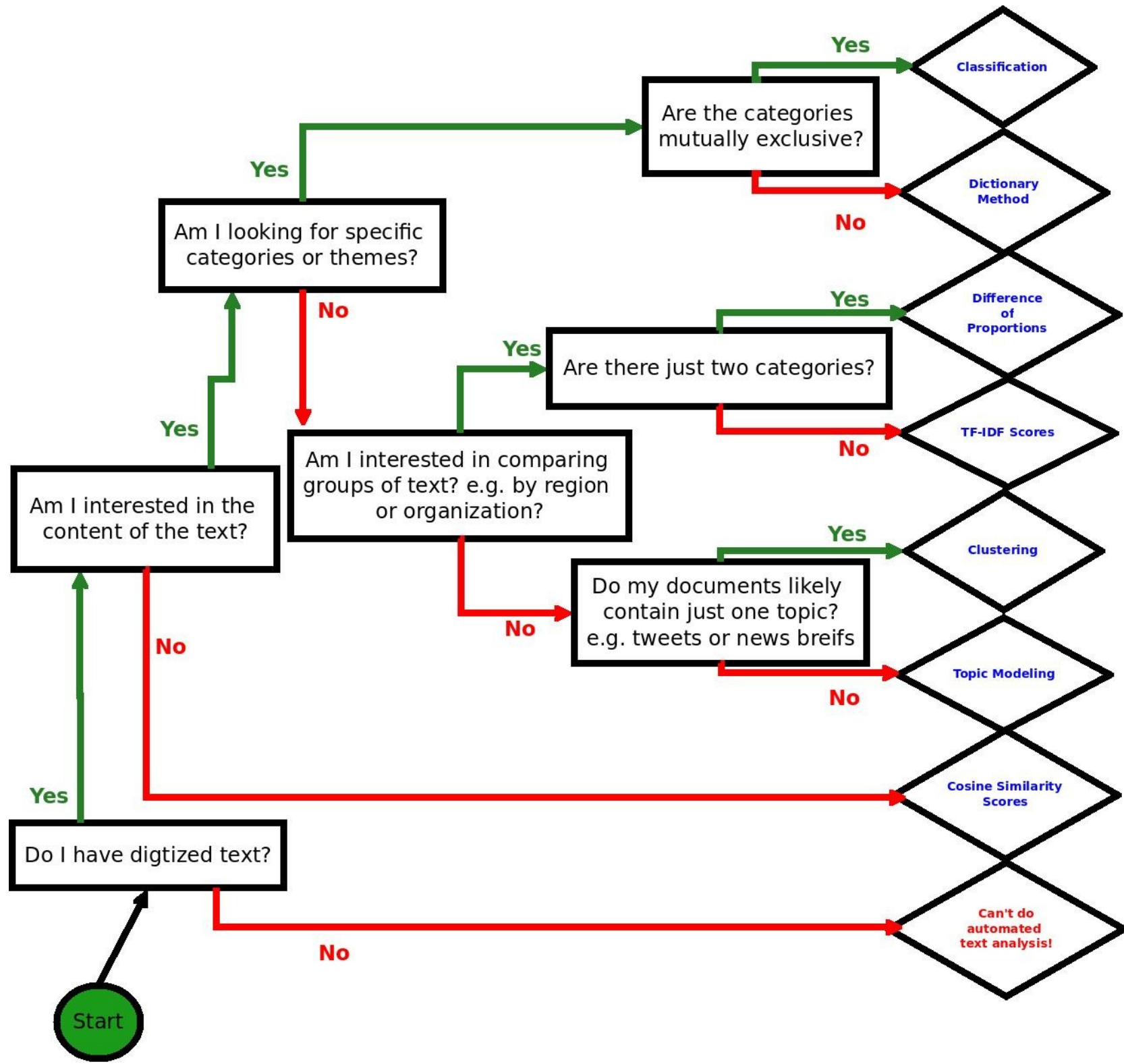
# Important Terms

- Corpus
  - Collection of texts/documents
- Lexical: fancy name for word
- N-gram
  - poverty: uni-gram
  - poverty of: bi-gram
  - poverty of ambition: tri-gram



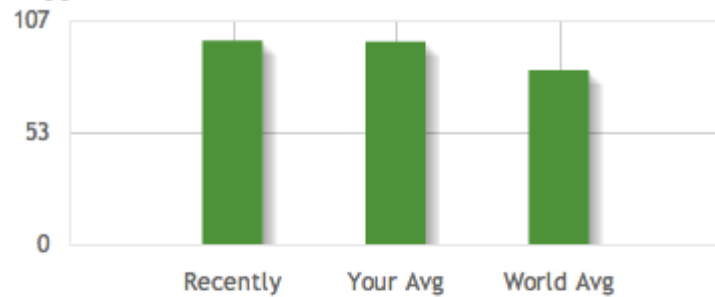
# Types of Automated Text Analysis

- Classification (deductive)
  - Dictionaries
  - Supervised Machine Learning
- Lexical selection (inductive)
  - Difference of proportions
  - Word scores, tf-idf
- Latent categorical analysis (inductive)
  - Clustering
  - Topic modeling

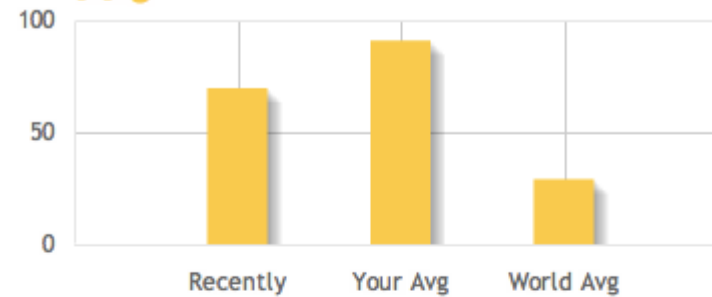


# Classification: Dictionaries

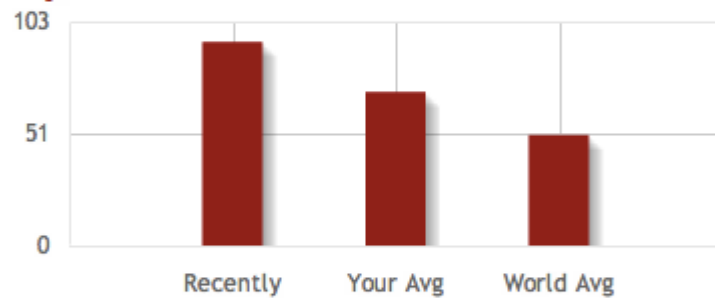
**Affectionate**



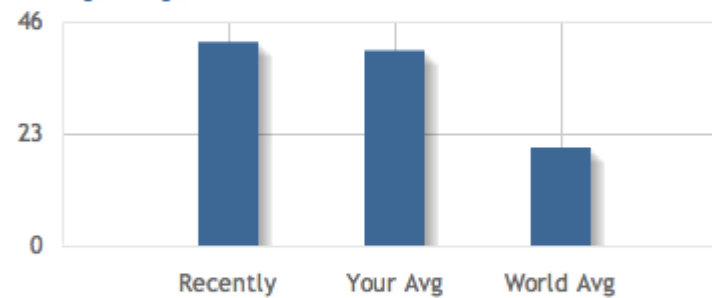
**Happy**



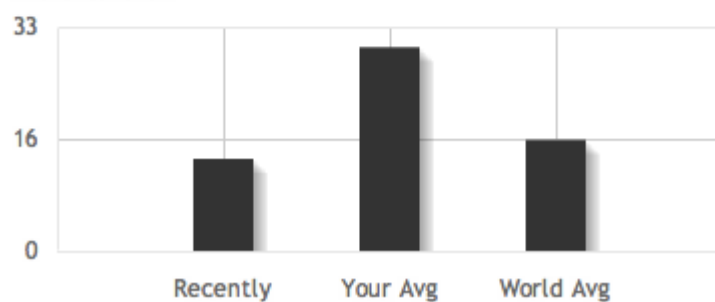
**Upset**



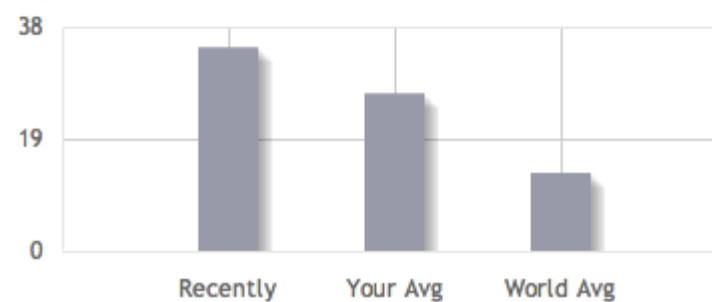
**Self-Expressive**



**Anxious**

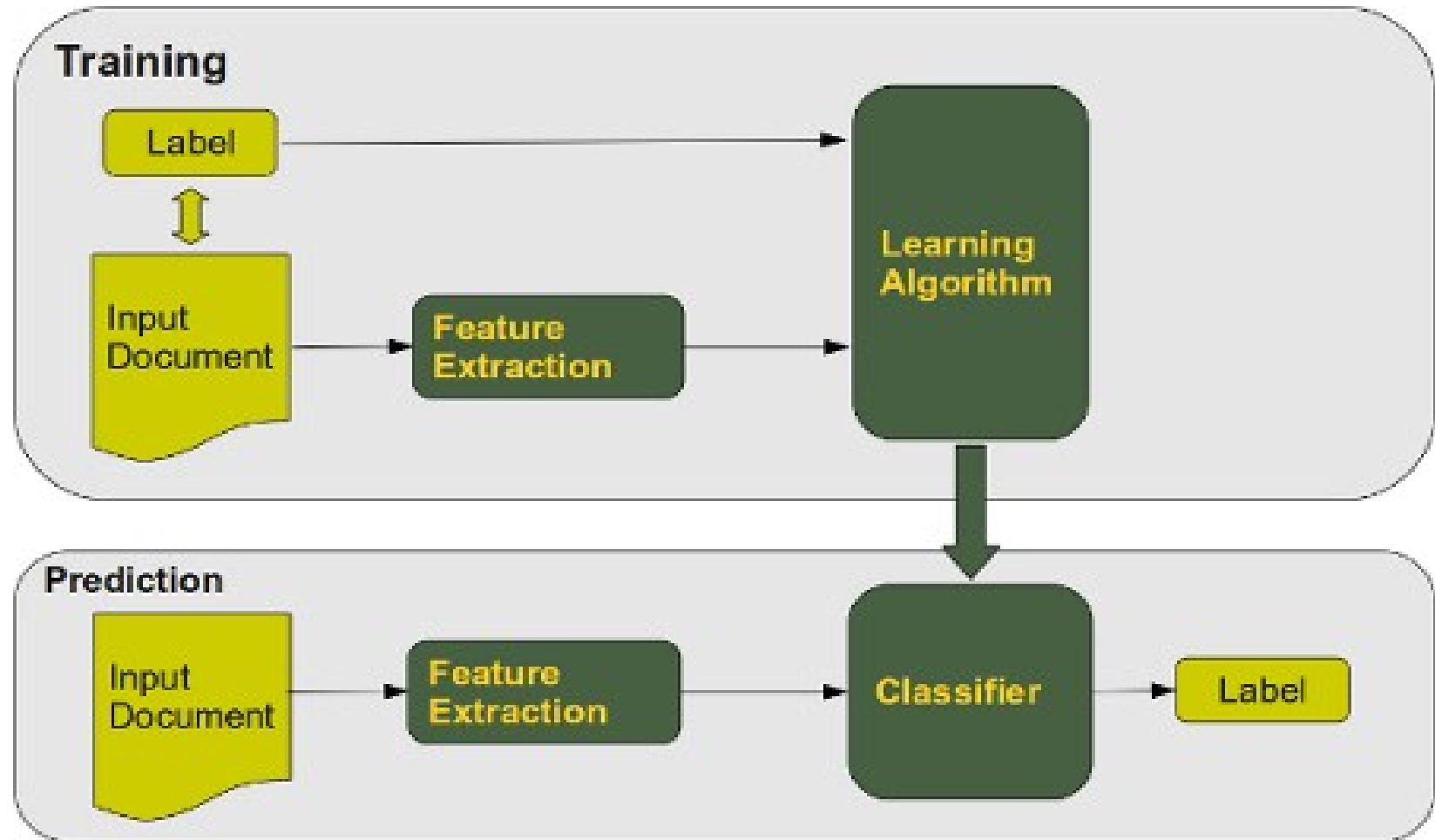


**Sad**



# Classification:

## Semi-Automated Machine Learning



# Classification:

## Semi-Automated Machine Learning

Label	Document Text	Function
Anti-materialist	Document 1 Text	Training
Pro-hard-work	Document 2 Text	Training
other	Document 3 Text	Training
Pro-hard-work	Document 4 Text	Test
Pro-hard-work	Document 5 Text	Test
Anti-materialist	Document 6 Text	Test
?	Document 7 Text	Unknown
?	Document 8 Text	Unknown
?	Document 9 Text	Unknown
?	Document 10 Text	Unknown

# Example: R

[github.com/lknelson/D-Lab--Text-Analysis-Workshop](https://github.com/lknelson/D-Lab--Text-Analysis-Workshop)



# Analyzing Output

$$\text{precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$

\*for each class, how many assigned to that class were actually in that class?

$$\text{recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

\*for each class, what percentage of the documents actually in that class assigned to that class?

# Analyzing Output

F score = harmonic mean of precision and recall

$$= \frac{\text{precision} * \text{recall}}{\text{precision} + \text{recall}} * 2$$

\*\*\*\*Tests accuracy

Problem of over-fitting: the model will perform poorly on unseen data.

Solution: Cross Validation

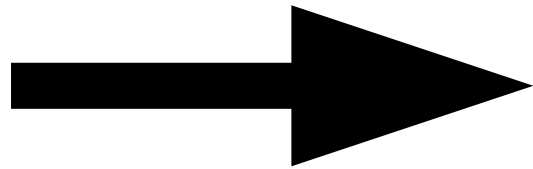
# Recap

- The goal: analyze natural languages, effectively *and* efficiently
- The approach: Utilize computational tools to speed up the process, make it more reliable, and possibly more valid
- First method, used for deductive analysis: Classification
- Today:
  - Lexical Selection
  - Latent Categorical Analysis, or Automated Machine Learning

# What is Machine Learning?

- Arthur Samuel defined machine learning as a "Field of study that gives computers the ability to learn without being explicitly programmed."
- Yesterday: computers *learn* how to recognize categories given by you (supervised machine learning, or *prediction*)
- Today: computers *learn* what categories arise from the text itself (unsupervised machine learning, or *discovery*)

# Automated Inductive Text Analysis: The Goal



Informative  
Groups of Words

# Pre-Processing

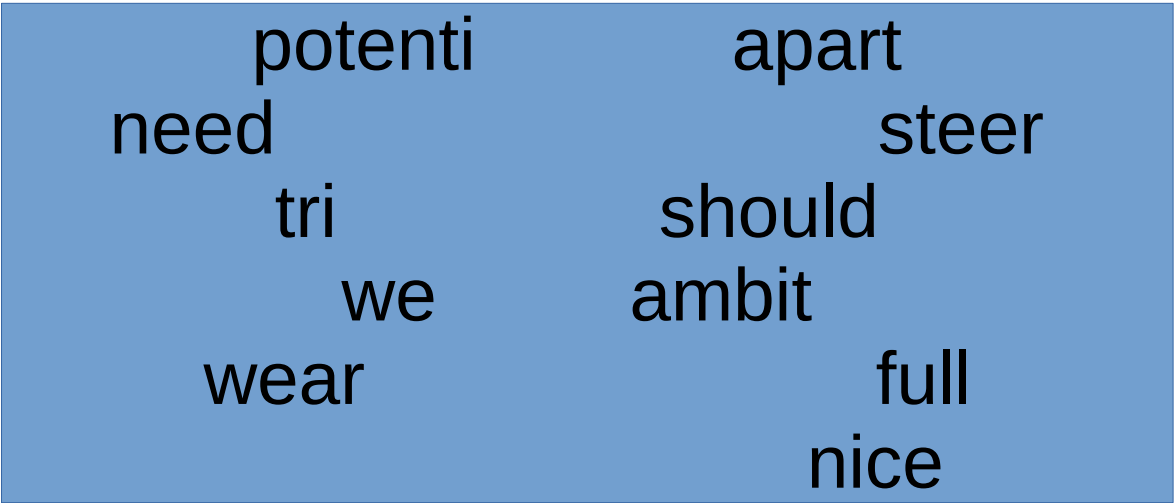
“We need to steer clear of this poverty of ambition, where people want to drive fancy cars and wear nice clothes and live in nice apartments but don't want to work hard to accomplish these things. Everyone should try to realize their full potential.”



# Pre-Processing

we need steer clear **poverti ambit**  
where **peopl** want drive **fanci** car wear  
nice **cloth** live nice **apart** want work  
hard accomplish **thing everyon** should  
**tri realiz** full **potenti**

# Pre-Processing



potenti      apart  
need      steer  
tri      should  
we      ambit  
wear      full  
         nice

# Document-Term Matrix

	ambit	poverti	peopl	full
Document1	4	2	0	0
Document2	1	3	7	0
Document3	2	0	0	0
Document4	9	1	4	0
Document5	0	0	0	6

# Lexical Selection: Difference of Proportions

	ambit	poverti	peopl	full	<i>Total</i>
Document1	.57	.29	0	.14	7
Document2	.09	.27	.64	0	11
Diff of Prop	.48	.02	-.64	.14	

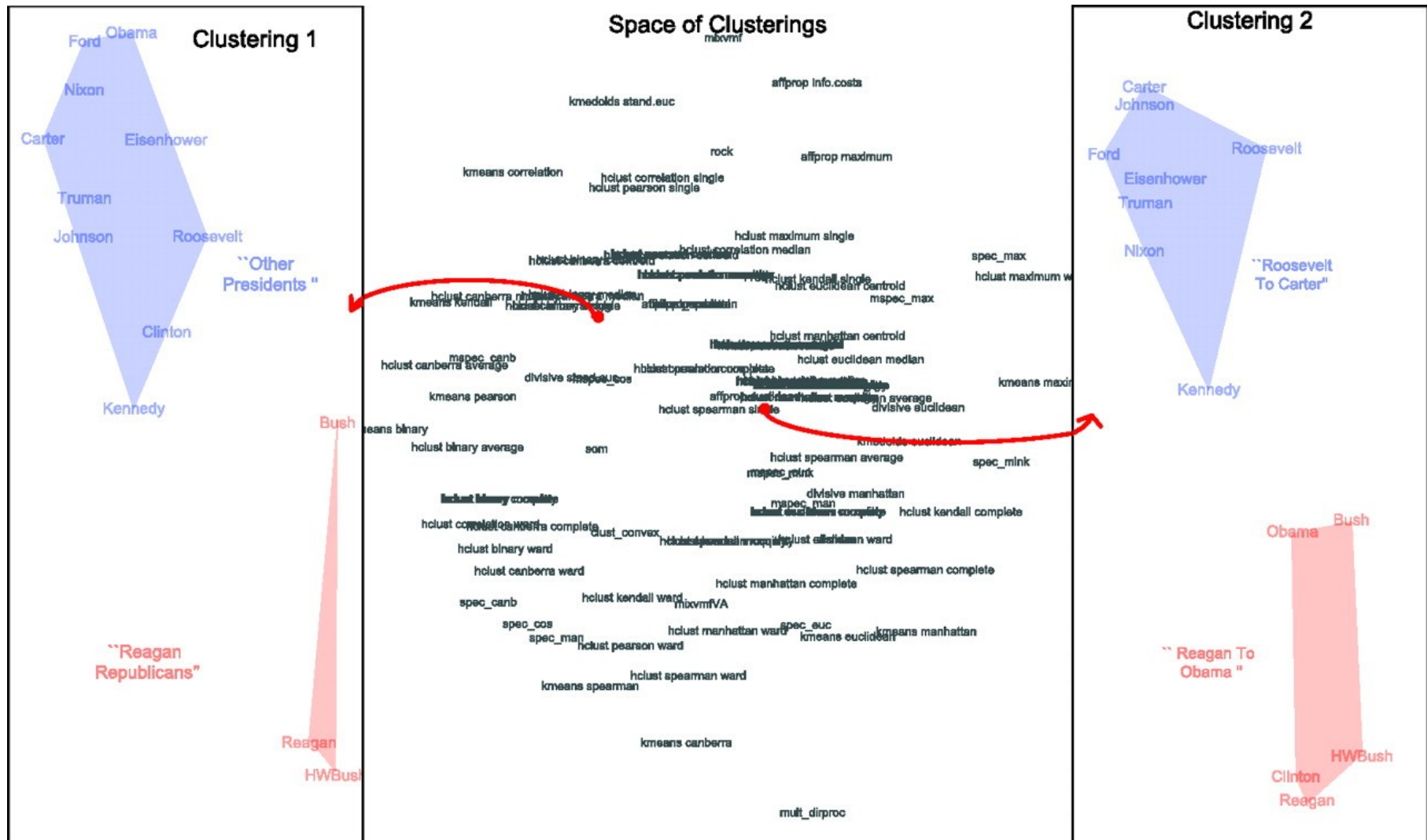
Document1: ambit

Document2: peopl

# Example: R

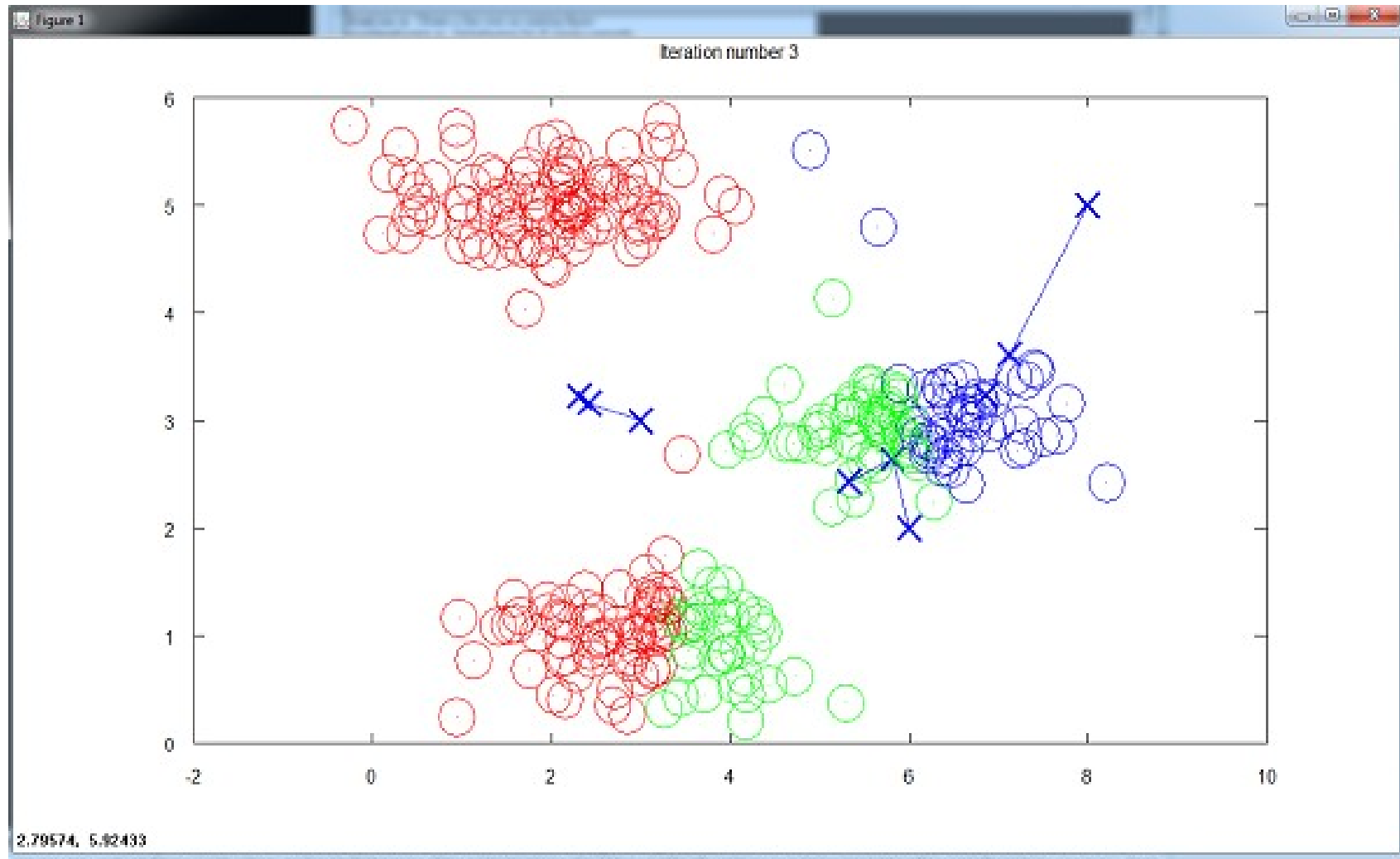
[github.com/lknelson/D-Lab--Text-Analysis-Workshop](https://github.com/lknelson/D-Lab--Text-Analysis-Workshop)

# Latent Categorical Analysis





# Clustering



# Topic Modeling

## It does:

- Allow categories to arise inductively
- Find latent categories
- Find patterns across text
- Handle large and diverse corpora
- Find key differences between categories

## It does not:

- Find the “one” best way to categorize text
- Capture the categories *you* want
- Tell you who does what to whom
- *Magically* reveal meaning

# Latent Dirichlet Allocation

## Topics

gene 0.04  
dna 0.02  
genetic 0.01  
...

life 0.02  
evolve 0.01  
organism 0.01  
...

brain 0.04  
neuron 0.02  
nerve 0.01  
...

data 0.02  
number 0.02  
computer 0.01  
...

## Documents

### Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here,\* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those predictions

"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game, particularly as more and more genomes are completely mapped and sequenced. "It may be a way of organizing any newly sequenced genome," explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an

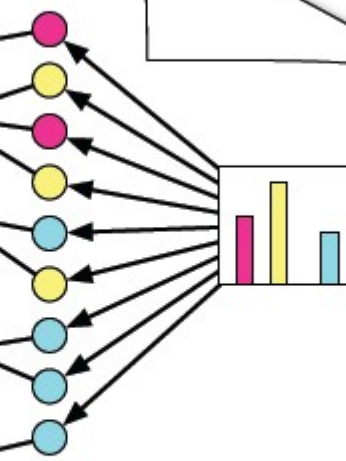


\* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

Stripping down. Computer analysis yields an estimate of the minimum modern and ancient genomes.

SCIENCE • VOL. 272 • 24 MAY 1996

## Topic proportions and assignments



# Latent Dirichlet Allocation

Each *topic* is a distribution over ALL words.

Topic1	topic1_weight	Topic2	topic2_weight	Topic3	topic3_weight
gene	0.5	genetic	0.4	dna	0.7
dna	0.3	dna	0.2	genetic	0.2
genetic	0.2	gene	0.2	gene	0.1
<b>Sum</b>	<b>1.0</b>		<b>1.0</b>		<b>1.0</b>

Each *document* is a distribution over ALL topics.

Doc1	doc1_weight	Doc2	doc2_weight	Doc3	doc3_weight
Topic1	0.6	Topic2	0.8	Topic3	0.5
Topic2	0.3	Topic3	0.1	Topic1	0.3
Topic3	0.1	Topic1	0.1	Topic2	0.2
<b>Sum</b>	<b>1.0</b>		<b>1.0</b>		<b>1.0</b>

# Example: R

[github.com/lknelson/D-Lab--Text-Analysis-Workshop](https://github.com/lknelson/D-Lab--Text-Analysis-Workshop)

# Example

- Four women's organizations in two cities, looking at difference between cities.
- Method: automated text analysis--LDA and difference of proportions
- Interested in differences between organizations, within cities, and over time.
- Combined methods to come to conclusion.



# Results: Difference of Proportions

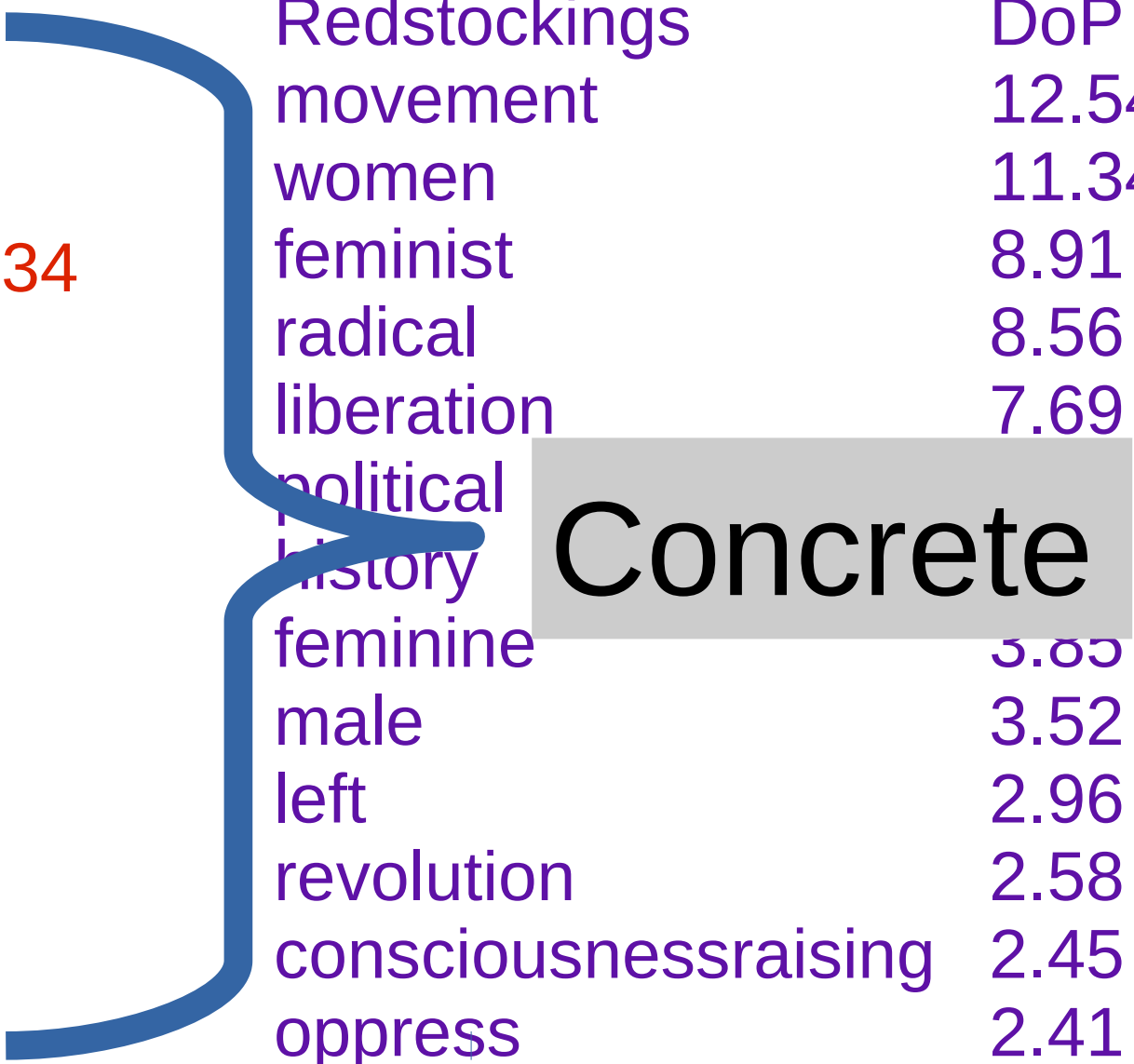
CWLU	DoP	Redstockings	DoP
chicago	5.31	movement	12.54
children	4.59	women	11.34
center	4.34	feminist	8.91
union	3.61	radical	8.56
school	3.48	liberation	7.69
abort	3.19	political	5.81
nixon	2.93	history	5.68
day	2.86	feminine	3.85
vietnam	2.57	male	3.52
people	2.50	left	2.96
city	2.44	revolution	2.58
hospital	2.38	consciousnessraising	2.45
cwlu	2.37	oppress	2.41

# Results: Difference of Proportions

CWLU	DoP	Redstockings	DoP
chicago	5.31	movement	12.54
children	4.59	women	11.34
center	4.34	feminist	8.91
union	3.61	radical	8.56
school	3.48	liberation	7.69
<b>Abstract</b>	3.19	political	5.81
day	2.93	history	5.68
vietnam	2.86	feminine	3.85
people	2.57	male	3.52
city	2.50	left	2.96
hospital	2.44	revolution	2.58
cwlu	2.38	consciousnessraising	2.45
	2.37	oppress	2.41

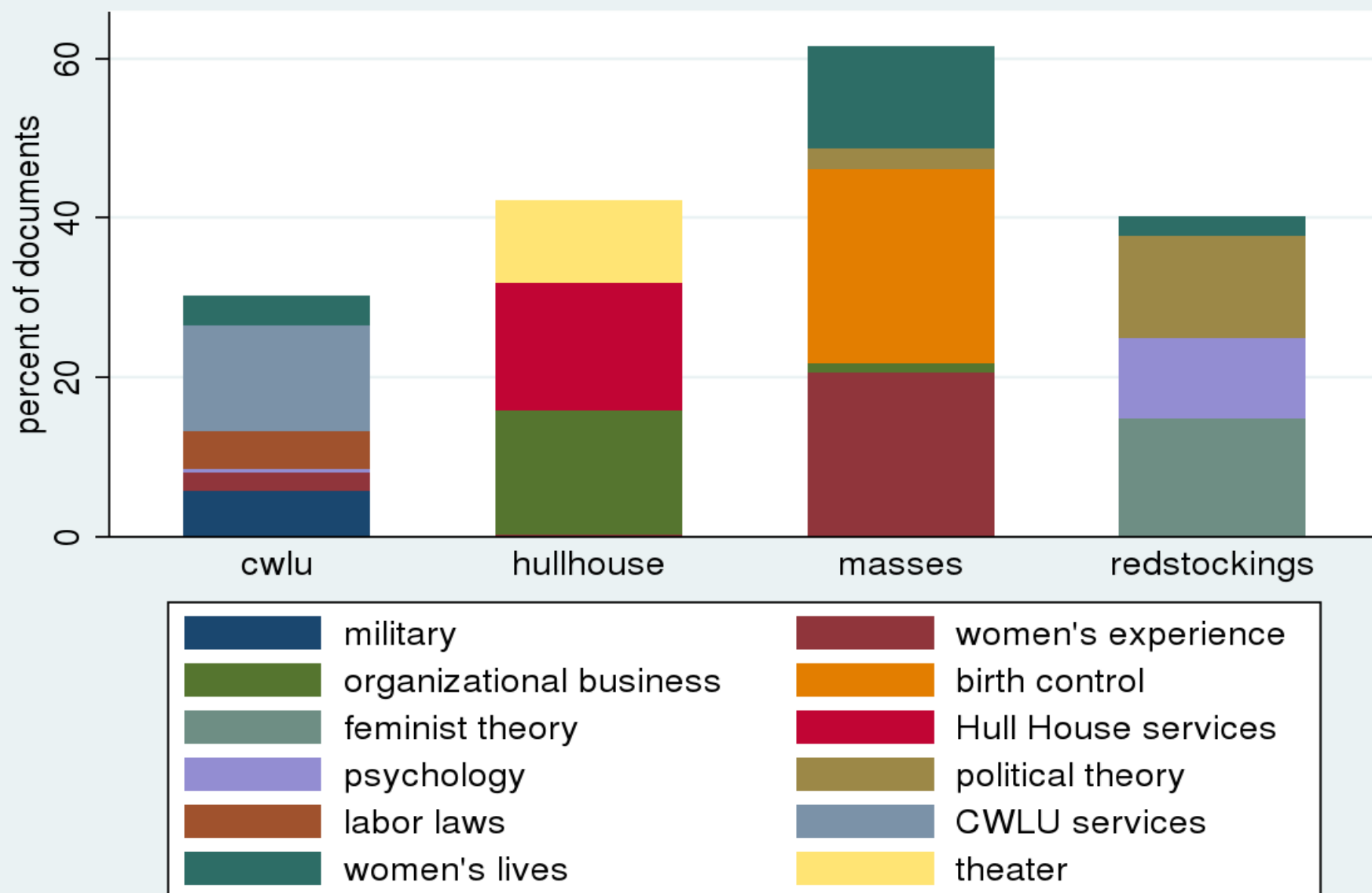
# Results: Difference of Proportions

CWLU	DoP	Redstockings	DoP
chicago	5.31	movement	12.54
children	4.59	women	11.34
center	4.34	feminist	8.91
union	3.61	radical	8.56
school	3.48	liberation	7.69
abort	3.19	political	
nixon	2.93	history	
day	2.86	feminine	5.85
vietnam	2.57	male	3.52
people	2.50	left	2.96
city	2.44	revolution	2.58
hospital	2.38	consciousnessraising	2.45
cwlu	2.37	oppress	2.41



Concrete

# Results: LDA, Top 12 Topics



# Results: LDA, Aggregated Topics



# Qualitative Check

Women attorneys in private practice who are members of the Chicago Women's Liberation Union legal clinic will counsel at the YWCA-Loop Center on everything from domestic relations to criminal law. Among the most frequently asked questions at these legal clinics, said one of the attorneys, are those concerning a woman's rights in marriage, ownership rights, property rights, rights in business, and labor union problems.

---*Womankind*, 1973

It was 1969 when she became pregnant again. This time she wanted the baby. The birth in the middle of the night at a public hospital ward affected her badly. And when the doctor finally arrived at 9:30 a.m., the local anesthesia had worn off so that she was stitched up without it. When she got over the exhaustion following all this, she had to go to work in a mens clothing factory (which today is under investigation for unsanitary working conditions).

---*Feminist Revolution*, 1979

# Summary of Findings

	Quantitative Results	Qualitative Analysis	Political Logic
New York City	Words and categories that are <b>abstract</b> and <b>general</b> dominate.	The authors used <b>stories</b> and <b>narratives</b> to make <b>generalizable claims</b> about the social world that are applicable to all women, <b>abstracting</b> these stories to claims about social structure.	Social change happens through <b>individuals</b>  Strategy: Change individual <b>consciousness</b>
Chicago	Words and categories that are <b>concrete</b> and <b>particular</b> dominate.	The documents outline each organization's attempts to identify <b>concrete needs</b> of the community and their efforts to take <b>practical steps</b> toward meeting those needs.	Social change happens through <b>institutions</b>  Strategy: short-term goals winning <b>concrete changes</b>

# Conclusion

- Computers can be helpful, but...
- Still requires qualitative choices
  - Stem words? Remove stop words? Which stop words?  
Remove frequent/infrequent words? How many clusters?  
How many topics?
- They will not do the interpretive work for you!
- Choose methods based on your question
- Combine different methods
- Validate, validate, validate!