

# Extracting Semantic Themes with Topic Models

Mark Steyvers  
Department of Cognitive Sciences  
University of California, Irvine

Collaborators:

Tom Griffiths, UC Berkeley

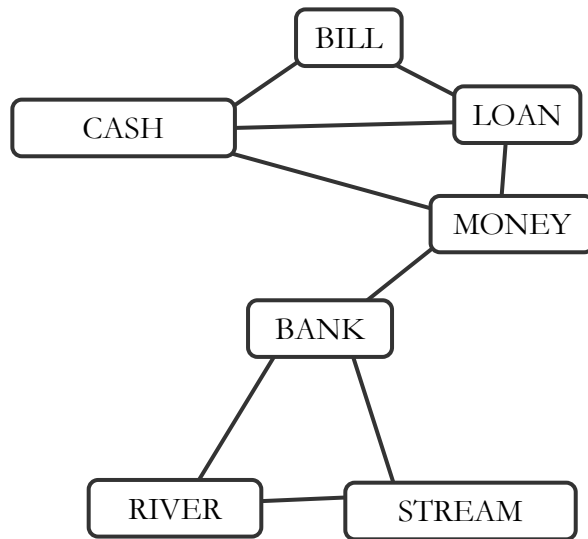
Padhraic Smyth, UC Irvine

# Research on Semantic Representations

- Cognitive Science:
  - How is information represented in semantic memory?
  - How do we retrieve relevant semantic memories?
- Text-mining/ Information Retrieval/ Machine Learning:
  - How can computers represent semantic information?
  - How can computers automatically extract semantic information from text?

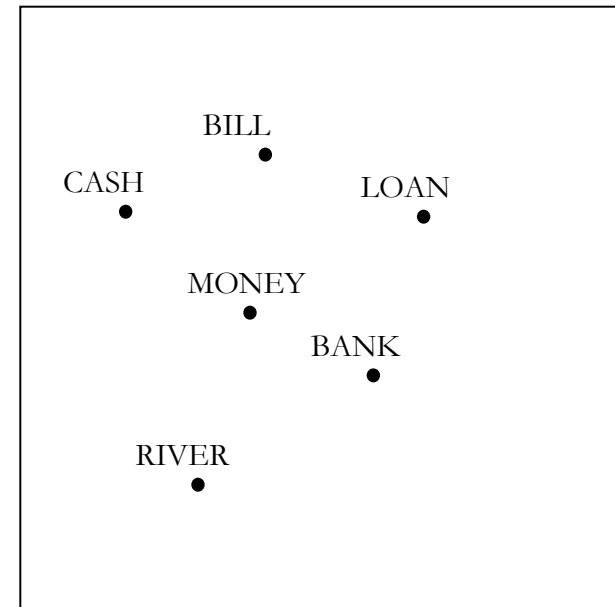
# Two approaches to semantic representation

## Semantic networks



*E.g. Spreading activation*  
*Collins & Quillian (1969)*  
*Collins & Loftus (1975)*

## Semantic Spaces



*Latent Semantic Analysis*  
*(Landauer & Dumais, 1997)*

New approach:  
Probabilistic Topic Models

# Topic Models

- Widely used model in machine learning and text mining
  - aka Latent Dirichlet Allocation (LDA) and pLSI
- *Automatic and unsupervised* extraction of semantic themes from large text collections.
- Essentially a Bayesian analysis of co-occurrence statistics.  
How are words correlated with other words across contexts?

# Topics provide quick summary of content

- What is in this corpus?
- What is in this document?
- What does this person/group of people write about?
- What are the topical trends over time?

Example:

# Topic Analysis of Search Queries

# AOL dataset

- Dataset:
  - 20,000,000+ web queries
  - 650,000+ users
- Users were given “anonymous” user-id
- Publicly available (...but controversial)



# Example query log from user #2178

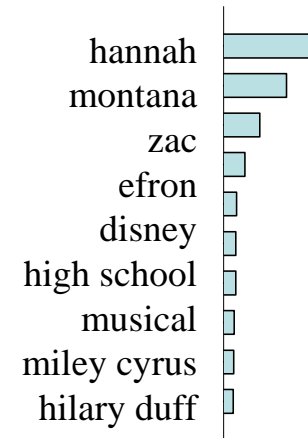
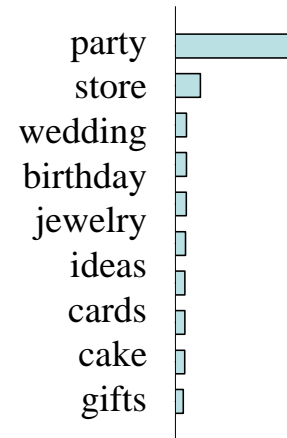
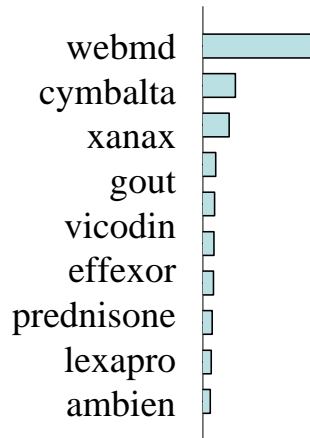
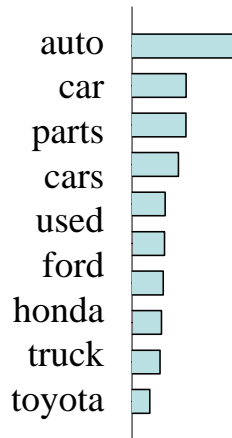
ID	Query	Date/Time	URL clicked
2178	dog eats uncooked pasta	2006-05-26 15:31:56	
2178	inducing dog vomiting	2006-05-26 15:32:46	<a href="http://www.twodogpress.com">http://www.twodogpress.com</a>
2178	inducing dog vomiting	2006-05-26 15:32:46	<a href="http://www.canismajor.com">http://www.canismajor.com</a>
2178	inducing dog vomiting	2006-05-26 15:32:46	<a href="http://kitchen.robبيهaf.com">http://kitchen.robبيهaf.com</a>
2178	inducing dog vomiting	2006-05-26 15:32:46	<a href="http://www.dog-first-aid-101.com">http://www.dog-first-aid-101.com</a>
2178	inducing dog vomiting	2006-05-26 15:38:36	
2178	walmart	2006-05-12 12:39:52	<a href="http://www.walmart.com">http://www.walmart.com</a>
2178	sears	2006-05-12 12:44:22	<a href="http://www.sears.com">http://www.sears.com</a>
2178	target	2006-05-12 17:05:36	<a href="http://www.target.com">http://www.target.com</a>
2178	babycenter.com	2006-05-12 17:43:59	<a href="http://www.babycenter.com">http://www.babycenter.com</a>
2178	google	2006-05-16 10:54:39	<a href="http://www.google.com">http://www.google.com</a>
2178	fit pregnancy	2006-05-16 15:34:23	
2178	baby center	2006-05-16 15:37:22	
2178	yahoo.com	2006-05-18 17:11:05	<a href="http://www.yahoo.com">http://www.yahoo.com</a>
2178	applebee's carside	2006-05-19 19:21:08	<a href="http://www.applebees.com">http://www.applebees.com</a>
2178	baby names	2006-05-20 15:02:38	<a href="http://www.babynames.com">http://www.babynames.com</a>
2178	baby names	2006-05-20 15:02:38	<a href="http://www.babynamesworld.com">http://www.babynamesworld.com</a>
2178	baby names	2006-05-20 15:02:38	<a href="http://www.thinkbabynames.com">http://www.thinkbabynames.com</a>
2178	mortgage calculator	2006-05-24 14:39:05	<a href="http://www.bankrate.com">http://www.bankrate.com</a>
2178	us zip codes	2006-05-25 21:26:47	<a href="http://www.usps.com">http://www.usps.com</a>
2178	us zip codes	2006-05-25 21:26:47	<a href="http://www.usps.com">http://www.usps.com</a>

# Query Topic Model

- Each *user* searches for a *mixture of topics*
- Each *topic* is a probability distribution over query words
- Simultaneously solve for all unknowns using efficient stochastic search (MCMC and Gibbs sampling)

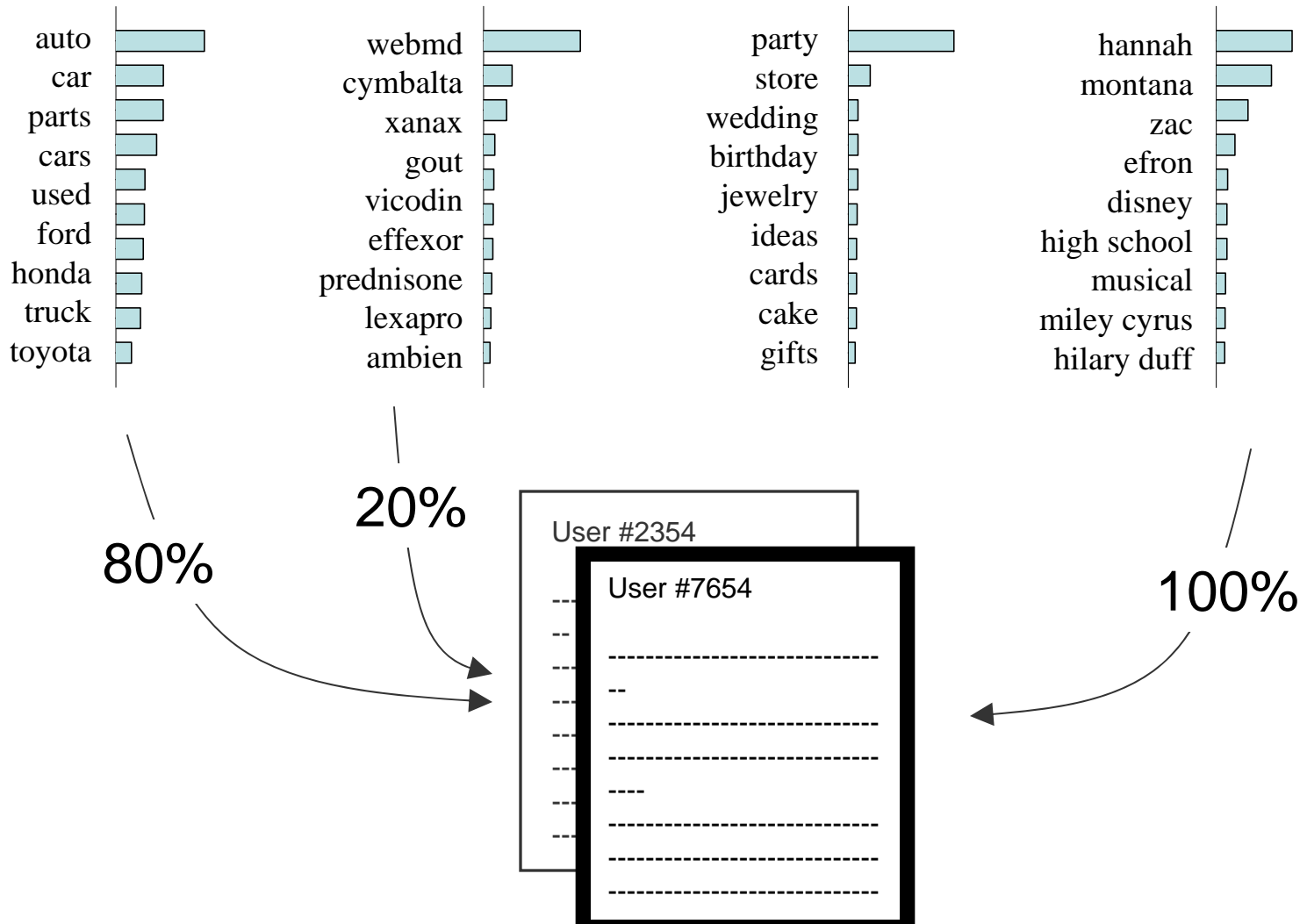
# Topic = distribution over words

Four example topics (out of 200)



Probability distribution  
over words. Most likely  
words listed at the top

# User = mixture of topics

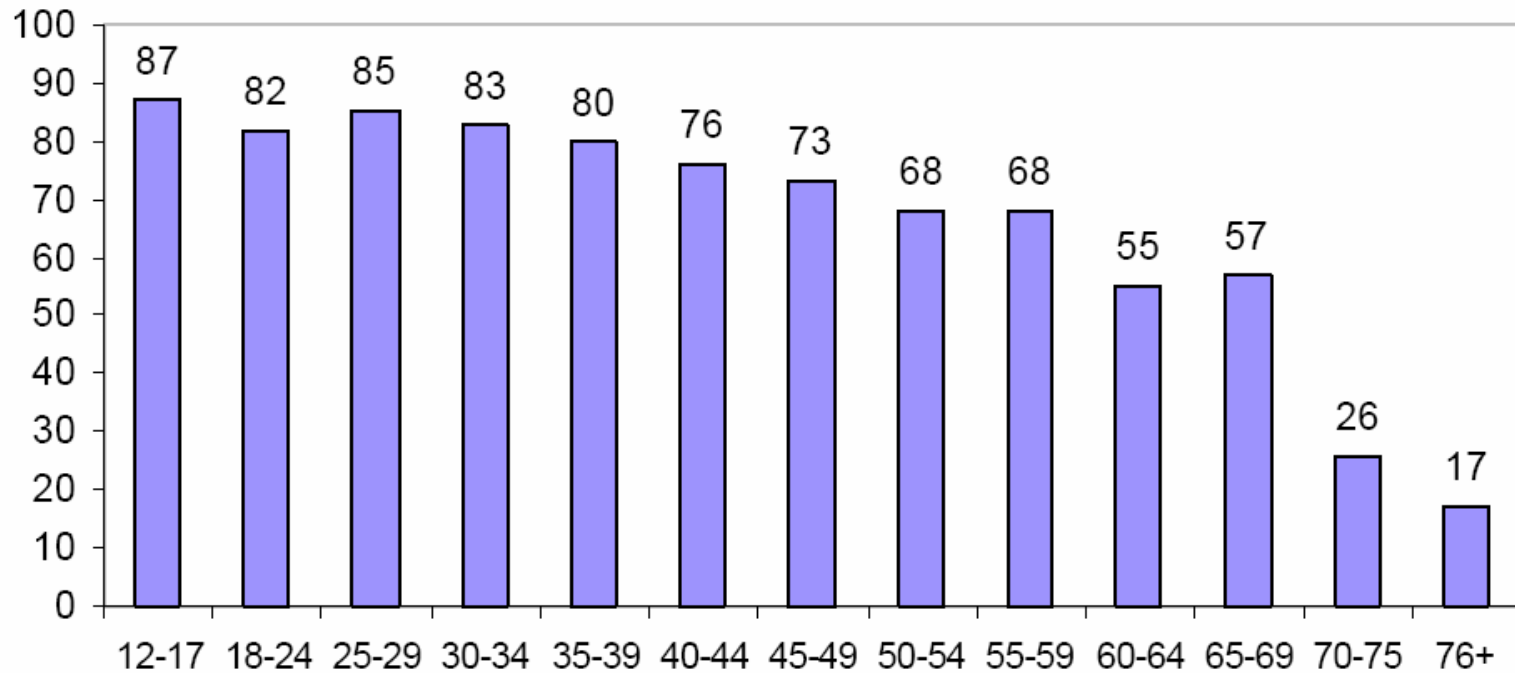


# Topic Model

- Model requires no background knowledge:
  - no dictionaries
  - no thesauri
  - no ontologies, etc

# Internet search behavior for different age groups

# Share of Americans online by age



Source: *PEW Internet & American Life Project* (2005)

# Survey of Online Activities by Age Group

Generational Differences in Online Activities								
	Online Teens <sup>a</sup> (12-17)	Gen Y (18-28)	Gen X (29-40)	Trailing Boomers (41-50)	Leading Boomers (51-59)	Matures (60-69)	After Work (70+)	All Online Adults <sup>b</sup>
Go online	87%	84%	87%	79%	75%	54%	21%	72%
<i>Teens and Gen Y are more likely to engage in the following activities compared with older users:</i>								
Online games	81	54	37	29	25	25	32	36
School research	*	73	60	61	48	33	14	57
Instant message	75	66	52	38	42	33	25	47
Text message	*	60	44	29	15	11	8	35
Get info about a school	57	59	42	50	40	30	14	45
Download music	51	45	28	16	14	8	5	25
Read blogs	38	41	30	20	21	19	16	27
Download video	31	27	22	14	8	8	1	18
Create a blog	19	20	9	3	9	3	4	9

Source: *PEW Internet & American Life Project* (2005)



# Alternative to Survey Analysis

- PEW surveys presuppose you already know the questions to ask for
- Alternative:
  - automatically topic analyze the search terms people enter into a search engine
  - Gives a more direct picture of the user goals and needs

(will incorporate result slides during presentation)

# Conclusions

- Topic modeling coupled with demographic analysis yields “windows” into the minds of different age groups
- Other potential applications:
  - clinical data, e.g. therapy discussions
  - open question surveys
  - internet behavior, e.g. chatting, SMS, email

# Open Questions

- How can research community get access to data from internet providers/ search engines?
- Privacy issues

# Additional Slides

# Example Topics from New York Times

## Terrorism

SEPT\_11  
WAR  
SECURITY  
IRAQ  
TERRORISM  
NATION  
KILLED  
AFGHANISTAN  
ATTACKS  
OSAMA\_BIN\_LADEN  
AMERICAN  
ATTACK  
NEW\_YORK\_REGION  
NEW  
MILITARY  
NEW\_YORK  
WORLD  
NATIONAL  
QAEDA  
TERRORIST\_ATTACKS

## Wall Street Firms

WALL\_STREET  
ANALYSTS  
INVESTORS  
FIRM  
GOLDMAN\_SACHS  
FIRMS  
INVESTMENT  
MERRILL\_LYNCH  
COMPANIES  
SECURITIES  
RESEARCH  
STOCK  
BUSINESS  
ANALYST  
WALL\_STREET\_FIRMS  
SALOMON\_SMITH\_BARNEY  
CLIENTS  
INVESTMENT\_BANKING  
INVESTMENT\_BANKERS  
INVESTMENT\_BANKS

## Stock Market

WEEK  
DOW\_JONES  
POINTS  
10\_YR\_TREASURY\_YIELD  
PERCENT  
CLOSE  
NASDAQ\_COMPOSITE  
STANDARD\_POOR  
CHANGE  
FRIDAY  
DOW\_INDUSTRIALS  
GRAPH\_TRACKS  
EXPECTED  
BILLION  
NASDAQ\_COMPOSITE\_INDEX  
EST\_02  
PHOTO\_YESTERDAY  
YEN  
10  
500\_STOCK\_INDEX

## Bankruptcy

BANKRUPTCY  
CREDITORS  
BANKRUPTCY\_PROTECTION  
ASSETS  
COMPANY  
FILED  
BANKRUPTCY\_FILING  
ENRON  
BANKRUPTCY\_COURT  
KMART  
CHAPTER\_11  
FILING  
COOPER  
BILLIONS  
COMPANIES  
BANKRUPTCY\_PROCEEDINGS  
DEBTS  
RESTRUCTURING  
CASE  
GROUP

# Algorithm input/output

**INPUT:** word-document counts (word order is irrelevant)

**OUTPUT:**

topic assignments to each word	$P( z_i )$
likely words in each topic	$P( w \mid z )$
likely topics in each document (“gist”)	$P( \theta \mid d )$

# Generative Process

- For each document, choose a mixture of topics

$$\theta \sim \text{Dirichlet}(\alpha)$$

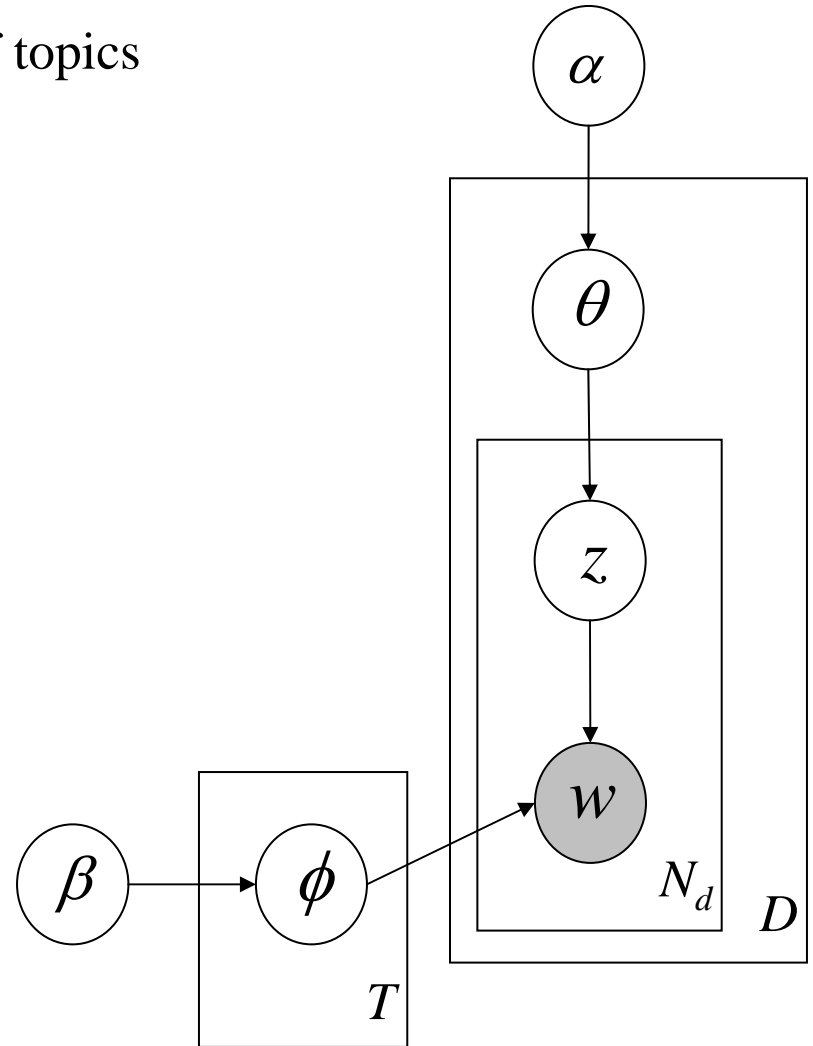
- Sample a topic [1..T] from the mixture

$$z \sim \text{Multinomial}(\theta)$$

- Sample a word from the topic

$$w \sim \text{Multinomial}(\phi^{(z)})$$

$$\phi \sim \text{Dirichlet}(\beta)$$

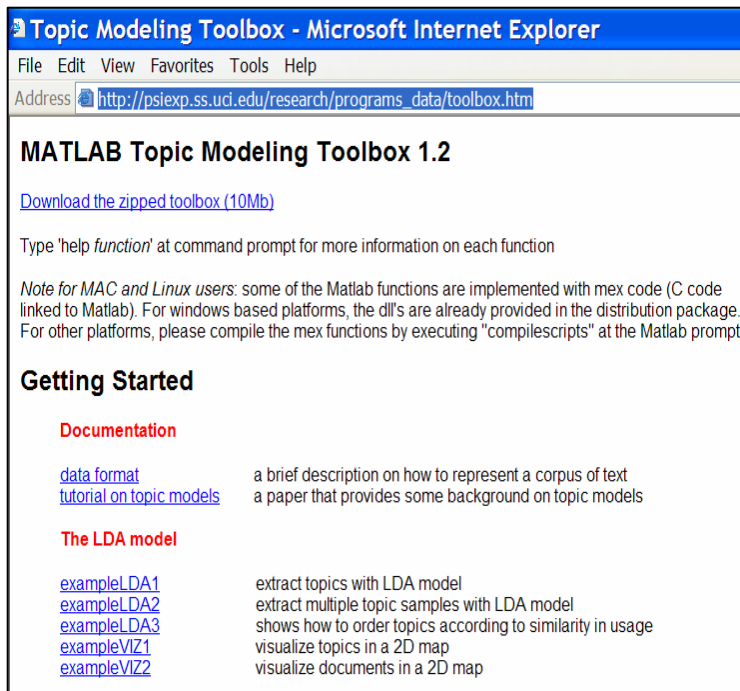




# Software

Public-domain MATLAB toolbox for topic modeling on the Web:

[http://psiexp.ss.uci.edu/research/programs\\_data/toolbox.htm](http://psiexp.ss.uci.edu/research/programs_data/toolbox.htm)



The screenshot shows a web browser window titled "Topic Modeling Toolbox - Microsoft Internet Explorer". The address bar displays the URL [http://psiexp.ss.uci.edu/research/programs\\_data/toolbox.htm](http://psiexp.ss.uci.edu/research/programs_data/toolbox.htm). The main content area features the heading "MATLAB Topic Modeling Toolbox 1.2" and a link to "Download the zipped toolbox (10Mb)". Below this, it instructs users to type 'help function' at the command prompt. A note for MAC and Linux users explains that some functions are implemented with mex code. The "Getting Started" section includes links to documentation and the LDA model, with brief descriptions of each.

**Topic Modeling Toolbox - Microsoft Internet Explorer**

File Edit View Favorites Tools Help

Address [http://psiexp.ss.uci.edu/research/programs\\_data/toolbox.htm](http://psiexp.ss.uci.edu/research/programs_data/toolbox.htm)

## MATLAB Topic Modeling Toolbox 1.2

[Download the zipped toolbox \(10Mb\)](#)

Type 'help *function*' at command prompt for more information on each function

*Note for MAC and Linux users:* some of the Matlab functions are implemented with mex code (C code linked to Matlab). For windows based platforms, the dll's are already provided in the distribution package. For other platforms, please compile the mex functions by executing "compilescripts" at the Matlab prompt

### Getting Started

**Documentation**

<a href="#">data format</a>	a brief description on how to represent a corpus of text
<a href="#">tutorial on topic models</a>	a paper that provides some background on topic models

**The LDA model**

<a href="#">exampleLDA1</a>	extract topics with LDA model
<a href="#">exampleLDA2</a>	extract multiple topic samples with LDA model
<a href="#">exampleLDA3</a>	shows how to order topics according to similarity in usage
<a href="#">exampleVIZ1</a>	visualize topics in a 2D map
<a href="#">exampleVIZ2</a>	visualize documents in a 2D map

# Choosing number of topics

- Bayesian model selection
- Generalization test
  - e.g., perplexity on out-of-sample data
- Non-parametric Bayesian approach
  - Number of topics grows with size of data
  - E.g. Hierarchical Dirichlet Processes (HDP)

# Polysemy

PRINTING  
PAPER  
PRINT  
PRINTED  
TYPE  
PROCESS  
INK  
PRESS  
IMAGE  
PRINTER  
PRINTS  
PRINTERS  
COPY  
COPIES  
FORM  
OFFSET  
GRAPHIC  
SURFACE  
PRODUCED  
**CHARACTERS**

**PLAY**  
PLAYS  
STAGE  
AUDIENCE  
THEATER  
ACTORS  
DRAMA  
SHAKESPEARE  
ACTOR  
THEATRE  
PLAYWRIGHT  
PERFORMANCE  
DRAMATIC  
COSTUMES  
COMEDY  
TRAGEDY  
**CHARACTERS**  
SCENES  
OPERA  
PERFORMED

TEAM  
GAME  
BASKETBALL  
PLAYERS  
PLAYER  
**PLAY**  
PLAYING  
SOCCER  
PLAYED  
BALL  
TEAMS  
BASKET  
FOOTBALL  
SCORE  
**COURT**  
GAMES  
TRY  
COACH  
GYM  
SHOT

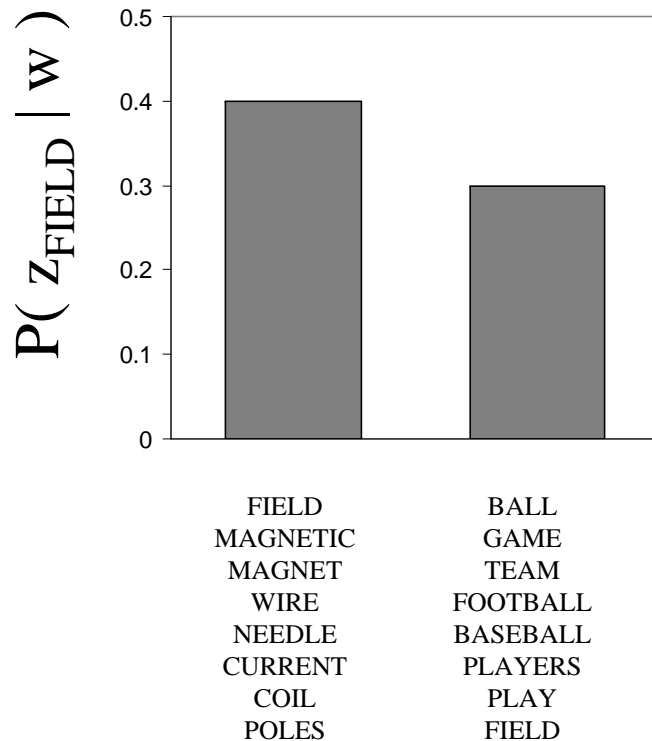
JUDGE  
TRIAL  
**COURT**  
CASE  
JURY  
ACCUSED  
GUILTY  
DEFENDANT  
JUSTICE  
**EVIDENCE**  
WITNESSES  
CRIME  
LAWYER  
WITNESS  
ATTORNEY  
HEARING  
INNOCENT  
DEFENSE  
CHARGE  
CRIMINAL

HYPOTHESIS  
EXPERIMENT  
SCIENTIFIC  
OBSERVATIONS  
SCIENTISTS  
EXPERIMENTS  
SCIENTIST  
EXPERIMENTAL  
**TEST**  
METHOD  
HYPOTHESES  
TESTED  
**EVIDENCE**  
BASED  
OBSERVATION  
SCIENCE  
FACTS  
DATA  
RESULTS  
EXPLANATION

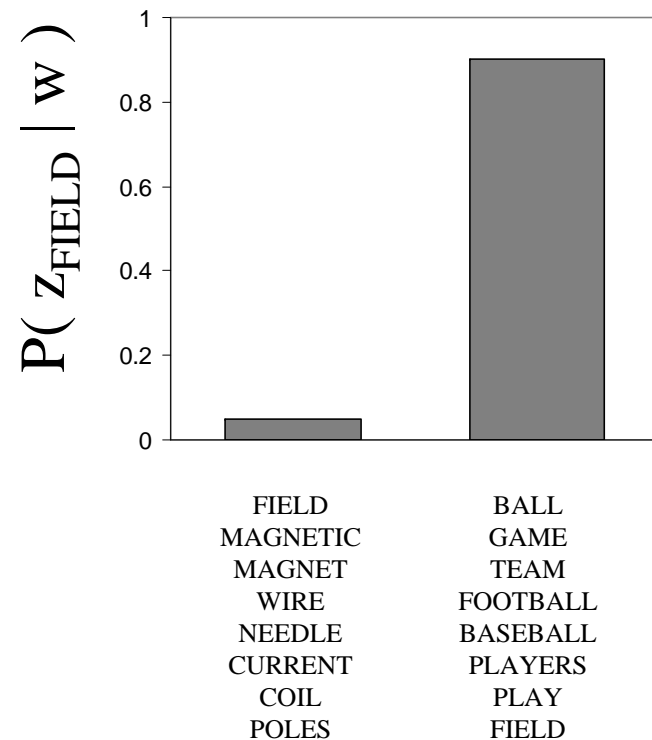
STUDY  
**TEST**  
STUDYING  
HOMEWORK  
NEED  
CLASS  
MATH  
TRY  
TEACHER  
WRITE  
PLAN  
ARITHMETIC  
ASSIGNMENT  
PLACE  
STUDIED  
CAREFULLY  
DECIDE  
IMPORTANT  
NOTEBOOK  
REVIEW

# Disambiguation

“FIELD”



“FOOTBALL FIELD”



# Recent Papers

- Steyvers, M., Griffiths, T.L., & Dennis, S. (2006). Probabilistic inference in human semantic memory. *Trends in Cognitive Sciences*, 10(7), 327-334.
- Griffiths, T.L., Steyvers, M., & Tenenbaum, J.B.T. (2007). Topics in Semantic Representation. *Psychological Review*, 114(2), 211-244.
- Griffiths, T.L., Steyvers, M., & Firl, A. (in press). Google and the mind: Predicting fluency with PageRank. *Psychological Science*.
- Steyvers, M. & Griffiths, T.L. (2008). Rational Analysis as a Link between Human Memory and Information Retrieval. In N. Chater and M Oaksford (Eds.) *The Probabilistic Mind: Prospects from Rational Models of Cognition*. Oxford University Press.
- Chemudugunta, C., Smyth, P., & Steyvers, M. (2007). Modeling General and Specific Aspects of Documents with a Probabilistic Topic Model. In: *Advances in Neural Information Processing Systems*, 19.