

CS 277, Data Mining

Topic Modeling for Text Documents

Padhraic Smyth

Department of Computer Science

Bren School of Information and Computer Sciences

University of California, Irvine

Project Progress Report 1

- Grade distribution
 - A few 9's: excellent reports
 - Many 8's: good reports, nothing outstanding
 - Many 7's: good, but could be improved, e.g., missing key information
 - Below 7: needs improvement

Project Progress Report 1

- Grade distribution
 - A few 9's: excellent reports
 - Many 8's: good reports, nothing outstanding
 - Many 7's: good, but could be improved, e.g., missing key information
 - Below 7: needs improvement
- What is required to get a top score?
 1. Technical progress on your project
 - Not just summarizing well-known knowledge (e.g., known algorithms)
 - Insights, figures, details: show me what you have learned so far
 - Clear that you have spent time on this project
(note that spending time alone is not sufficient to get a high score)
 2. Clearly written report
 - Well organized, each section builds on the preceding section
 - Good use of figures and graphs
 - Goals, methods, results are clearly explained: reader is not guessing what you did
 - Appropriate level of detail (not too much...but some details are good)

TimeTable (Updated)

- Progress Report 2: now due Wednesday Feb 26th (1 week away)
 - Hand in hardcopy in class
 - AND submit electronic copy online via EEE
 - You can re-use a small amount of figures/text that you think is necessary/helpful from your earlier report, but clearly point this out. At least 80% of what you report should be new.
- Final Report
 - Due noon Friday March 14th (electronically to EEE)
 - Will discuss the format and expectations in more detail later in the quarter

Basis of Final Grade

- Weighted combination of Assignments and Reports
 - 10% for Assignment 1
 - 20% for each of the Progress Reports
 - 30% for the Final Report

Topic Models for Text

Examples of Large Text Corpora



NYT: 1.5 million news articles



NIH: 1 million grant abstracts

From: PGE News
To: ALL PGE EMPLOYEES
Date: 8/14/01 2:54PM
Subject: Jeff Dilling resigns as CEO of Enron

PGE News August 14, 2001

Jeff Dilling resigns as CEO of Enron

Enron today announced that President and CEO Jeff Dilling has resigned, effective immediately, and that the Enron Board of Directors has elected Ken Lay to resume his role as Chairman and CEO.

"Shir Horner called this afternoon to inform me of Jeff's decision to step down for personal reasons," says PGE CEO and President Peggy Fowler Horner, CEO of Enron Transportation, a Fowler's executive connection to the Enron team. "He wanted to let me know that Mr. Dilling's departure will not in any way impact Enron's ongoing strategy for success and we should expect no near-term dramatic organizational change."

"Clearly, Enron will continue to focus on increasing the company's stock value," Fowler added. "PGE can help in this effort by continuing committed to our investment goals and operational success."

Below is the letter Ken Lay is sending to Enron employees this afternoon announcing the decision:

To: Enron Employees Worldwide
From: Ken Lay

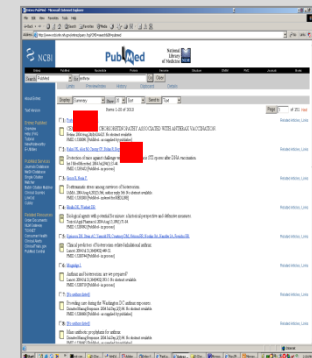
It is with regret that I have to announce that Jeff Dilling is leaving Enron. Today, the Board of Directors accepted his resignation as President and CEO of Enron. Jeff is resigning for personal reasons and his departure is voluntary. I regret the decision, but I understand it. I have worked closely with Jeff for more than 15 years, including 11 years at Enron, and have had the privilege of professional relationships that value them. I am pleased to say that he has agreed to enter into a consulting arrangement with the company to advise me and the Board of Directors.

Now it's time to look forward.

With Jeff leaving, the Board has asked me to resume the responsibilities of President and CEO in addition to my role as Chairman of the Board. I have agreed. I want to assure you that I have never felt better about the prospects for the company. All of you know that our stock price has suffered substantially over the last few months. One of my top priorities will be to restore a significant amount of the stock value we have lost as far as possible. Our performance has never been stronger; our business model has never been more robust; our growth has never been more certain; and, most importantly, we have never had a better and more professional management team. We have the finest organization in American business today. Together, we will make Enron the world's leading company.

CC: Kathy & George Wyatt, Kathy Wyatt

Enron:
250,000 emails



PubMed: 20 million biomedical articles

Unsupervised Learning from Text

- Large collections of unlabeled documents..
 - Web
 - Digital libraries
 - Email archives, etc
- Often wish to organize/summarize/index/tag these documents automatically
- We will look at probabilistic techniques for clustering and topic extraction from sets of documents

Problems of Interest

- What topics do these documents “span”?
- Which documents are about a particular topic?
- How have topics changed over time?
- What does author X write about?
- Who is likely to write about topic Y?
- Who wrote this specific document?
- and so on.....

“Bag of Words” Matrix for Documents

	water	rights	cattle	hunting	land	use	dust	maize	chumash	navajo	reservation	war	disease
doc1	1	1			2								
doc2			1	4	1								1
doc3	2	1	1		1	1	1						
doc4			2		1						1		
doc5	1	1	1	1	3	1	1						
doc6		1		2									
doc7									1	1	1	3	
doc8					1			1	1				
doc9										1	1		
doc10									1		1		1

Statistical Topic Models for Count Data

- Simple hypothetical “generative” models for sparse counts
 - A description of how the data might have been generated
 - Simple in nature (“all models are wrong but some are useful”)
 - Can handle counts, metadata, etc
- Learning the parameters given the data
 - Generative model = $P(D | \theta)$: how likely data D are given the parameters θ
 - Use Bayes rule to get $P(\theta | D)$: how likely parameters θ are given data D

Statistical Topic Models for Count Data

- Simple hypothetical “generative” models for sparse counts
 - A description of how the data might have been generated
 - Simple in nature (“all models are wrong but some are useful”)
 - Can handle counts, metadata, etc
- Learning the parameters given the data
 - Generative model = $P(D | \theta)$: how likely data D are given the parameters θ
 - Use Bayes rule to get $P(\theta | D)$: how likely parameters θ are given data D
- Key Features
 - Multimembership: rows can “belong” to multiple factors
 - Leverage sparsity: computational advantages
 - Can build in dependence on metadata (e.g., document authors)

Modeling Word Frequencies given Count Data

Tossing a die: 6 sides, equally likely, memoryless

Parameters of a “model” for a die:

A vector of 6 probabilities $\theta_1, \dots, \theta_6$ sum to 1, $\sum \theta = 1$

Modeling Word Frequencies given Count Data

Tossing a die: 6 sides, equally likely, memoryless

Parameters of a “model” for a die:

A vector of 6 probabilities $\theta_1, \dots, \theta_6$ sum to 1, $\sum \theta = 1$

Same model for text?

Now we have a K-sided die, where K could be 100,000

A vector of K probabilities $\underline{\theta}$, sum to 1, $\sum \theta = 1$

Modeling Word Frequencies given Count Data

Tossing a die: 6 sides, equally likely, memoryless

Parameters of a “model” for a die:

A vector of 6 probabilities $\theta_1, \dots, \theta_6$ sum to 1, $\sum \theta = 1$

Same model for text?

Now we have a K-sided die, where K could be 100,000

A vector of K probabilities $\underline{\theta}$, sum to 1, $\sum \theta = 1$

Can learn these probabilities from a corpus, via smoothed frequency counts

Applications? detecting shifts in language usage in scientific literature, in customer complaints, in social media, etc

Topic = “Focused” Probability Distribution over Words

Word	Probability
president	0.129
roosevelt	0.032
congress	0.030
johnson	0.026
office	0.021
wilson	0.021
nixon	0.020
reagan	0.018
kennedy	0.018
...	...

Different Topics for Different Semantic Concepts

Word	Probability
red	0.202
blue	0.099
green	0.096
yellow	0.073
white	0.048
color	0.030
bright	0.029
colors	0.027
brown	0.027
....

Word	Probability
president	0.129
roosevelt	0.032
congress	0.030
johnson	0.026
office	0.021
wilson	0.021
nixon	0.020
reagan	0.018
kennedy	0.018
....

Key Idea: Documents as Mixtures of Topics

Topic 1: search_query (0.4), precision (0.3), retrieval (0.3)

Topic 2: classification (0.5), neural_network (0.3), labels (0.2)

Topic 3: experiment (0.6), result (0.2), significance (0.2)

Key Idea: Documents as Mixtures of Topics

Topic 1: search_query (0.4), precision (0.3), retrieval (0.3)

Topic 2: classification (0.5), neural_network (0.3), labels (0.2)

Topic 3: experiment (0.6), result (0.2), significance (0.2)

Topic model: documents = convex combinations of Topics 1, 2, 3: e.g.,

$P(\text{Words}) \text{ for Doc 1} = 0.4 * \text{Topic 1} + 0.4 * \text{Topic 2} + 0.2 * \text{Topic 3}$

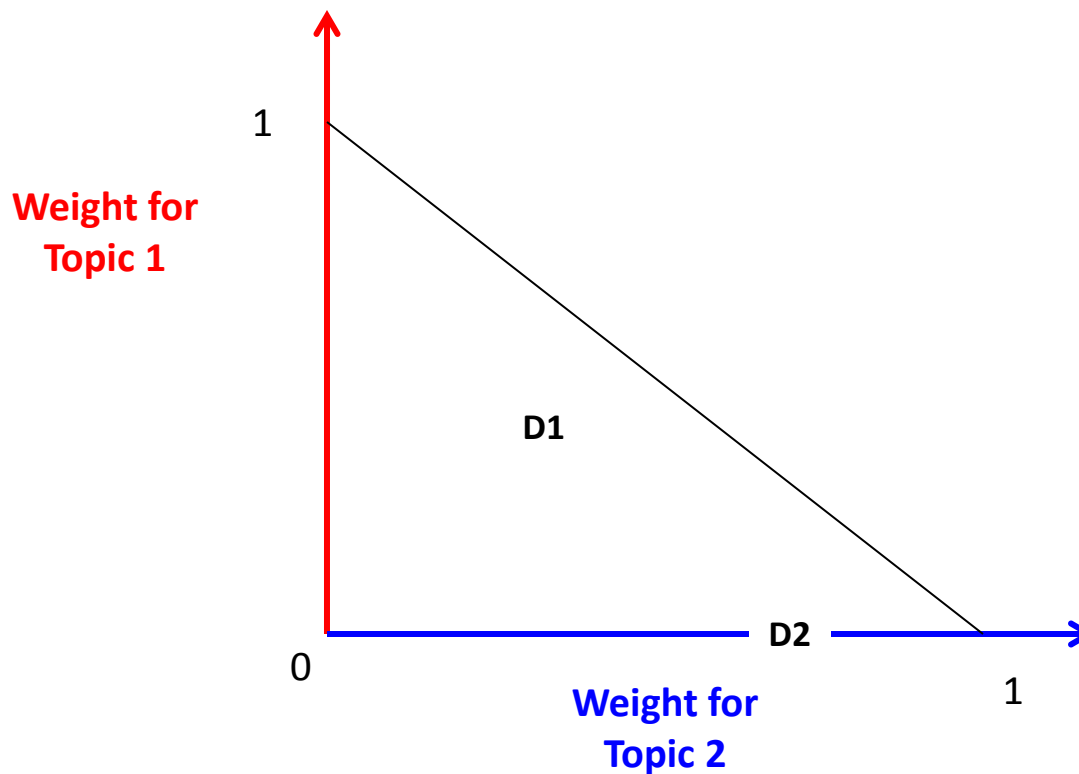
$P(\text{Words}) \text{ for Doc 2} = 0.0 * \text{Topic 1} + 0.8 * \text{Topic 2} + 0.2 * \text{Topic 3}$

Key Idea: Documents as Mixtures of Topics

Topic 1: search_query (0.4), precision (0.3), retrieval (0.3)

Topic 2: classification (0.5), neural_network (0.3), labels (0.2)

Topic 3: experiment (0.6), result (0.2), significance (0.2)



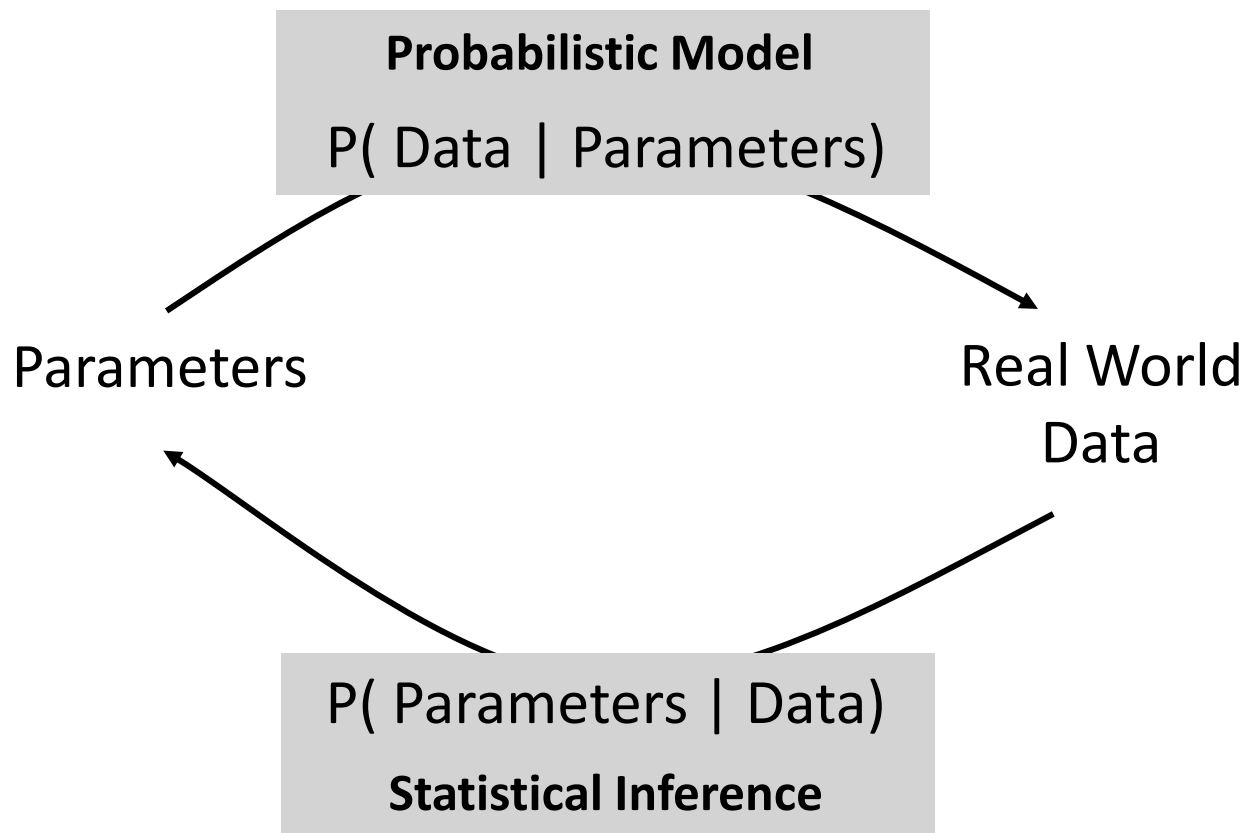
A Generative Model for Documents

- Our topic model is a simple “forward generative” model for observed data (i.e., counts of words in documents)
- We can think of it as a simulator (with pseudocode)
 - For each document in our corpus
 - For each word in our document
 - Sample a topic from $P(\text{topics} \mid \text{document})$
 - Given the topic, sample a word from $P(\text{words} \mid \text{topic})$
 - End
 - End

A Generative Model for Documents

- Our topic model is a simple “forward generative” model for observed data (i.e., counts of words in documents)
- We can think of it as a simulator (with pseudocode)
 - For each document in our corpus
 - For each word in our document
 - Sample a topic from $P(\text{topics} \mid \text{document})$
 - Given the topic, sample a word from $P(\text{words} \mid \text{topic})$
 - End
 - End
- Key points:
 - Many useful statistical models in machine learning have a very simple “forward mechanism” (a few lines of pseudocode)
 - Learning this model is essentially “inverting” this code via Bayes rule
 - The reverse “inference” step is usually much harder than the forward step!

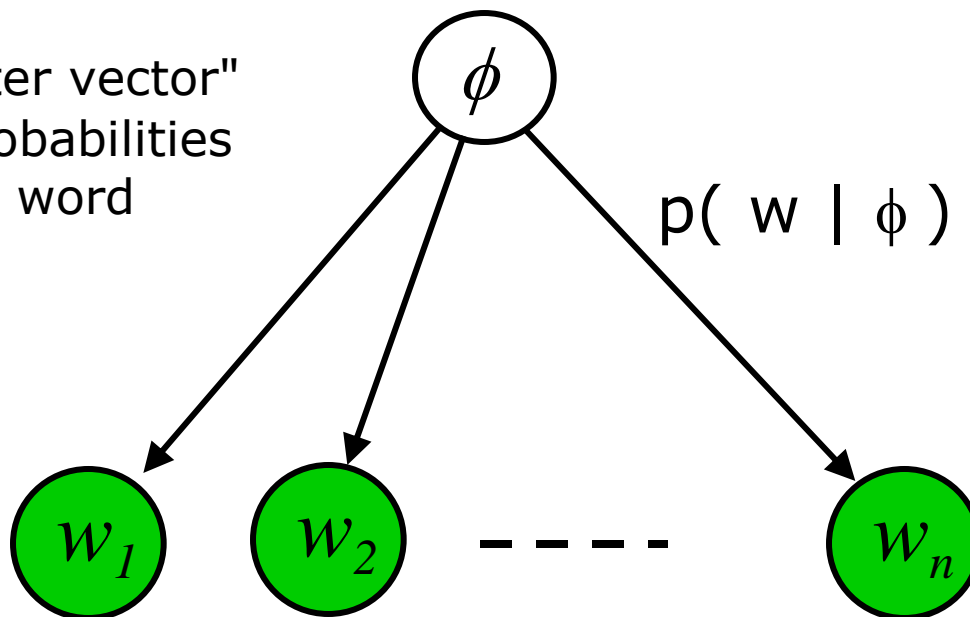
Generative Statistical Models for Data



A Graphical Model

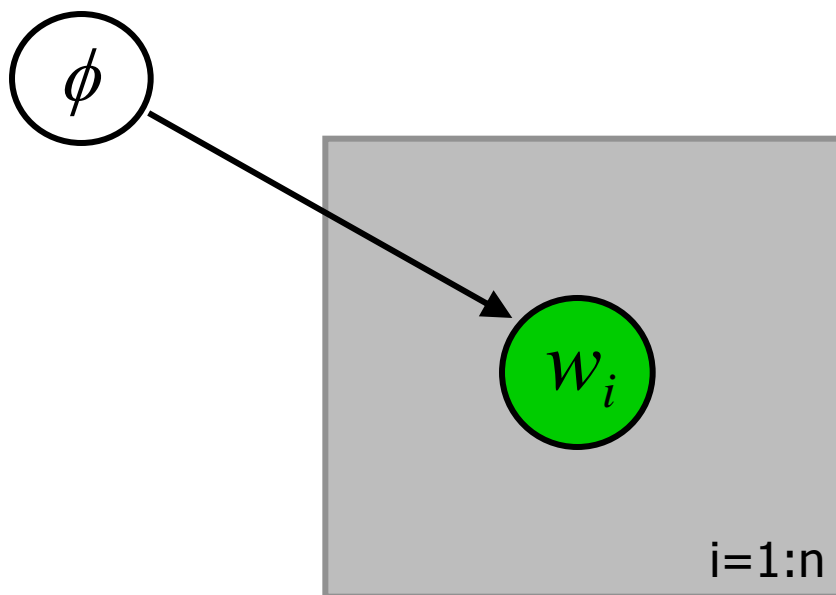
$$p(\text{doc} \mid \phi) = \prod p(w_i \mid \phi)$$

ϕ = "parameter vector"
= set of probabilities
one per word



Another view....

$$p(\text{doc} \mid \phi) = \prod p(w_i \mid \phi)$$

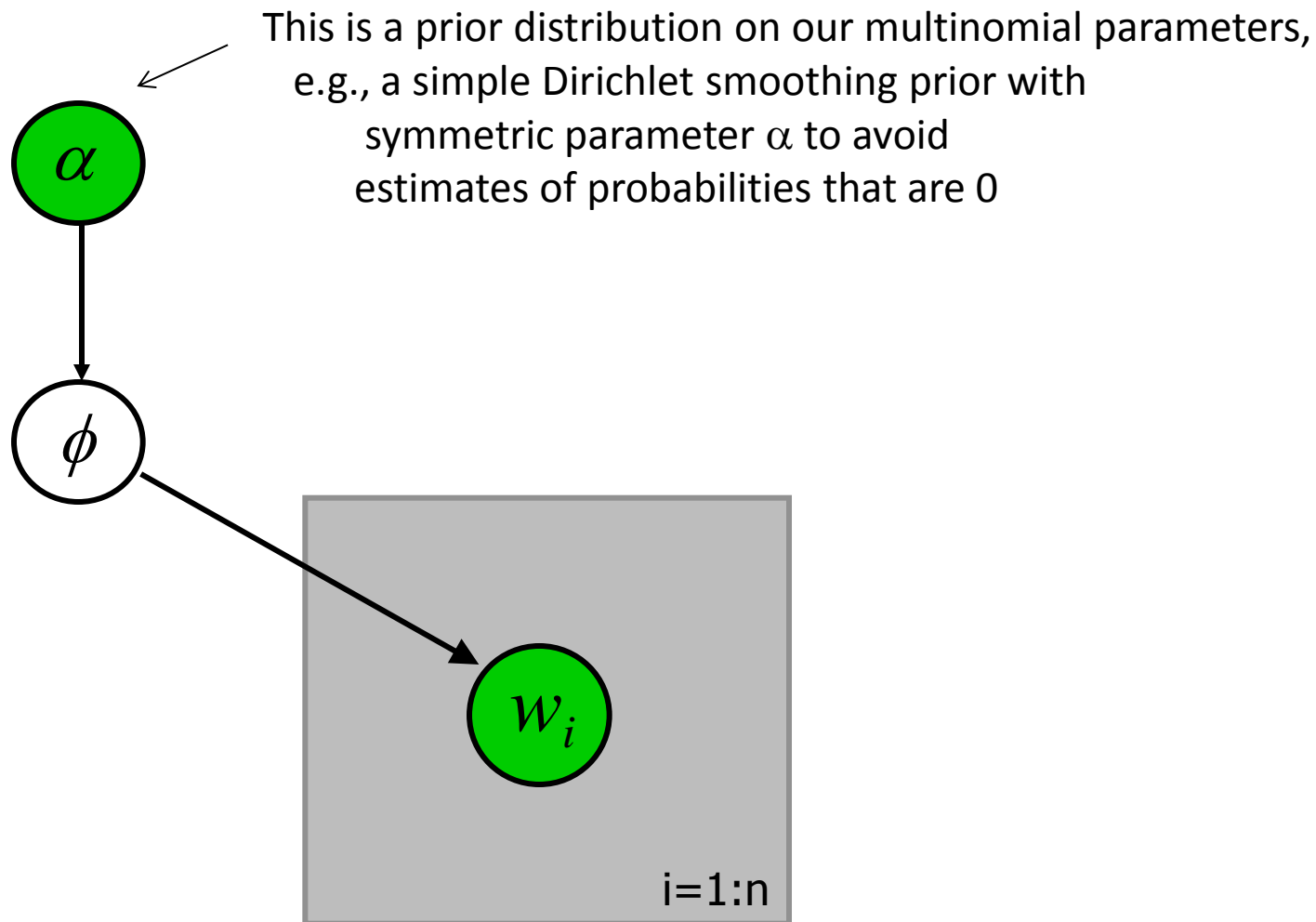


This is “plate notation”

Items inside the plate
are conditionally independent
given the variable outside
the plate

There are “n” conditionally
independent replicates
represented by the plate

Being Bayesian....

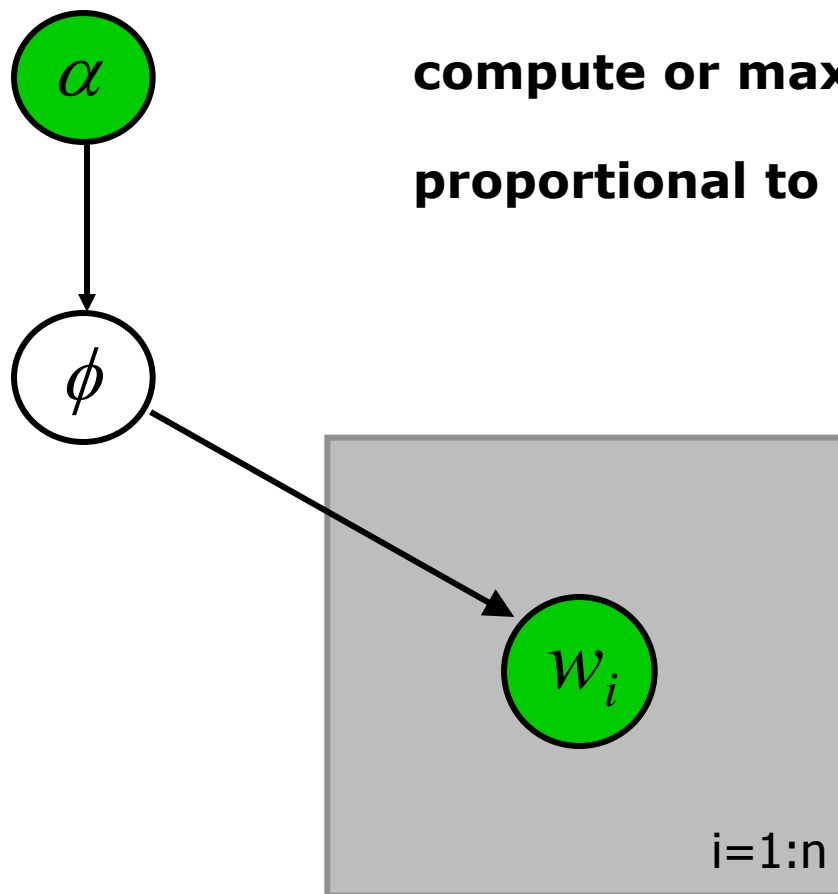


Being Bayesian....

Learning:

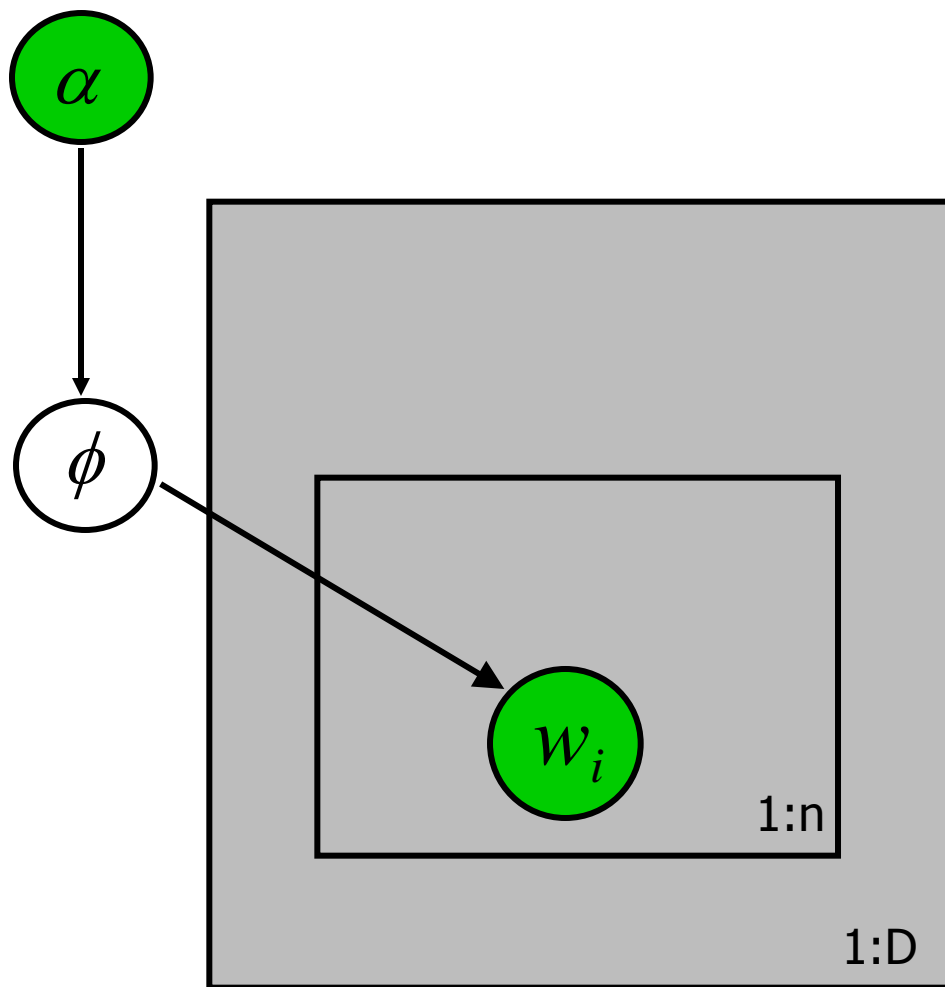
compute or maximize $p(\phi \mid \text{words}, \alpha)$

proportional to $p(\text{words} \mid \phi) p(\phi \mid \alpha)$



Multiple Documents

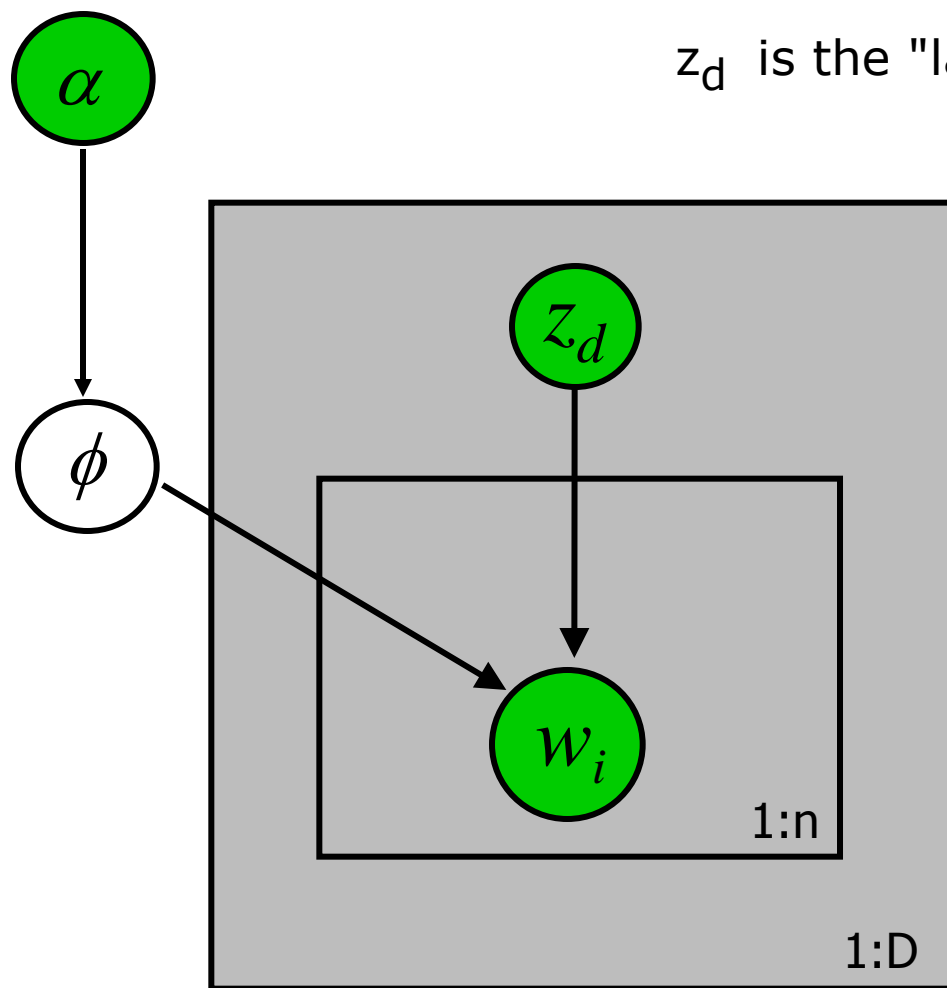
$$p(\text{corpus} \mid \phi) = \prod p(\text{doc} \mid \phi)$$



Different Document Types

$p(w \mid \phi, z_d)$ is a multinomial over words

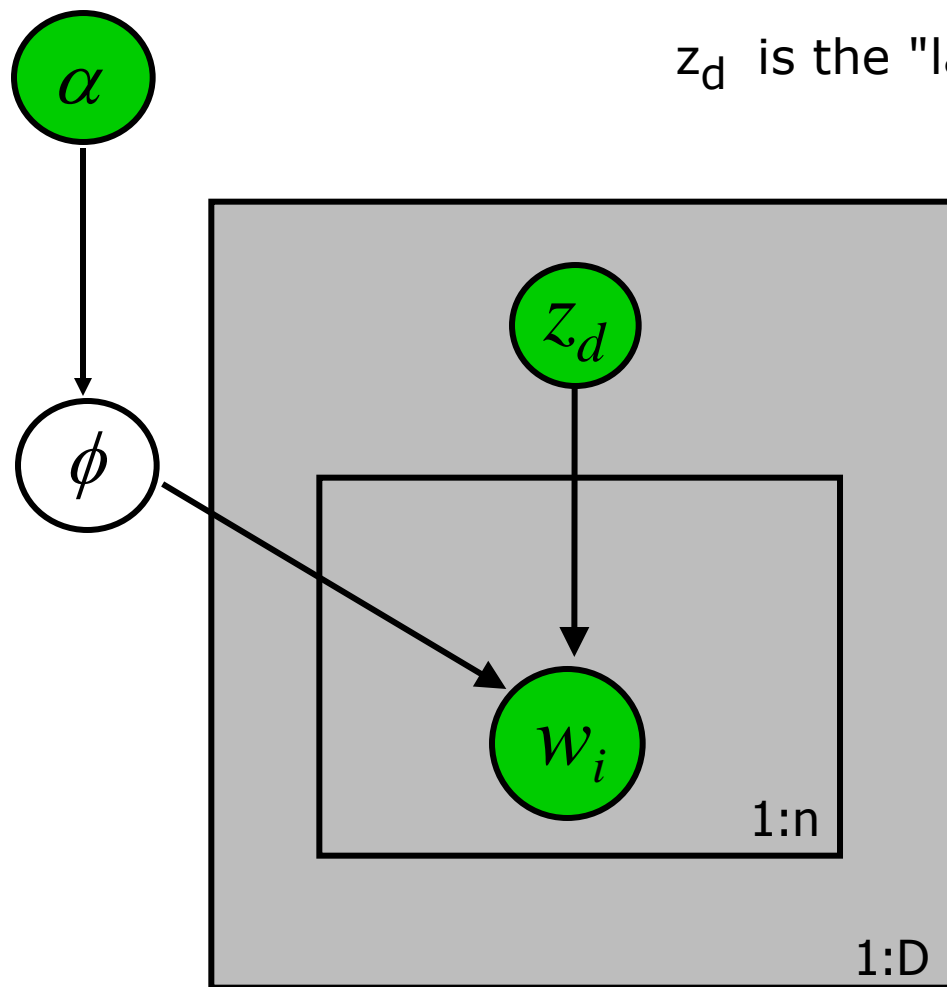
z_d is the "label" for each doc



Different Document Types

$p(w | \phi, z_d)$ is a multinomial over words

z_d is the "label" for each doc

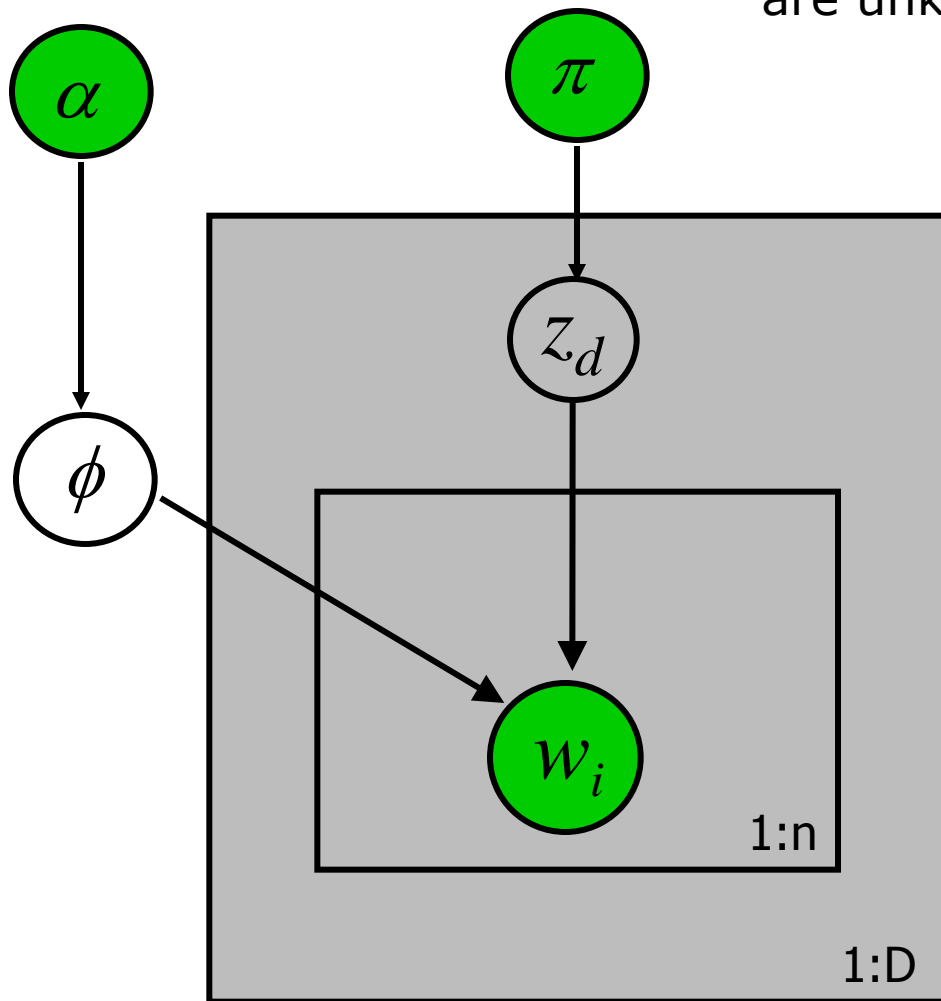


Different multinomials,
depending on the
value of z_d (discrete)

ϕ now represents $|z|$ different
multinomials

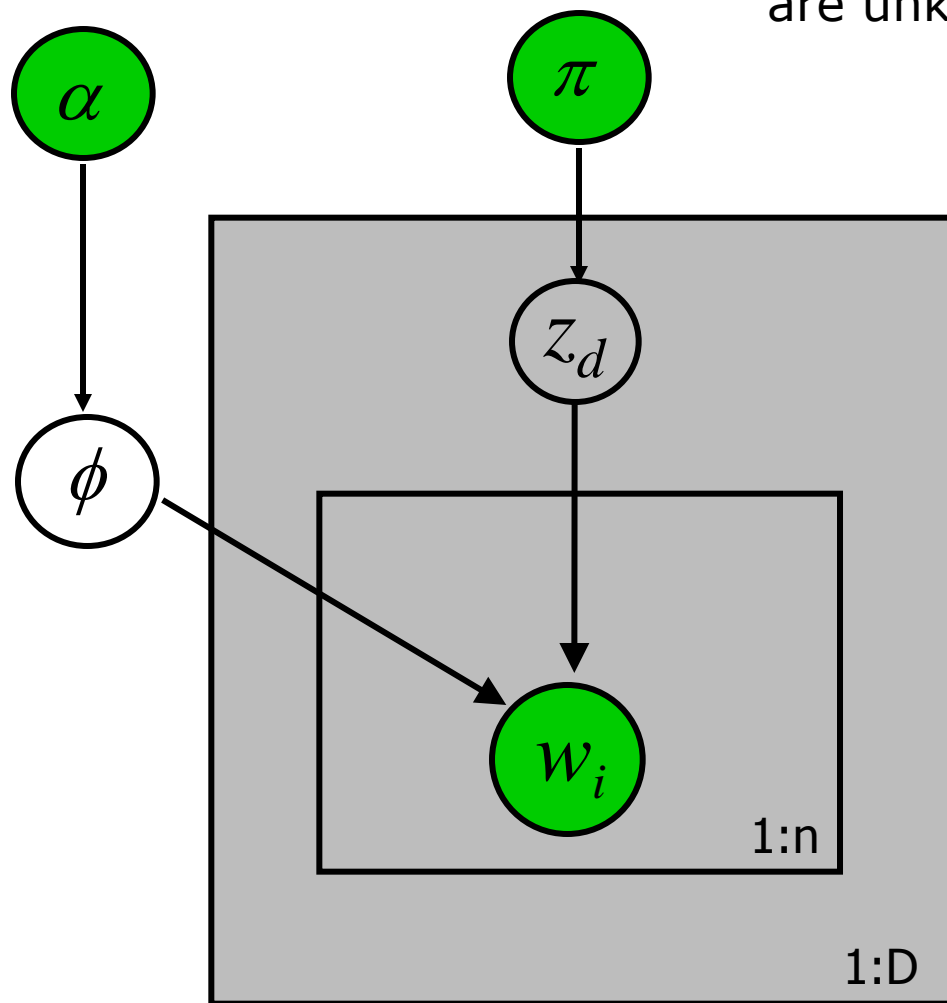
Unknown Document Types

Now the values of z for each document are unknown - hopeless?



Unknown Document Types

Now the values of z for each document are unknown - hopeless?



Not hopeless :)

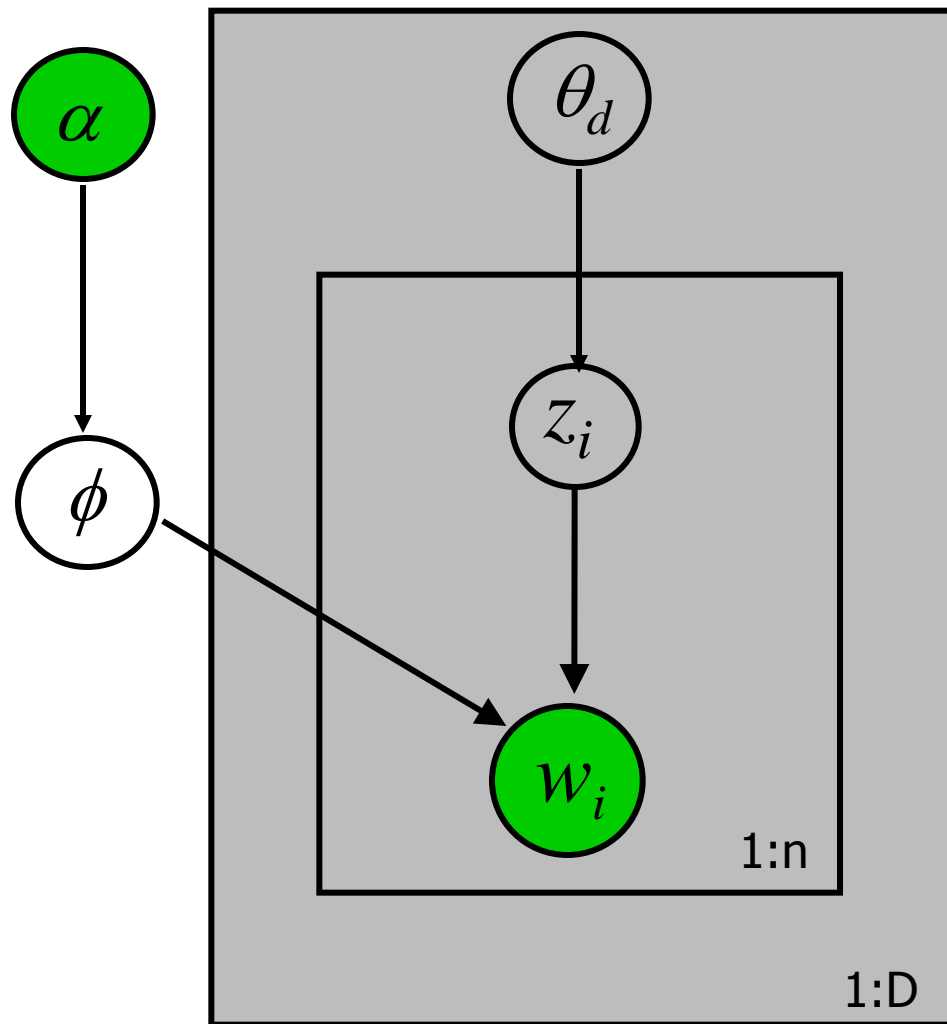
Can learn about both z and θ

e.g., EM algorithm

This gives probabilistic clustering

$p(w \mid z=k, \phi)$ is the k th multinomial over words

Topic Model



z_i is a "label" for each word

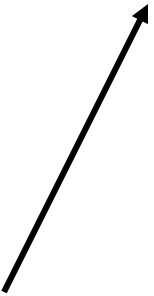
$p(w \mid \phi, z_i = k)$
 = multinomial over words
 = a "topic"

$p(z_i \mid \theta_d)$ = distribution
 over topics that is
 document specific


Mixture Model Equation for Topic Models

$$p(w_i|d) = \sum_{j=1}^T p(w_i|z_j)p(z_j|d)$$

Multinomial over words
for topic z
(the ϕ 's)



Multinomial over topics
for document d
(the θ 's)



Key Features of Topic Models

- Generative model for documents in form of bags of words
- Allows a document to be composed of multiple topics
 - More flexible than clustering which assumes 1 cluster for each document
- Completely unsupervised
 - Topics learned directly from data
 - Leverages strong dependencies at word level AND large data sets
- Learning algorithm
 - Collapsed Gibbs sampling is the method of choice
- Scalable
 - Linear in number of word tokens
 - Can be run on millions of documents

Learning the Model

- Three sets of latent variables we can learn
 - topic-word distributions ϕ
 - document-topic distributions ϑ
 - topic assignments for each word z
- Options:
 - EM algorithm to find point estimates of ϕ and θ
 - e.g., Chien and Wu, IEEE Trans ASLP, 2008
 - Gibbs sampling
 - Find $p(\phi \mid \text{data})$, $p(\theta \mid \text{data})$, $p(z \mid \text{data})$
 - Can be slow to converge
 - Collapsed Gibbs sampling
 - Most widely used method

[See also Asuncion, Welling, Smyth, Teh, UAI 2009 for additional discussion]

Gibbs Sampling

- Say we have 3 parameters x, y, z , and some data
- Bayesian learning:
 - We want to compute $p(x, y, z \mid \text{data})$
 - But frequently it is impossible to compute this exactly
 - However, often we can compute conditionals for individual variables, e.g.,
 $p(x \mid y, z, \text{data})$
 - Not clear how this is useful yet, since it assumes y and z are known (i.e., we condition on them).

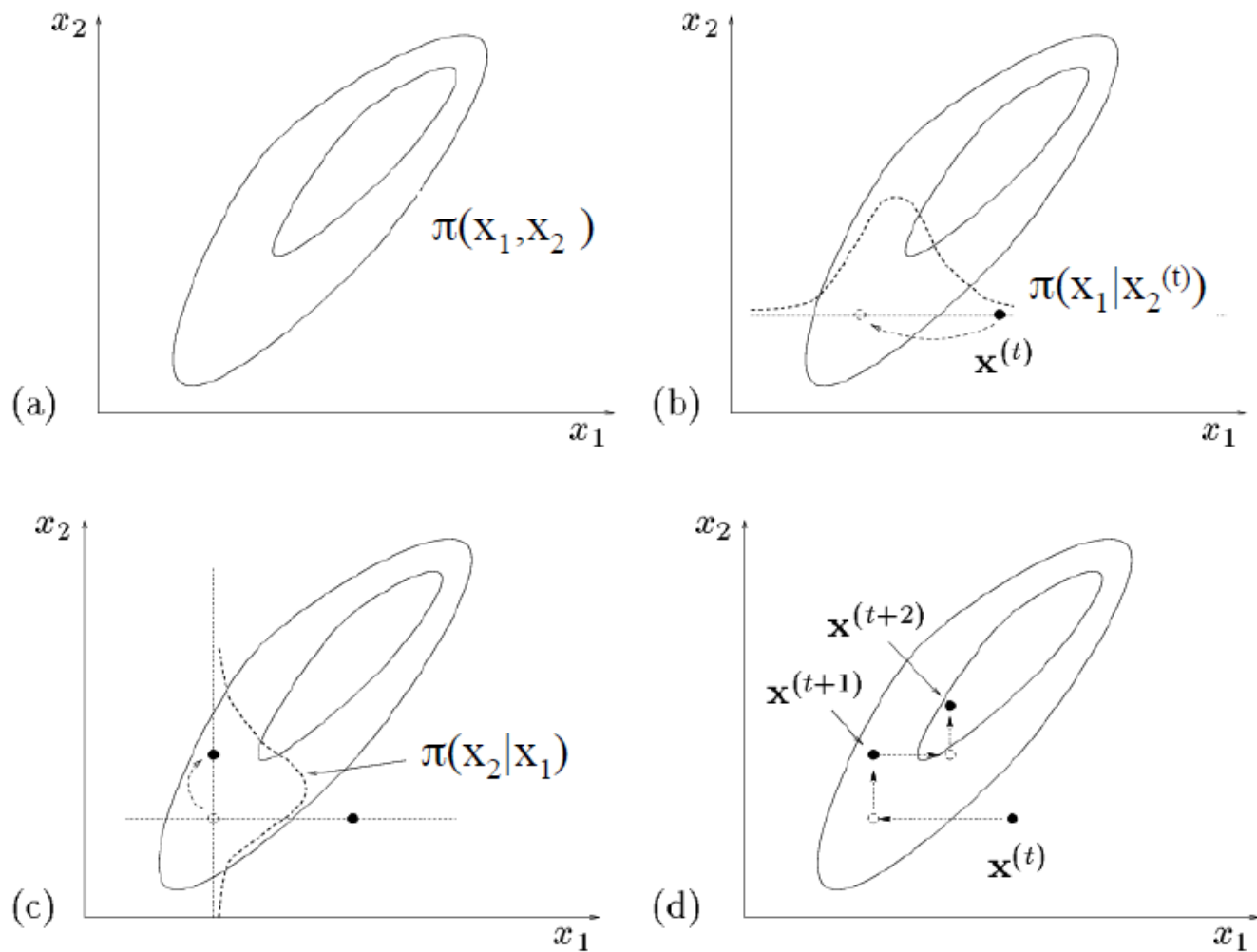
Gibbs Sampling 2

- Example of Gibbs sampling:
 - Initialize with x', y', z' (e.g., randomly)
 - Iterate:
 - Sample new $x' \sim P(x \mid y', z', \text{data})$
 - Sample new $y' \sim P(y \mid x', z', \text{data})$
 - Sample new $z' \sim P(z \mid x', y', \text{data})$
 - Continue for some number (large) of iterations
 - Each iteration consists of a sweep through the hidden variables or parameters (here, x, y , and z)
 - Gibbs = a Markov Chain Monte Carlo (MCMC) method

In the limit, the samples x', y', z' will be samples from the true joint distribution $P(x, y, z \mid \text{data})$

This gives us an *empirical estimate* of $P(x, y, z \mid \text{data})$

Example of Gibbs Sampling in 2d



From online MCMC tutorial notes by Frank Dellaert, Georgia Tech

Gibbs Sampling for the Topic Model

- Recall: 3 sets of latent variables we can learn
 - topic-word distributions ϕ
 - document-topic distributions ϑ
 - topic assignments for each word z
- Gibbs sampling algorithm
 - Initialize all the z 's randomly to a topic, z_1, \dots, z_N
 - Iteration
 - For $i = 1, \dots, N$
 - Sample $z_i \sim p(z_i \mid \text{all other } z\text{'s, data})$
 - Continue for a fixed number of iterations or convergence
 - Note that this is collapsed Gibbs sampling
 - Sample from $p(z_1, \dots, z_N \mid \text{data})$, “collapsing” over ϕ and θ

Review of Topic Models.....

- A topic model represents:
 - Documents as probability distributions over topics
 - Topics as probability distributions over words
- The model is learned from bag-of-words data (docs x words count matrix) using unsupervised learning
- The most commonly used learning algorithm is based on a technique called Gibbs sampling

Topic Model Learning Algorithm

- Input:
 - N documents as count vectors, number of topics T
- Output:
 - T topics, i.e., T topic-word probability vectors
 - N sets of topic weights, one per document
 - Assignment of each word in each document to 1 of T topics
- Algorithm (based on Gibbs Sampling)
 - Randomly initialize all word tokens to a topic (a number from 1 to T)
 - For each iteration
 - Sample a new topic for each word token, keeping all other assignments fixed
 - Iterate through all word tokens in all documents
 - Typically converges quickly (20 to 100 iterations)
 - Each iteration is linear in the number of word tokens

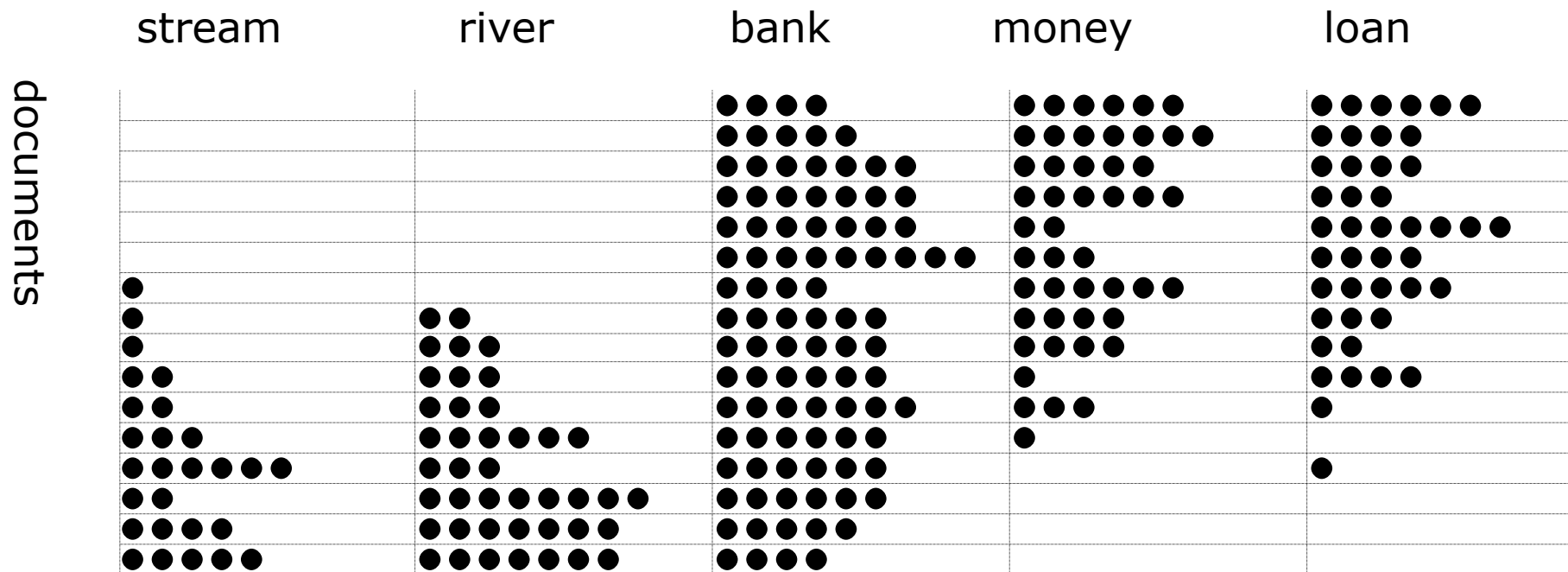
Collapsed Gibbs Sampling Equation for Topic Modeling

The diagram shows the collapsed Gibbs sampling equation for topic modeling, with three callouts explaining the variables:

- count of topic t assigned to doc d** : This callout points to the $n_{td}^{-i} + \alpha$ term in the numerator of the first fraction.
- count of word w assigned to topic t** : This callout points to the $n_{wt}^{-i} + \beta$ term in the numerator of the second fraction.
- probability that word i is assigned to topic t** : This callout points to the entire equation, which represents the probability $p(z_i = t | z_{-i})$.

$$p(z_i = t | z_{-i}) = \frac{n_{td}^{-i} + \alpha}{\sum_{t'} n_{t'd}^{-i} + T\alpha} \times \frac{n_{wt}^{-i} + \beta}{\sum_{w'} n_{w't}^{-i} + W\beta}$$

Word/Document counts for 16 Artificial Documents



Can we recover the original topics and topic mixtures from this data?

Example of Collapsed Gibbs Sampling

- Assign word tokens randomly to 2 topics:

stream	river	bank	money	loan
		○ ○ ○ ○	● ○ ○ ○ ● ○	● ● ○ ● ○ ○
		○ ○ ● ○ ○	● ● ● ● ● ● ○	● ○ ○ ●
		○ ○ ○ ● ○ ○ ○	○ ● ○ ● ○	● ○ ○ ○
		● ● ● ○ ● ○ ○	○ ● ● ○ ○ ○	○ ○ ○
		● ● ○ ● ○ ● ○	● ○	○ ● ○ ○ ○ ○ ○
		○ ● ● ○ ● ● ● ● ●	○ ● ○	○ ○ ● ●
○		○ ● ● ●	● ● ○ ○ ● ○	○ ● ● ● ○
●	○ ●	○ ○ ● ● ● ●	○ ● ● ○	● ● ○
●	○ ○ ●	○ ○ ○ ○ ○ ●	● ○ ● ●	○ ●
● ○	● ● ○	● ○ ○ ○ ○ ○	●	● ○ ○ ●
○ ●	○ ● ●	○ ○ ○ ● ● ○ ○	● ● ●	●
○ ○ ○	○ ○ ○ ○ ● ○	● ○ ● ● ○ ●	○	
○ ○ ○ ● ● ●	○ ● ○	● ○ ○ ○ ● ●		○
○ ○	● ● ○ ○ ○ ● ● ●	● ● ○ ● ○ ○		
○ ● ● ●	● ● ● ○ ○ ● ○	● ○ ● ○ ●		
● ○ ● ● ○	● ● ○ ○ ○ ○ ●	● ● ● ○		

After 1 iteration

stream	river	bank	money	loan
		● ● ○ ○	○ ○ ○ ○ ○ ●	● ○ ○ ○ ○ ○
		● ○ ○ ○ ○	○ ● ● ● ● ● ○	○ ○ ○ ●
		○ ○ ○ ○ ○ ○ ●	○ ○ ○ ○ ●	○ ○ ● ○
		○ ○ ○ ○ ○ ○ ○	● ○ ○ ○ ○ ○ ○	○ ○ ○
		● ● ● ● ● ● ○	● ●	● ○ ● ○ ● ● ●
		● ○ ● ○ ● ● ● ○ ●	● ● ●	● ○ ● ●
●		● ● ● ●	● ● ● ● ● ● ●	● ● ● ● ● ●
○	● ○	● ● ● ● ● ●	● ● ● ●	● ● ●
●	○ ● ●	○ ○ ○ ● ○ ○	○ ● ● ●	● ●
○ ●	○ ○ ○	○ ○ ○ ○ ● ○	○	● ○ ○ ○
○ ●	● ● ○	● ● ● ● ● ● ○	○ ○ ●	○
○ ● ●	○ ○ ● ○ ○ ●	○ ○ ○ ○ ○ ○	○	
● ● ● ● ● ○	○ ● ○	○ ○ ○ ● ○ ○		●
● ●	○ ○ ○ ● ○ ○ ○ ○	○ ● ● ● ○ ●		
● ○ ○ ○	○ ○ ○ ○ ● ○ ○	○ ○ ○ ● ○		
● ● ● ● ●	○ ● ○ ● ○ ● ●	● ○ ● ●		

After 4 iterations

stream	river	bank	money	loan
		● ● ● ●	● ● ● ● ● ●	● ● ● ● ● ●
		● ○ ○ ● ○	● ○ ● ● ● ● ●	● ● ● ●
		○ ○ ● ○ ● ● ●	● ● ● ○ ●	○ ○ ○ ●
		○ ○ ○ ○ ○ ○ ○	○ ● ● ● ● ○	○ ● ○
		● ○ ● ● ● ● ●	● ●	● ● ● ● ● ● ●
		● ● ● ● ● ● ● ● ●	● ● ●	● ● ● ●
●		● ● ● ●	● ● ● ● ● ● ●	● ● ● ● ●
○	○ ○	○ ● ○ ○ ● ○	● ○ ● ●	● ● ●
●	● ○ ●	● ● ● ● ● ○	● ● ● ●	● ●
● ○	○ ○ ○	○ ● ● ○ ○ ●	○	● ● ○ ○
○ ○	○ ○ ○	○ ● ○ ● ○ ○ ○	● ○ ○	○
○ ○ ○	○ ○ ○ ○ ○ ○	● ○ ○ ○ ○ ○	○	
○ ○ ○ ○ ○ ○ ○	○ ○ ○	○ ● ● ○ ○ ○		●
○ ○	○ ○ ○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○		
○ ○ ○ ○	○ ○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○		
○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○ ○	○ ○ ○ ○		

After 32 iterations

topic 1	
stream	.40
bank	.35
river	.25

topic 2	
bank	.39
money	.32
loan	.29

stream

river

bank

money

loan

		● ● ● ●	● ● ● ● ● ●	● ● ● ● ● ●
		● ● ● ○ ●	● ● ● ● ● ● ●	● ● ● ●
		● ● ● ● ● ● ●	● ● ● ● ●	● ● ● ●
		● ● ● ● ● ● ●	● ● ● ● ● ●	● ● ●
		● ● ● ○ ● ● ●	● ●	● ● ● ● ● ● ●
		● ● ● ● ● ● ● ● ●	● ● ●	● ● ● ●
○		● ● ● ○	● ● ● ● ● ● ●	● ● ● ● ●
○	○ ○	● ○ ○ ○ ● ○	● ● ● ●	● ● ●
○	○ ○ ○	● ○ ● ● ● ●	● ● ● ●	● ●
○ ○	○ ○ ○	○ ● ● ● ● ●	●	● ● ● ●
○ ○	○ ○ ○	○ ○ ○ ● ○ ○ ○	● ● ●	●
○ ○ ○	○ ○ ○ ○ ○ ○	○ ● ○ ● ● ●	●	
○ ○ ○ ○ ○ ○	○ ○ ○	● ○ ○ ○ ○ ○		●
○ ○	○ ○ ○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○		
○ ○ ○ ○	○ ○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○		
○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○ ○	○ ○ ○ ○		

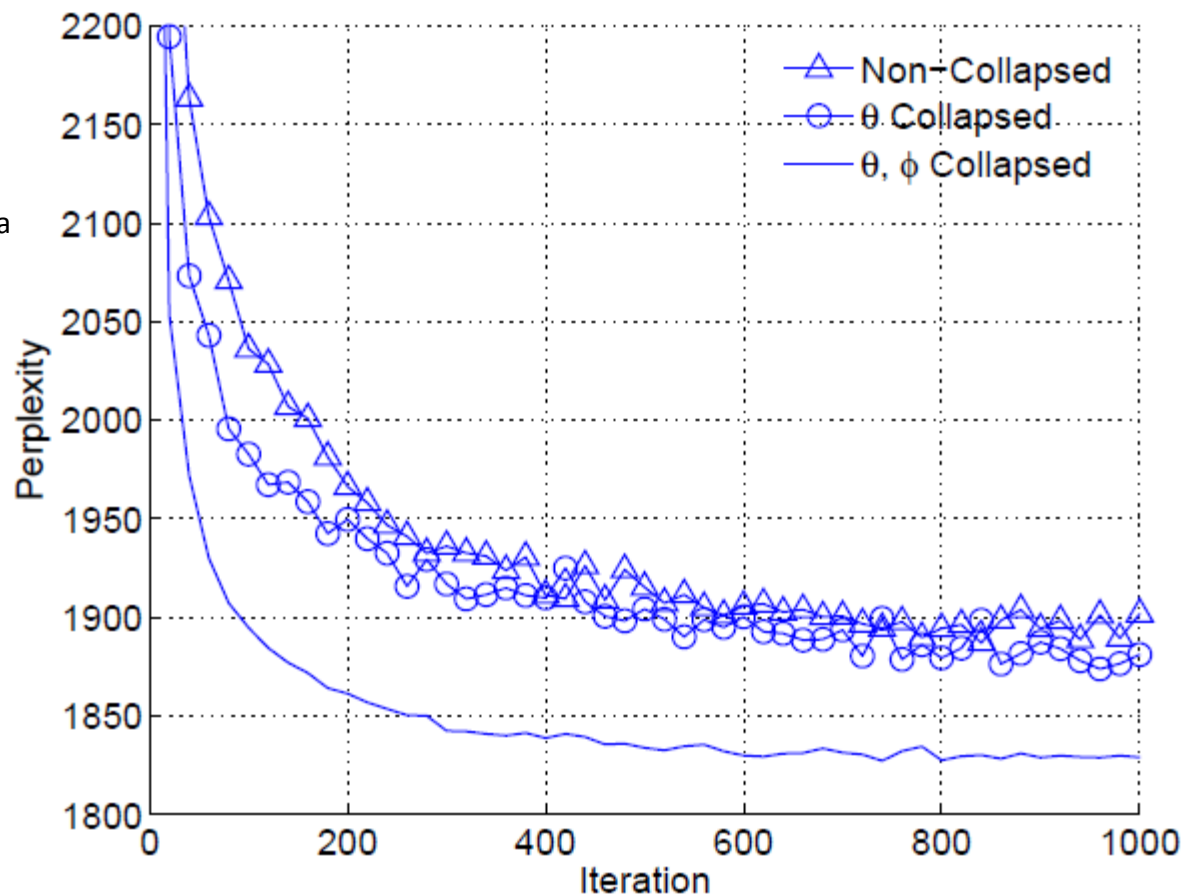
Computational Aspects

- Convergence
 - In the limit, samples x', y', z' are from $P(x, y, z \mid \text{data})$
 - How many iterations are needed?
 - Cannot be computed ahead of time
 - Early iterations are discarded (“burn-in”)
 - Typically monitor some quantities of interest to monitor convergence
 - Detecting convergence in Gibbs/MCMC is a tricky issue!
- Complexity per iteration
 - Linear in number of hidden variables and parameters
 - Times the complexity of generating a sample each time

Convergence Example

(from Newman et al, JMLR, 2009)

A measure of how well
the model is fitting the data



The New York Times

music band songs rock album jazz pop song singer night	book life novel story books man stories love children family	art museum show exhibition artist artists paintings painting century works	game knicks nets points team season play games night coach	show film television movie series says life man character know
theater play production show stage street broadway director musical directed	clinton bush campaign gore political republican dole presidential senator house	stock market percent fund investors funds companies stocks investment trading	restaurant sauce menu food dishes street dining dinner chicken served	budget tax governor county mayor billion taxes plan legislature fiscal

Each box contains the high-probability words from each topic

Note that this is entirely unsupervised – no human labeling, just the words

Example from Hoffman et al, 2013

3 of 300 example topics (from TASA Corpus)

Example from Mark Steyvers

TOPIC 82	
WORD	PROB.
PLAY	0.0601
PLAYS	0.0362
STAGE	0.0305
MOVIE	0.0288
SCENE	0.0253
ROLE	0.0245
AUDIENCE	0.0197
THEATER	0.0186
PART	0.0178
FILM	0.0148
ACTORS	0.0145
DRAMA	0.0136
REAL	0.0128
CHARACTER	0.0122
ACTOR	0.0116
ACT	0.0114
MOVIES	0.0114
ACTION	0.0101
SET	0.0097
SCENES	0.0094

TOPIC 77	
WORD	PROB.
MUSIC	0.0903
DANCE	0.0345
SONG	0.0329
PLAY	0.0301
SING	0.0265
SINGING	0.0264
BAND	0.0260
PLAYED	0.0229
SANG	0.0224
SONGS	0.0208
DANCING	0.0198
PIANO	0.0169
PLAYING	0.0159
RHYTHM	0.0145
ALBERT	0.0134
MUSICAL	0.0134
DRUM	0.0129
GUITAR	0.0098
BEAT	0.0097
BALLET	0.0096

TOPIC 166	
WORD	PROB.
PLAY	0.1358
BALL	0.1288
GAME	0.0654
PLAYING	0.0418
HIT	0.0324
PLAYED	0.0312
BASEBALL	0.0274
GAMES	0.0250
BAT	0.0193
RUN	0.0186
THROW	0.0158
BALLS	0.0154
TENNIS	0.0107
HOME	0.0099
CATCH	0.0098
FIELD	0.0097
PLAYER	0.0096
FUN	0.0092
THROWING	0.0083
PITCHER	0.0080

Topic Modeling on different text sources ...

Collection	# docs	Description
New York Times	1,500,000	News articles from New York Times
Austen	1,400	The six Jane Austen novels, broken up into 100-line sections
Blogs	4,000	Blog entries harvested from Daily Kos
Bible	1,200	Chapters in the bible (KJV)
Police Reports	250,000	Police accident reports from North Carolina
CiteSeer	750,000	Abstracts from research publications in computer science and engineering
Search Queries	1,000,000	Queries issued to web search engine
Enron	250,000	Enron emails seized by the US Government for the federal case against the company

... sample topics

Collection	Sample Topic
New York Times	[WMD] IRAQ iraqi weapon war SADDAM_HUSSEIN SADDAM resolution UNITED_STATES military inspector U_N UNITED_NATION BAGHDAD inspection action SECURITY_COUNCIL
Austen	[SENTIMENT] felt comfort feeling feel spirit mind heart ill evil fear impossible hope poor distress end loss relief suffering concern dreadful misery unhappy
Blogs	[ELECTIONS] november poll house electoral governor polls account ground republicans trouble
Bible	[COMMANDS] thou thy thee shalt thine lord god hast unto not shall
Police Reports	[RAN OFF ROAD] v1 off road ran came rest ditch traveling struck side shoulder tree overturned control lost
CiteSeer	[GRAPH THEORY] graph edge vertices edges vertex number directed connected degree coloring subgraph set drawing
Search Queries	[CREDIT] credit card loans bill loan report bad visa debt score
Enron	[ENERGY CRISIS] state california power electricity utilities davis energy prices generators edison public deregulation billion governor federal consumers commission plants companies electric wholesale crisis summer

Enron Email Topics

TOPIC 36	
WORD	PROB.
FEEDBACK	0.0781
PERFORMANCE	0.0462
PROCESS	0.0455
PEP	0.0446
MANAGEMENT	0.03
COMPLETE	0.0205
QUESTIONS	0.0203
SELECTED	0.0187
COMPLETED	0.0146
SYSTEM	0.0146

TOPIC 72	
WORD	PROB.
PROJECT	0.0514
PLANT	0.028
COST	0.0182
CONSTRUCTION	0.0169
UNIT	0.0166
FACILITY	0.0165
SITE	0.0136
PROJECTS	0.0117
CONTRACT	0.011
UNITS	0.0106

TOPIC 54	
WORD	PROB.
FERC	0.0554
MARKET	0.0328
ISO	0.0226
COMMISSION	0.0215
ORDER	0.0212
FILING	0.0149
COMMENTS	0.0116
PRICE	0.0116
CALIFORNIA	0.0110
FILED	0.0110

TOPIC 23	
WORD	PROB.
ENVIRONMENTAL	0.0291
AIR	0.0232
MTBE	0.019
EMISSIONS	0.017
CLEAN	0.0143
EPA	0.0133
PENDING	0.0129
SAFETY	0.0104
WATER	0.0092
GASOLINE	0.0086

“Personal” Topics...

TOPIC 66	
WORD	PROB.
HOLIDAY	0.0857
PARTY	0.0368
YEAR	0.0316
SEASON	0.0305
COMPANY	0.0255
CELEBRATION	0.0199
ENRON	0.0198
TIME	0.0194
RECOGNIZE	0.019
MONTH	0.018

TOPIC 182	
WORD	PROB.
TEXANS	0.0145
WIN	0.0143
FOOTBALL	0.0137
FANTASY	0.0129
SPORTSLINE	0.0129
PLAY	0.0123
TEAM	0.0114
GAME	0.0112
SPORTS	0.011
GAMES	0.0109

TOPIC 113	
WORD	PROB.
GOD	0.0357
LIFE	0.0272
MAN	0.0116
PEOPLE	0.0103
CHRIST	0.0092
FAITH	0.0083
LORD	0.0079
JESUS	0.0075
SPIRITUAL	0.0066
VISIT	0.0065

TOPIC 109	
WORD	PROB.
AMAZON	0.0312
GIFT	0.0226
CLICK	0.0193
SAVE	0.0147
SHOPPING	0.0140
OFFER	0.0124
HOLIDAY	0.0122
RECEIVE	0.0102
SHIPPING	0.0100
FLOWERS	0.0099

Political Topics

TOPIC 18	
WORD	PROB.
POWER	0.0915
CALIFORNIA	0.0756
ELECTRICITY	0.0331
UTILITIES	0.0253
PRICES	0.0249
MARKET	0.0244
PRICE	0.0207
UTILITY	0.0140
CUSTOMERS	0.0134
ELECTRIC	0.0120

TOPIC 22	
WORD	PROB.
STATE	0.0253
PLAN	0.0245
CALIFORNIA	0.0137
POLITICIAN Y	0.0137
RATE	0.0131
BANKRUPTCY	0.0126
SOCAL	0.0119
POWER	0.0114
BONDS	0.0109
MOU	0.0107

TOPIC 114	
WORD	PROB.
COMMITTEE	0.0197
BILL	0.0189
HOUSE	0.0169
WASHINGTON	0.0140
SENATE	0.0135
POLITICIAN X	0.0114
CONGRESS	0.0112
PRESIDENT	0.0105
LEGISLATION	0.0099
DC	0.0093

TOPIC 194	
WORD	PROB.
LAW	0.0380
TESTIMONY	0.0201
ATTORNEY	0.0164
SETTLEMENT	0.0131
LEGAL	0.0100
EXHIBIT	0.0098
CLE	0.0093
SOCALGAS	0.0093
METALS	0.0091
PERSON Z	0.0083

Automated Tagging of Words

(numbers & colors → topic assignments)

Example from Mark Steyvers

A **Play**⁰⁸² is written⁰⁸² to be performed⁰⁸² on a stage⁰⁸² before a live⁰⁹³ audience⁰⁸² or before motion²⁷⁰ picture⁰⁰⁴ or television⁰⁰⁴ cameras⁰⁰⁴ (for later⁰⁵⁴ viewing⁰⁰⁴ by large²⁰² audiences⁰⁸²). A **Play**⁰⁸² is written⁰⁸² because playwrights⁰⁸² have something

He was listening⁰⁷⁷ to music⁰⁷⁷ coming⁰⁰⁹ from a passing⁰⁴³ riverboat. The music⁰⁷⁷ had already captured⁰⁰⁶ his heart¹⁵⁷ as well as his ear¹¹⁹. It was jazz⁰⁷⁷. Bix beiderbecke had already had music⁰⁷⁷ lessons⁰⁷⁷. He wanted²⁶⁸ to **play**⁰⁷⁷ the cornet. And he wanted²⁶⁸ to **play**⁰⁷⁷ jazz⁰⁷⁷

Jim²⁹⁶ **plays**¹⁶⁶ the game¹⁶⁶. Jim²⁹⁶ likes⁰⁸¹ the game¹⁶⁶ for one. The game¹⁶⁶ book²⁵⁴ helps⁰⁸¹ jim²⁹⁶. Don¹⁸⁰ comes⁰⁴⁰ into the house⁰³⁸. Don¹⁸⁰ and jim²⁹⁶ read²⁵⁴ the game¹⁶⁶ book²⁵⁴. The boys⁰²⁰ see a game¹⁶⁶ for two. The two boys⁰²⁰ **play**¹⁶⁶ the game¹⁶⁶.

What is this paper about?

Empirical Bayes screening for multi-item associations

Bill DuMouchel and Daryl Pregibon, ACM SIGKDD 2001

Most likely topics according to the model are...

1. data, mining, discovery, association, attribute..
2. set, subset, maximal, minimal, complete,...
3. measurements, correlation, statistical, variation,
4. Bayesian, model, prior, data, mixture,.....

Original article

Most likely words from top topics



TECHVIEW: DNA SEQUENCING

Sequencing the Genome, Fast

James C. Mullikin and Amanda A. McMurtry

Genome sequencing projects reveal the genetic makeup of an organism by reading off the sequence of the DNA bases, which encodes all of the information necessary for the life of the organism. The base sequence contains four nucleotides—adenine, thymine, guanine, and cytosine—which are linked together into long double-helical chains. Over the last two decades, automated DNA sequencers have made the process of obtaining the base-by-base sequence of DNA easier. By application of an electric field across a gel matrix, these sequencers separate fluorescently labeled DNA molecules that differ in size by one base. As the molecules move past a given point in the gel, laser excitation of a fluorescent dye specific to the base at the end of the molecule yields a base-specific signal that can be automatically recorded.

The latest sequencer to be launched is Perkin-Elmer's much-anticipated ABI Prism 3700 DNA Analyzer which, like the Molecular Dynamics MegabACE 1000 launched last year, incorporates a capillary tube to hold the sequence gel rather than a traditional slab-gel apparatus. Extra interest in the ABI 3700 has been generated because Craig Venter of Celera Genomics Corporation anticipates that ~230 of these machines (1) will enable the company to produce raw sequence for the entire 3 gigabases (Gb) of the human genome in 3 years. The specifications of the ABI 3700 machine say that, with less than 1 hour of human labor per day, it can sequence 768 samples per day. Assuming that each sample gives an average of 400 base pairs (bp) of usable sequence data (its read length) and any section from the entire human genome is covered by an average of 10 overlapping independent reads (2), the 75 million samples that Celera must process will require ~300,000 ABI 3700 machine days. With ~230 machines, that works out to less than 2 years or about 434 days, which affords some margin of error for unexpected developments.

At the Sanger Centre, we have finished 146 Mb of genomic sequence from a vari-

ety of genomes, including 81 Mb of sequence from the human genome, the largest amount of any center so far (3). We are aiming to sequence 1 Gb of human sequence in rough-draft form by 2003, with a finished version by 2005. Our sequencing equipment includes 44 ABI 373XL, 41 ABI 377XL, and 31 ABI 377XL-96 slab gel sequencers from Perkin-Elmer plus 6 Molecular Dynamics MegabACE 1000 capillary sequencers, allowing a maximum throughput of 32,000 samples per day. Two ABI 3700 capillary sequencers—delivered

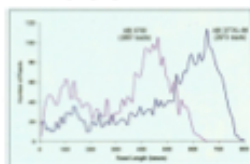


Fig. 5. Comparison of read length histograms for sequencing collected with the ABI 3700 capillary machine and the ABI 3770-96 slab gel machine. The capillary machine underperforms the slab gel machine by about 200 bases. Both sets of reads are from runs with 400 bp Big Dye Terminator chemistry. Read length is computed as the number of bases per read where the predicted error rate is less than or equal to 1.0% (2). The "gtred" Q value was recalculated for each type of read.

to the Sanger Centre in December 1998—are in our Research and Development department for evaluation. Then, the ABI 3700 will ultimately be added to our present capacity to reach our goal.

The ABI 3700 DNA sequencer is built into a floor-standing cabinet, which contains in its base all the reagents required for its operation. The major components are readily accessible for replacement, which is required every day under high-throughput operation. At bench height within the cabinet is a four-position bed, on which microtiter plates of DNA samples are located. The operator places the prepared plates in the position, closes the front of the machine and programs it by using a personal computer. A robotic arm transfers DNA sam-

ples from the plates into wells that open into the capillaries. This and the rest of the sequencing operation is fully automatic. The machine can currently process four 96-well plates of DNA samples unattended, taking approximately 16 hours before operator intervention is required. This rate falls short of the design specification of four 96-well plates in 12 hours.

The main innovation of the ABI 3700 is the use of a sheath flow fluorescence detection system (4). Detection of the DNA fragments occurs 300 µm past the end of the capillary within a fused silica sheath. A laminar fluid flows over the ends of the capillaries, drawing the DNA fragments as they emerge from the capillaries through a fused silica sheath simultaneously interacting with all of the samples. The emitted fluorescence is detected with a spectral CCD (charge-coupled device) detector. This arrangement means that there are no moving parts in the detection system, other than a shutter in front of the CCD detector.

We have evaluated these machines for their performance, operation, ease of use, and reliability in comparison to the more commonly used slab gel sequencing machines. In automated sequencers, there are two methods for containing the gel matrix. One is to polymerize a gel matrix between two finely suspended glass plates (0.4 mm or less)—the slab gel method. The other is to inject a polymer matrix into a capillary (internal diameter ~0.2 mm). Most sequencing facilities use the slab gel method, because multicapillary sequencers have only recently become commercially available.

With either type of system, the aim is to read as many bases as possible for a given sample of DNA—that is, long read lengths are desirable. In fact, a system that could read twice as many bases but at half the speed of another system is preferable, if both systems cost the same. This is because assembling relatively fewer long-sequenced fragments is easier than assembling many short ones. So, read length is an important parameter when evaluating new sequencing technologies.

We have directly compared the ABI 3700 sequencer to the ABI 377XL slab gel sequencer by evaluating the sequence data obtained from both machines with human DNA samples. These samples were subcloned into plasmid or in 1) phage and prepared and sequenced with our standard protocols for Perkin-Elmer Big Dye Terminator chemistry.

sequence
genome
genes
sequences
human
gene
dna
sequencing
chromosome
regions
analysis
data
genomic
number

devices
device
materials
current
high
gate
light
silicon
material
technology
electrical
fiber
power
based

data
information
network
web
computer
language
networks
time
software
system
words
algorithm
number
internet

The authors are at The Sanger Centre, Wellcome Trust Genome Campus, Hinxton, Cambs, CB10 1TA, UK. E-mail: jcm@sanger.ac.uk

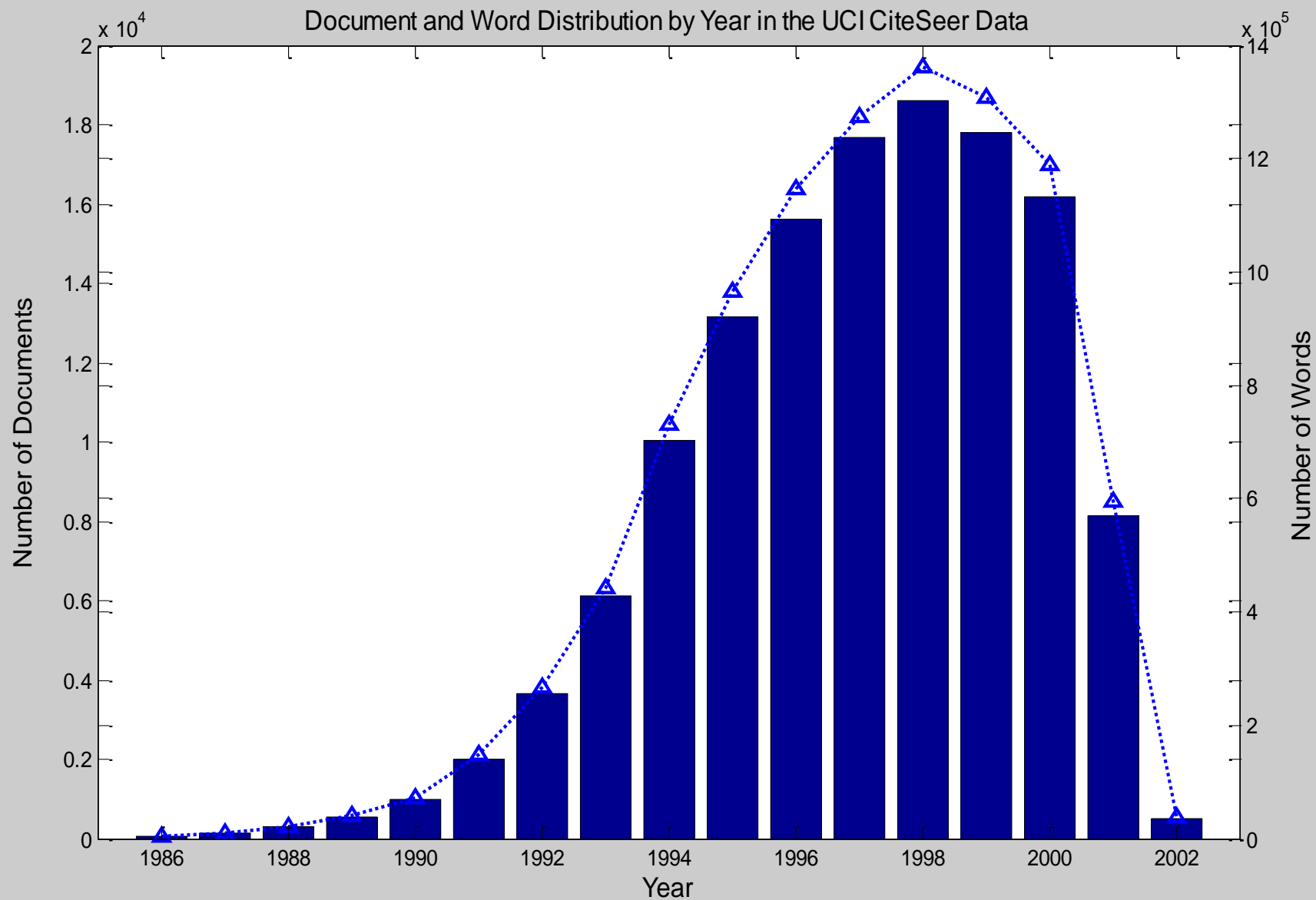
www.sciencemag.org SCIENCE VOL 283 19 MARCH 1999

1867

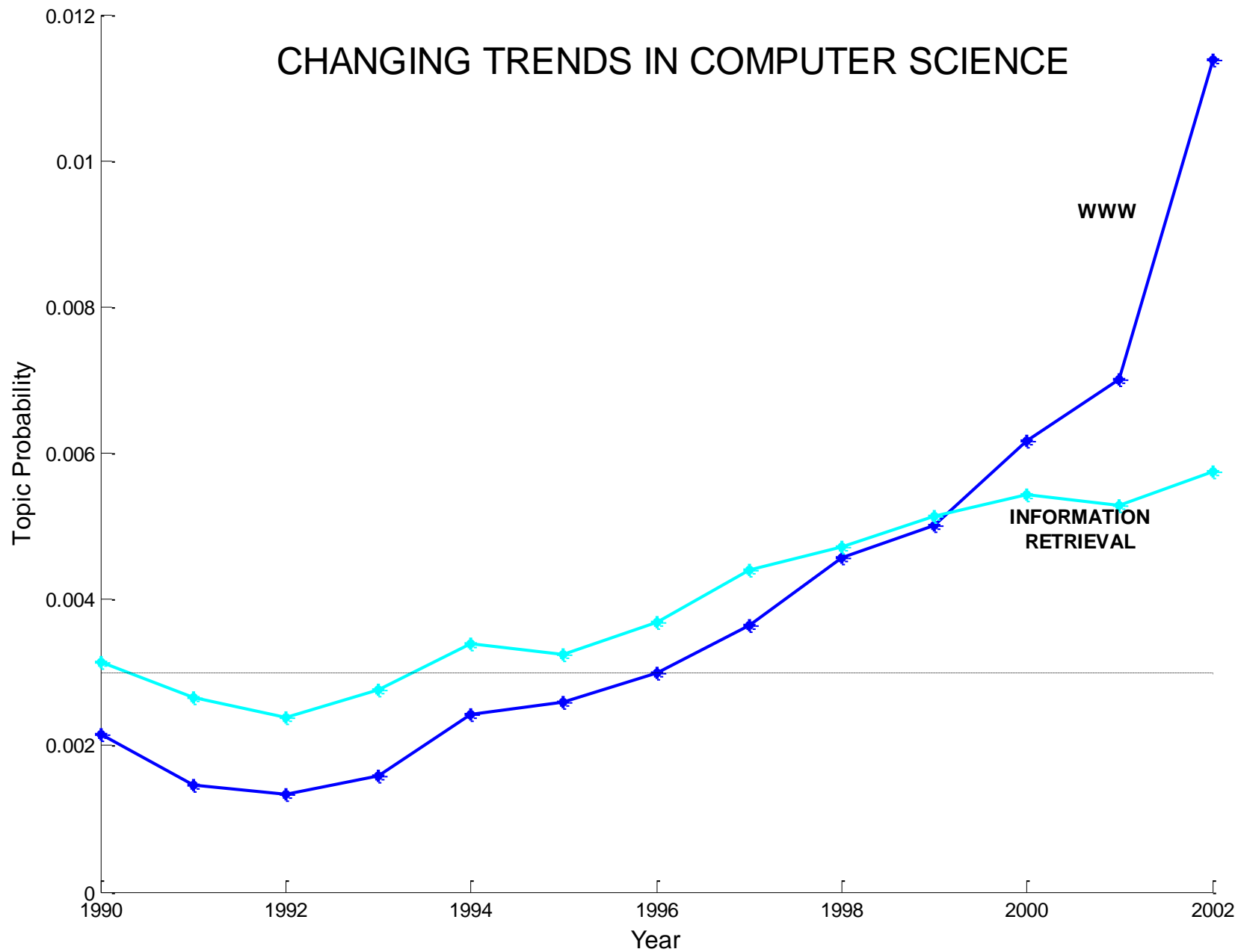
Example courtesy of David Blei, Princeton

Temporal patterns in Topics: Hot and Cold Topics

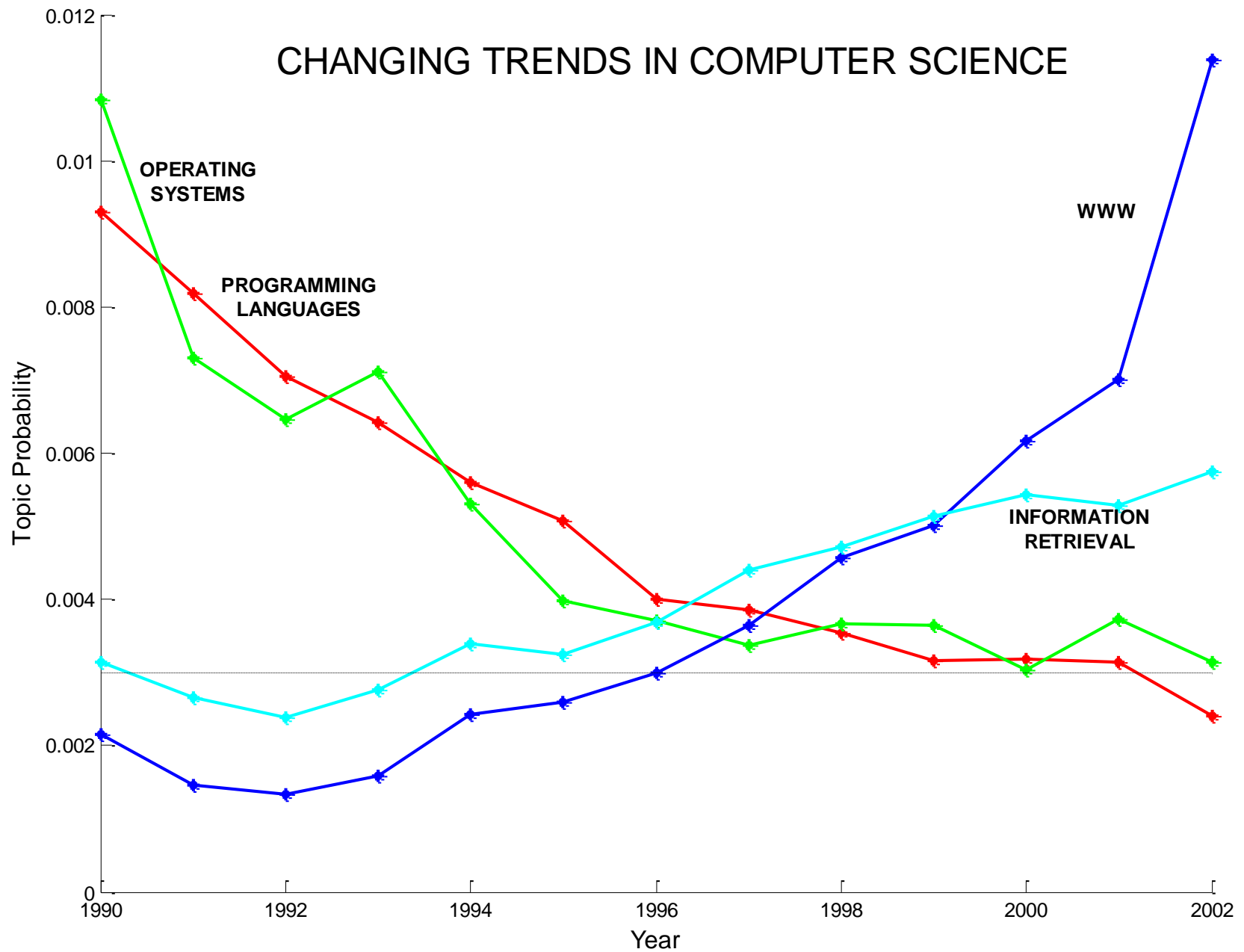
- CiteSeer papers from 1986-2002, about 200k papers
- For each year, calculate the fraction of words assigned to each topic
- This gives us time-series for topics
 - Hot topics become more prevalent
 - Cold topics become less prevalent



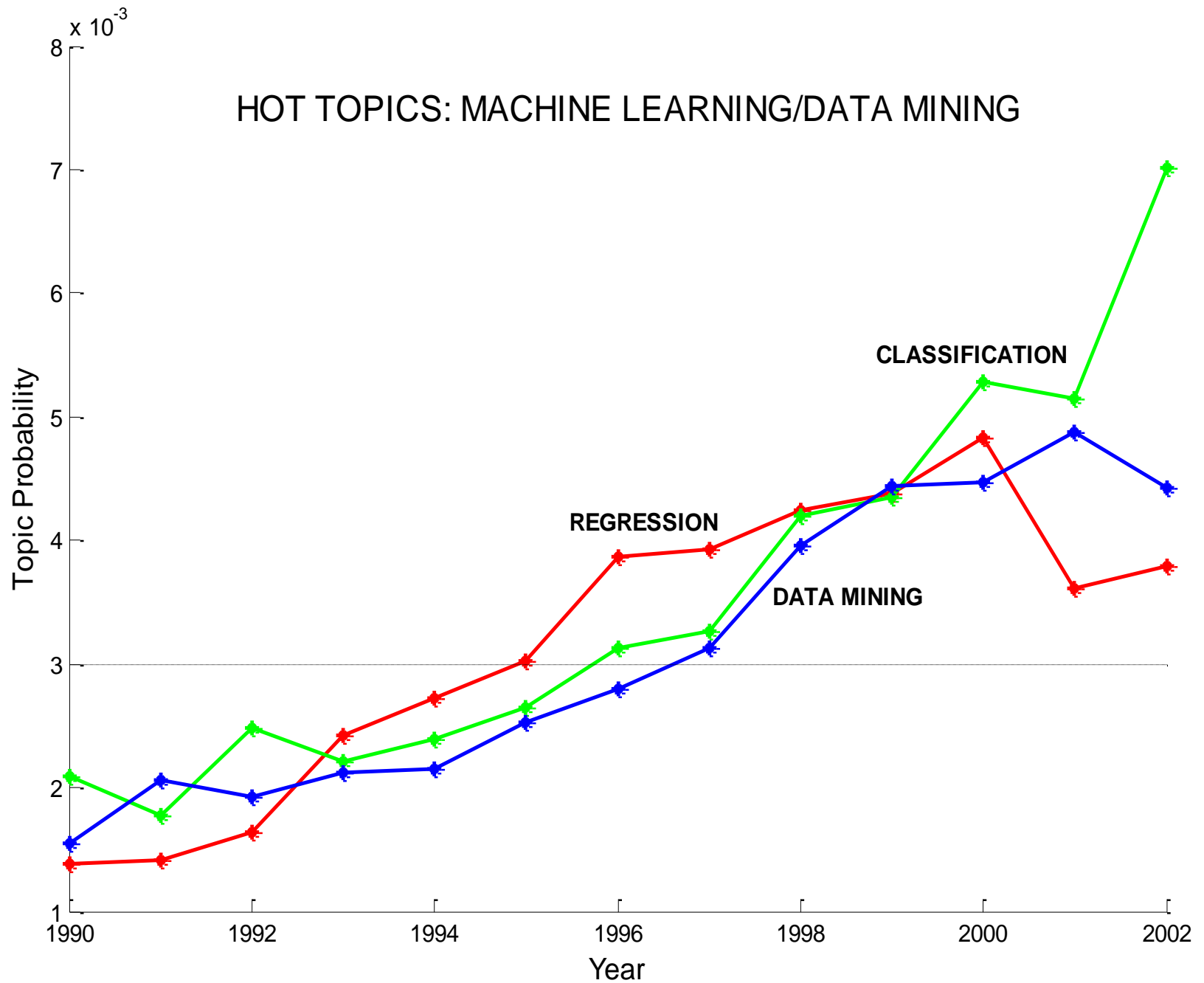
CHANGING TRENDS IN COMPUTER SCIENCE

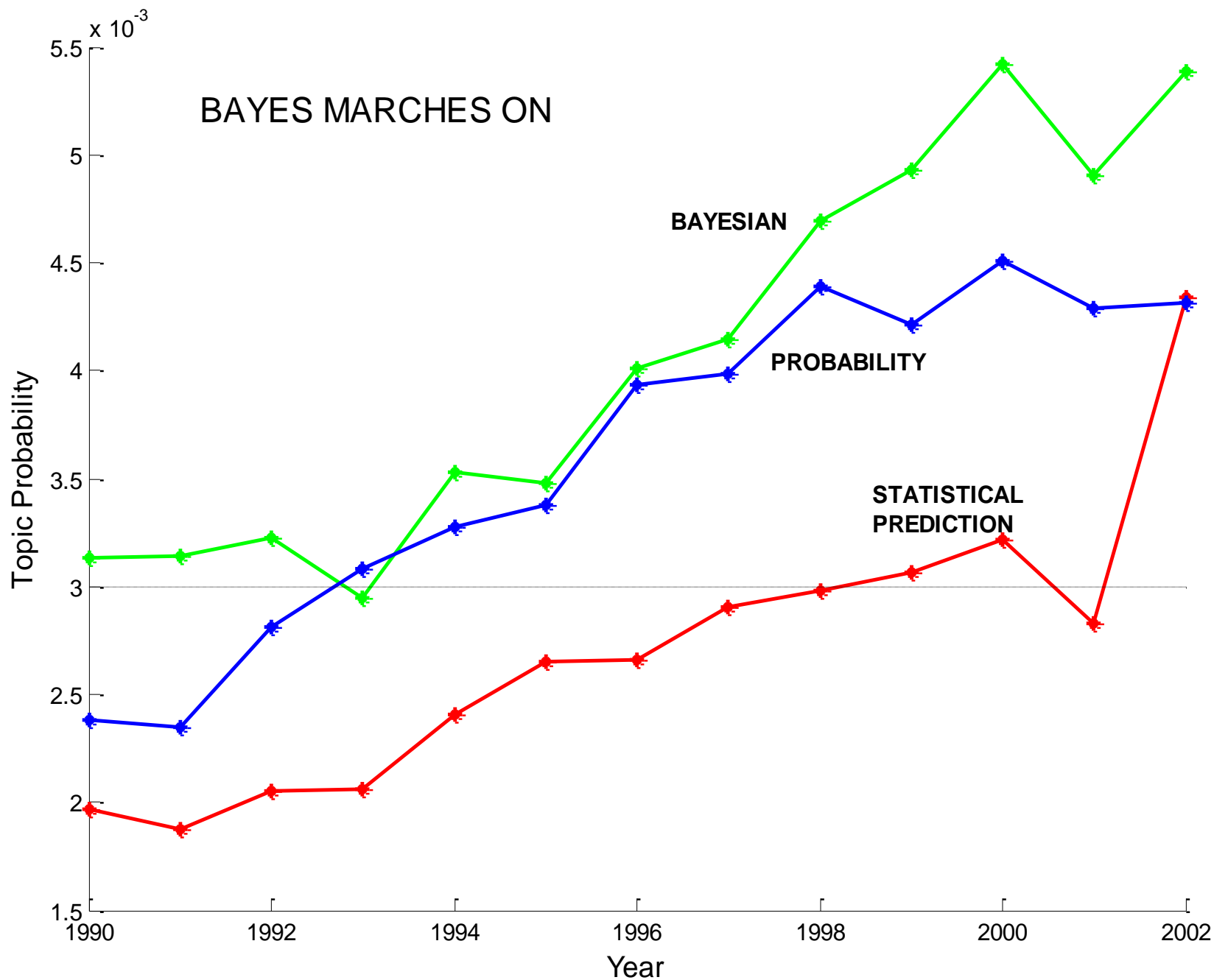


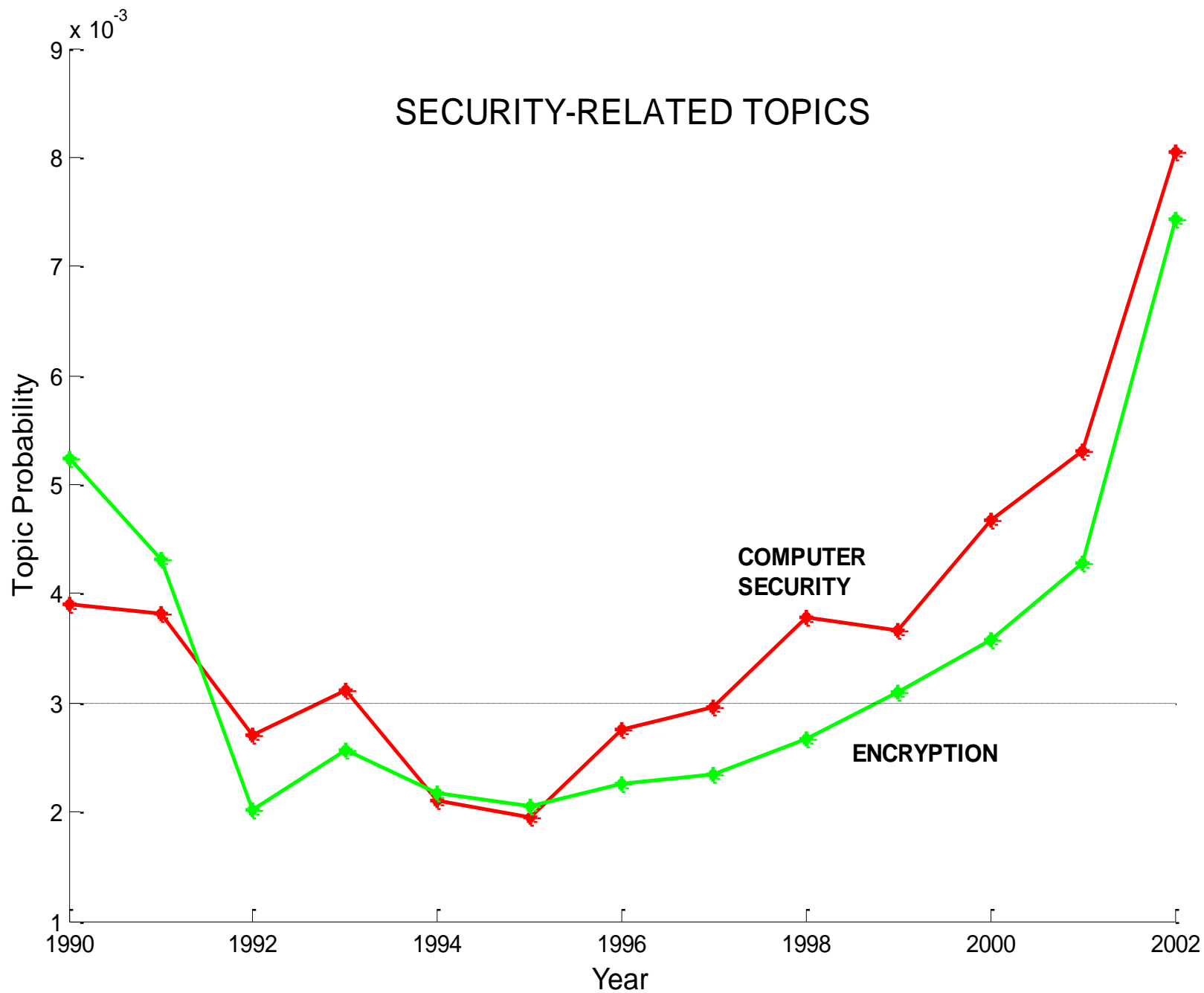
CHANGING TRENDS IN COMPUTER SCIENCE

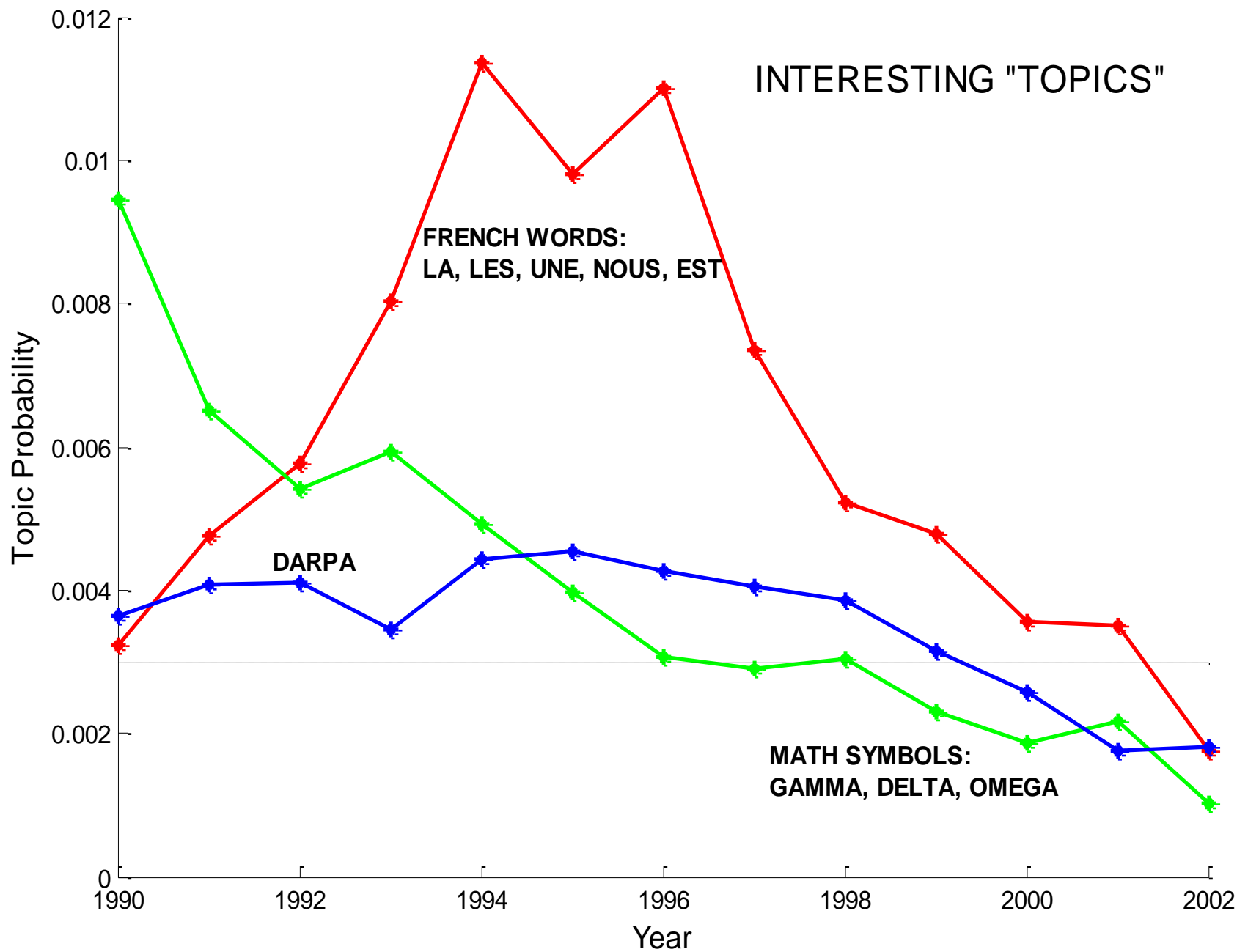


HOT TOPICS: MACHINE LEARNING/DATA MINING

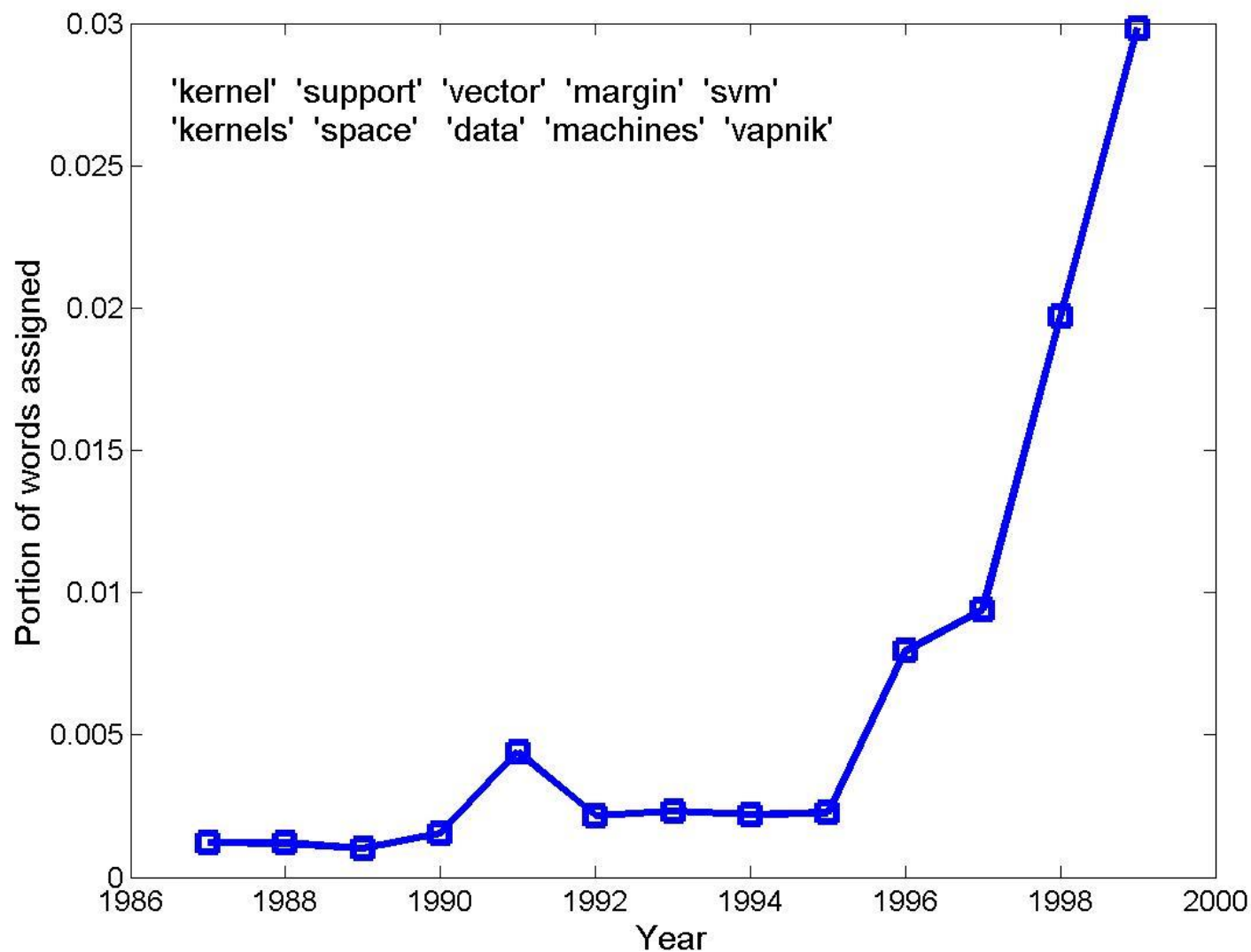




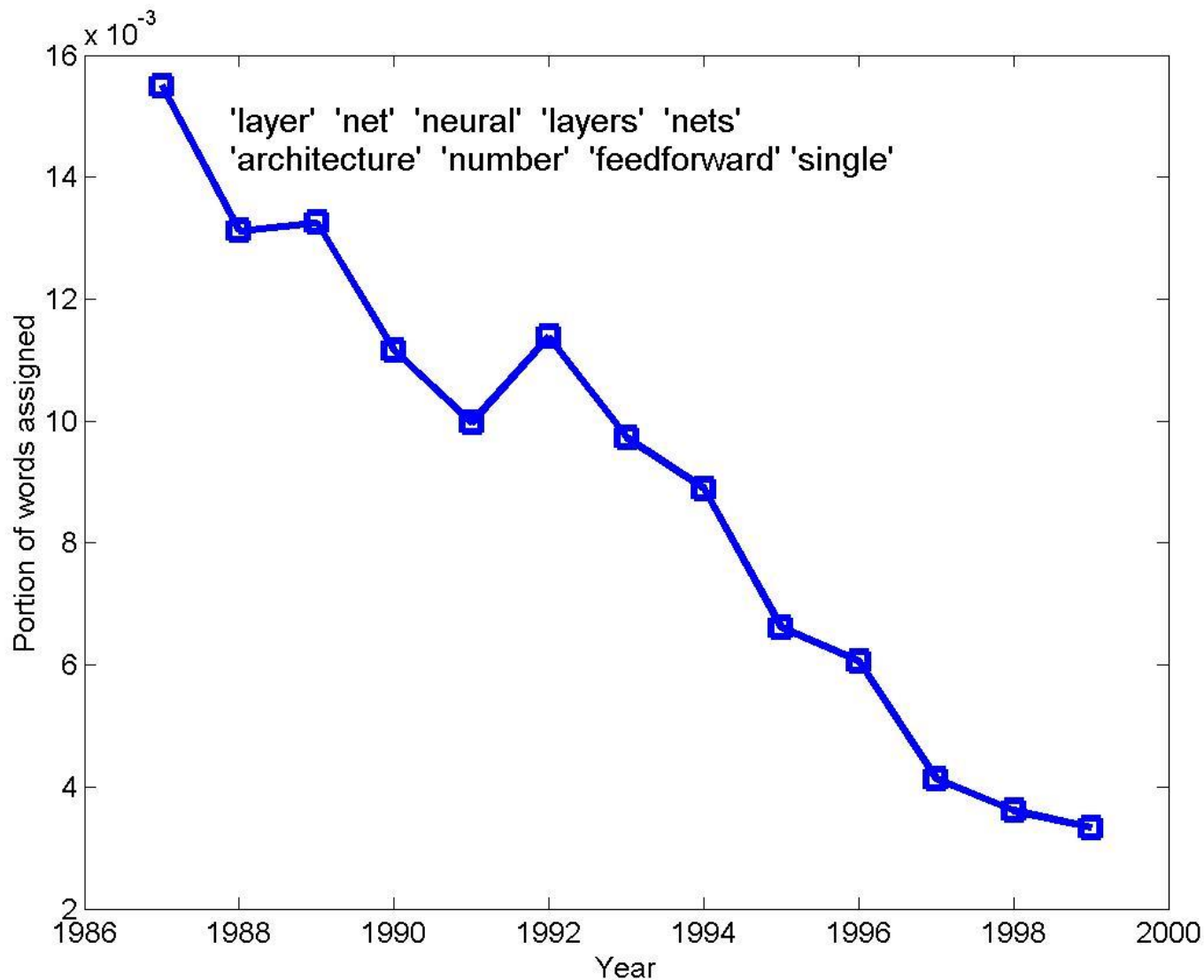




NIPS: SVM Topic



NIPS: neural network topic



Topics over Time for NIPS Proceedings

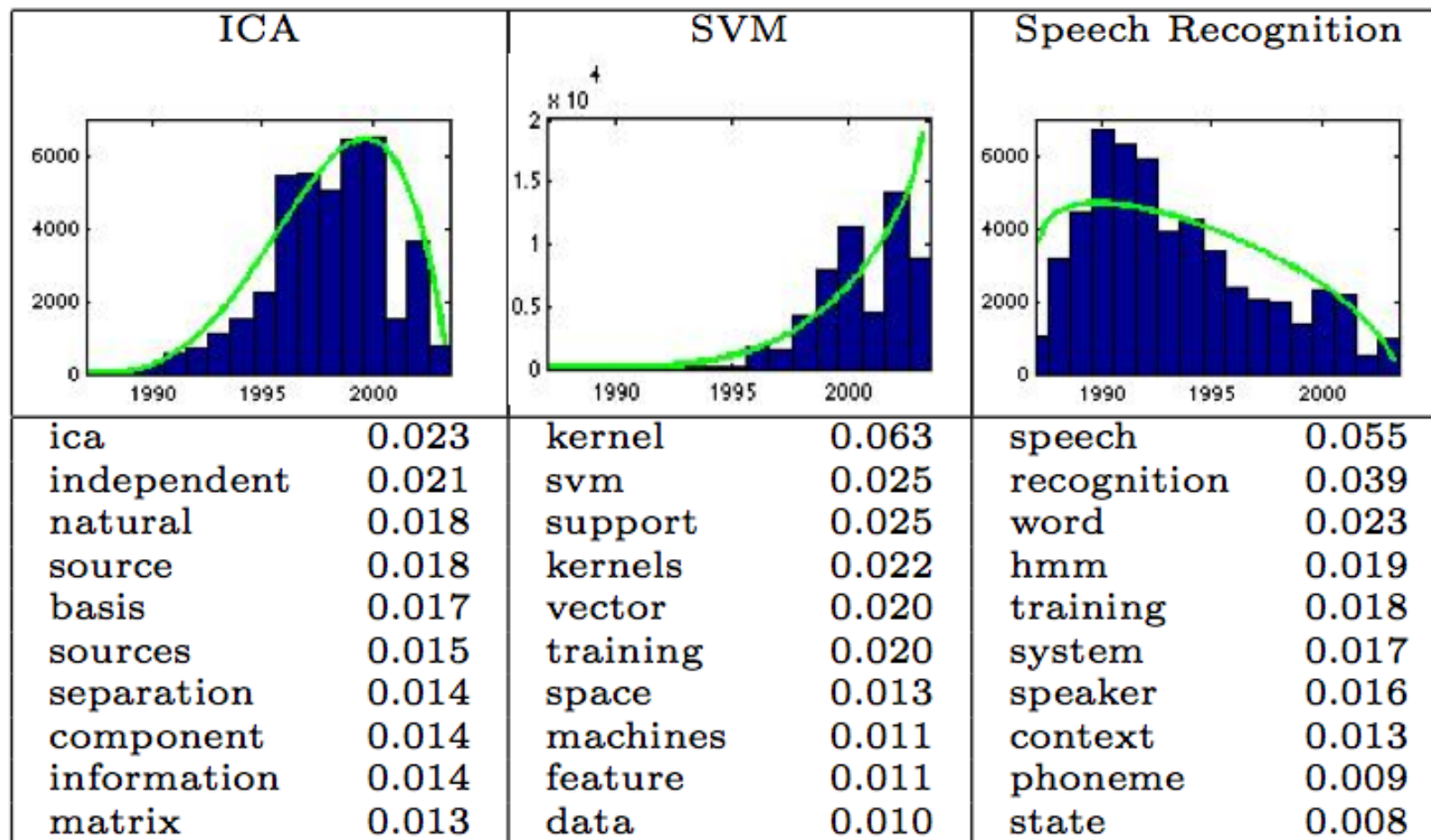


Figure courtesy of Xuerie Wang and Andrew McCallum, U Mass Amherst

All NIPS Topics over Time

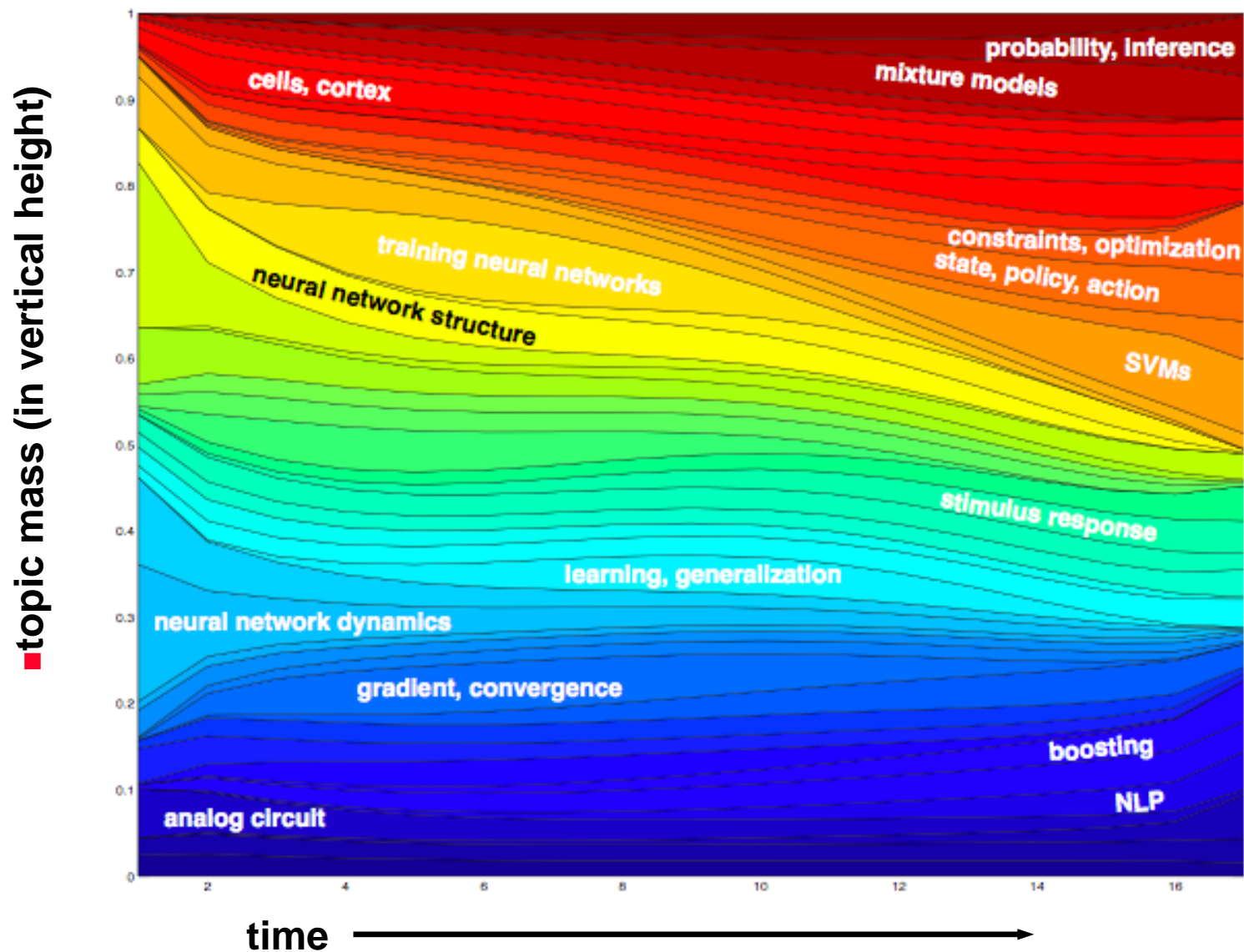


Figure courtesy of Xuerie Wang and Andrew McCallum, U Mass Amherst

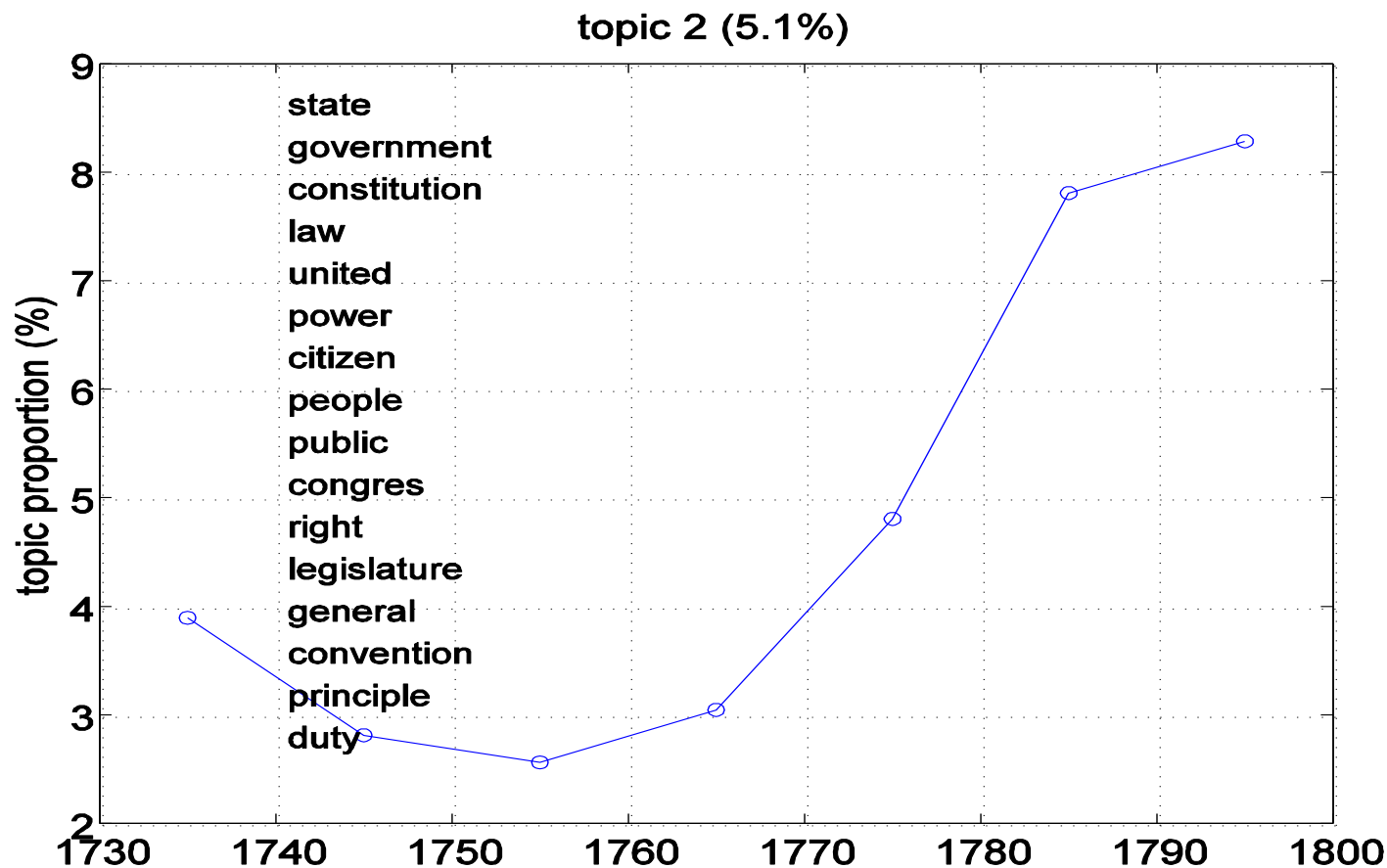
Topics from the Pennsylvania Gazette



Size	Most likely words in topic
6%	away reward servant named feet jacket high paid hair coat run inches master
5%	state government constitution law united power citizen people public congress
5%	good house acre sold land meadow mile premise plantation stone mill dwelling
4%	silk cotton ditto white black linen cloth women blue worsted fine thread plain
2%	church life god society great friend christian good virtue religion minister rev

(from Dave Newman and Sharon Block, UC Irvine)

Historical Trends in Pennsylvania Gazette Data



Topics from DNA Microarray Literature

49,000 PubMed abstracts related to DNA Microarrays

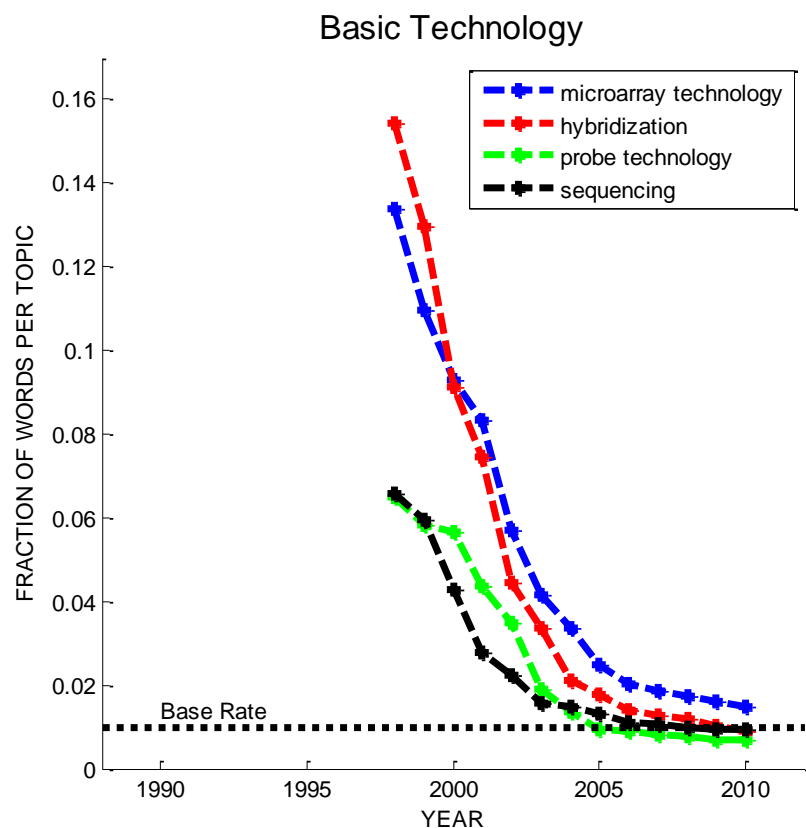
Displayed below are the top 5 highest probability words for 5 selected topics

Microarray Chip Technology	Classification Methods	Databases and Annotation	Regulatory Networks	Cancer
detection	classification	databases	network	patient
surface	selection	tool	regulatory	tumor
fluorescence	cancer	annotation	pathway	cancer
hybridization	algorithm	data set	interaction	survival
array	feature	web	transcriptional	prognostic

From basic technology to applications

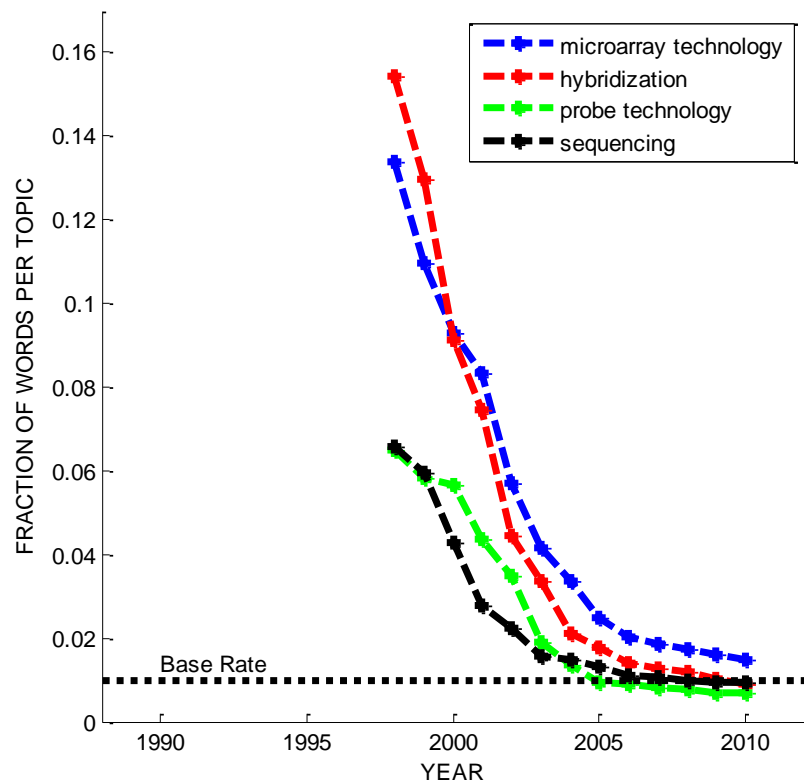


Technology v. Application Topics

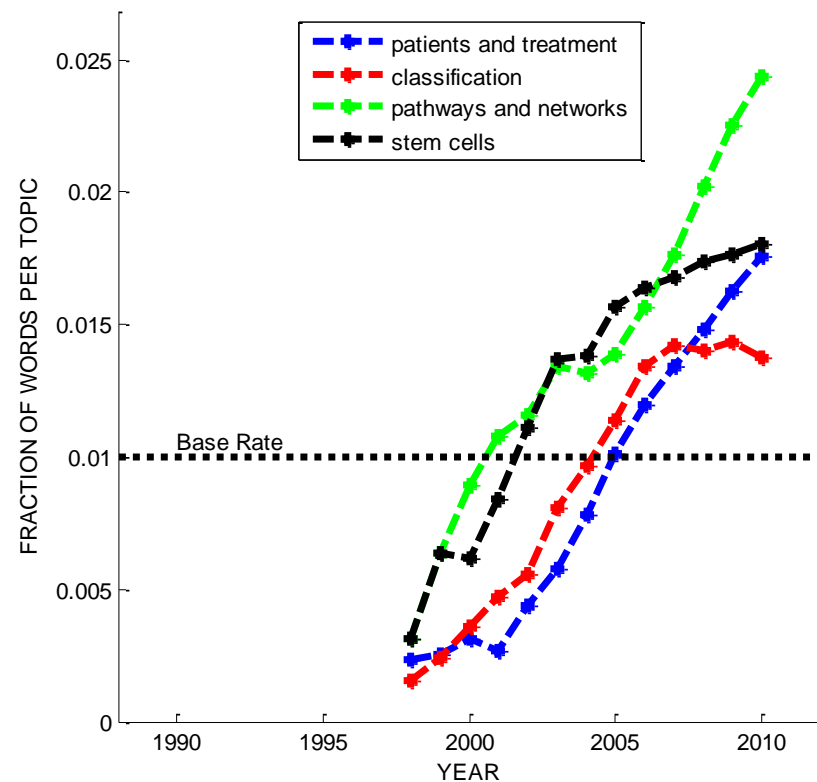


Technology v. Application Topics

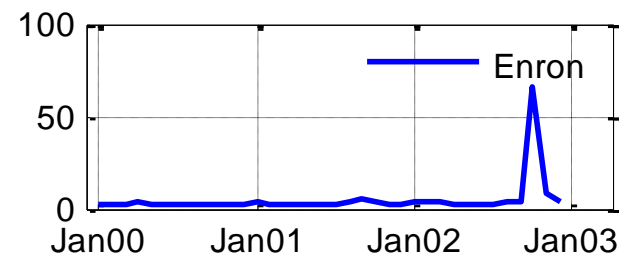
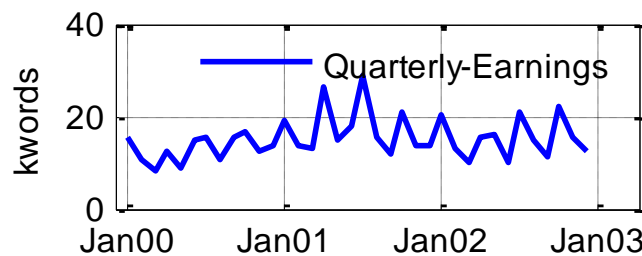
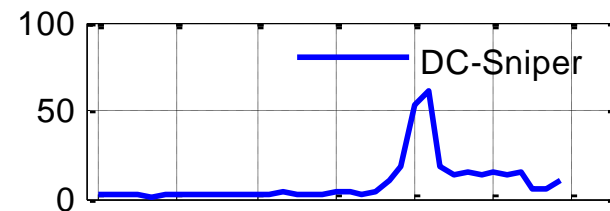
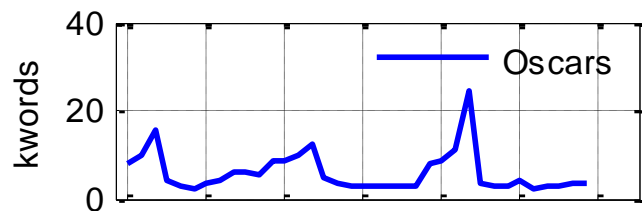
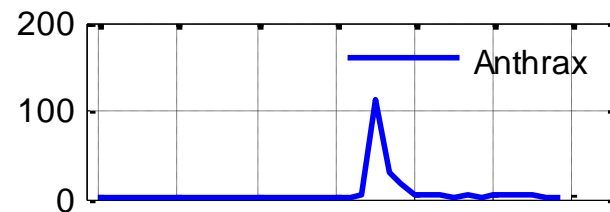
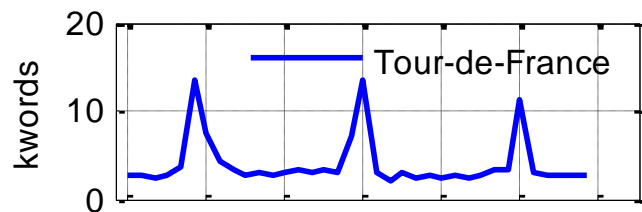
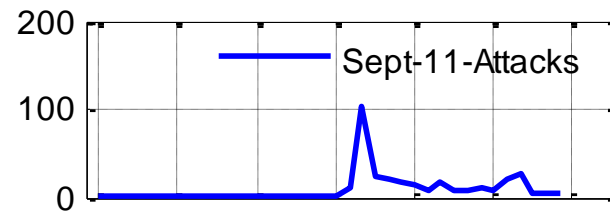
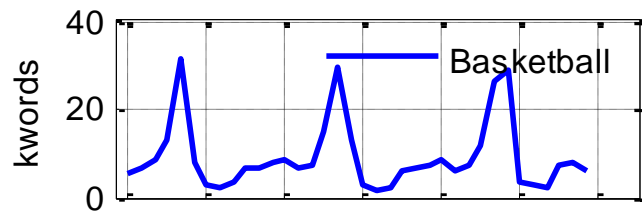
Basic Technology



Applications

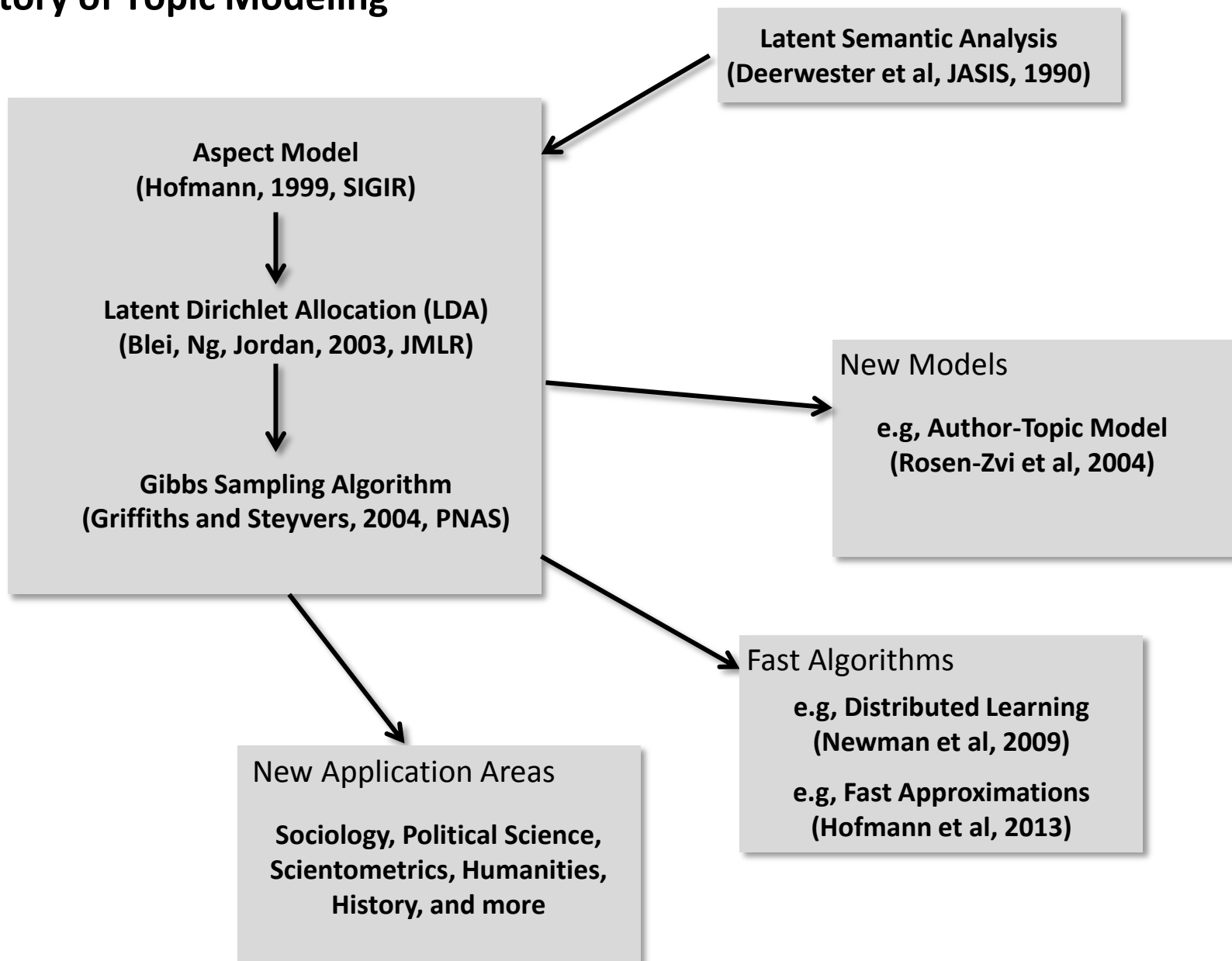


Topic Trends (New York Times articles)

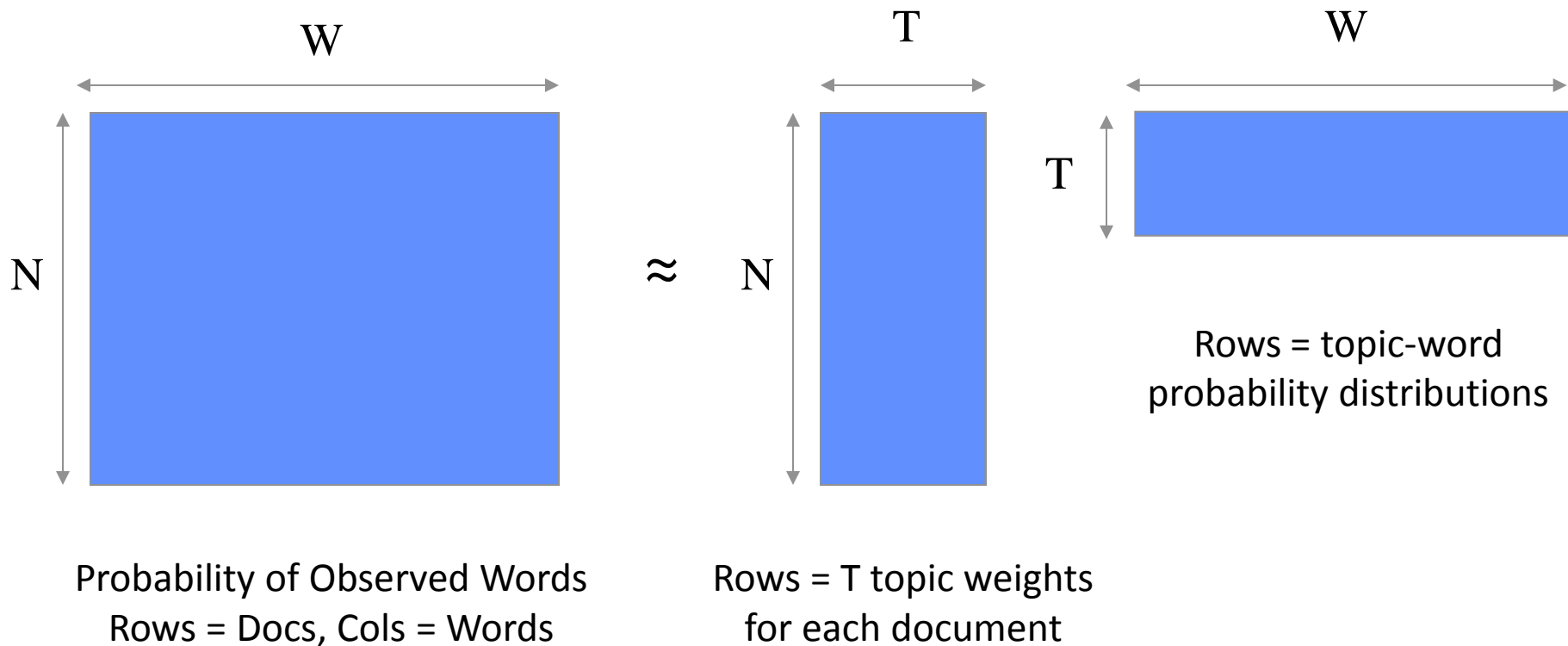


Topic Models and Related Approaches

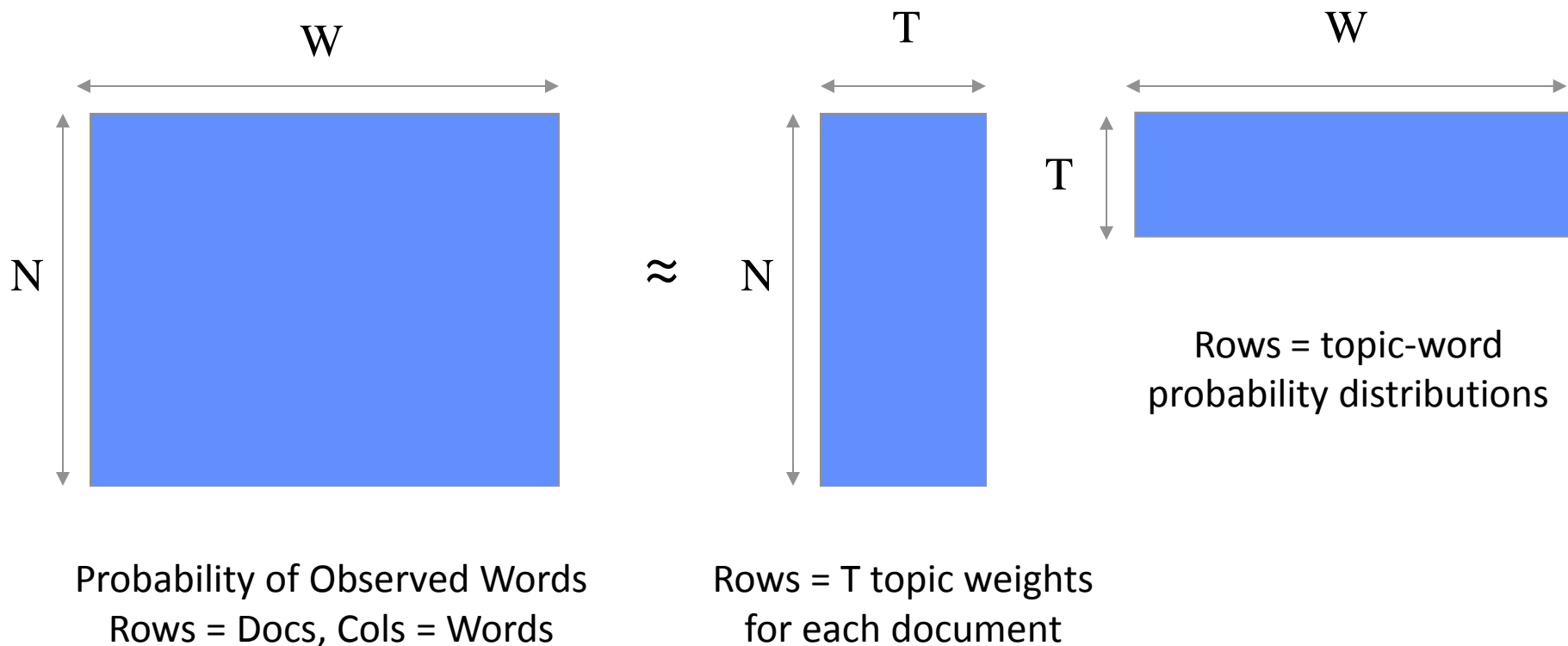
A Brief History of Topic Modeling



Topics as Matrix Factorization



Topics as Matrix Factorization



Directly analogous to principal components and factor models:

data matrix \approx weights * basis functions

Clusters v. Topics

Original Document

Hidden Markov Models in Molecular Biology: New Algorithms and Applications

Pierre Baldi, Yves C Hauvin, Tim Hunkapiller, Marcella A. McClure

Hidden Markov Models (HMMs) can be applied to several important problems in molecular biology. We introduce a new convergent learning algorithm for HMMs that, unlike the classical Baum-Welch algorithm is smooth and can be applied on-line or in batch mode, with or without the usual Viterbi most likely path approximation. Left-right HMMs with insertion and deletion states are then trained to represent several protein families including immunoglobulins and kinases. In all cases, the models derived capture all the important statistical properties of the families and can be used efficiently in a number of important tasks such as multiple alignment, motif detection, and classification.

Clusters v. Topics

Original Document

Hidden Markov Models in Molecular Biology: New Algorithms and Applications

Pierre Baldi, Yves C Hauvin, Tim Hunkapiller, Marcella A. McClure

Hidden Markov Models (HMMs) can be applied to several important problems in molecular biology. We introduce a new convergent learning algorithm for HMMs that, unlike the classical Baum-Welch algorithm is smooth and can be applied on-line or in batch mode, with or without the usual Viterbi most likely path approximation. Left-right HMMs with insertion and deletion states are then trained to represent several protein families including immunoglobulins and kinases. In all cases, the models derived capture all the important statistical properties of the families and can be used efficiently in a number of important tasks such as multiple alignment, motif detection, and classification.

One Cluster

[cluster 88]

model data
models time
neural figure state
learning set
parameters
network
probability
number networks
training function
system algorithm
hidden markov

Clusters v. Topics

Original Document

Hidden Markov Models in Molecular Biology: New Algorithms and Applications

Pierre Baldi, Yves C Hauvin, Tim Hunkapiller, Marcella A. McClure

Hidden Markov Models (HMMs) can be applied to several important problems in molecular biology. We introduce a new convergent learning algorithm for HMMs that, unlike the classical Baum-Welch algorithm is smooth and can be applied on-line or in batch mode, with or without the usual Viterbi most likely path approximation. Left-right HMMs with insertion and deletion states are then trained to represent several protein families including immunoglobulins and kinases. In all cases, the models derived capture all the important statistical properties of the families and can be used efficiently in a number of important tasks such as multiple alignment, motif detection, and classification.

One Cluster

[cluster 88]
model data
models time
neural figure state
learning set
parameters
network
probability
number networks
training function
system algorithm
hidden markov

Multiple Topics

[topic 10] state hmm markov
sequence models hidden
states probabilities sequences
parameters transition
probability training hmms
hybrid model likelihood
modeling

[topic 37] genetic structure
chain protein population
region algorithms human
mouse selection fitness
proteins search evolution
generation function sequence
sequences genes

Methodological Issues

- Selection of smoothing parameters/priors
- Convergence of the sampler?
- Consistency of results across sampling runs?
- Selecting the numbers of topics
- Publicly available code?
 - Mallet from U Mass
 - See also David Blei's Web page

Methods for Evaluating Topic Models

- Human inspection
 - e.g., coherence of high probability words
- Log Probability of Test Documents
 - Better models assign higher probability to unseen documents
- Performance on specific tasks
 - Information retrieval
 - Document classification

Extensions to Topic Models

Adding Metadata to Topic Models?

- **Topic models with metadata? With time? With spatial information?**
 - Yes: can write down simple generative models and then “invert”

- **Example: Probabilistic Pseudocode for Author-Topic Model**

For each document in our corpus

For each word in our document

Randomly select an author of the document

Sample a topic from $P(\text{topics} \mid \text{author})$

Given the topic, sample a word from $P(\text{words} \mid \text{topic})$

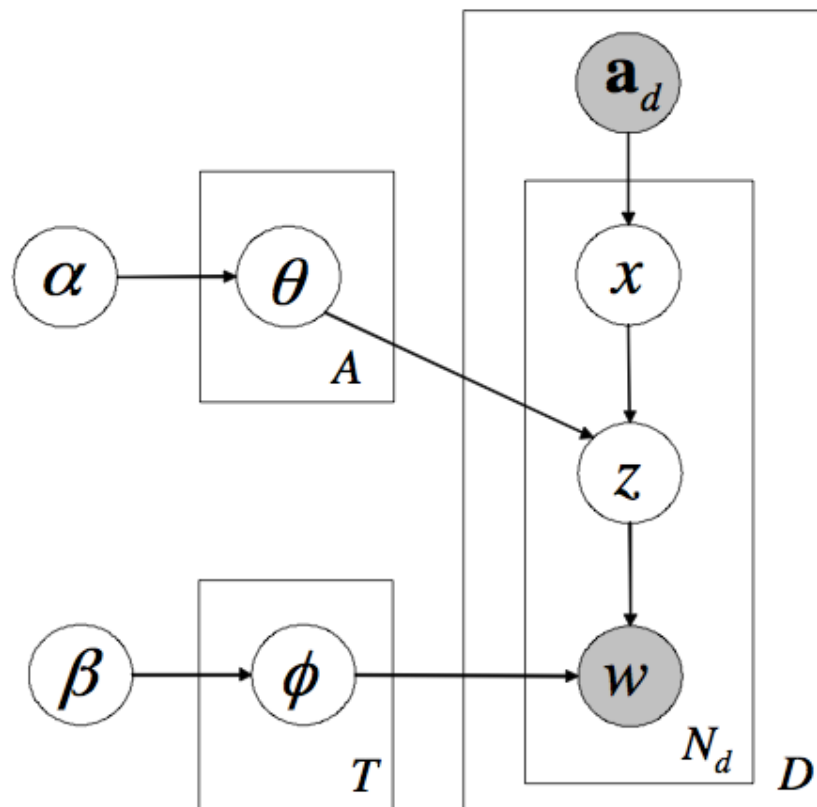
End

End

- **We can write down many interesting models in a similar way**
 - Much harder to do inference, i.e., “invert” these models given data

The Author-Topic Model

Steyvers et al, 2004; Rosen-Zvi et al, 2010



Examples of Author-Topic Models

Learned using an author-topic model applied to 20 years of papers from the NIPS conference on machine learning

TOPIC 4	
WORD	PROB.
LIGHT	.0306
RESPONSE	.0282
INTENSITY	.0252
RETINA	.0241
OPTICAL	.0233
KOCH	.0190
BACKGROUND	.0162
CONTRAST	.0145
CENTER	.0124
FEEDBACK	.0118
AUTHOR	PROB.
Koch_C	.0903
Boahen_K	.0320
Skrzypek_J	.0283
Liu_S	.0250
Delbruck_T	.0232
Etienne-C_R	.0210
Bair_W	.0178
Bialek_W	.0133
Yasui_S	.0106
Hsu_K	.0103

TOPIC 13	
WORD	PROB.
RECOGNITION	.0500
CHARACTER	.0334
TANGENT	.0246
CHARACTERS	.0232
DISTANCE	.0197
HANDWRITTEN	.0166
DIGITS	.0154
SEGMENTATION	.0142
DIGIT	.0124
IMAGE	.0111
AUTHOR	PROB.
Simard_P	.0602
Martin_G	.0340
LeCun_Y	.0339
Henderson_D	.0289
Denker_J	.0245
Revow_M	.0206
Rashid_M	.0205
Rumelhart_D	.0185
Sackinger_E	.0181
Flann_N	.0142

TOPIC 28	
WORD	PROB.
KERNEL	.0547
VECTOR	.0293
SUPPORT	.0293
MARGIN	.0239
SVM	.0196
DATA	.0165
SPACE	.0161
KERNELS	.0160
SET	.0146
MACHINES	.0132
AUTHOR	PROB.
Scholkopf_B	.0774
Smola_A	.0685
Vapnik_V	.0487
Burges_C	.0411
Ratsch_G	.0296
Mason_L	.0232
Platt_J	.0225
Cristianini_N	.0179
Laskov_P	.0160
Chapelle_O	.0152

TOPIC 9	
WORD	PROB.
SOURCE	.0389
INDEPENDENT	.0376
SOURCES	.0344
SEPARATION	.0322
INFORMATION	.0319
ICA	.0276
BLIND	.0227
COMPONENT	.0226
SEJNOWSKI	.0224
NATURAL	.0183
AUTHOR	PROB.
Sejnowski_T	.0627
Bell_A	.0378
Yang_H	.0349
Lee_T	.0348
Attias_H	.0290
Parra_L	.0271
Cichocki_A	.0262
Hyvarinen_A	.0242
Amari_S	.0160
Oja_E	.0143

Author-Topic Models for CiteSeer

TOPIC 205		TOPIC 209		TOPIC 289		TOPIC 10	
WORD	PROB.	WORD	PROB.	WORD	PROB.	WORD	PROB.
DATA	0.1563	PROBABILISTIC	0.0778	RETRIEVAL	0.1179	QUERY	0.1848
MINING	0.0674	BAYESIAN	0.0671	TEXT	0.0853	QUERIES	0.1367
ATTRIBUTES	0.0462	PROBABILITY	0.0532	DOCUMENTS	0.0527	INDEX	0.0488
DISCOVERY	0.0401	CARLO	0.0309	INFORMATION	0.0504	DATA	0.0368
ASSOCIATION	0.0335	MONTE	0.0308	DOCUMENT	0.0441	JOIN	0.0260
LARGE	0.0280	DISTRIBUTION	0.0257	CONTENT	0.0242	INDEXING	0.0180
KNOWLEDGE	0.0260	INFERENCE	0.0253	INDEXING	0.0205	PROCESSING	0.0113
DATABASES	0.0210	PROBABILITIES	0.0253	RELEVANCE	0.0159	AGGREGATE	0.0110
ATTRIBUTE	0.0188	CONDITIONAL	0.0229	COLLECTION	0.0146	ACCESS	0.0102
DATASETS	0.0165	PRIOR	0.0219	RELEVANT	0.0136	PRESENT	0.0095
AUTHOR	PROB.	AUTHOR	PROB.	AUTHOR	PROB.	AUTHOR	PROB.
Han_J	0.0196	Friedman_N	0.0094	Oard_D	0.0110	Suciu_D	0.0102
Rastogi_R	0.0094	Heckerman_D	0.0067	Croft_W	0.0056	Naughton_J	0.0095
Zaki_M	0.0084	Ghahramani_Z	0.0062	Jones_K	0.0053	Levy_A	0.0071
Shim_K	0.0077	Koller_D	0.0062	Schauble_P	0.0051	DeWitt_D	0.0068
Ng_R	0.0060	Jordan_M	0.0059	Voorhees_E	0.0050	Wong_L	0.0067
Liu_B	0.0058	Neal_R	0.0055	Singhal_A	0.0048	Chakrabarti_K	0.0064
Mannila_H	0.0056	Rafer_A	0.0054	Hawking_D	0.0048	Ross_K	0.0061
Brin_S	0.0054	Lukasiewicz_T	0.0053	Merkel_D	0.0042	Hellerstein_J	0.0059
Liu_H	0.0047	Halpern_J	0.0052	Allan_J	0.0040	Lenzerini_M	0.0054
Holder_L	0.0044	Muller_P	0.0048	Doermann_D	0.0039	Moerkotte_G	0.0053

Author-Profiles

Author = Andrew McCallum, U Mass:

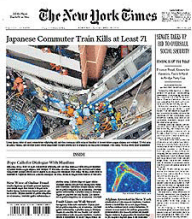
- Topic 1: classification, training, generalization, decision, data,...
- Topic 2: learning, machine, examples, reinforcement, inductive,....
- Topic 3: retrieval, text, document, information, content,...

Author = Hector Garcia-Molina, Stanford:

- Topic 1: query, index, data, join, processing, aggregate....
- Topic 2: transaction, concurrency, copy, permission, distributed....
- Topic 3: source, separation, paper, heterogeneous, merging.....

Author = Jerry Friedman, Stanford:

- Topic 1: regression, estimate, variance, data, series,...
- Topic 2: classification, training, accuracy, decision, data,....
- Topic 3: distance, metric, similarity, measure, nearest,...



Learning Word Combinations

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

Reinforcement Learning Topic with Topical N-grams

LDA

state
learning
policy
action
reinforcement
states
time
optimal
actions
function
algorithm
reward
step
dynamic
control
sutton
rl
decision
algorithms
agent

Topical N-grams (2+)

reinforcement learning
optimal policy
dynamic programming
optimal control
function approximator
prioritized sweeping
finite-state controller
learning system
reinforcement learning RL
function approximators
markov decision problems
markov decision processes
local search
state-action pair
markov decision process
belief states
stochastic policy
action selection
upright position
reinforcement learning methods

Figure courtesy of Xuerie Wang and Andrew McCallum, U Mass Amherst

Support Vector Machine Topic with Topical N-grams

LDA

kernel
linear
vector
support
set
nonlinear
data
algorithm
space
pca
function
problem
margin
vectors
solution
training
svm
kernels
matrix
machines

Topical N-grams (2+)

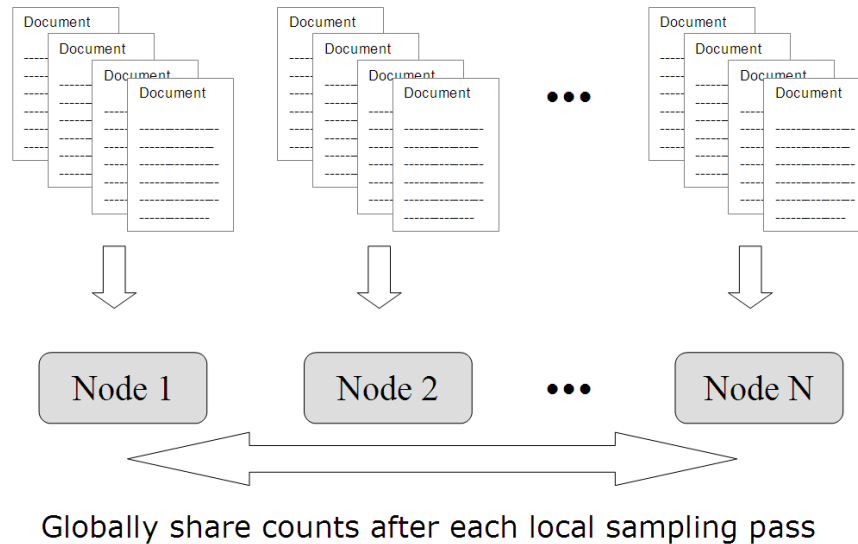
support vectors
test error
support vector machines
training error
feature space
training examples
decision function
cost functions
test inputs
kkt conditions
leave-one-out procedure
soft margin
bayesian transduction
training patterns
training points
maximum margin
strictly convex
regularization operators
base classifiers
convex optimization

Figure courtesy of Xuerie Wang and Andrew McCallum, U Mass Amherst

Scaling to Large Corpora

- Time complexity: linear in number of word tokens and topics
 - But even this can be slow on millions of documents
- Distributed algorithms
 - Distribute documents across multiple processors (Newman et al, 2009)
 - Approximate, but works very well in practice
- Fast sampling tricks
 - Inner sampling loop is computed billions of times
 - Can re-order operations to get order of magnitude speedup (Porteous et al, 2006)
- Stochastic gradient methods
 - Standard: 1 iteration = full sweep through all words in the corpus
 - Stochastic/online: update parameters after every few documents
 - Algorithm can converge even before a single iteration is complete!

Distributed Topic Modeling

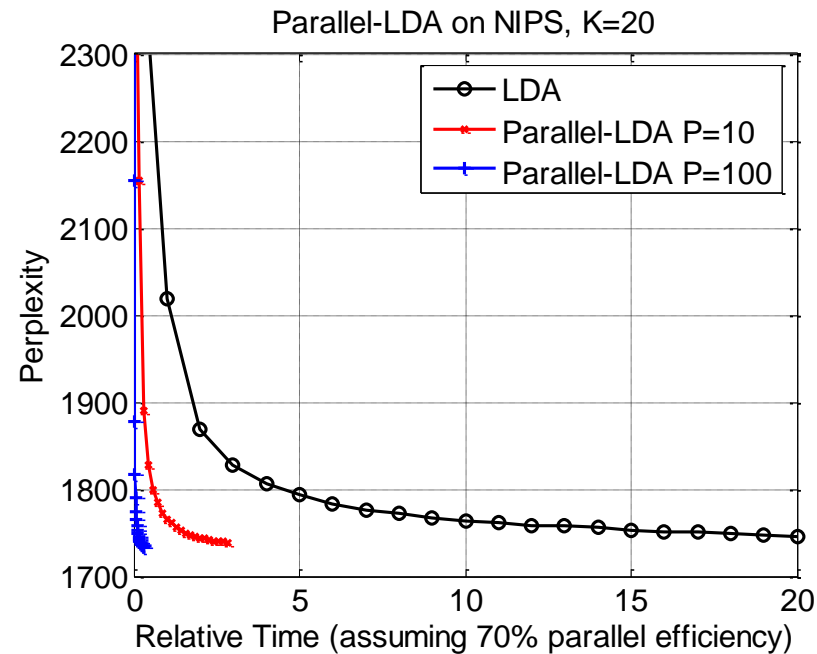
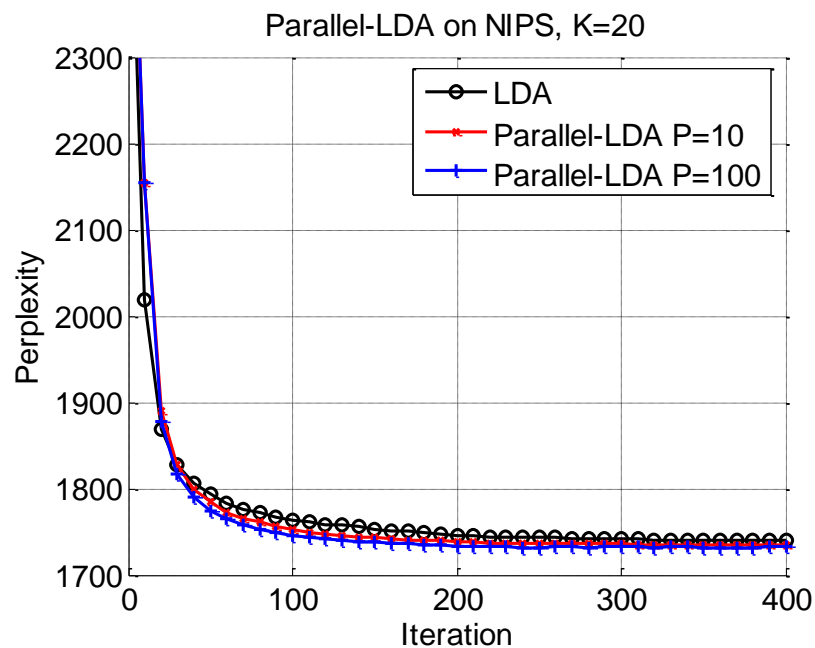


Parallel-LDA [Newman, Asuncion, Smyth, Welling, NIPS 2007, JMLR 2009]

- Each processor performs Gibbs sampling over local set of documents
- At the end of each iteration, all processors combine topic counts to create global model

Parallel-LDA Results

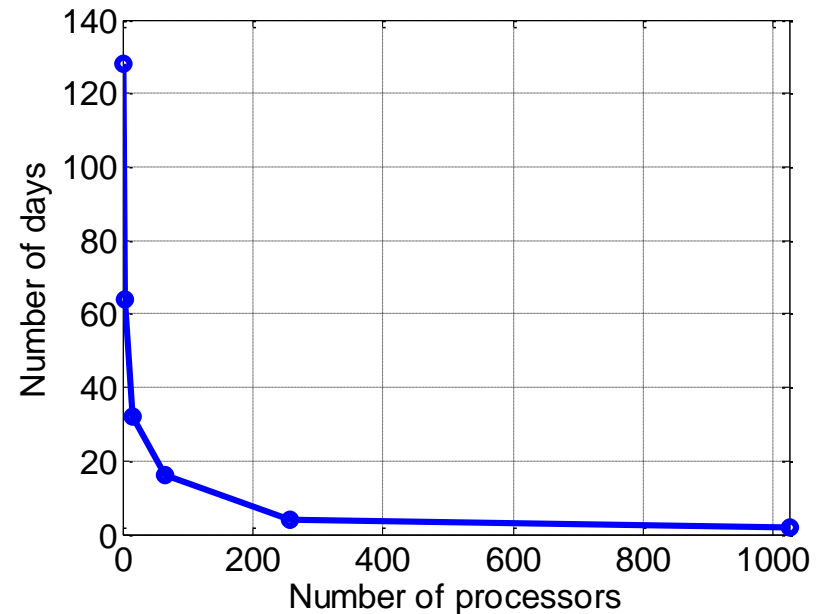
[Newman, Asuncion, Smyth, Welling, NIPS 2007, JMLR 2009]



Experiments with 8 Million Documents



All of MEDLINE
8 million abstracts
1 billion words
2000 topics

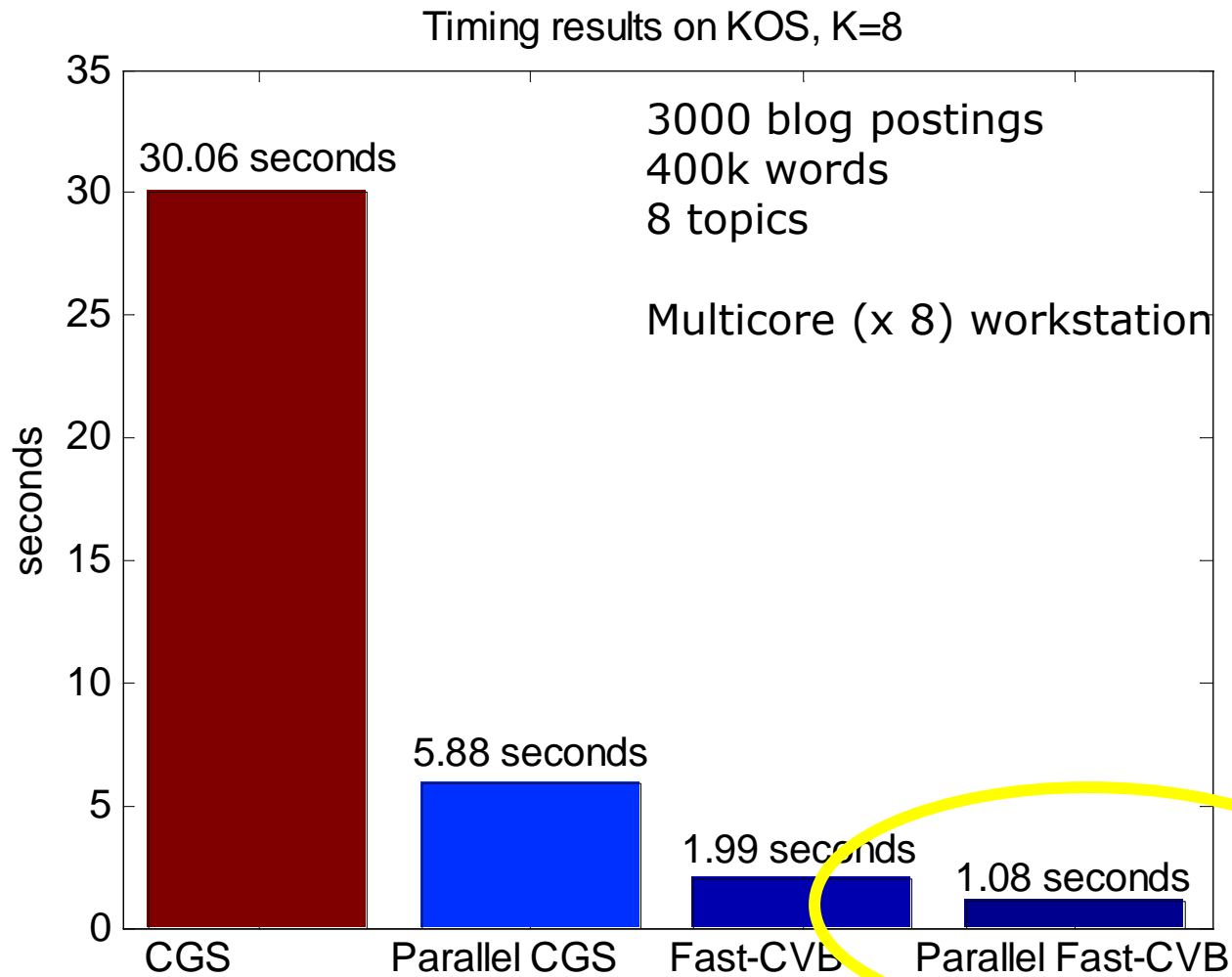


Experiments with 1000 processors
at the San Diego Supercomputing
Center (SDSC)

Time goes from 4 months to a few hours

Real-Time Topic Modeling

Asuncion et al, 2009



Real-Time Topic Modeling of Search Results

TASER

☐ Web ☒ News

Theme-Assisted Search Engine in Real-time

Themes found:

T1: mouse school
cartilage skin cell team
mice worship pyramid
messages
[\[order by\]](#) [\[zoom\]](#)

T2: mouse ces multitouch
show consumer microsoft
electronic scanner vegas
las
[\[order by\]](#) [\[zoom\]](#)

T3: model cell stem gene
blood disorder inherited
children embryonic lines
[\[order by\]](#) [\[zoom\]](#)

T4: mouse haifa prweb
disney big amusement
park january mickey
home
[\[order by\]](#) [\[zoom\]](#)

T5: mouse click computer
control surface coming
xbox super cost world
[\[order by\]](#) [\[zoom\]](#)

T6: mouse touch
microsoft gestures based
looking research window
top june
[\[order by\]](#) [\[zoom\]](#)

T7: mouse magic
microsoft apple touch air
multi function finally
sensitive
[\[order by\]](#) [\[zoom\]](#)

[The \\$130 TRON mouse and mouse surface: beautiful overkill](#)

It's going to take a lot of convincing to get us to give our blessing to a mouse and surface that cost \$130 . It is covered with TRON branding and style, so fans of the franchise have a little extra reason to pick it up, but that's a very expensive mouse, accessory, collectible, or whatever it is. The packaging is suitably high-end, as the mouse and the surface are displayed in heavy cardboard ...

http://arstechnica.com/gaming/reviews/2011/01/the-130-tron-mouse-and-mouse-surface-beautiful-overkill.ars?utm_source=rss&utm_medium=m=rss&utm_campaign=rss -- 0kb -- 2011/01/10


[Similar Pages](#)
[Thematic Markup](#)

[Microsoft Touch Mouse Announced for June Release](#)

Microsoft is looking to enhance your Window 7 navigating experience with its new gesture-based Touch Mouse.

<http://www.pcmag.com/article2/0,2817,2375277,00.asp?k=PCRSS03069TX1K0001121> -- 0kb -- 2011/01/06


[Similar Pages](#)
[Thematic Markup](#)

[Microsoft Touch Mouse due midyear](#)

Microsoft has announced the Touch Mouse, the company's answer to Apple's Magic Mouse. Designed for use with Windows 7, it uses capacitive multitouch technology.

<http://www.itwire.com/business-it-news/technology/44261-microsoft-touch-mouse-due-midyear> -- 0kb -- 2011/01/09


[Similar Pages](#)
[Thematic Markup](#)

[See how they run \(amok\): Mouse calls are up this year](#)

Big snowdrifts and recent mild winters are being blamed for more mouse calls from Twin Cities residents. The best mouse trap? Prevention.

<http://www.startribune.com/local/113177959.html> -- 0kb -- 2011/01/10


[Similar Pages](#)
[Thematic Markup](#)

[Julia Gillard's 'Mouse Pack' and other dumb stuff](#)

Julia Gillard is "not well informed" and is part of a Melbourne-based gang called The Mouse Pack, ,while Tony Abbott has "good manners", is "formidable" and possessed of a "first-class mind".

<http://www.brisbanetimes.com.au/entertainment/books/julia-gillards-mouse-pack-and-other-dumb-stuff-20110110-19kbe.html> -- 0kb -- 2011/01/10

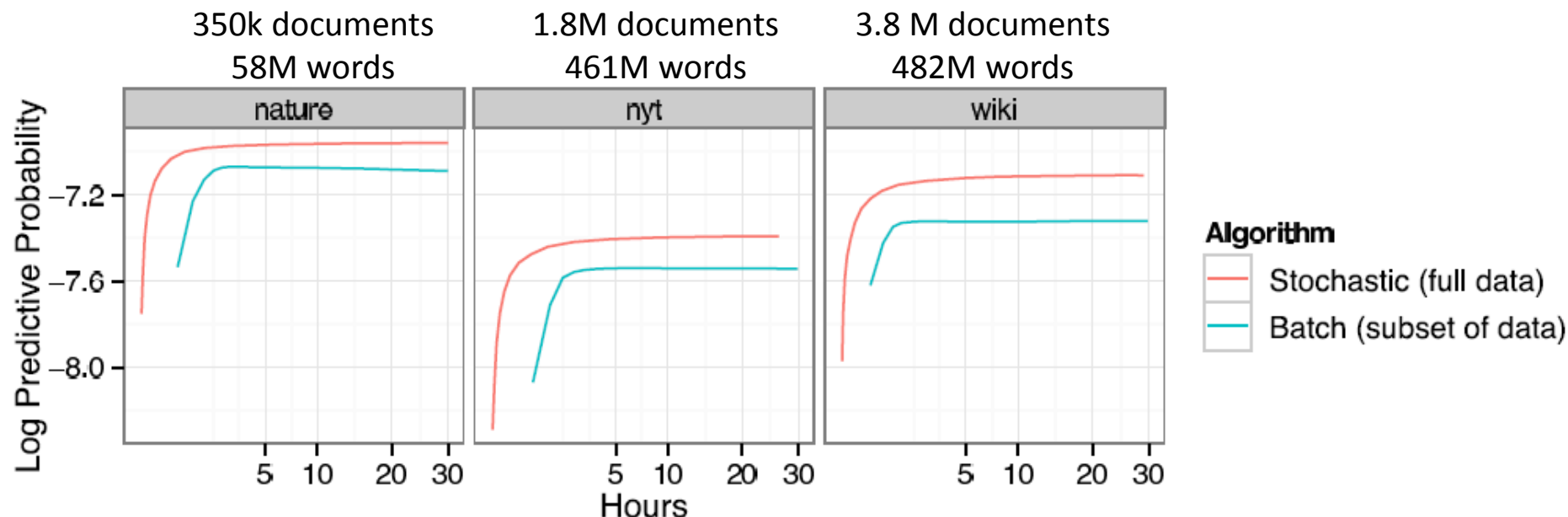

[Similar Pages](#)
[Thematic Markup](#)

Topic Mixtures

Learned
Topics

Fast Topic Learning with Stochastic Methods

(From Hoffmann et al, 2013)



Y-axis = predictive log probability on test documents: higher is better

X-axis = time on a log-scale

Stochastic algorithm learns a better model in minutes than batch algorithm does in hours

Applications of Topic Modeling

Application: Calit2 Research Browser

- System crawled UCI/UCSD faculty websites
- Browser built on topic model learned from faculty papers
- Query-answering = computation of conditional probabilities

Topic Modeling of Researchers and Research at UCSD and UCI

After automatically collecting 12,000 publications from 460 UCSD and UCI faculty, we used our probabilistic topic model to characterize the nature of each researcher's work and find researchers with similar interests.

Researchers (more...)



Topics (more...)



one
topic

neural network models and algorithms

network input unit learning output training pattern neural_network representation weight grammar class structure connectionist learn net performance simple prediction connection elman classes experiment features architecture modeling training_set recognition initial vowel mit_press chaotic epoch mapping rules dynamical feature label

Other researchers in neural network models and algorithms (UCSD, UCI):

(19%) DE SA, VIRGINIA
 (11%) COTTRELL, GARRISON
 (11%) ELMAN, JEFFREY L.
 (5%) MJOLSNESS, ERIC D.
 (4%) BELEW, RICHARD K.
 (4%) YOUSEFIZADEH, HOMAYOUN
 (3%) GRANGER, RICHARD H.
 (3%) BALDI, PIERRE F.
 (2%) WELLING, MAX
 (2%) ABARBANEL, HENRY D.
 (2%) BORK, ALFRED
 (1%) KIBLER, DENNIS F.
 (1%) CHANCE, FRANCES S.
 (1%) TRIESCH, JOCHEN
 (1%) STEYVERS, MARK
 (1%) TODOROV, EMANUEL
 (1%) BATALI, JOHN D.
 (1%) ESKIN, ELEAZAR

most prolific
researchers
for this topic

one
researcher

COTTRELL, GARRISON

COG SCI
DIVISION OF SOCIAL SCIENCES
UCSD

email: gary@ucsd.edu

publications URL: <http://www-cse.ucsd.edu/users/gary/> (53 papers collected)

Research topics:

topics this
researcher
works on

(28%) [[neural network models and algorithms](#)] network input unit learning output
(14%) [[image and vision modeling](#)] image images face recognition pixel features
(7%) [[information retrieval](#)] query retrieval feature image user document system
(7%) [[cognitive experiments](#)] subject word memory experiment task participant
(4%) [[data analysis](#)] data correlation analysis sample average estimates parameter
(4%) [[cognition and EEG](#)] word erp processing brain sentence language semantics
(4%) [[language modeling](#)] language verb theory sense structure word meaning
(4%) [[human learning and development](#)] children word development learning age
(3%) [[modeling](#)] model simulation parameter modeling process

Related researchers ([UCSD](#), [UCI](#)) :

other researchers
with similar topical
interests

(0.9) [DE SA, VIRGINIA](#)
(0.7) [ELMAN, JEFFREY L.](#)
(0.6) [MJOLSNESS, ERIC D.](#)
(0.5) [BELONGIE, SERGE J.](#)
(0.5) [VASCONCELOS, NUNO](#)
(0.5) [BELEW, RICHARD K.](#)
(0.5) [TRIESCH, JOCHEN](#)
(0.4) [KREIGMAN, DAVID](#)
(0.4) [WELLING, MAX](#)
(0.3) [STEYVERS, MARK](#)
(0.3) [ESKIN, ELEAZAR](#)
(0.3) [KIRSH, DAVID J.](#)
(0.3) [BROWN, SCOTT D.](#)
(0.3) [GRANGER, RICHARD H.](#)
(0.3) [JAIN, RAMESH CHANDRA](#)

VII: Aligned topics in English (AJP), in dark gray extending upwards, and German (Hermes), in light gray extending downwards. Lines are at 1920, 1940, 1960, 1980 and 2000.

case languages language example nominative english genitive means cases object dative subject grammatical accusative article finnish like noun feminine

sanskrit self vedic mind buddhist three nature buddhism hindu type means through buddha reality meditation practice veda seeing compound

linguistics context words use specific type example particular meanings linguistic same terms different lexical analysis sound function literal within

rhetoric demosthenes cicero against speeches isocrates speech orator others style public along should teacher important private orators ten funeral law legal court case criminal laws under right rule principle will act crime civil cases decision parties judge courts

she daughter wife married mother sister marriage husband child children became queen gave bore woman named father birth herself

women woman men man female young male considered famous seen than sexual feminine beautiful role make reports sex parallel

sprachen sprache kasus beispiel fall akkusativ haus genitiv grammatik deutschen dativ latein beispiele kennt deutsch person ausdrücken präposition regel

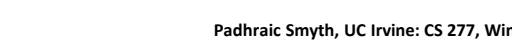
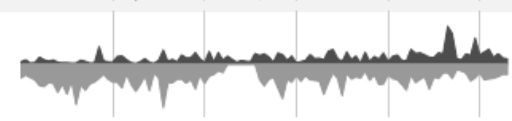
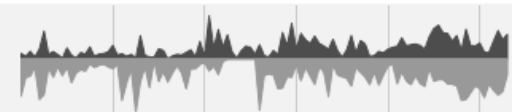
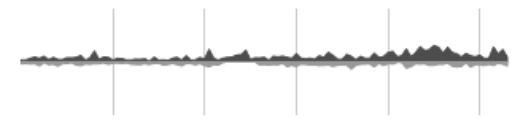
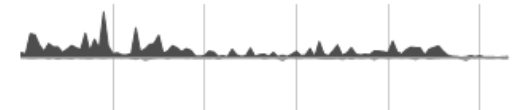
sanskrit drei vier buddhismus indischen lehre bedeutet prakriti purusha dukkha wesen hinduismus zustand buddhistischen buddha meditation yoga

bedeutung zwischen beispiel verschiedene linguistik also beispielsweise wort verschiedenen zeichen lassen wörtern kontext sprachwissenschaft englisch gegenstand siehe unterscheidung spricht

cicero rede rhetorik reden redner allem ersten demosthenes ciceros bedeutung tod prozess seinem letzte wobei damit rhetor auseinandersetzung meidias recht gesetz wenn also lat law keine ohne wegen liegt tat non sog strafe gemäß deutschen grundsatz vertrag deutschland

tochter mutter ihr frau schwester vater ihrem ihre ihrer ihres heiratete ehe gattin kind verheiratet ihren kinder ehfrau geboren

frauen männer gehört mädchen männlichen frau anderen personen jungen junge ihre weiblichen mann männern bestand tritt bringen weiblicher männliche



From Nguyen et al., Modeling topic control to detect influence in conversations,
Machine Learning Journal, 2013

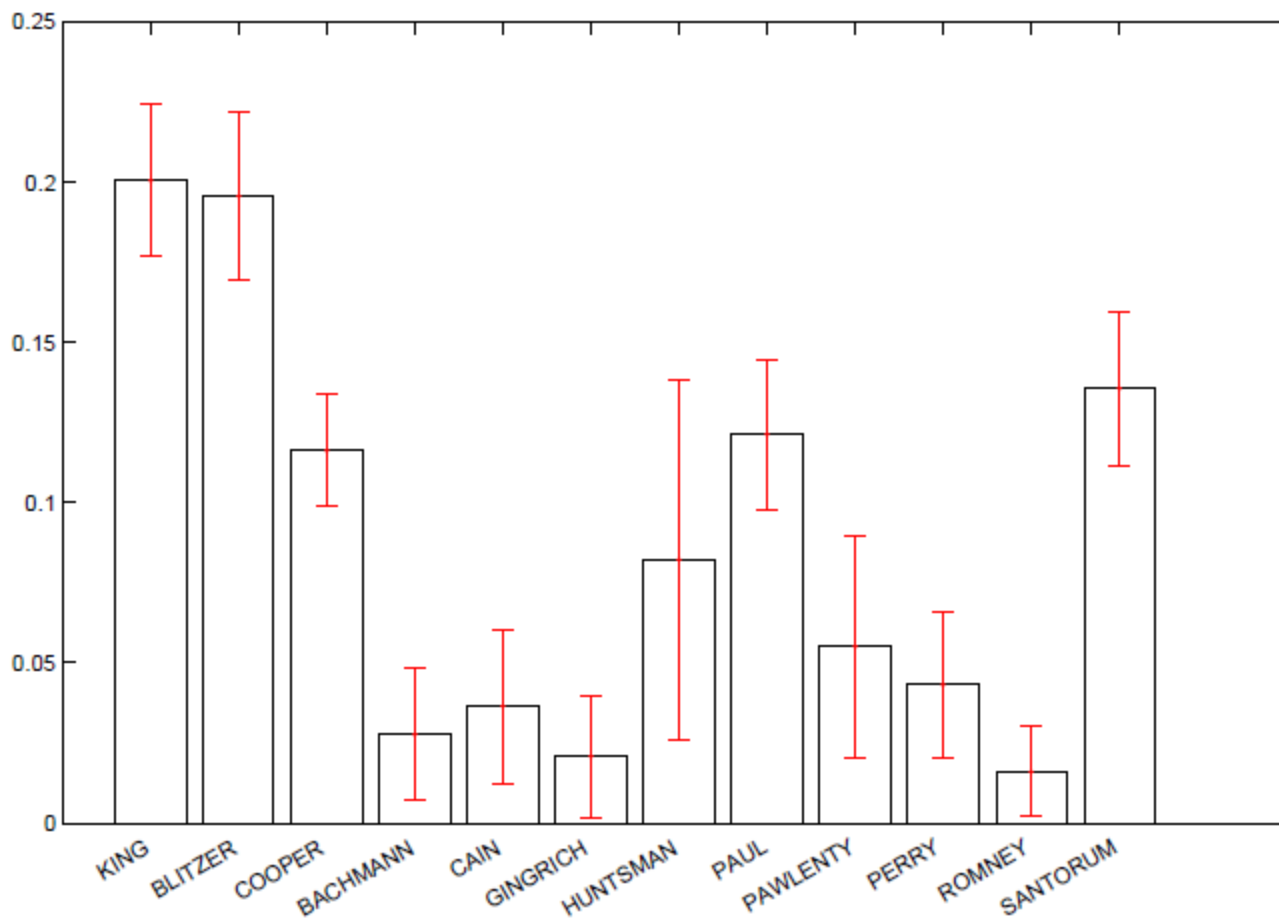


Fig. 6: Topic shift tendency π of speakers in the 2012 Republican Primary Debates (larger means greater tendency). KING, BLITZER and COOPER are moderators in these debates; the rest are candidates.

Analyzing Psychotherapy Transcripts

(joint work with Mark Steyvers, Cognitive Science, UCI)

Each transcript treated as a document

“Subject” and “Symptom” labels manually assigned to some transcripts

Can build in sequential dependence, talk turns, etc

Label Type	Label	High Probability Words
Subject	Medications	dose, mg, medicine, need, wellbutrin, lamictal, mood, sleep, medicines, prescription, medication, use, xanax, klonopin, lexapro, morning, blood, zoloft
	Spousal relationships	wife, married, marriage, home, husband, children, relationship, situation, talked, love, guy, problems, course, work, accept, divorce, couple, meet, girl, attitude, happy, type
Symptom	Fatigue	sleep, blood, depression, energy, tired, low, thyroid, pressure, hormone, fatigue, problems, pituitary, months, growth, level, exercise, treatment, depakote, pharmacy
	Depression	depressed, depression, doctor, sleep, worse, pain, problems, upset, tired, afraid, care, sad, crying, medication, feels, help, sorry, left, understand, hurt, lexapro, remember

Analyzing Psychotherapy Transcripts

(joint work with Mark Steyvers, Cognitive Science, UCI)

Each transcript treated as a document

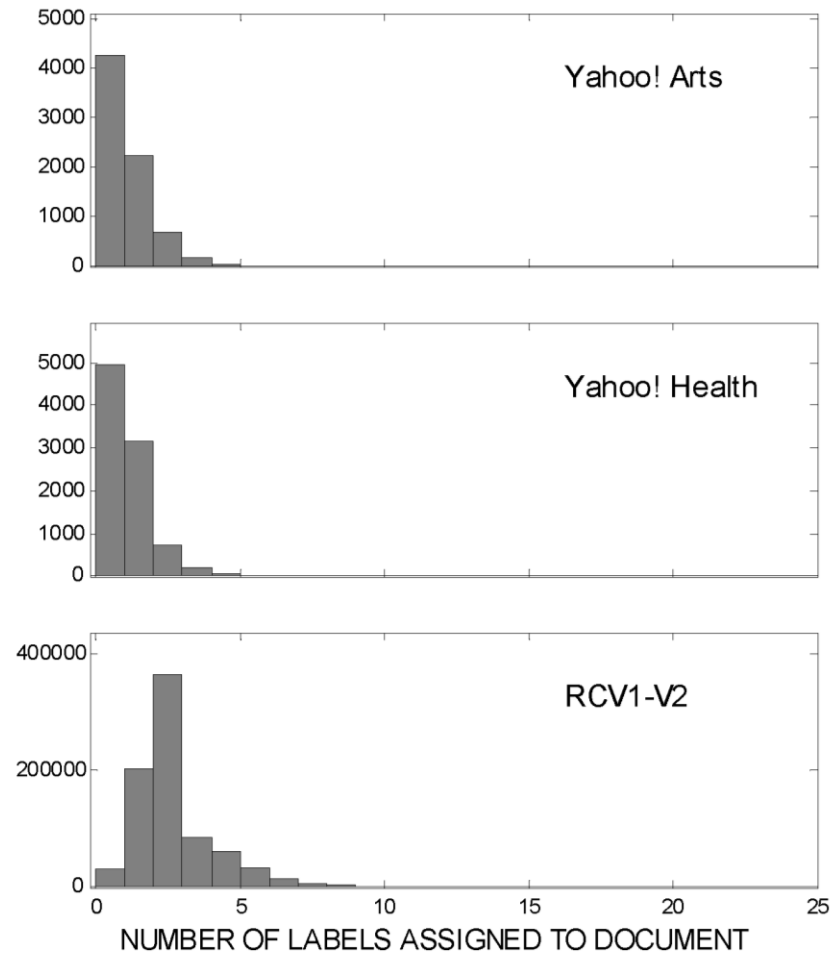
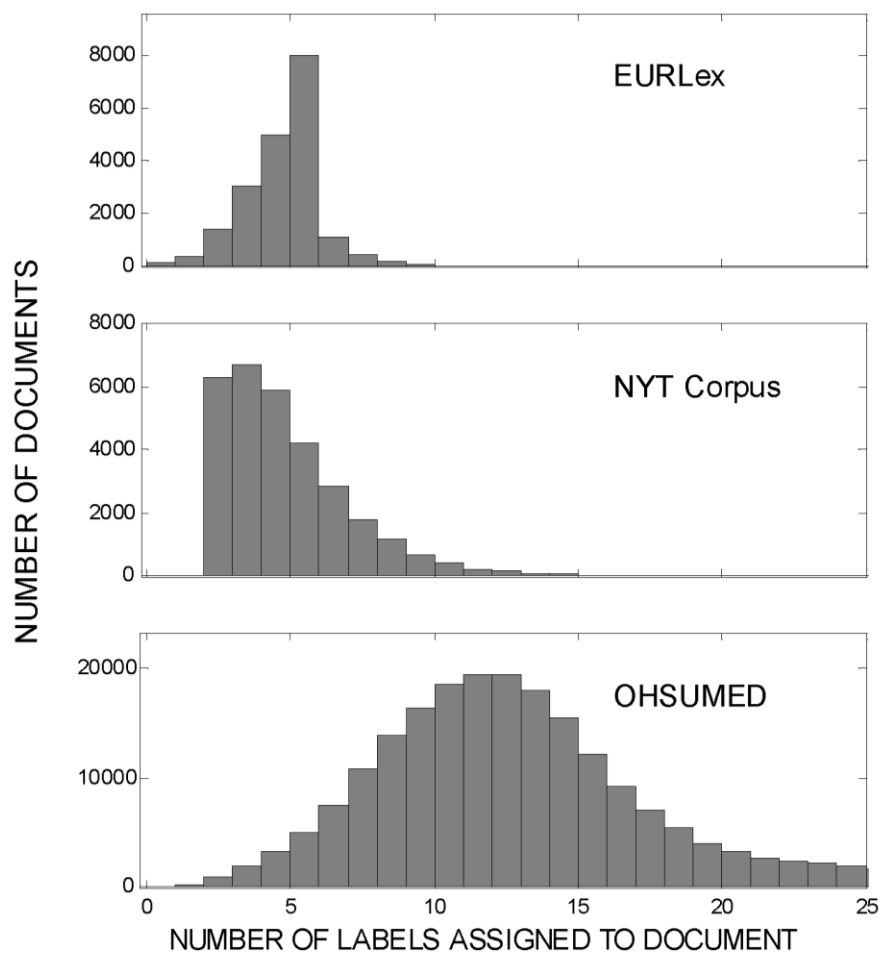
“Subject” and “Symptom” labels manually assigned to some transcripts

Can build in sequential dependence, talk turns, etc

Label Type	Label	High Probability Words	Example Talk Turn Assigned by Model
Subject	Medications	dose, mg, medicine, need, wellbutrin, lamictal, mood, sleep, medicines, prescription, medication, use, xanax, klonopin, lexapro, morning, blood, zoloft	[THERAPIST] so risperdal, geodon, i like geodon, you haven't been on that. so, we got risperdal, geodon and invega.
	Spousal relationships	wife, married, marriage, home, husband, children, relationship, situation, talked, love, guy, problems, course, work, accept, divorce, couple, meet, girl, attitude, happy, type	[PATIENT] see, because by saying like, `` yes, i'll go home with my wife and i'll stay home, " then i reaffirm my way of life and my relationship with my wife.
Symptom	Fatigue	sleep, blood, depression, energy, tired, low, thyroid, pressure, hormone, fatigue, problems, pituitary, months, growth, level, exercise, treatment, depakote, pharmacy	[PATIENT] pretty good, i'm still having such trouble with being able to have energy. my energy level is down and i am just sleepy, sleepy all day long.
	Depression	depressed, depression, doctor, sleep, worse, pain, problems, upset, tired, afraid, care, sad, crying, medication, feels, help, sorry, left, understand, hurt, lexapro, remember	[PATIENT] and it's sort of - i don't know. i just really don't want to be stuck, like i really don't want to wake up one morning and realize that i hate my life completely. that i'm really miserable where i am or where i decided to go.

Topic Models for Supervised Learning with Multilabel Documents

Multilabel Document Data Sets



Real-World Multilabel Data Sets



- Annotated New York Times
 - 1.5 million news articles
 - Thousands of topics
 - Tagged by team of librarians



- Open Directory Project (ODP)
 - Science subtree
 - 10,817 categories
 - 78k Web pages
 - 11-level hierarchy



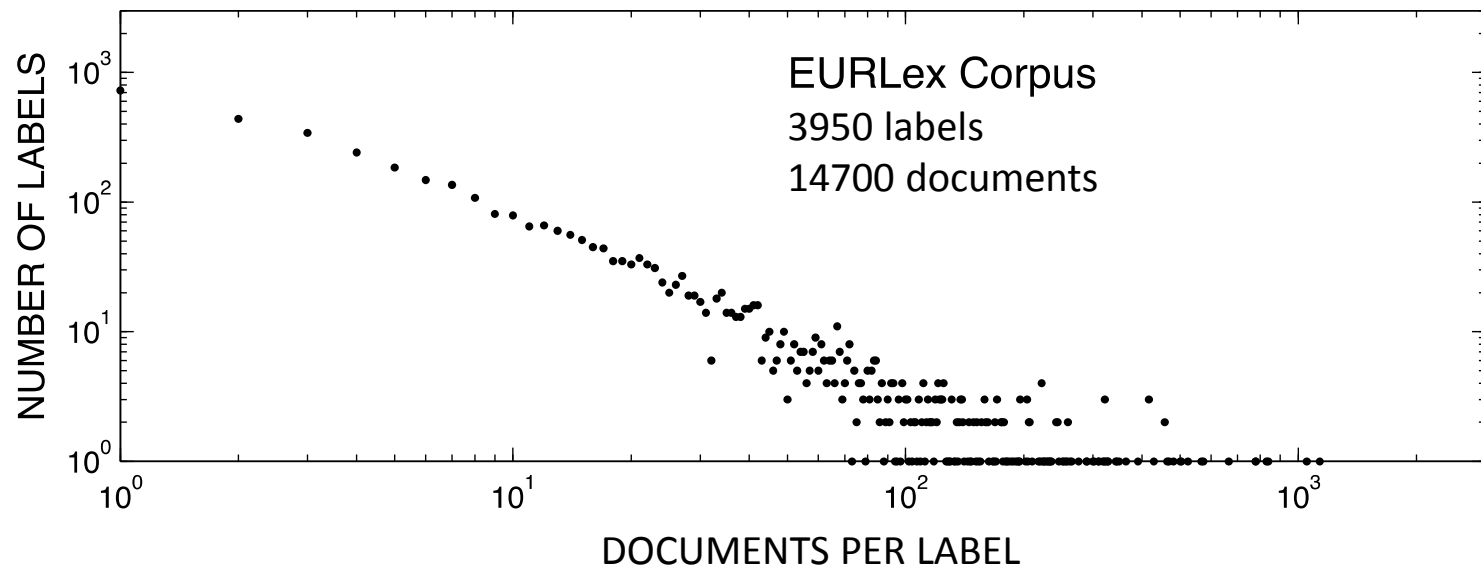
- Wikipedia
 - $O(100k)$ categories
 - Network of relationships
 - Wiki pages tagged with categories

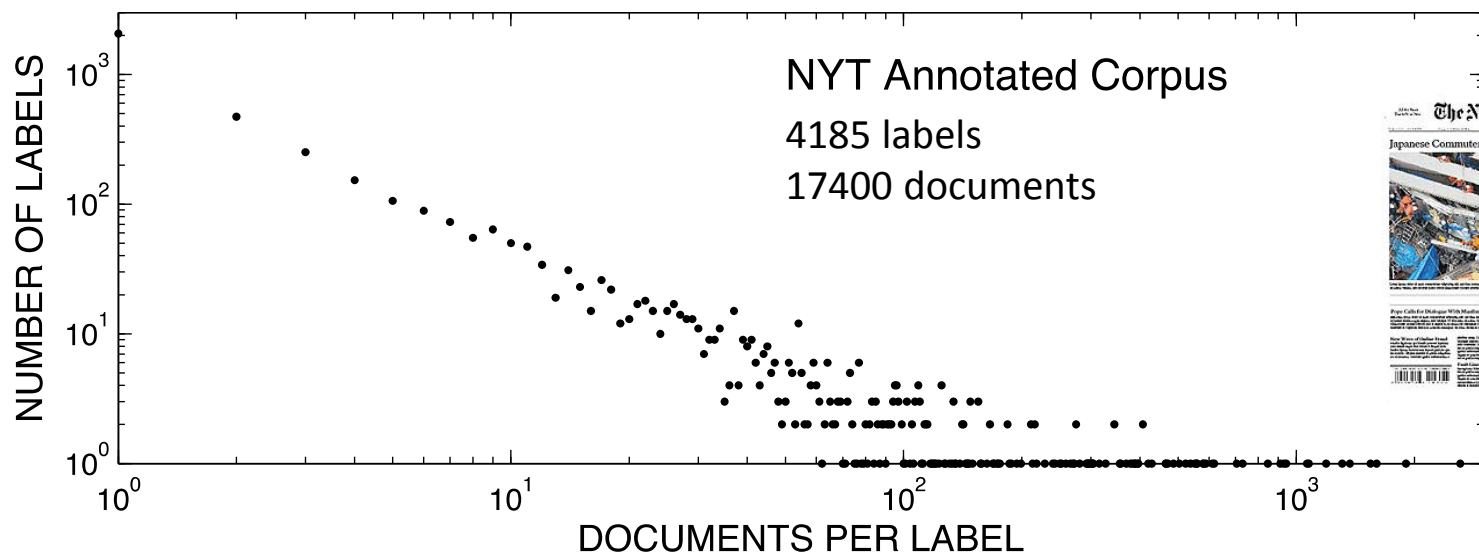
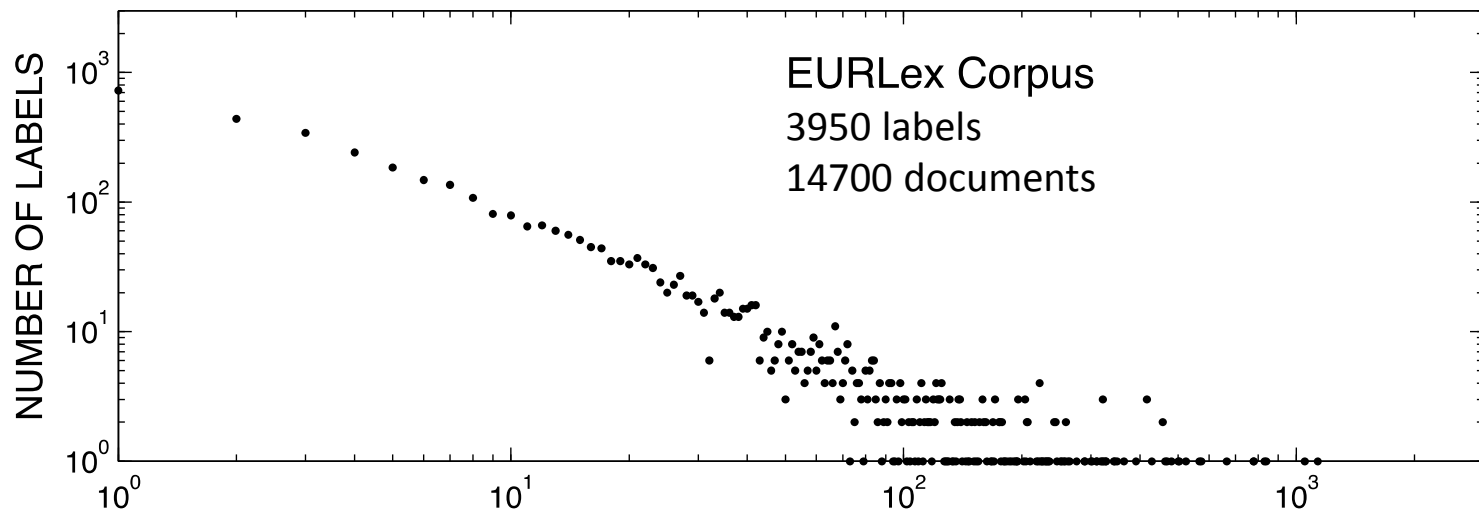
MultiLabel Document Data Sets

Data Set	Number of Unique Labels	Median Number of Documents per Label
RCV1-V2	103	7410
Yahoo! Arts	14	530
Yahoo! Health	19	500

MultiLabel Document Data Sets

Data Set	Number of Unique Labels	Median Number of Documents per Label
RCV1-V2	103	7410
Yahoo! Arts	14	530
Yahoo! Health	19	500
EUR-Lex	3993	6
New York Times	4185	3





Prior Work in Multilabel Text Classification

Unrealistically small numbers of labels are often used, e.g.,

- Yahoo! directories: ~20 labels
- RCV-1 ~ 100 labels
- OHSUMED ~ 100 labels

SVMs do well on common labels, but do poorly when there are few documents per label

- See study by Liu et al (*SIGKDD Explorations*, 2005)

Room for new approaches for data sets with many labels

Discriminative Learning

- “One-versus-all” discriminative learning
 - Learn a binary classifier for each label
 - e.g., SVMs, logistic regression
- Potential limitations of discriminative approach
 - documents with many labels
 - labels with few documents
- Differences with generative (topic) model
 - Discriminative: assigns labels at the document level
 - Generative: assigns labels at the word level

Applying Topic Models to Multilabel Classification

Rubin, Chambers, Smyth, Steyvers, MLJ 2012

- Simple idea:
 - Associate each label with a topic (see Ramage et al, EMNLP, 2009)
 - During learning, restrict the sampler to the known labels for the document
 - Algorithm learns a distribution over words for each label
 - Key difference with discriminative methods: labels are assigned per word, not per document
- Modeling label dependencies
 - Extend standard LDA to allow for label (topic) dependencies – significantly improves performance

Topic Modeling with Labels

- Say our documents can have multiple labels (supervised data)
- Simple observation:
 - Labels and topics: 1-1 correspondence
 - When sampling with Gibbs sampler, we can restrict sampling to “topics” (labels) assigned to that document
- Algorithm learns
 - Which labels are associated with each word within a document
 - Probability distribution over words for each label
- Probabilistic basis for multi-label document classification

NY Times Article

Document Labels	Label Freq.	SVM (weight)
ANTITRUST ACTIONS AND LAWS	19	nintendo
SUITS AND LITIGATION	67	mcgowan
VIDEO GAMES	1	futuristic
		compatible
		illusion
		shrewd
		inception
		truthful
		profiles
		billionayear
		suing
		infringement
		architecture
		handheld
		tantamount
		payoff

Document Excerpt

A flurry of lawsuits, started by a small American software developer, now surrounds the Nintendo Entertainment System, the best-selling toy in the United States last year...Atari Games argues that Nintendo's high degree of control is tantamount to monopoly, and is suing Nintendo for antitrust violations...

NY Times Article

Models for VIDEO GAMES

Document Labels**Label Freq.****SVM (weight)****LDA (prob.)**

ANTITRUST ACTIONS AND LAWS

19

nintendo

nintendo

SUITS AND LITIGATION

67

mcgowan

games

VIDEO GAMES

1

futuristic

software

compatible

video

illusion

system

shrewd

game

inception

chip

truthful

control

profiles

market

billionayear

home

suing

computer

infringement

shortage

architecture

say

handheld

buy

tantamount

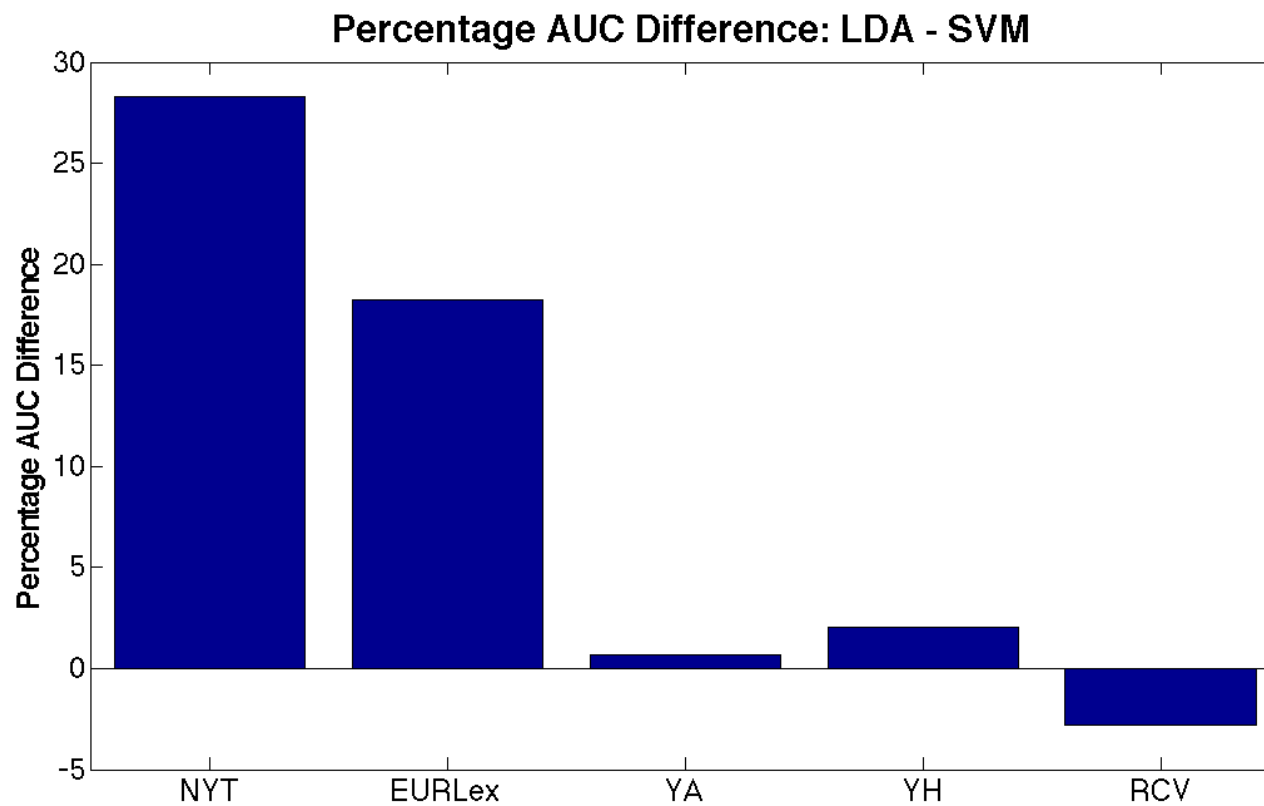
demand

payoff

developer

Document Excerpt

A flurry of lawsuits, started by a small American software developer, now surrounds the Nintendo Entertainment System, the best-selling toy in the United States last year...Atari Games argues that Nintendo's high degree of control is tantamount to monopoly, and is suing Nintendo for antitrust violations...



From Rubin, Chambers, Smyth, Steyvers, MLJ 2012

Data Set	Median Number of Documents per Label	Metrics where Topics were better	Metrics where SVMs were better
RCV1-V2	7410	1	24
Yahoo! Arts	530	11	13
Yahoo! Health	500	13	12
EUR-Lex	6	18	6
New York Times	3	22	1

From Rubin, Chambers, Smyth, Steyvers, MLJ 2012

Background Reading on Topic Models

David Blei's Topic Modeling Web page:

<https://www.cs.princeton.edu/~blei/topicmodeling.html>

See introductory papers and slides from various tutorials

See also code, browsers, visualizations, discussion list, etc

Original paper on topic modeling,

Latent Dirichlet allocation, David Blei, Andrew Y. Ng and Michael Jordan. *Journal of Machine Learning Research*, 3:993-1022, 2003.

Probabilistic topic models, Steyvers, M. & Griffiths, T. (2006). In T. Landauer, D. McNamara, S. Dennis, and W. Kintsch (eds), *Latent Semantic Analysis: A Road to Meaning*. Laurence Erlbaum

Distributed algorithms for topic models, D. Newman, A. Asuncion, P. Smyth, and M. Welling, *Journal of Machine Learning Research*, 10(Aug):1801-1828, 2009.