

**10-701/15-781, Spring 2008**

**Eric Xing**

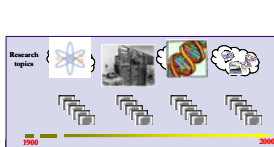
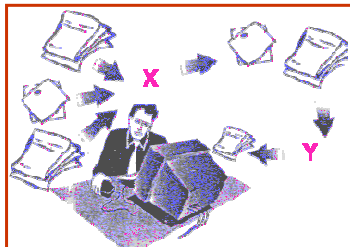
## Lecture 25, April 21, 2008

**Reading: see class homepage**

[illegible]

We want:

- Semantic-based search
- infer topics and categorize documents
- Multimedia inference
- Automatic translation
- Predict how topics evolve
- ...



## Modeling document collections

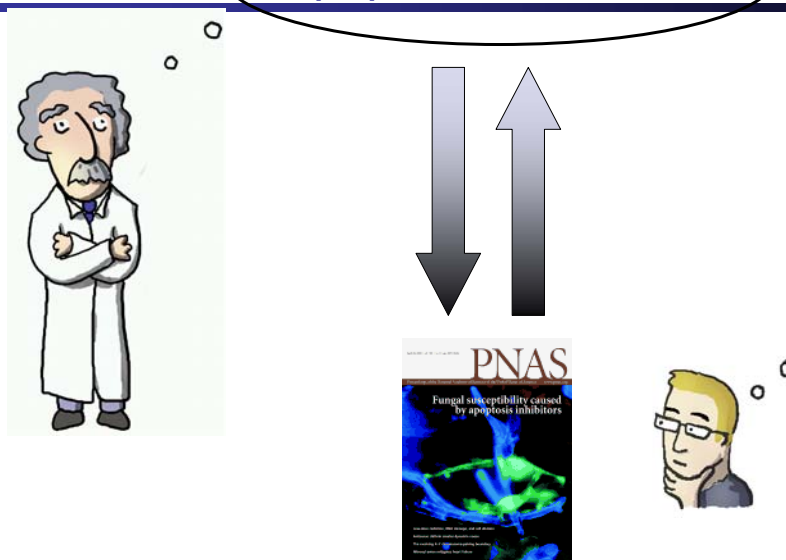


- A document collection is a dataset where each data point is itself a collection of simpler data.
  - Text documents are collections of words.
  - Segmented images are collections of regions.
  - User histories are collections of purchased items.
- Many modern problems ask questions of such data.
  - Is this text document relevant to my query?
  - 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?
  - Which category is this image in? Create a caption for this image.
  - What movies would I probably like?
  - and so on.....

Eric Xing

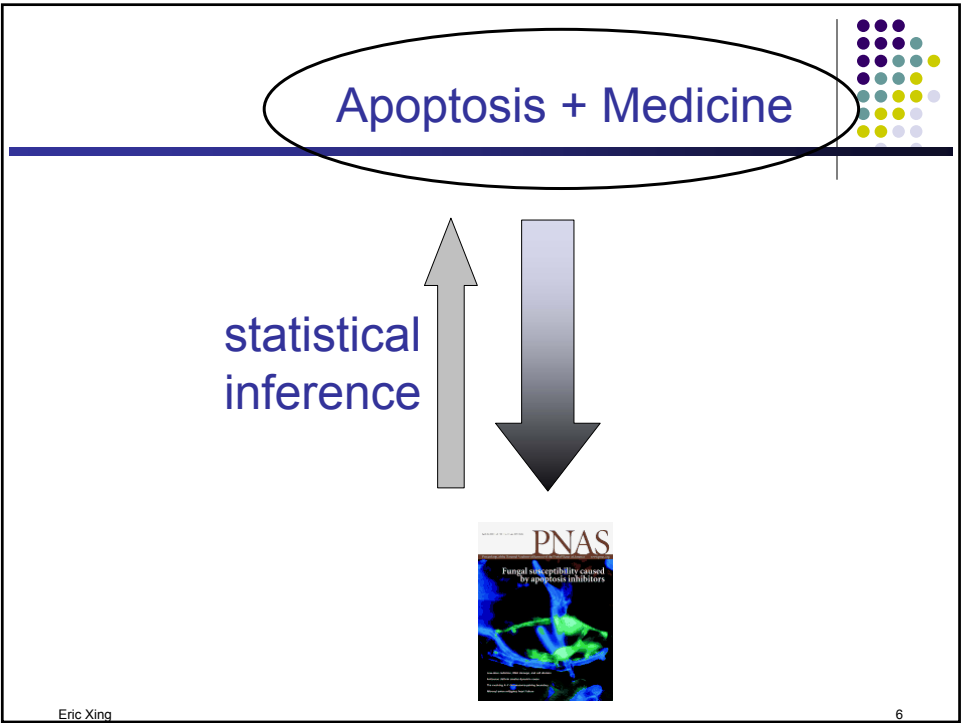
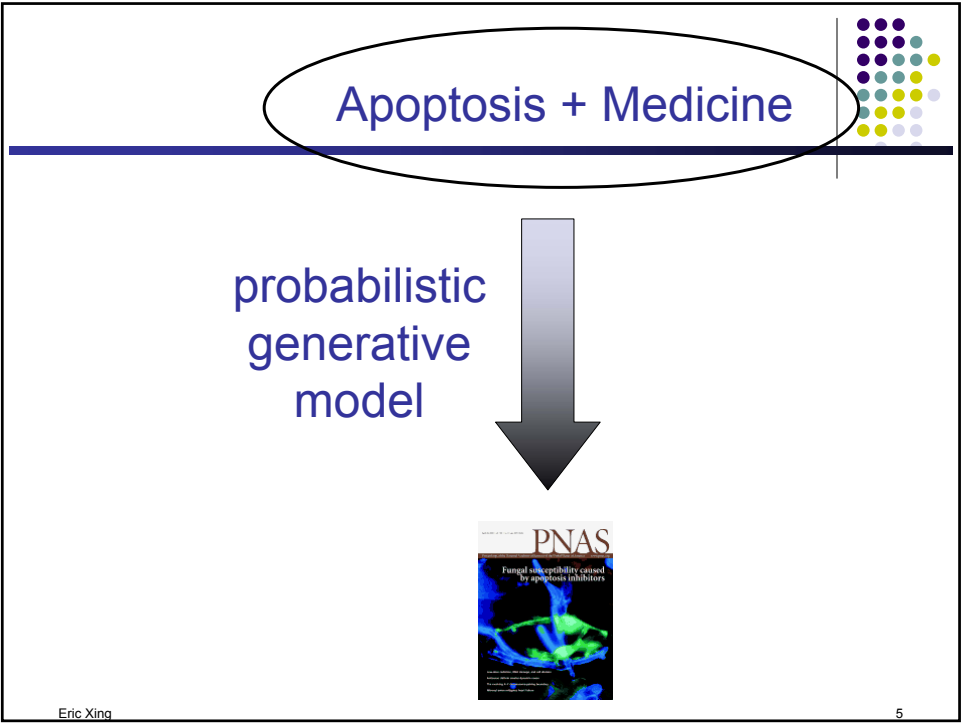
3

## Apoptosis + Medicine

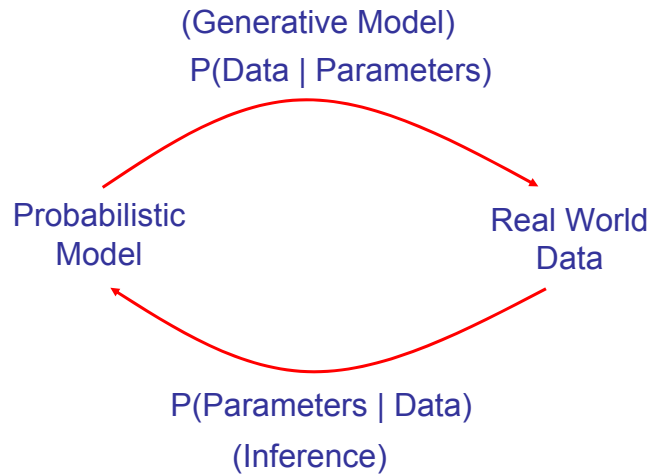


Eric Xing

4



## Connecting Probability Models to Data



Eric Xing

7

## Motivation for modeling latent topical aspects



- Dimensionality reduction
  - A VSM lives in a very high-dimensional feature space (usually larger vocabulary,  $V$ )
  - Sparse representation of documents ( $|V| \gg$  actual number of appeared words in any given document) --- often too spurious for many IR tasks
- Semantic analysis and comprehension
  - A need to define conceptual closeness,
  - to capture relation between features,
  - to distinguish and infer features from heterogeneous sources ...

Eric Xing

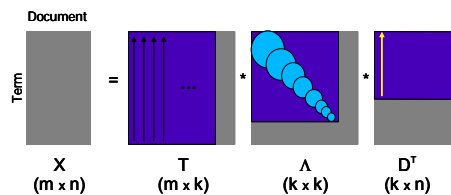
8

# Latent Semantic Indexing

(Deerwester et al., 1990)



- Classic attempt at solving this problem in information retrieval

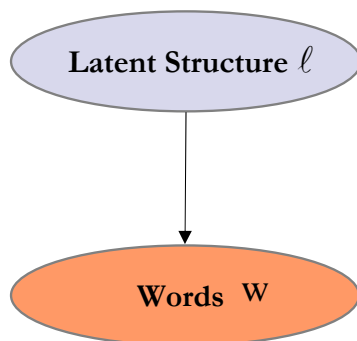


- Uses SVD to reduce document representations
- Models synonymy and polysemy
- Computing SVD is slow
- Non-probabilistic model

Eric Xing

9

# Latent Semantic Structure



Distribution over words

$$P(w) = \sum_{\ell} P(w, \ell)$$

Inferring latent structure

$$P(\ell | w) = \frac{P(w | \ell)P(\ell)}{P(w)}$$

Prediction

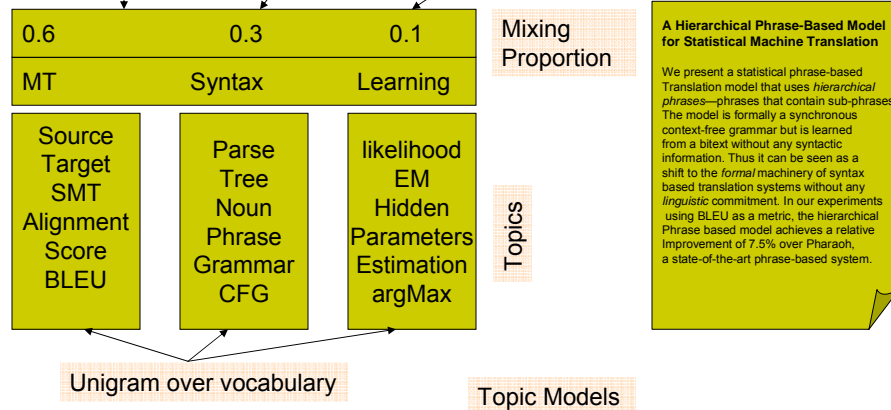
$$P(w_{n+1} | w) = \dots$$

Eric Xing

10

# How to Model Semantic?

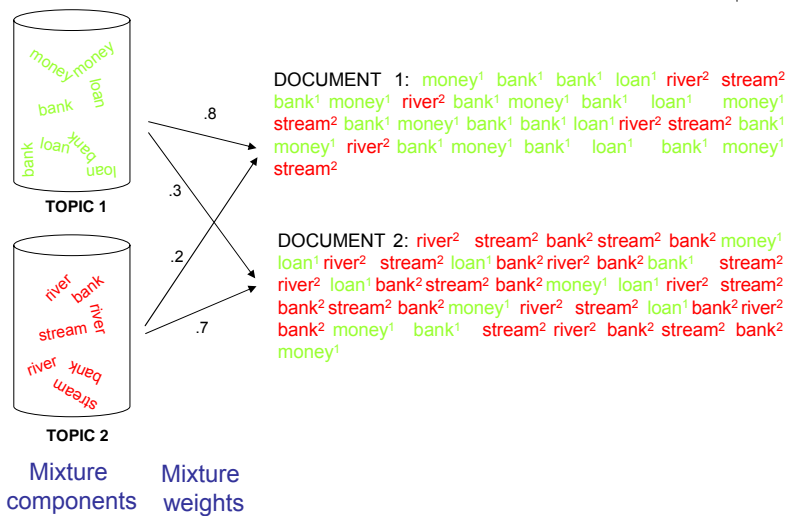
- Q: What is it about?
- A: Mainly MT, with syntax, some learning



Eric Xing

11

# GENERATIVE PROCESS



Eric Xing

12

## Why this is Useful?

- Q: What is it about?
- A: Mainly MT, with syntax, some learning

0.6	0.3	0.1
MT	Syntax	Learning

Mixing  
Proportion

### A Hierarchical Phrase-Based Model for Statistical Machine Translation

We present a statistical phrase-based Translation model that uses *hierarchical phrases*—phrases that contain sub-phrases. The model is formally a synchronous context-free grammar but is learned from a bitext without any syntactic information. Thus it can be seen as a shift to the *formal* machinery of syntax based translation systems without any *linguistic* commitment. In our experiments using BLEU as a metric, the hierarchical Phrase based model achieves a relative Improvement of 7.5% over Pharaoh, a state-of-the-art phrase-based system.

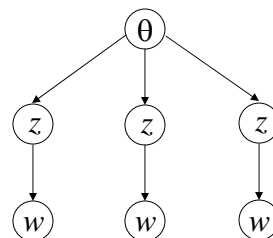
- Q: give me similar document?
  - Structured way of browsing the collection
- Other tasks
  - Dimensionality reduction
    - TF-IDF vs. topic mixing proportion
    - Classification, clustering, and more ...

Eric Xing

13

## A generative model for documents

$$P(w_i) = \sum_{j=1}^T P(w_i | z_i = j) P(z_i = j)$$

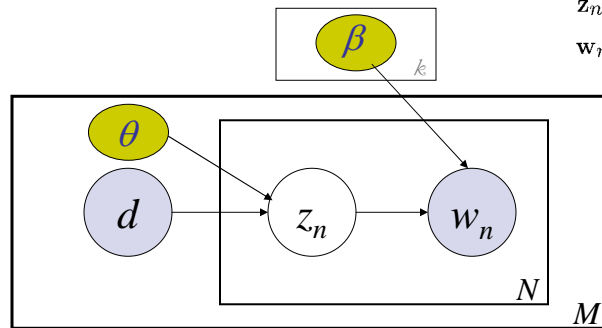


Eric Xing

14

# Probabilistic LSI

Hoffman (1999)



$$z_n \sim \text{Mult}(\theta)$$

$$w_n \sim p(w_n | z_n, \beta)$$

$$p(d, w_n) = p(d) \sum_{\mathbf{z}} \left( \prod_{n=1}^N p(w_n | z_n) p(z_n | d) \right)$$

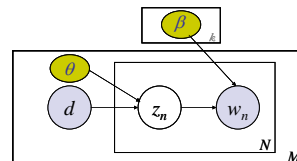
Eric Xing

15

# Probabilistic LSI



- A "generative" model
- Models each word in a document as a sample from a mixture model.
- Each word is generated from a single topic, different words in the document may be generated from different topics.
- A topic is characterized by a distribution over words.
- Each document is represented as a list of mixing proportions for the components (i.e. topic vector  $\theta$ ).



Eric Xing

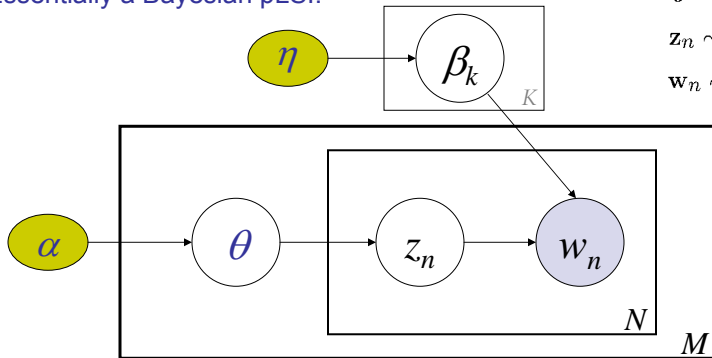
16

# Latent Dirichlet Allocation

Blei, Ng and Jordan (2003)



Essentially a Bayesian pLSI:



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

$$z_n \sim \text{Mult}(\theta)$$

$$w_n \sim p(w_n | z_n, \beta)$$

$$p(\mathbf{w}) = \sum_{\mathbf{z}} \int p(\theta) p(\beta) \left( \prod_{n=1}^N p(z_n | \theta) p(w_n | \beta_{z_n}) \right) d\theta d\beta$$

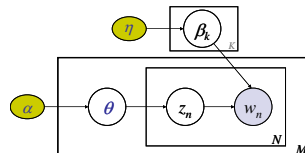
Eric Xing

17

# LDA



- Generative model
- Models each word in a document as a sample from a mixture model.
- Each word is generated from a single topic, different words in the document may be generated from different topics.
- A topic is characterized by a distribution over words.
- Each document is represented as a list of mixing proportions for the components (i.e. topic vector).
- The topic vectors and the word rates each follows a Dirichlet prior --- essentially a Bayesian pLSI

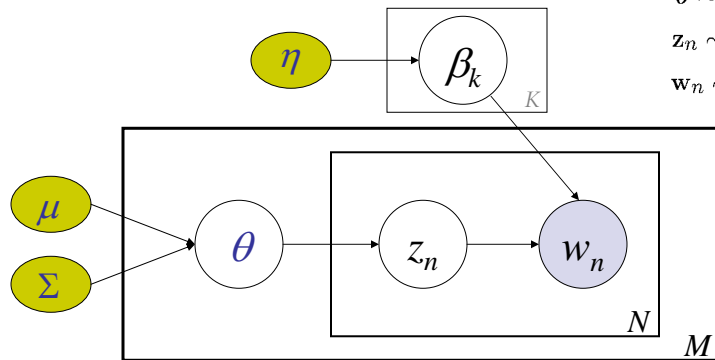


Eric Xing

18

# Correlated Topic Model

Blei & Lafferty (2005)



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

$$z_n \sim \text{Mult}(\theta)$$

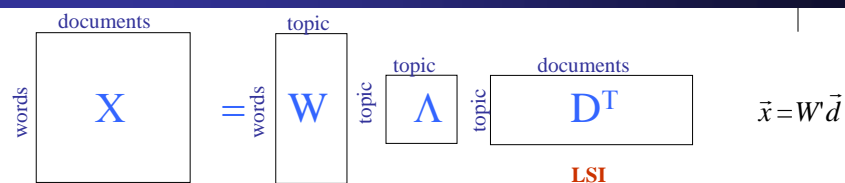
$$w_n \sim p(w_n | z_n, \beta)$$

$$p(\mathbf{w}) = \sum_{\mathbf{z}} \int p(\theta) p(\beta) \left( \prod_{n=1}^N p(z_n | \theta) p(w_n | \beta_{z_n}) \right) d\theta d\beta$$

Eric Xing

19

# Comparison of model semantics



LSI

Topic models

Topic-Mixing is via repeated word labeling

Eric Xing

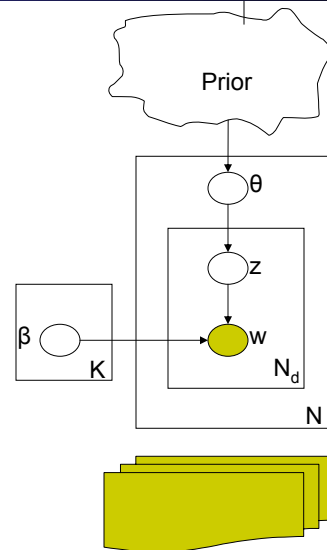
20

# Topic Models = Mixture Membership Models

## Generating a document

- Draw  $\theta$  from the prior
- For each word  $n$ 
  - Draw  $z_n$  from  $\text{multinomial} l(\theta)$
  - Draw  $w_n | z_n, \{\beta_{1:k}\}$  from  $\text{multinomial} l(\beta_{z_n})$

Which prior to use?

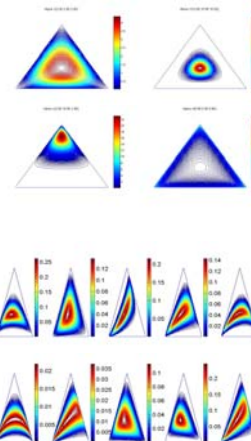


Eric Xing

21

## Prior Comparison

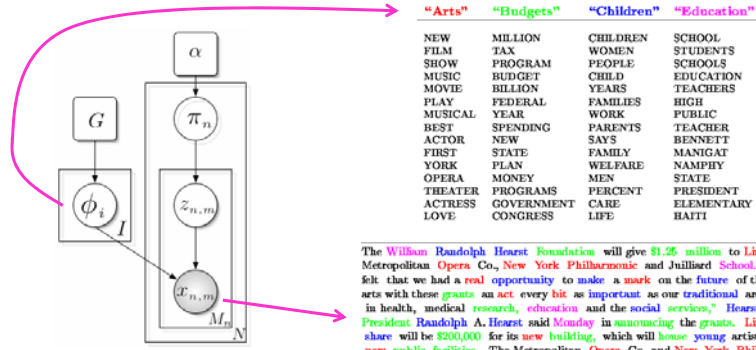
- Dirichlet (LDA) (Blei et al. 2003)
  - Conjugate prior means efficient inference
  - Can **only** capture variations in each topic's intensity **independently**
- Logistic Normal (CTM=LoNTAM) (Blei & Lafferty 2005, Ahmed & Xing 2006)
  - Capture the intuition that some topics are highly correlated and can rise up in intensity together
  - Not a conjugate prior implies **hard** inference



Eric Xing

22

# Inference Tasks



The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. "Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services," Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center's share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.

Eric Xing

23

# Bayesian inference

- A possible query:

$$p(\pi_n | D) = ?$$

$$p(z_{n,m} | D) = ?$$

- Close form solution?

$$p(\pi_n | D) = \frac{p(\pi_n, D)}{p(D)}$$

$$= \frac{\sum_{\{z_{n,m}\}} \int \left( \prod_n \left( \prod_m p(x_{n,m} | \phi_{z_n}) p(z_{n,m} | \pi_n) \right) p(\pi_n | \alpha) \right) p(\phi | G) d\pi_{-i} d\phi}{p(D)}$$

$$p(D) = \sum_{\{z_{n,m}\}} \int \cdots \int \left( \prod_n \left( \prod_m p(x_{n,m} | \phi_{z_n}) p(z_{n,m} | \pi_n) \right) p(\pi_n | \alpha) \right) p(\phi | G) d\pi_1 \cdots d\pi_N d\phi$$

- Sum in the denominator over  $T^n$  terms, and integrate over  $n$   $k$ -dimensional topic vectors

Eric Xing

24

# Approximate Inference

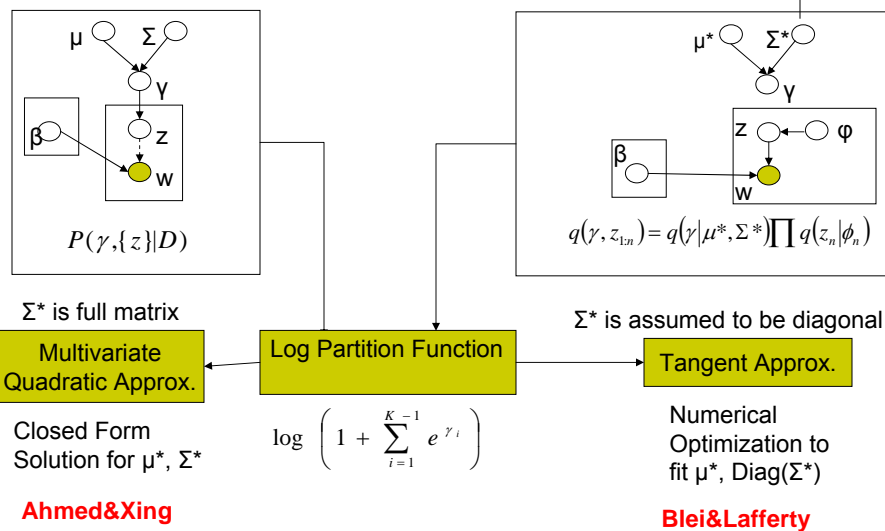
- Variational Inference
  - Mean field approximation (Blei et al)
  - Expectation propagation (Minka et al)
  - Variational 2<sup>nd</sup>-order Taylor approximation (Xing)
- Markov Chain Monte Carlo
  - Gibbs sampling (Griffiths et al)

Eric Xing

25

# Variational Inference

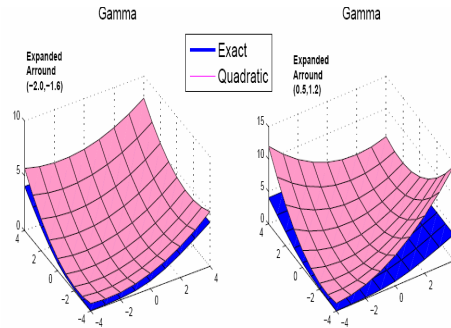
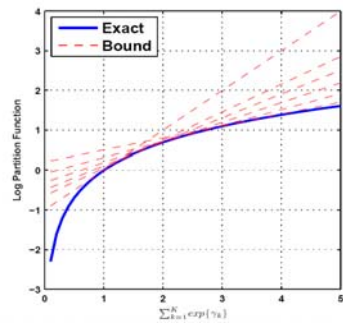
(e.g., MF, Jordan et al 1999, GMF, Xing et al 2004)



Eric Xing

26

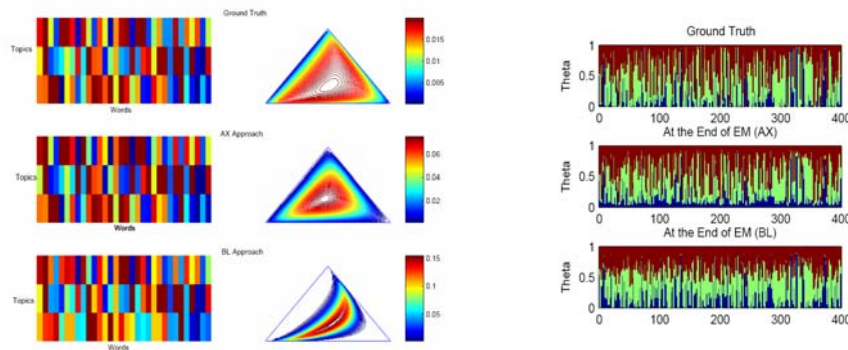
# Tangent Approximation



Eric Xing

27

# Test on Synthetic Text



Eric Xing

28

## Comparison: accuracy and speed

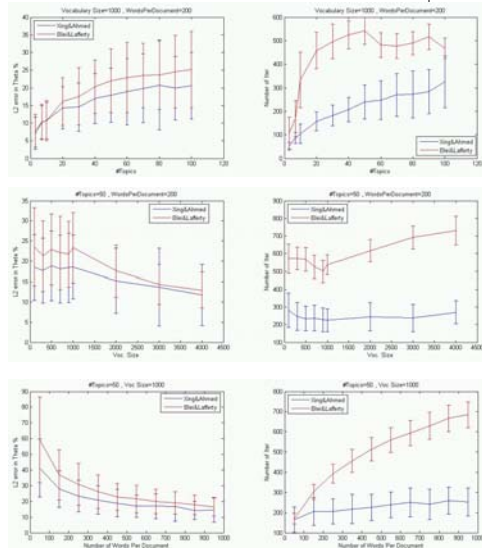


L2 error in topic vector est.  
and # of iterations

- Varying Num. of Topics

- Varying Voc. Size

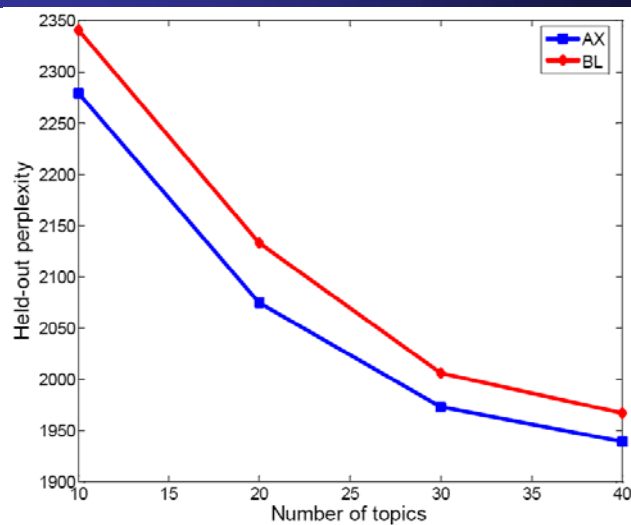
- Varying Num. Words Per Document



Eric Xing

29

## Comparison: perplexity



Eric Xing

30

# Collapsed Gibbs sampling

(Tom Griffiths & Mark Steyvers)

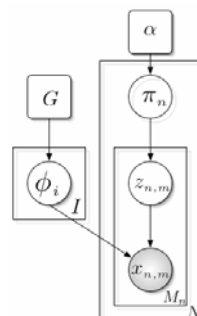
- Collapsed Gibbs sampling

- Integrate out  $\pi$

For variables  $\mathbf{z} = z_1, z_2, \dots, z_n$

Draw  $z_i^{(t+1)}$  from  $P(z_i | \mathbf{z}_{-i}, \mathbf{w})$

$\mathbf{z}_{-i} = z_1^{(t+1)}, z_2^{(t+1)}, \dots, z_{i-1}^{(t+1)}, z_{i+1}^{(t)}, \dots, z_n^{(t)}$



Eric Xing

31

# Gibbs sampling

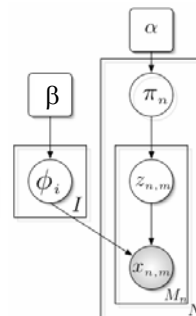
- Need full conditional distributions for variables
- Since we only sample  $z$  we need

$$P(z_i = j | \mathbf{z}_{-i}, \mathbf{w}) \propto P(w_i | z_i = j, \mathbf{z}_{-i}, \mathbf{w}_{-i}) P(z_i = j | \mathbf{z}_{-i})$$

$$= \frac{n_{-i,j}^{(w_i)} + \beta}{n_{-i,j}^{(\cdot)} + W\beta} \frac{n_{-i,j}^{(d_i)} + \alpha}{n_{-i,j}^{(d)} + T\alpha}$$

$n_j^{(w)}$  number of times word  $w$  assigned to topic  $j$

$n_j^{(d)}$  number of times topic  $j$  used in document  $d$



Eric Xing

32

# Gibbs sampling



$i$	$w_i$	$d_i$	iteration	
			1	
1	MATHEMATICS	1	2	
2	KNOWLEDGE	1	2	
3	RESEARCH	1	1	
4	WORK	1	2	
5	MATHEMATICS	1	1	
6	RESEARCH	1	2	
7	WORK	1	2	
8	SCIENTIFIC	1	1	
9	MATHEMATICS	1	2	
10	WORK	1	1	
11	SCIENTIFIC	2	1	
12	KNOWLEDGE	2	1	
.	.	.	.	
.	.	.	.	
.	.	.	.	
50	JOY	5	2	

Eric Xing

33

# Gibbs sampling



$i$	$w_i$	$d_i$	iteration	
			1	2
1	MATHEMATICS	1	2	?
2	KNOWLEDGE	1	2	
3	RESEARCH	1	1	
4	WORK	1	2	
5	MATHEMATICS	1	1	
6	RESEARCH	1	2	
7	WORK	1	2	
8	SCIENTIFIC	1	1	
9	MATHEMATICS	1	2	
10	WORK	1	1	
11	SCIENTIFIC	2	1	
12	KNOWLEDGE	2	1	
.	.	.	.	
.	.	.	.	
.	.	.	.	
50	JOY	5	2	

Eric Xing

34

# Gibbs sampling



$i$	$w_i$	$d_i$	iteration	
			1	2
1	MATHEMATICS	1	$z_i$	$z_i$
2	KNOWLEDGE	1	2	?
3	RESEARCH	1	1	
4	WORK	1	2	
5	MATHEMATICS	1	1	
6	RESEARCH	1	2	
7	WORK	1	2	
8	SCIENTIFIC	1	1	
9	MATHEMATICS	1	2	
10	WORK	1	1	
11	SCIENTIFIC	2	1	
12	KNOWLEDGE	2	1	
.	.	.	.	
.	.	.	.	
.	.	.	.	
50	JOY	5	2	

$$P(z_i = j | \mathbf{z}_{-i}, \mathbf{w}) \propto \frac{n_{-i,j}^{(w_i)} + \beta}{n_{-i,j}^{(\cdot)} + W\beta} \frac{n_{-i,j}^{(d_i)} + \alpha}{n_{-i,\cdot}^{(d_i)} + T\alpha}$$

Eric Xing

35

# Gibbs sampling



$i$	$w_i$	$d_i$	iteration	
			1	2
1	MATHEMATICS	1	2	?
2	KNOWLEDGE	1	2	
3	RESEARCH	1	1	
4	WORK	1	2	
5	MATHEMATICS	1	1	
6	RESEARCH	1	2	
7	WORK	1	2	
8	SCIENTIFIC	1	1	
9	MATHEMATICS	1	2	
10	WORK	1	1	
11	SCIENTIFIC	2	1	
12	KNOWLEDGE	2	1	
.	.	.	.	
.	.	.	.	
.	.	.	.	
50	JOY	5	2	

$$P(z_i = j | \mathbf{z}_{-i}, \mathbf{w}) \propto \frac{n_{-i,j}^{(w_i)} + \beta}{n_{-i,j}^{(\cdot)} + W\beta} \frac{n_{-i,j}^{(d_i)} + \alpha}{n_{-i,\cdot}^{(d_i)} + T\alpha}$$

Eric Xing

36

# Gibbs sampling



$i$	$w_i$	$d_i$	iteration	
			1	2
1	MATHEMATICS	1	2	2
2	KNOWLEDGE	1	2	?
3	RESEARCH	1	1	
4	WORK	1	2	
5	MATHEMATICS	1	1	
6	RESEARCH	1	2	
7	WORK	1	2	
8	SCIENTIFIC	1	1	
9	MATHEMATICS	1	2	
10	WORK	1	1	
11	SCIENTIFIC	2	1	
12	KNOWLEDGE	2	1	
.	.	.	.	
.	.	.	.	
50	JOY	5	2	

$$P(z_i = j | \mathbf{z}_{-i}, \mathbf{w}) \propto \frac{n_{-i,j}^{(w_i)} + \beta}{n_{-i,j}^{(\cdot)} + W\beta} \frac{n_{-i,j}^{(d_i)} + \alpha}{n_{-i,j}^{(d_i)} + T\alpha}$$

Eric Xing

37

# Gibbs sampling



$i$	$w_i$	$d_i$	iteration	
			1	2
1	MATHEMATICS	1	2	2
2	KNOWLEDGE	1	2	1
3	RESEARCH	1	1	?
4	WORK	1	2	
5	MATHEMATICS	1	1	
6	RESEARCH	1	2	
7	WORK	1	2	
8	SCIENTIFIC	1	1	
9	MATHEMATICS	1	2	
10	WORK	1	1	
11	SCIENTIFIC	2	1	
12	KNOWLEDGE	2	1	
.	.	.	.	
.	.	.	.	
50	JOY	5	2	

$$P(z_i = j | \mathbf{z}_{-i}, \mathbf{w}) \propto \frac{n_{-i,j}^{(w_i)} + \beta}{n_{-i,j}^{(\cdot)} + W\beta} \frac{n_{-i,j}^{(d_i)} + \alpha}{n_{-i,j}^{(d_i)} + T\alpha}$$

Eric Xing

38

# Gibbs sampling



$i$	$w_i$	$d_i$	iteration	
			1	2
1	MATHEMATICS	1	2	2
2	KNOWLEDGE	1	2	1
3	RESEARCH	1	1	1
4	WORK	1	2	?
5	MATHEMATICS	1	1	
6	RESEARCH	1	2	
7	WORK	1	2	
8	SCIENTIFIC	1	1	
9	MATHEMATICS	1	2	
10	WORK	1	1	
11	SCIENTIFIC	2	1	
12	KNOWLEDGE	2	1	
.	.	.	.	
.	.	.	.	
50	JOY	5	2	

$$P(z_i = j | \mathbf{z}_{-i}, \mathbf{w}) \propto \frac{n_{-i,j}^{(w_i)} + \beta}{n_{-i,j}^{(\cdot)} + W\beta} \frac{n_{-i,j}^{(d_i)} + \alpha}{n_{-i,j}^{(d_i)} + T\alpha}$$

Eric Xing

39

# Gibbs sampling



$i$	$w_i$	$d_i$	iteration	
			1	2
1	MATHEMATICS	1	2	2
2	KNOWLEDGE	1	2	1
3	RESEARCH	1	1	1
4	WORK	1	2	2
5	MATHEMATICS	1	1	?
6	RESEARCH	1	2	
7	WORK	1	2	
8	SCIENTIFIC	1	1	
9	MATHEMATICS	1	2	
10	WORK	1	1	
11	SCIENTIFIC	2	1	
12	KNOWLEDGE	2	1	
.	.	.	.	
.	.	.	.	
50	JOY	5	2	

$$P(z_i = j | \mathbf{z}_{-i}, \mathbf{w}) \propto \frac{n_{-i,j}^{(w_i)} + \beta}{n_{-i,j}^{(\cdot)} + W\beta} \frac{n_{-i,j}^{(d_i)} + \alpha}{n_{-i,j}^{(d_i)} + T\alpha}$$

Eric Xing

40

# Gibbs sampling



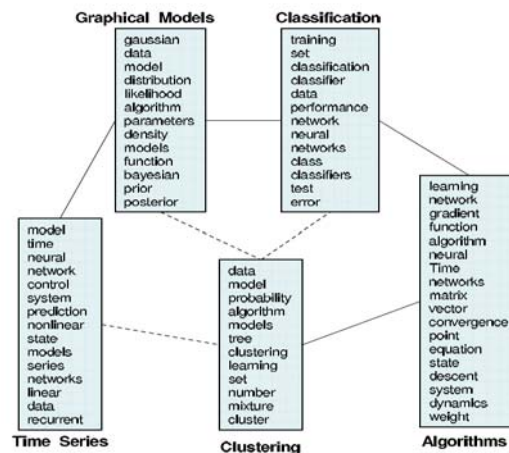
$i$	$w_i$	$d_i$	iteration			
			1	2	...	1000
1	MATHEMATICS	1	2	2		2
2	KNOWLEDGE	1	2	1		2
3	RESEARCH	1	1	1		2
4	WORK	1	2	2		1
5	MATHEMATICS	1	1	2		2
6	RESEARCH	1	2	2		2
7	WORK	1	2	2		2
8	SCIENTIFIC	1	1	1	...	1
9	MATHEMATICS	1	2	2		2
10	WORK	1	1	2		2
11	SCIENTIFIC	2	1	1		2
12	KNOWLEDGE	2	1	2		2
.	.	.	.	.		.
.	.	.	.	.		.
50	JOY	5	2	1		1

$$P(z_i = j | \mathbf{z}_{-i}, \mathbf{w}) \propto \frac{n_{-i,j}^{(w_i)} + \beta}{n_{-i,j}^{(\cdot)} + W\beta} \frac{n_{-i,j}^{(d_i)} + \alpha}{n_{-i,j}^{(d_i)} + T\alpha}$$

Eric Xing

41

# Topics and topic graphs



Eric Xing

42

## Result on PNAS collection

- PNAS abstracts from 1997-2002
  - 2500 documents
  - Average of 170 words per document
- Fitted 40-topics model using both approaches
- Use low dimensional representation to predict the abstract category
  - Use SVM classifier
  - 85% for training and 15% for testing

Classification Accuracy

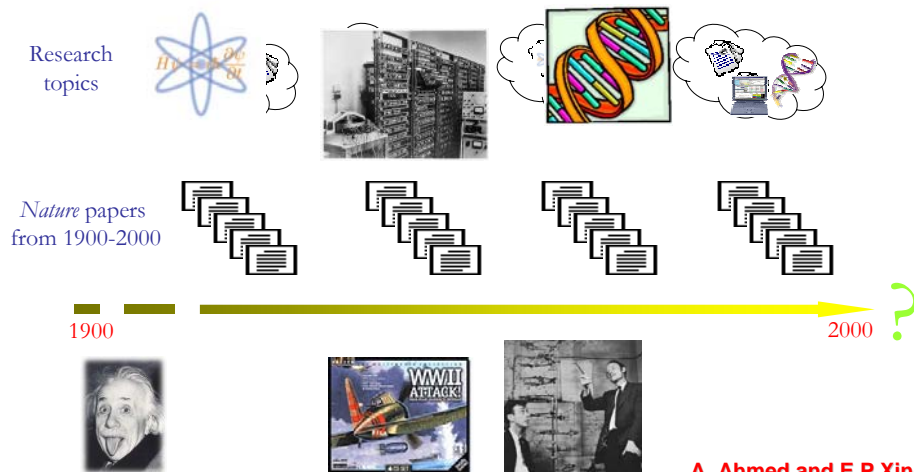
Category	Doc	BL	AX
Genetics	21	61.9	61.9
Biochemistry	86	65.1	77.9
Immunology	24	70.8	66.6
Biophysics	15	53.3	66.6
Total	146	64.3	72.6

-Notable Difference  
-Examine the low dimensional representations below

Eric Xing

43

## Extension 1: topic evolution?



Eric Xing

44

# How to Model Topic Evolution

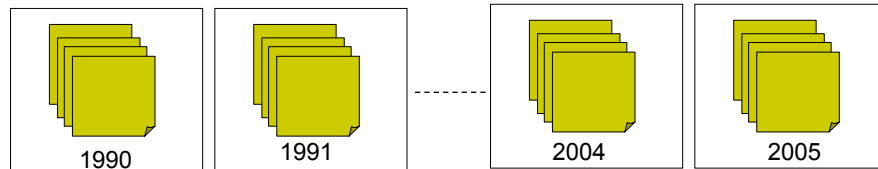
Topic Trends

Topic Keywords

Topic correlations

~~Number of topics~~

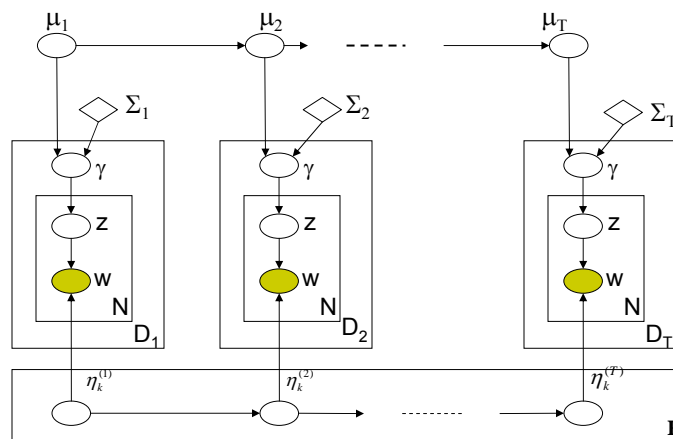
The Dynamic Correlated Topic model



Eric Xing

45

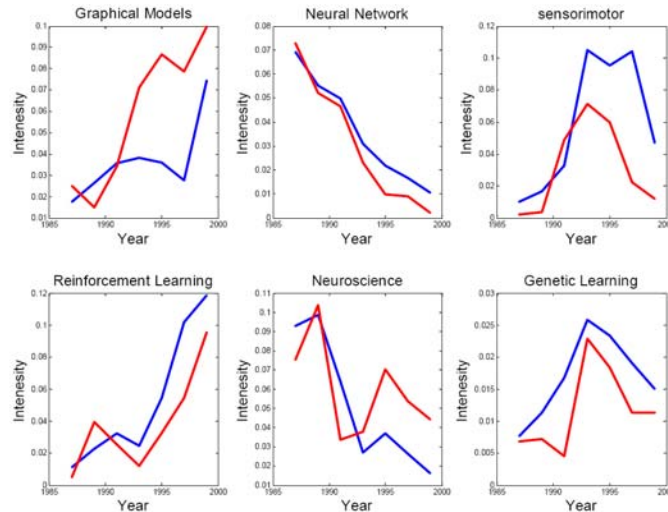
## The Dynamic CTM



Eric Xing

46

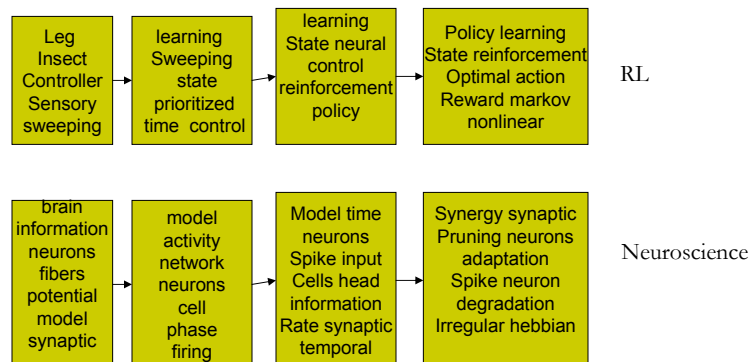
## Topic Trends



Eric Xing

47

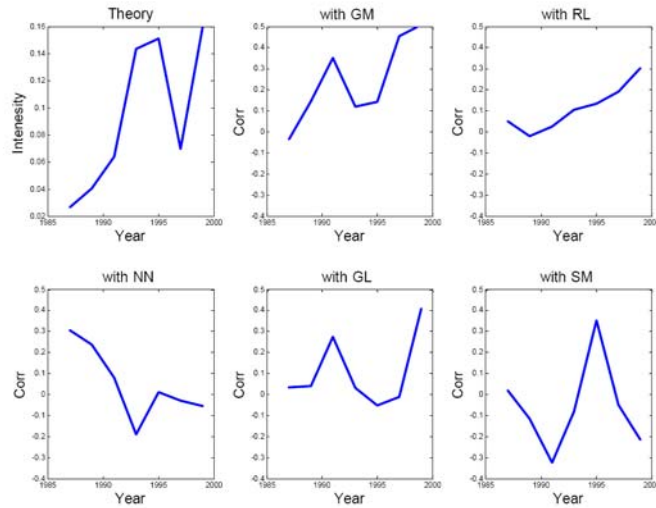
## Topic Words over Time



Eric Xing

48

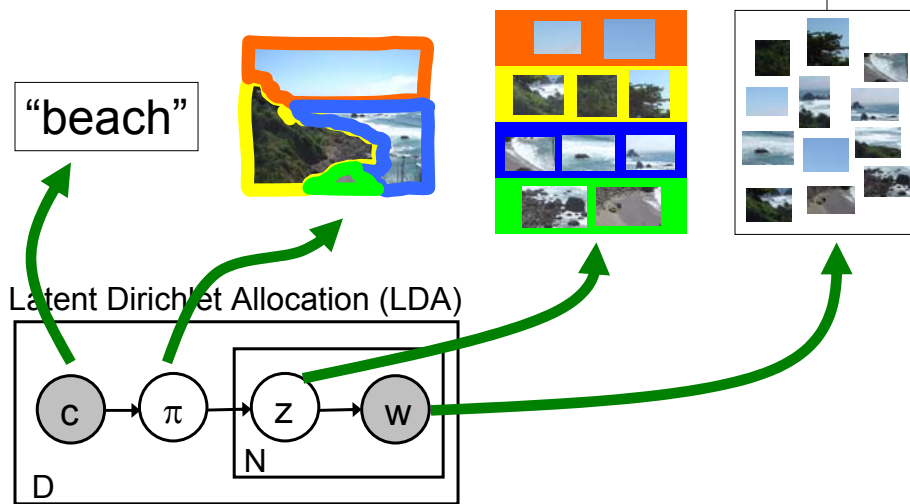
## Topic Correlations Over Time



Eric Xing

49

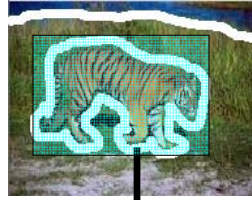
## Extension 2: Topic Models for Images



Eric Xing

Fei-Fei et al. ICCV 2005

## Image Representation



cat, grass, tiger, water

$$[r_{11} \cdots r_{1d}], [w_1 \cdots w_{|V|}]$$

**representation vector** : **annotation vector**  
 (real, 1 per image segment) : (binary, same for each segment)

$$[r_{n1} \cdots r_{nd}], [w_1 \cdots w_{|V|}]$$

Eric Xing

51

## To Generate an Image ...

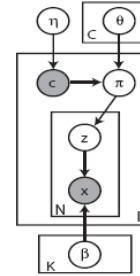
$$p(x, z, \pi, c | \theta, \eta, \beta) = p(c | \eta) p(\pi | c, \theta) \cdot \prod_{n=1}^N p(z_n | \pi) p(x_n | z_n, \beta)$$

$$p(c | \eta) = \text{Mult}(c | \eta)$$

$$p(\pi | c, \theta) = \prod_{j=1}^C \text{Dir}(\pi | \theta_{j \cdot})^{\delta(c, j)}$$

$$p(z_n | \pi) = \text{Mult}(z_n | \pi)$$

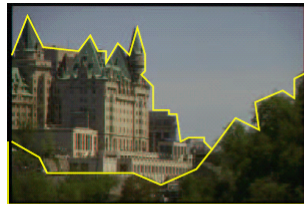
$$p(x_n | z_n, \beta) = \prod_{k=1}^K p(x_n | \beta_{k \cdot})^{\delta(z_n^k, 1)}$$



Eric Xing

52

## Annotated images



This cozy place is nestled in the heart of the Mission. Easy access to bars, restaurants, and BART.

$\{9.32, 2.44, 0.02, 3.23\}$   
 $\{4.35, 3.12, -0.23, 9.41\}$   
 $\{6.65, 2.11, 1.02, 2.31\}$

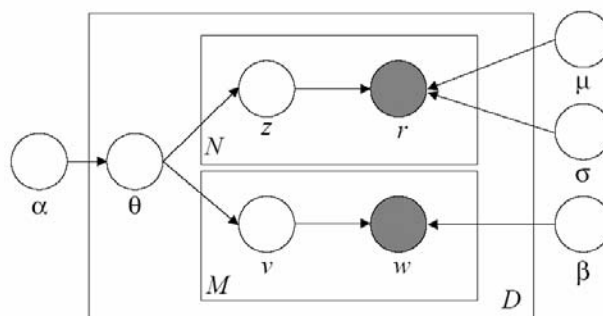
This, cozy, place, is, nestled, in, the, heart, of, the, Mission. Easy, access, to, bars, restaurants, and, BART

- Forsyth et. al. (2001): images as documents where region-specific feature vectors are like visual words.
- A captioned image can be thought of as annotated data: two documents, one of which describes the other.

Eric Xing

53

## Gaussian-multinomial LDA



- A natural next step is to glue two LDA models together.
- Bottom: a traditional LDA model on captions
- Top: a Gaussian-LDA model on images
  - each region is a multivariate Gaussian
- Does not work well

Eric Xing

54

## Automatic annotation



**True caption**

birds tree

**Corr-LDA**

birds nest leaves branch tree

**GM-LDA**

water birds nest tree sky

**GM-Mixture**

tree ocean fungus mushrooms coral



**True caption**

fish reefs water

**Corr-LDA**

fish water ocean tree coral

**GM-LDA**

water sky vegetables tree people

**GM-Mixture**

fungus mushrooms tree flowers leaves

Eric Xing

55

## Conclusion



- GM-based topic models are cool
  - Flexible
  - Modular
  - Interactive
- There are many ways of implementing topic models
  - Directed
  - Undirected
- Efficient Inference/learning algorithms
  - GMF, with Laplace approx. for non-conjugate dist.
  - MCMC
- Many applications
  - ...
  - Word-sense disambiguation (with WeiHao Lin and Alex Hauptman)
  - Word-net (with Amr)
  - Network inference (with Fan Guo and Steve Fienberg)

Eric Xing

56