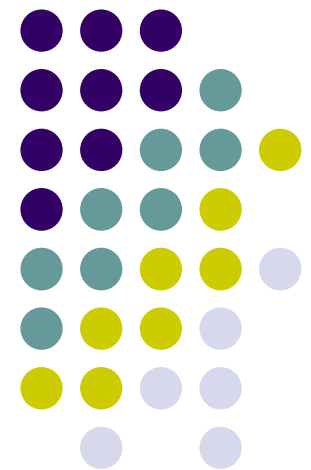**School of Computer Science**
**Carnegie Mellon**

# Probabilistic Graphical Models

## Mean Fiend Approximation

## &

## Topic Models

**Eric Xing**

**Lecture 15, March 5, 2014**

**Reading:** See class website

1

# Variational Principle

- Exact variational formulation

$$A(\theta) = \sup_{\mu \in \mathcal{M}} \left\{ \theta^T \mu - A^*(\mu) \right\}$$

  - $\mathcal{M}$: the marginal polytope, difficult to characterize
  - $A^*$: the negative entropy function, no explicit form

- Mean field method: non-convex inner bound and exact form of entropy

- Bethe approximation and loopy belief propagation: polyhedral outer bound and non-convex Bethe approximation

# Mean Field Approximation

# Mean Field Methods

- For a given tractable subgraph F, a subset of canonical parameters is

$$\mathcal{M}(F; \phi) := \{\tau \in \mathbb{R}^d \mid \tau = \mathbb{E}_\theta[\phi(X)] \text{ for some } \theta \in \Omega(F)\}$$

- Inner approximation

$$\mathcal{M}(F; \phi)^o \subseteq \mathcal{M}(G; \phi)^o$$

- Mean field solves the relaxed problem

$$\max_{\tau \in \mathcal{M}_F(G)} \{\langle \tau, \theta \rangle - A_F^*(\tau)\}$$

- $A_F^* = A^*\big|_{\mathcal{M}_F(G)}$ is the exact dual function restricted to $\mathcal{M}_F(G)$
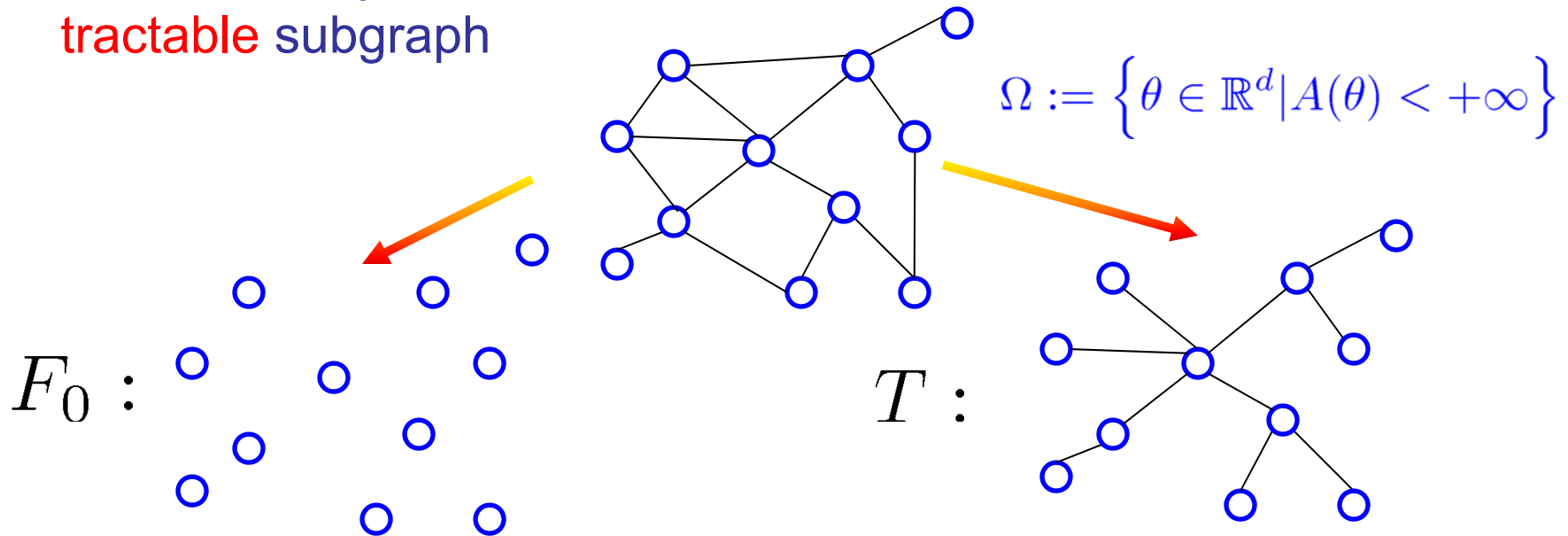
# Tractable Subgraphs

- For an exponential family with sufficient statistics $\phi$ defined on graph G, the set of realizable mean parameter set

$$\mathcal{M}(G;\phi) := \{\mu \in \mathbb{R}^d \mid \exists p \text{ s.t. } \mathbb{E}_p[\phi(X)] = \mu\}$$

- Idea: restrict $p$ to a subset of distributions associated with a tractable subgraph

$$\Omega := \left\{\theta \in \mathbb{R}^d \mid A(\theta) < +\infty\right\}$$

$$F_0 : \qquad T :$$

$$\Omega(F_0) := \left\{\theta \in \Omega \mid \theta_{(s,t)} = 0 \ \forall \ (s,t) \in E\right\}. \quad \Omega(T) := \left\{\theta \in \Omega \mid \theta_{(s,t)} = 0 \ \forall \ (s,t) \notin E(T)\right\}.$$

# Example: Naïve Mean Field for Ising Model

- Ising model in {0,1} representation

$$p(x) \propto \exp\left\{\sum_{s \in V} x_s \theta_s + \sum_{(s,t) \in E} x_s x_t \theta_{st}\right\}$$

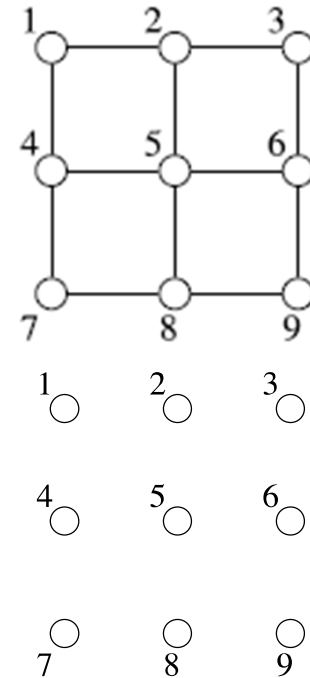- Mean parameters

  $\mu_s = E_p[X_s] = P[X_s = 1]$  for all $s \in V$, and

  $\mu_{st} = E_p[X_s X_t] = P[(X_s, X_t) = (1,1)]$  for all $(s,t) \in E$.

- For fully disconnected graph F,

$$\mathcal{M}_F(G) := \{\tau \in \mathbb{R}^{|V|+|E|} \mid 0 \le \tau_s \le 1, \forall s \in V, \tau_{st} = \tau_s \tau_t, \forall (s,t) \in E\}$$

- The dual decomposes into sum, one for each node

$$A_F^*(\tau) = \sum_{s \in V}[\tau_s \log \tau_s + (1 - \tau_s)\log(1 - \tau_s)]$$
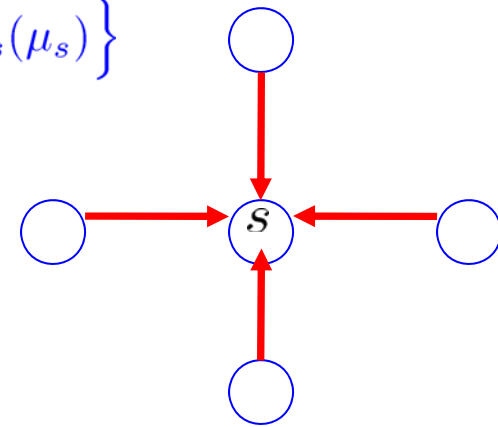
# Naïve Mean Field for Ising Model

- Optimization Problem

$$\max_{\mu \in [0,1]^m} \left\{ \sum_{s \in V} \theta_s \mu_s + \sum_{(s,t) \in E} \theta_{st} \mu_s \mu_t + \sum_{s \in V} H_s(\mu_s) \right\}$$

- Update Rule

$$\mu_s \leftarrow \sigma\left(\theta_s + \sum_{t \in N(s)} \theta_{st} \mu_t\right)$$

- $\mu_t = p(X_t = 1) = \mathbb{E}_p[X_t]$ resembles "message" sent from node $t$ to $s$

- $\{\mathbb{E}_p[X_t], t \in N(s)\}$ forms the "mean field" applied to $s$ from its neighborhood

- Also yields lower bound on log partition function

$$KL(Q \| P) = -H_Q(X) - \sum_{f_a \in F} E_Q \log f_a(X_a) + \log Z$$

# Geometry of Mean Field

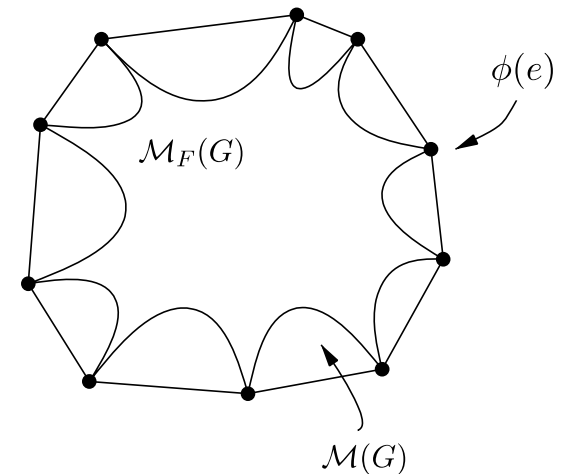- Mean field optimization is always <span style="color:red">non-convex</span> for any exponential family in which the state space $\mathcal{X}^m$ is finite

- Recall the marginal polytope is a convex hull

$$\mathcal{M}(G) = \text{conv}\{\phi(e); e \in \mathcal{X}^m\}$$



- $\mathcal{M}_F(G)$ contains all the extreme points

  - If it is a <span style="color:red">strict</span> subset, then it must be non-convex

- Example: two-node Ising model

$$\mathcal{M}_F(G) = \{0 \le \tau_1 \le 1, 0 \le \tau_2 \le 1, \tau_{12} = \tau_1\tau_2\}$$

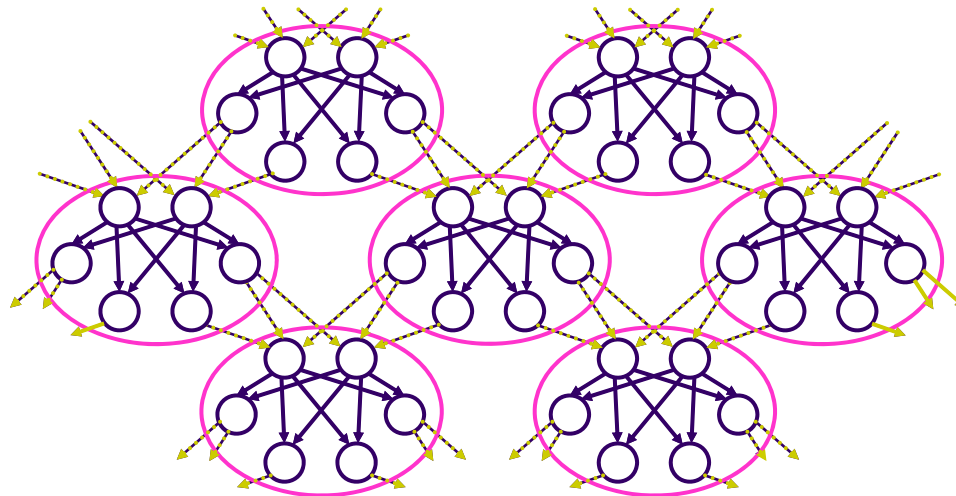  - It has a parabolic cross section along $\tau_1 = \tau_2$ , hence non-convex

# Cluster-based approx. to the Gibbs free energy

(Wiegerinck 2001, Xing *et al* 03,04)

Exact: $G[p(X)]$ *(intractable)*

Clusters: $G[\{q_c(X_c)\}]$

# Mean field approx. to Gibbs free energy

- Given a disjoint clustering, $\{C_1, \ldots, C_l\}$, of all variables
- Let

$$q(\mathbf{X}) = \prod_i q_i(\mathbf{X}_{Ci}),$$

- Mean-field free energy

$$G_{\mathrm{MF}} = \sum_i \sum_{\mathbf{x}_{C_i}} \prod_i q_i(\mathbf{x}_{C_i}) E(\mathbf{x}_{C_i}) + \sum_i \sum_{\mathbf{x}_{C_i}} q_i(\mathbf{x}_{C_i}) \ln q_i(\mathbf{x}_{C_i})$$

e.g., $\quad G_{\mathrm{MF}} = \sum_{i<j} \sum_{x_i x_j} q(x_i) q(x_j) \phi(x_i x_j) + \sum_i \sum_{x_i} q(x_i) \phi(x_i) + \sum_i \sum_{x_i} q(x_i) \ln q(x_i) \quad$ (naïve mean field)

- - Will **never** equal to the exact Gibbs free energy no matter what clustering is used, but it does **always** define a lower bound of the likelihood

- Optimize each $q_i(x_c)$'s.
  - Variational calculus …
  - Do inference in each $q_i(x_c)$ using any tractable algorithm
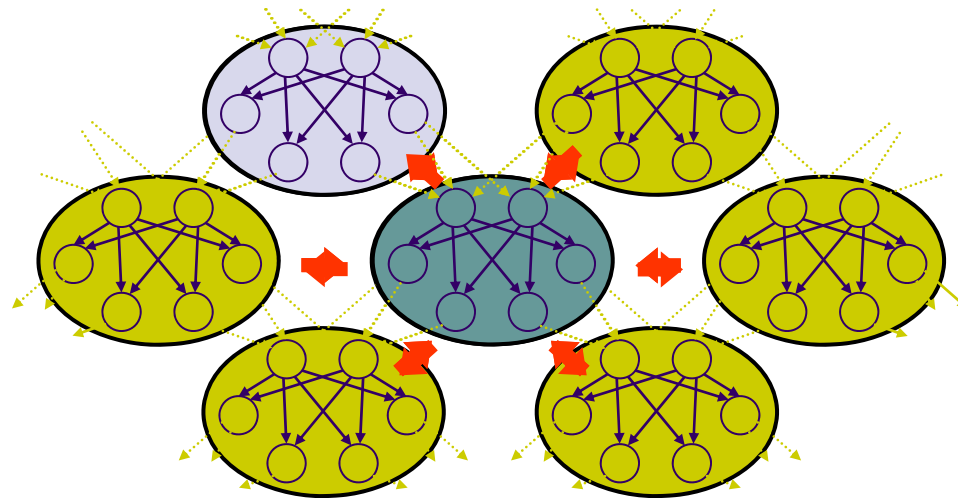
# The Generalized Mean Field theorem

**Theorem:** The optimum GMF approximation to the cluster marginal is isomorphic to the cluster posterior of the original distribution given internal evidence and its generalized mean fields:

$$q_i^*(\mathbf{X}_{H,C_i}) = p(\mathbf{X}_{H,C_i} \mid \mathbf{x}_{E,C_i}, \langle \mathbf{X}_{H,MB_i} \rangle_{q_{j \neq i}})$$
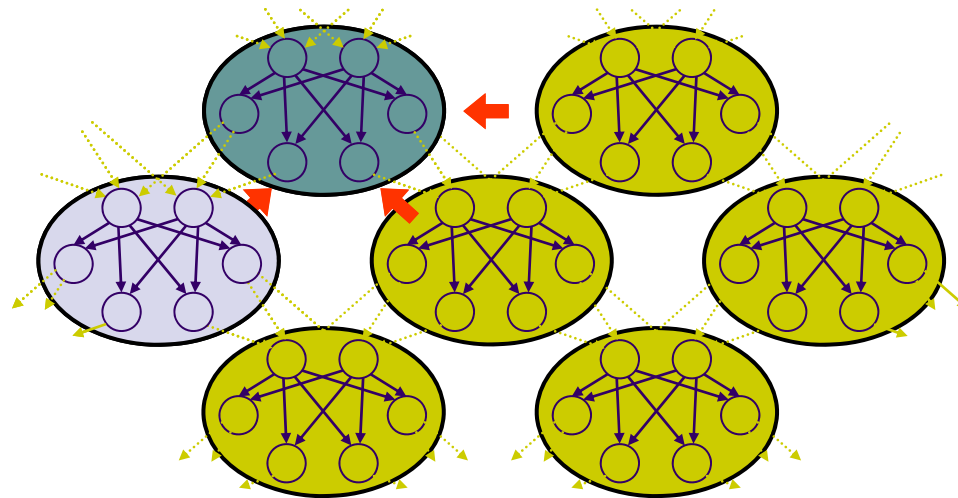
GMF algorithm: Iterate over each $q_i$

# A generalized mean field algorithm [xing *et al*. UAI 2003]

# A generalized mean field algorithm [xing *et al*. UAI 2003]

# Convergence theorem

**Theorem:** The GMF algorithm is guaranteed to converge to a local optimum, and provides a lower bound for the likelihood of evidence (or partition function) the model.
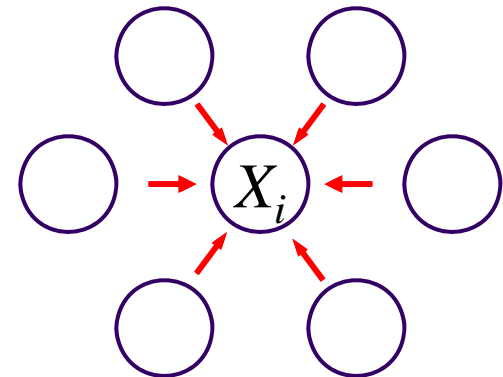
# The naive mean field approximation

- Approximate $p(\mathbf{X})$ by fully factorized $q(\mathbf{X}) = P_i q_i(X_i)$

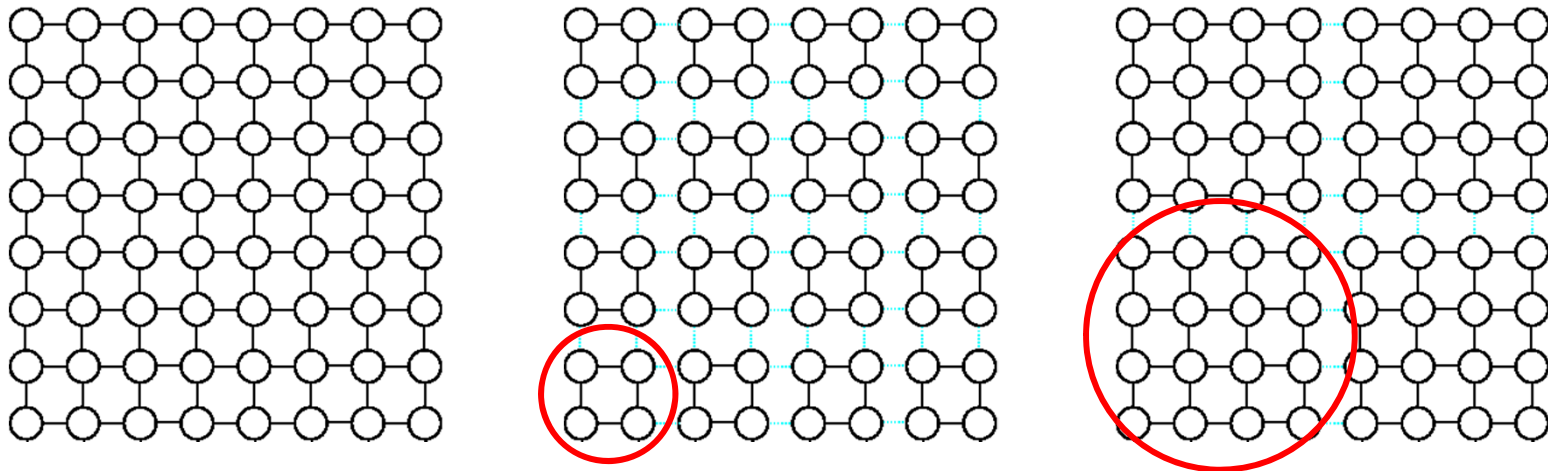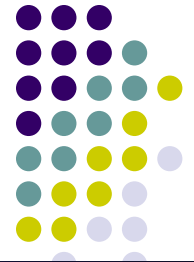- For Boltzmann distribution $p(X) = \exp\{\sum_{i<j} q_{ij} X_i X_j + q_{io} X_i\}/Z$ :

mean field equation:

$$q_i(X_i) = \exp\left\{ \theta_{i0} X_i + \sum_{j \in \mathcal{N}_i} \theta_{ij} X_i \langle X_j \rangle_{q_j} + A_i \right\}$$

$$= p(X_i \mid \{\langle X_j \rangle_{q_j} : j \in \mathcal{N}_i\})$$



- $\langle X_j \rangle_{q_j}$ resembles a "message" sent from node $j$ to $i$

- $\{\langle X_j \rangle_{q_j} : j \in \mathcal{N}_i\}$ forms the "mean field" applied to $X_i$ from its neighborhood

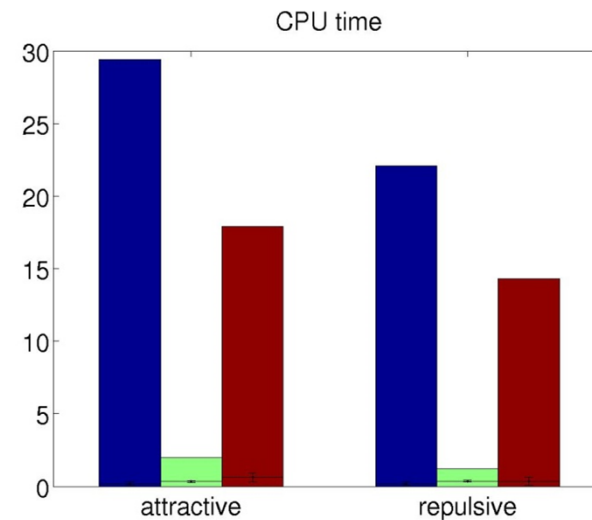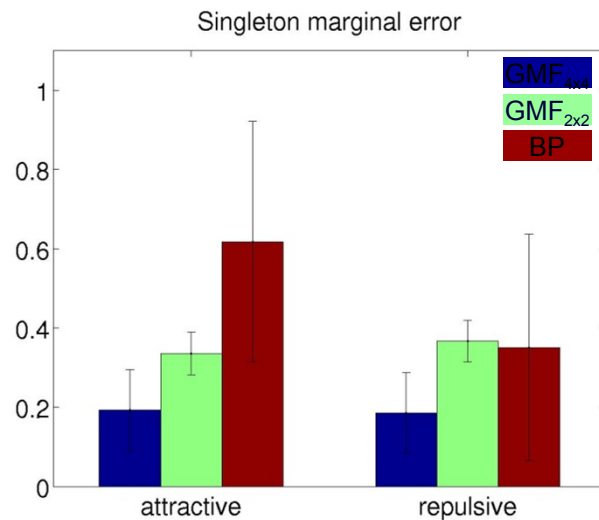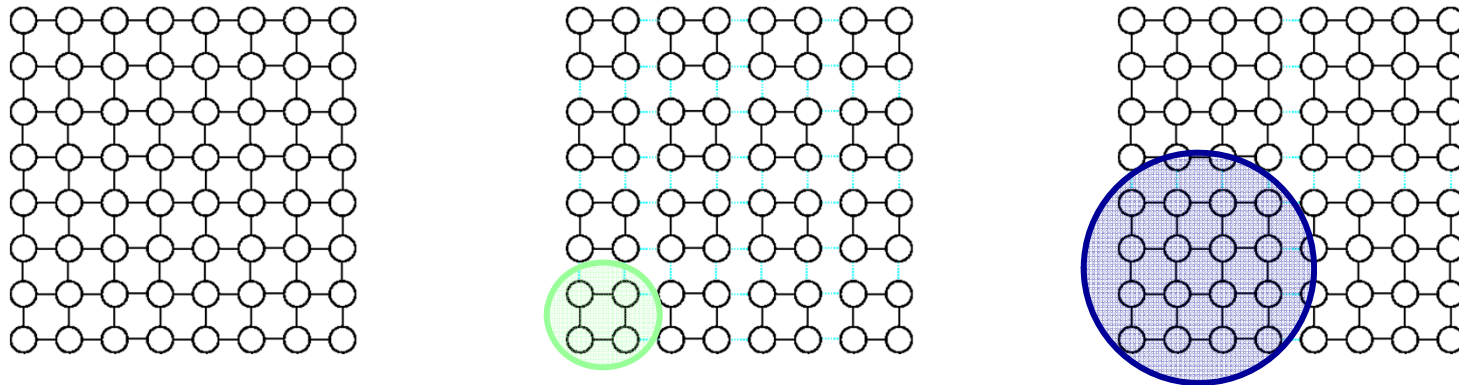# Example 1: Generalized MF approximations to Ising models

Cluster marginal of a square block $C_k$:

$$q(X_{C_k}) \propto \exp\left\{ \sum_{i,j \in C_k} \theta_{ij} X_i X_j + \sum_{i \in C_k} \theta_{i0} X_i + \sum_{\substack{i \in C_k, j \in MB_k, \\ k' \in MBC_k}} \theta_{ij} X_i \left\langle X_j \right\rangle_{q(X_{C_{k'}})} \right\}$$
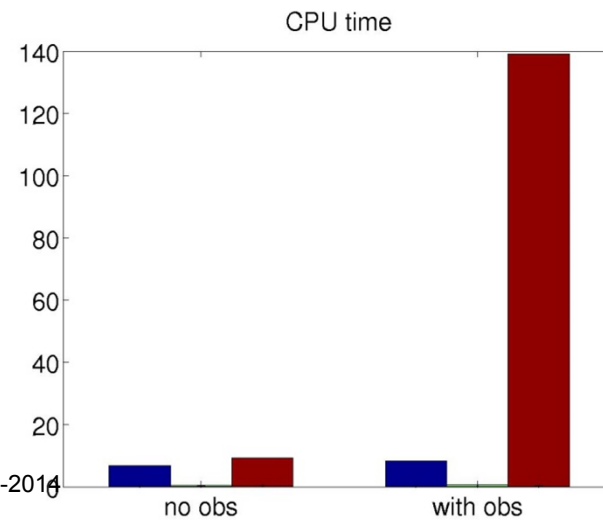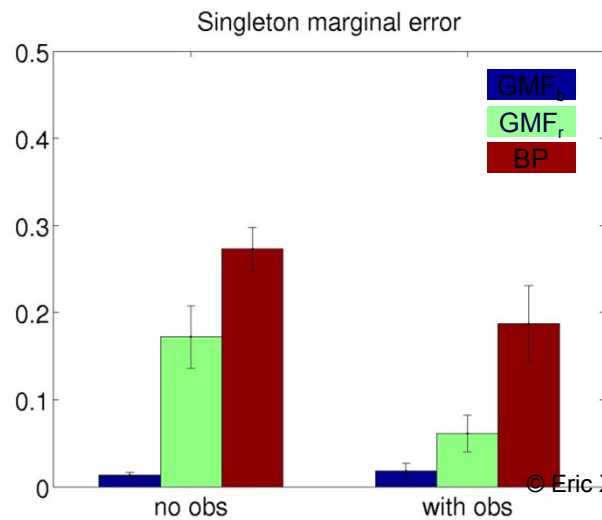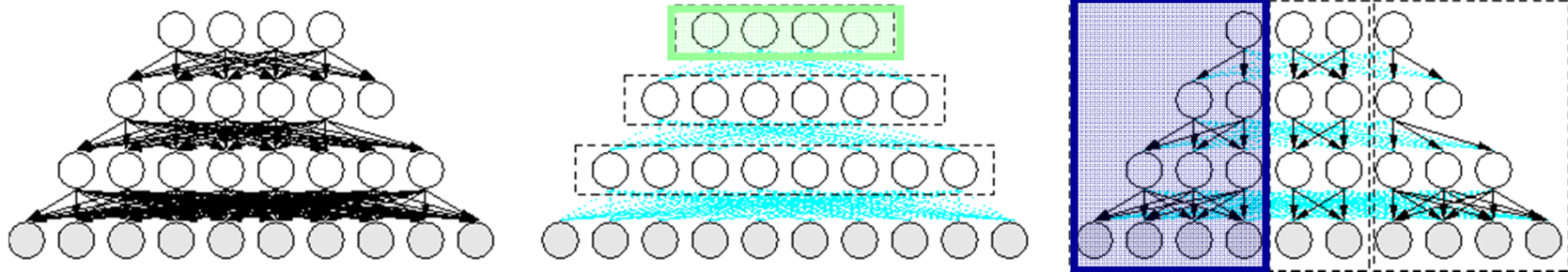
Virtually a reparameterized Ising model of small size.

# GMF approximation to Ising models



Singleton marginal error

CPU time

- GMF$_{4x4}$
- GMF$_{2x2}$
- BP

Attractive coupling: positively weighted
Repulsive coupling: negatively weighted

# Example 2: Sigmoid belief network

# Example 3: Factorial HMM

19

# Automatic Variational Inference



fHMM        Mean field approx.        Structured variational approx.

- Currently for each new model we have to
  - derive the variational update equations
  - write application-specific code to find the solution

- Each can be time consuming and error prone

- Can we build a general-purpose inference engine which automates these procedures?

# Probabilistic Topic Models



- Humans cannot afford to deal with (e.g., search, browse, or measure similarity) a huge number of text documents
- We need computers to help out …

# How to get started?

- **Here are some important elements to consider before you start:**
  - Task:
    - Embedding? Classification? Clustering? Topic extraction? …
  - Data representation:
    - Input and output (e.g., continuous, binary, counts, …)
  - Model:
    - BN? MRF? Regression? SVM?
  - Inference:
    - Exact inference? MCMC? Variational?
  - Learning:
    - MLE? MCLE? Max margin?
  - Evaluation:
    - Visualization? Human interpretability? Perperlexity? Predictive accuracy?

- **It is better to consider one element at a time!**

# Tasks: document embedding

- Say, we want to have a mapping ..., so that

⇒

- Compare similarity
- Classify contents
- Cluster/group/categorizing
- Distill semantics and perspectives
- ..

# Summarizing the data using topics

| Bayesian modeling | Visual cortex | Education | Market |
|---|---|---|---|
| Bayesian | cortex | students | market |
| model | cortical | education | economic |
| inference | areas | learning | financial |
| models | visual | educational | economics |
| probability | area | teaching | markets |
| probabilistic | primary | school | returns |
| Markov | connections | student | price |
| prior | ventral | skills | stock |
| hidden | cerebral | teacher | value |
| approach | sensory | academic | investment |

# See how data changes over time

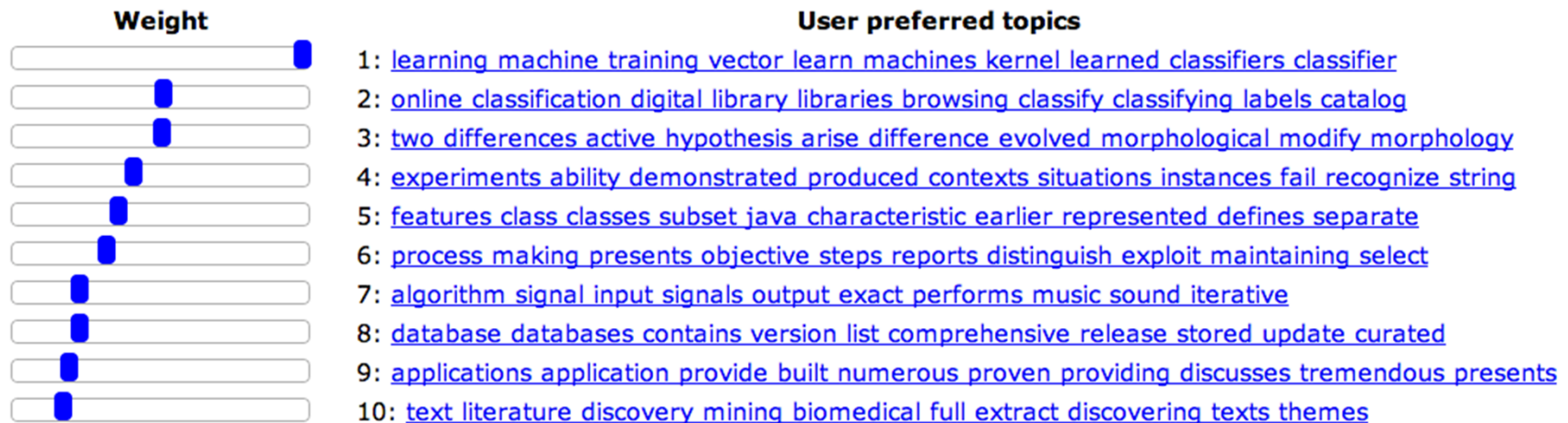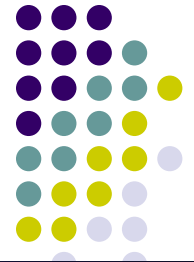| 2/27/2007 | 4/24/2007 | 6/26/2007 | 8/28/2007 | 10/23/2007 | 12/25/2007 | 2/19/2008 |
|---|---|---|---|---|---|---|
| healthcare | healthcare | healthcare | healthcare | kucinich | obama | obama |
| abc | abc | wisconsin | wisconsin | ron | clinton | clinton |
| wisconsin | wisconsin | vegas | vegas | obama | paul | hillary |
| vegas | vegas | superdelegate | superdelegate | healthcare | ron | barack |
| superdelegate | superdelegate | nevada | kucinich | paul | kucinich | campaign |
| nevada | nevada | abc | nevada | wisconsin | hillary | democratic |
| delegate | delegate | fundraising | fundraising | vegas | iowa | iowa |
| civil | civil | delegate | delegate | superdelegate | campaign | kucinich |
| recount | fundraising | civil | florida | iowa | new | paul |
| florida | recount | florida | civil | nevada | barack | ron |

# User interest modeling using topics

**User interest profile (adjustable with sliders---Changing these changes recommendations.)**

**Weight**

**User preferred topics**

1: learning machine training vector learn machines kernel learned classifiers classifier

2: online classification digital library libraries browsing classify classifying labels catalog

3: two differences active hypothesis arise difference evolved morphological modify morphology

4: experiments ability demonstrated produced contexts situations instances fail recognize string

5: features class classes subset java characteristic earlier represented defines separate

6: process making presents objective steps reports distinguish exploit maintaining select

7: algorithm signal input signals output exact performs music sound iterative

8: database databases contains version list comprehensive release stored update curated

9: applications application provide built numerous proven providing discusses tremendous presents

10: text literature discovery mining biomedical full extract discovering texts themes

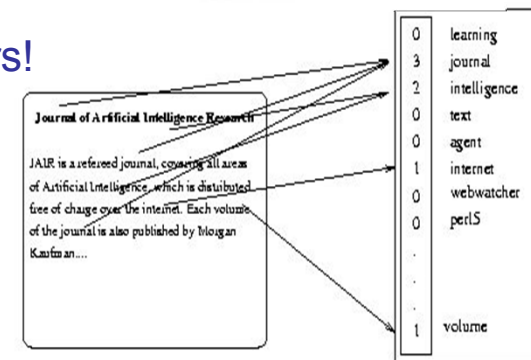**http://cogito-demos.ml.cmu.edu/cgi-bin/recommendation.cgi**

# Representation:

- ## Data: **Bag of Words Representation**

As for the Arabian and Palestinean voices that are against the current negotiations and the so-called peace process, they are not against peace per se, but rather for their well-founded predictions that Israel would NOT give an inch of the West bank (and most probably the same for Golan Heights) back to the Arabs. An 18 months of "negotiations" in Madrid, and Washington proved these predictions. Now many will jump on me saying why are you blaming israelis for no-result negotiations. I would say why would the Arabs stall the negotiations, what do they have to loose ?

Arabian
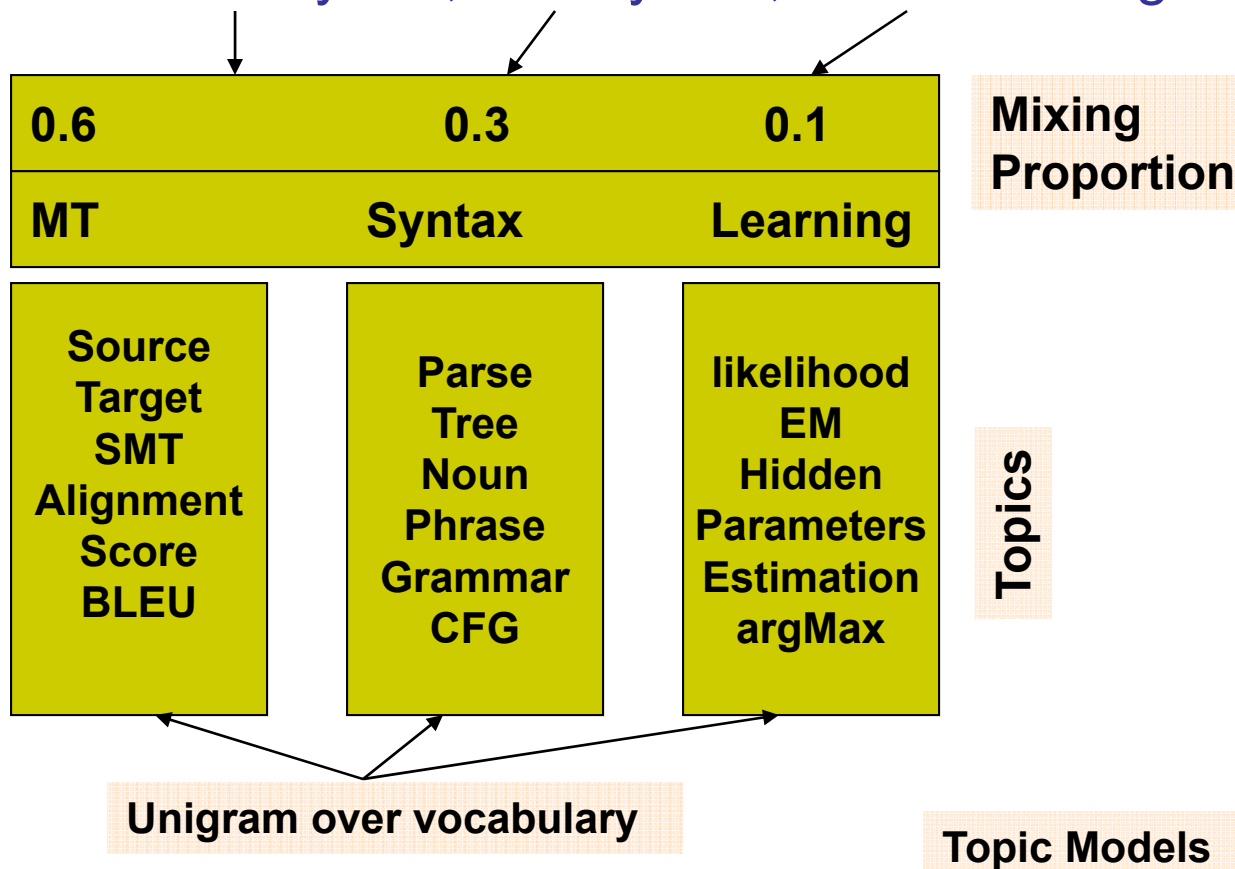negotiations
against
peace
Israel
Arabs
blaming

- Each document is a vector in the word space
- Ignore the order of words in a document. Only count matters!

- A high-dimensional and sparse representation $(|V| \gg D)$
  - Not efficient text processing tasks, e.g., search, document classification, or similarity measure
  - Not effective for browsing

| 0 | learning |
| 3 | journal |
| 2 | intelligence |
| 0 | text |
| 0 | agent |
| 1 | internet |
| 0 | webwatcher |
| 0 | perlS |
| . | |
| . | |
| 1 | volume |

# How to Model Semantic?

- Q: What is it about?

- A: Mainly MT, with syntax, some learning

| 0.6 | 0.3 | 0.1 | Mixing Proportion |
|---|---|---|---|
| MT | Syntax | Learning | |
| **Source Target SMT Alignment Score BLEU** | **Parse Tree Noun Phrase Grammar CFG** | **likelihood EM Hidden Parameters Estimation argMax** | **Topics** |

**A Hierarchical Phrase-Based Model for Statistical Machine Translation**
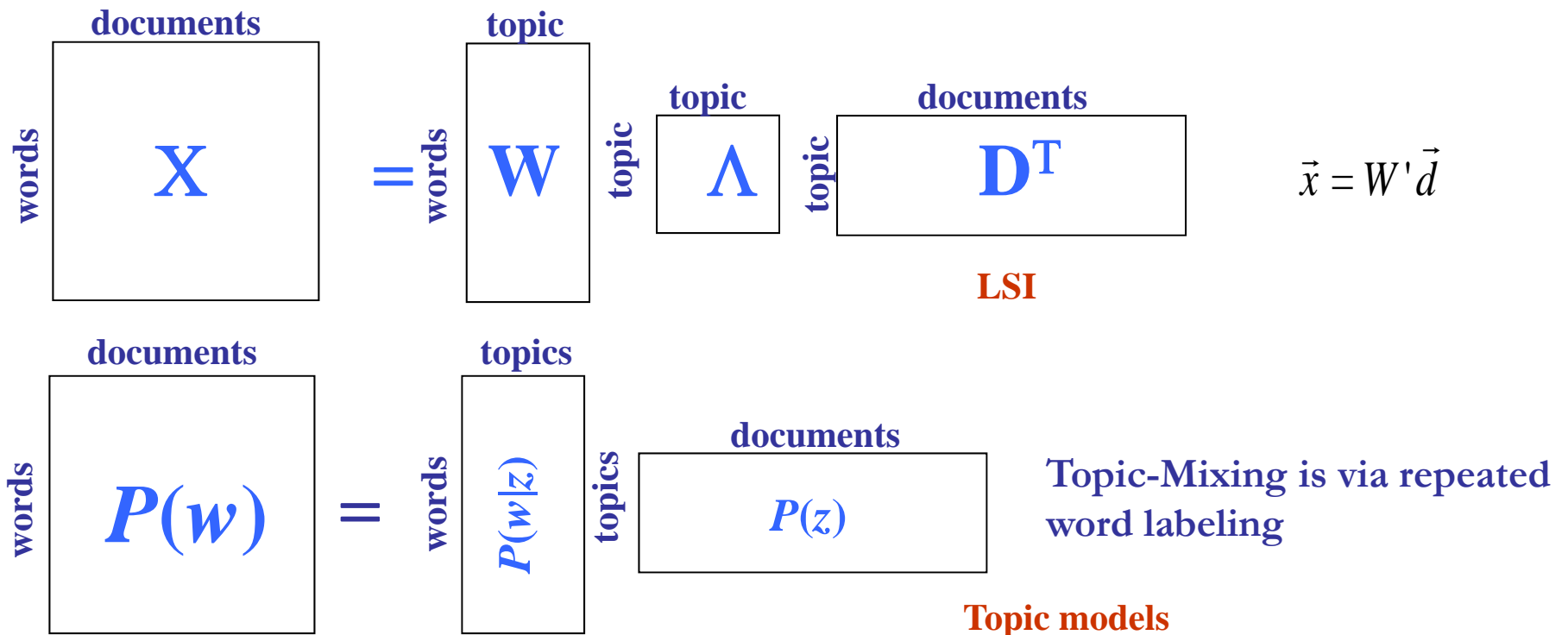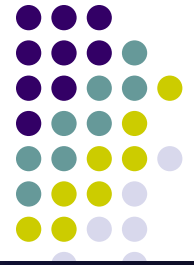
We present a statistical phrase-based Translation model that uses *hierarchical phrases*—phrases that contain sub-phrases. The model is formally a synchronous context-free grammar but is learned from a bitext without any syntactic information. Thus it can be seen as a shift to the *formal* machinery 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.
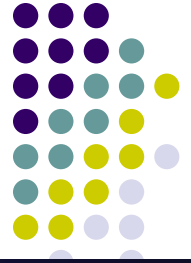
**Unigram over vocabulary**

**Topic Models**

# 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 …

# Topic Models: The Big Picture

Unstructured Collection

Structured Topic Network

Topic Discovery

Word Simplex

$W_1$

$W_n$

$W_2$

Dimensionality Reduction

Topic Simplex

$T_1$

$T_k$

$T_2$

# LSI versus Topic Model (probabilistic LSI)



documents

words $\mathbf{X}$ $=$ words $\mathbf{W}$ topic $\quad$ topic $\mathbf{\Lambda}$ topic $\quad$ documents $\mathbf{D^T}$

topic

$$\vec{x} = W'\vec{d}$$

**LSI**

documents

words $P(w)$ $=$ words $P(w|z)$ topics $\quad$ topics documents $P(z)$

Topic-Mixing is via repeated word labeling

**Topic models**

# Words in Contexts

- " It was a nice **shot**. "

# Words in Contexts (con'd)

- the opposition Labor **Party** fared even worse, with a predicted 35 **seats**, seven less than last **election**.

# "Words" in Contexts (con'd)

Sivic et al. ICCV 2005

# Admixture Models

- **Objects are bags of elements**

- **Mixtures are distributions over element**

- **Objects have mixing vector $\theta$**
  - **Represents each mixtures' contributions**

- **Object is generated as follows:**
  - **Pick a mixture component from $\theta$**
  - **Pick an element from that component**

| 0.1 | 0.1 | ….. | 0.5 |
| 0.1 | 0.5 | ….. | 0.1 |
| 0.5 | 0.1 | ….. | 0.1 |

# Topic Models

## Generating a document

– *Draw $\theta$ from the prior*

For each word $n$

   - Draw $z_n$ from *multinomia l$(\theta)$*

   - Draw $w_n \mid z_n, \{\beta_{1:k}\}$ from *multinomia l$(\beta_{z_n})$*

**Prior**

$\theta$

$z$

$\beta$    $\rightarrow$ w

K

$N_d$

N

**Which prior to use?**

# Choices of Priors

- ## 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

# Generative Semantic of LoNTAM

Generating a document

– *Draw θ from the prior*

For each word *n*

  - Draw $z_n$ from *multinomial* $(\theta)$

  - Draw $w_n \mid z_n, \{\beta_{1:k}\}$ from *multinomial* $(\beta_{z_n})$

$$\theta \sim LN_K(\mu, \Sigma)$$

$$\gamma \sim N_{K-1}(\mu, \Sigma) \qquad \gamma_K = 0$$

$$\theta_i = \exp\left\{\gamma_i - \log\left(1 + \sum_{i=1}^{K-1} e^{\gamma_i}\right)\right\}$$

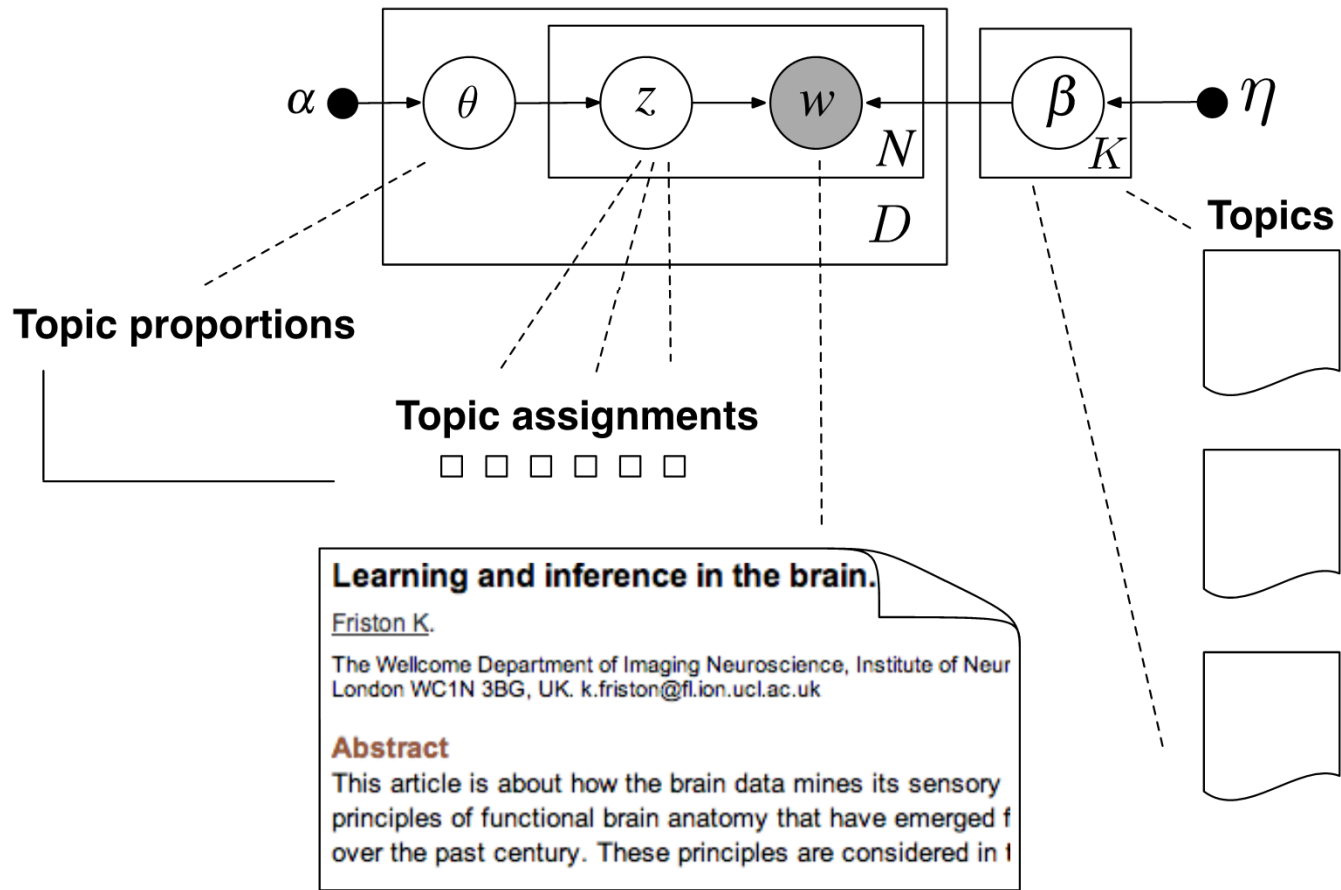$$C(\gamma) = \log\left(1 + \sum_{i=1}^{K-1} e^{\gamma_i}\right)$$

**Problem**

- **Log Partition Function**
- **Normalization Constant**

μ    Σ

γ

β   K

z

w

$N_d$

N

# Posterior inference



Topic proportions

Topic assignments

Learning and inference in the brain.

Friston K.

The Wellcome Department of Imaging Neuroscience, Institute of Neur
London WC1N 3BG, UK. k.friston@fl.ion.ucl.ac.uk

Abstract
This article is about how the brain data mines its sensory
principles of functional brain anatomy that have emerged f
over the past century. These principles are considered in

Topics

# Posterior inference results



© Eric Xing @ CMU, 2005-2013

# Joint likelihood of all variables

$$p(\beta, \theta, z, w) = \prod_{k=1}^{K} p(\beta_k | \eta) \prod_{d=1}^{D} p(\theta_d | \alpha) \prod_{n=1}^{N} p(z_{dn} | \theta_d) p(w_{dn} | z_{dn}, \beta)$$



**We are interested in computing the posterior, and the data likelihood!**

# Inference and Learning are both intractable

- A possible query:

$$p(\theta_n \mid D) = ?$$

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

- Close form solution?

$$p(\theta_n \mid D) = \frac{p(\theta_n, D)}{p(D)}$$

$$= \frac{\sum\limits_{\{z_{n,m}\}} \int \left( \prod\limits_n \left( \prod\limits_m p(w_{n,m} \mid \beta_{z_n}) p(z_{n,m} \mid \theta_n) \right) p(\theta_n \mid \alpha) \right) p(\beta \mid \eta) d\theta_{-i} d\beta}{p(D)}$$

$$p(D) = \sum\limits_{\{z_{n,m}\}} \int \cdots \int \left( \prod\limits_n \left( \prod\limits_m p(x_{n,m} \mid \beta_{z_n}) p(z_{n,m} \mid \theta_n) \right) p(\theta_n \mid \alpha) \right) p(\beta \mid \eta) d\theta_1 \cdots d\theta_N d\beta$$

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

- Learning: What to learn? What is the objective function?

# Approximate Inference

- Variational Inference

  - Mean field approximation (Blei et al)

  - Expectation propagation (Minka et al)

  - Variational $2^{nd}$-order Taylor approximation (Xing)

- Markov Chain Monte Carlo

  - Gibbs sampling (Griffiths et al)

# Mean-field assumption

- True posterior

$$p(\beta, \theta, \boldsymbol{z} | \boldsymbol{w}) = \frac{p(\beta, \theta, \boldsymbol{z}, \boldsymbol{w})}{p(\boldsymbol{w})}$$

- Break the dependency using the fully factorized distribution

$$q(\beta, \theta, \boldsymbol{z}) = \prod_k q(\beta_k) \prod_d q(\theta_d) \prod_n q(z_{dn})$$

- Mean-field family usually does NOT include the true posterior.

# Update each marginals

- Update

$$q(\theta_d) \propto \exp \left\{ \mathbb{E}_{\prod_n q(z_{dn})} \left[ \log p(\theta_d | \alpha) + \sum_n \log p(z_{dn} | \theta_d) \right] \right\}$$

- In LDA,

$$p(\theta_d | \alpha) \propto \exp \left\{ \sum_{k=1}^{K} (\alpha_k - 1) \log \theta_{dk} \right\} --\text{Dirichlet}$$

$$p(z_{dn} | \theta_d) = \exp \left\{ \sum_{k=1}^{K} 1[z_{dn} = k] \log \theta_{dk} \right\} --\text{Multinomial}$$

- We obtain

$$q(\theta_d) \propto \exp \left\{ \sum_{k=1}^{K} \left( \sum_{n=1}^{N} q(z_{dn} = k) + \alpha_k - 1 \right) \log \theta_{dk} \right\}$$

> **This is also a Dirichlet---the same as its prior!**

# Coordinate ascent algorithm for LDA

1: Initialize variational topics $q(\beta_k)$, $k = 1, ..., K$.

2: **repeat**

3:     **for** each document $d \in \{1, 2, ..., D\}$ **do**

4:         Initialize variational topic assigments $q(z_{dn})$, $n = 1, ..., N$

5:         **repeat**

6:             Update variational topic proportions $q(\theta_d)$

7:             Update variational topic assigments $q(z_{dn})$, $n = 1, ..., N$

8:         **until** Change of $q(\theta_d)$ is small enough

9:     **end for**

0:     Update variational topics $q(\beta_k)$, $k = 1, ..., K$.

1: **until** Lower bound $L(q)$ converges

# Choice of q() does matter



$$P(\gamma, \{z\} | D)$$

$$q(\gamma, z_{1:n}) = q(\gamma | \mu^*, \Sigma^*) \prod q(z_n | \phi_n)$$

**Σ* is full matrix**

**Σ* is assumed to be diagonal**

| Multivariate Quadratic Approx. | Log Partition Function | Tangent Approx. |
|---|---|---|

**Closed Form Solution for μ*, Σ***

$$\log \left( 1 + \sum_{i=1}^{K-1} e^{\gamma_i} \right)$$

**Numerical Optimization to fit μ*, Diag(Σ*)**

**Ahmed&Xing**

**Blei&Lafferty**

# Tangent Approximation

# How to evaluate?

- Empirical Visualization: e.g., topic discovery on New York Times

The 5 most frequent topics from the HDP on the *New York Times.*

| game | life | film | book | wine |
|------|------|------|------|------|
| season | know | movie | life | street |
| team | school | show | books | hotel |
| coach | street | life | novel | house |
| play | man | television | story | room |
| points | family | films | man | night |
| games | says | director | author | place |
| giants | house | man | house | restaurant |
| second | children | story | war | park |
| players | night | says | children | garden |

# How to evaluate?

- **Test on Synthetic Text where ground truth is known:**

# Comparison: accuracy and speed
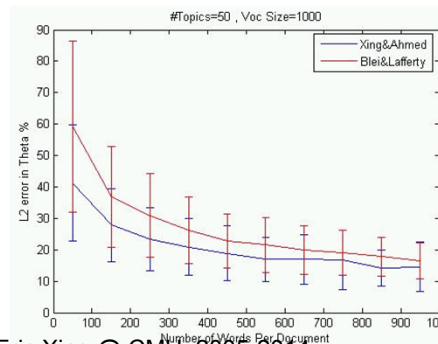
L2 error in topic vector est.
and # of iterations

- Varying Num. of Topics
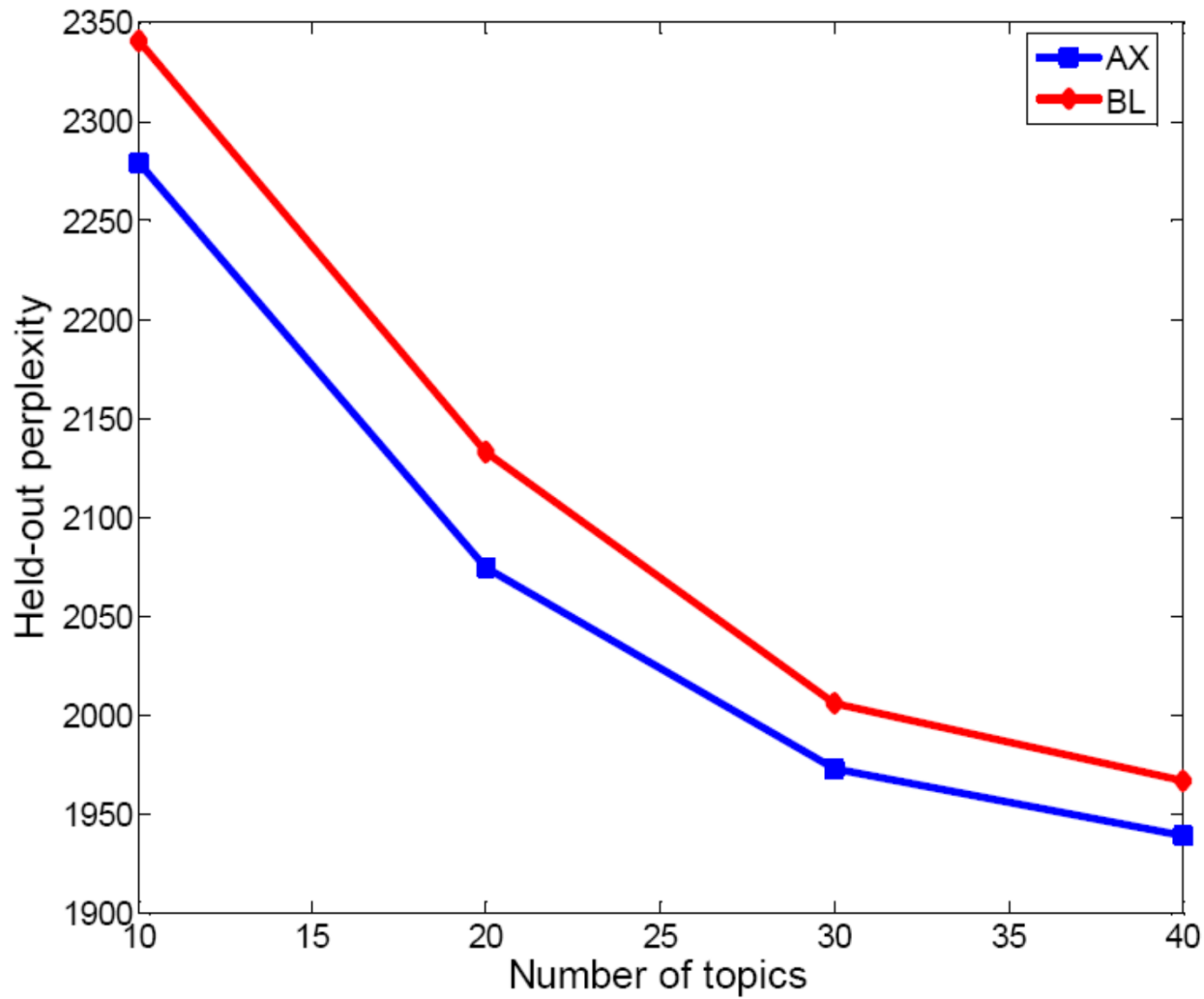
- Varying Voc. Size

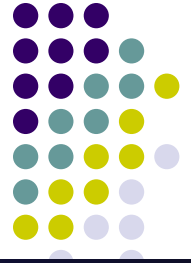- Varying Num. Words Per
Document

# Comparison: perplexity
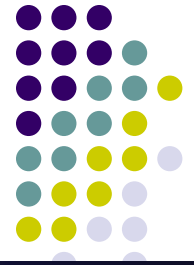
# Classification Result on PNAS collection

- 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

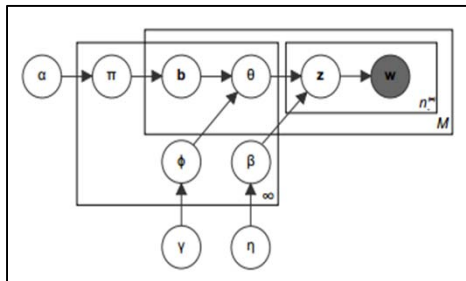| Category | Doc | BL | AX |
|---|---|---|---|
| Genetics | 21 | 61.9 | 61.9 |
| Biochemistry | 86 | 65.1 | 77.9 |
| Immunology | 24 | 70.8 | 66.6 |
| Biophysics | 15 | 53.3 | 66.6 |
| Total | 146 | 64.3 | 72.6 |

-Notable Difference
-Examine the low dimensional representations below
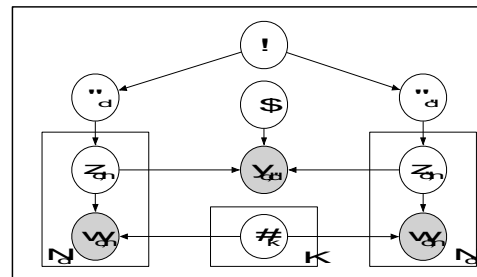
# What makes topic models useful -
# -- The Zoo of Topic Models!
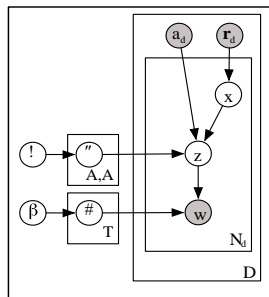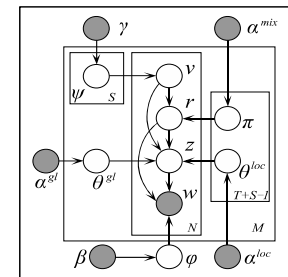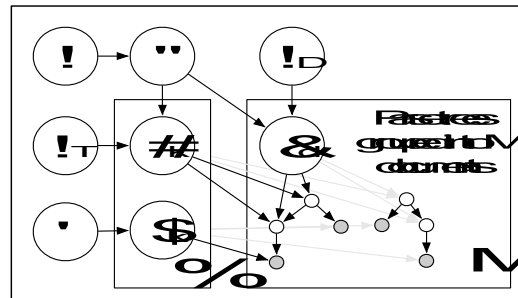
- It is a building block of many models.
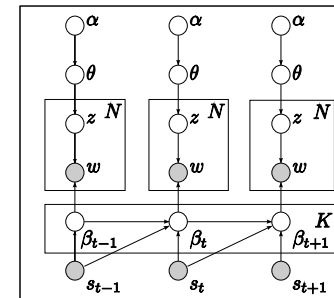
**Williamson et al. 2010**     **Chang & Blei, 2009**     **Titov & McDonald, 2008**



**McCallum et al. 2007**     **Boyd-Graber & Blei, 2008**     **Wang & Blei, 2008**

# Conclusion

- GM-based topic models are cool
  - Flexible
  - Modular
  - Interactive
- There are many ways of implementing topic models
  - unsupervised
  - supervised
- Efficient Inference/learning algorithms
  - GMF, with Laplace approx. for non-conjugate dist.
  - MCMC
- Many applications
  - …
  - Word-sense disambiguation
  - Image understanding
  - Network inference

# Summary on VI

- Variational methods in general turn inference into an optimization problem via exponential families and convex duality

- The exact variational principle is intractable to solve; there are two distinct components for approximations:
  - Either inner or outer bound to the marginal polytope
  - Various approximation to the entropy function

- <u>Mean field</u>: non-convex inner bound and exact form of entropy
- <u>BP</u>: polyhedral outer bound and non-convex Bethe approximation
- <u>Kikuchi and variants</u>: tighter polyhedral outer bounds and better entropy approximations (Yedidia et. al. 2002)