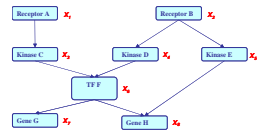


Approximate Inference: Mean Field Methods

Probabilistic Graphical Models (10-708)

Lecture 17, Nov 12, 2007



Eric Xing

Reading: KF-Chap. 12

1



- Questions????
 - Kalman Filters
 - Complex models
 - LBP-Bethe Minimization

Approximate Inference



Variational Methods



- For a distribution $p(X/\theta)$ associated with a complex graph, computing the marginal (or conditional) probability of arbitrary random variable(s) is intractable
- Variational methods
 - formulating probabilistic inference as an optimization problem:

$$e.g. \quad f^* = \arg \max_{f \in \mathcal{S}} \{ F(f) \}$$

f : a (tractable) probability distribution
or, solutions to certain probabilistic queries

Exponential Family



- Exponential representation of graphical models:

$$P(\mathbf{X}) = \frac{1}{Z} \prod_{c \in \mathcal{C}} \psi_c(\mathbf{X}_c) \Rightarrow p(\mathbf{X} | \boldsymbol{\theta}) = \exp \left\{ \sum_{\alpha} \theta_{\alpha} \phi_{\alpha}(\mathbf{X}_{D_{\alpha}}) - A(\boldsymbol{\theta}) \right\}$$

- Includes discrete models, Gaussian, Poisson, exponential, and many others

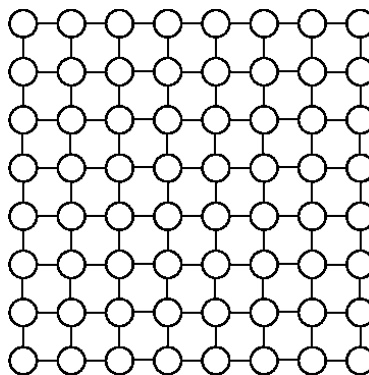
$E(\mathbf{X}) = -\sum_{\alpha} \theta_{\alpha} \phi_{\alpha}(\mathbf{X}_{D_{\alpha}})$ is referred to as the *energy* of state \mathbf{x}

$$\begin{aligned} \Rightarrow p(\mathbf{X} | \boldsymbol{\theta}) &= \exp \{ -E(\mathbf{X}) - A(\boldsymbol{\theta}) \} \\ &= \exp \{ -E(\mathbf{X}_H, \mathbf{x}_E) - A(\boldsymbol{\theta}, \mathbf{x}_E) \} \end{aligned}$$

Eric Xing

5

Example: the Boltzmann distribution on atomic lattice

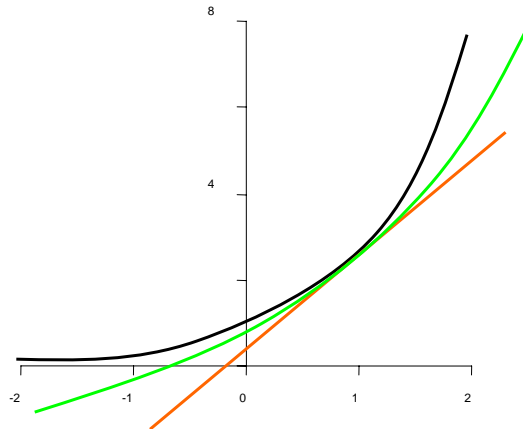


$$p(X) = \frac{1}{Z} \exp \left\{ \sum_{i < j} \theta_{ij} X_i X_j + \sum_i \theta_{i0} X_i \right\}$$

Eric Xing

6

Lower bounds of exponential functions



$$\exp(x) \geq \exp(\mu)(x - \mu + 1)$$

$$\exp(x) \geq \frac{1}{6} \exp(\mu) \left((x - \mu)^3 + 3(x - \mu)^2 + 6(x - \mu + 1) \right)$$

Eric Xing

7

Lower bounding likelihood



Representing $q(\mathbf{X}_H)$ by $\exp\{-E'(\mathbf{X}_H)\}$:

Lemma: Every marginal distribution $q(\mathbf{X}_H)$ defines a lower bound of likelihood:

$$p(\mathbf{x}_E) \geq \int d\mathbf{x}_H \exp\{-E'(\mathbf{x}_H)\} \\ (1 - A(\mathbf{x}_E) - (E(\mathbf{x}_H, \mathbf{x}_E) - E'(\mathbf{x}_H))),$$

where \mathbf{x}_E denotes observed variables (evidence).

Upgradeable to higher order bound [Leisink and Kappen, 2000]

Eric Xing

8

Lower bounding likelihood



Representing $q(\mathbf{X}_H)$ by $\exp\{-E'(\mathbf{X}_H)\}$:

Lemma: Every marginal distribution $q(\mathbf{X}_H)$ defines a lower bound of likelihood:

$$\begin{aligned} p(\mathbf{x}_E) &\geq C - \langle E(\mathbf{X}_H, \mathbf{x}_E) \rangle_{q(\mathbf{X}_H)} - \int d\mathbf{x}_H q(\mathbf{x}_H) \log q(\mathbf{x}_H) \\ &= C - \langle E \rangle_q + H_q, \end{aligned}$$

where \mathbf{x}_E denotes observed variables (evidence).

$\langle E \rangle_q$: expected energy $\langle E \rangle_q - H_q$: Gibbs free energy
 H_q : entropy

Eric Xing

9

KL and variational (Gibbs) free energy



- Kullback-Leibler Distance:

$$KL(q \parallel p) \equiv \sum_z q(z) \ln \frac{q(z)}{p(z)}$$

- “Boltzmann’s Law” (definition of “energy”):

$$p(z) = \frac{1}{C} \exp[-E(z)]$$

$$KL(q \parallel p) \equiv \underbrace{\sum_z q(z) E(z) + \sum_z q(z) \ln q(z) + \ln C}_{\text{Gibbs Free Energy } G(q)}$$

minimized when $q(Z) = p(Z)$

Eric Xing

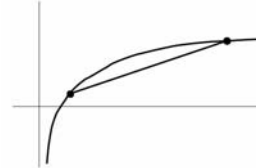
10

KL and Log Likelihood



- Jensen's inequality

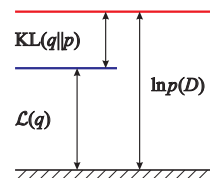
$$\begin{aligned} \mathcal{L}(\theta; x) &= \log p(x|\theta) \\ &= \log \sum_z p(x, z|\theta) \\ &= \log \sum_z q(z|x) \frac{p(x, z|\theta)}{q(z|x)} \\ &\geq \sum_z q(z|x) \log \frac{p(x, z|\theta)}{q(z|x)} \end{aligned}$$



$$\Rightarrow \mathcal{L}(\theta; x) \geq \langle \mathcal{L}_c(\theta; x, z) \rangle_q + H_q = \mathcal{L}(q)$$

- KL and Lower bound of likelihood

$$\begin{aligned} \mathcal{L}(\theta; x) &= \log p(x|\theta) = \log \frac{p(x, z|\theta)}{p(z|x, \theta)} = \sum_z q(z) \log \frac{p(x, z|\theta)}{p(z|x, \theta)} \\ &= \sum_z q(z) \log \frac{p(x, z|\theta)}{q(z)} - \sum_z q(z) \log \frac{q(z)}{p(z|x, \theta)} \\ &= \sum_z q(z) \log \frac{p(x, z|\theta)}{q(z)} + \sum_z q(z) \log \frac{q(z)}{p(z|x, \theta)} \end{aligned}$$



$$\Rightarrow \mathcal{L}(\theta; x) = \mathcal{L}(q) + KL(q \| p)$$

- Setting $q()=p(z|x)$ closes the gap (c.f. EM)

Eric Xing

11

A variational representation of probability distributions



$$\begin{aligned} q &= \arg \max_{q \in \mathcal{Q}} \{ -\langle E \rangle_q + H_q \} \\ &= \arg \min_{q \in \mathcal{Q}} \{ \langle E \rangle_q - H_q \} \end{aligned}$$

where \mathcal{Q} is the equivalent sets of realizable distributions, e.g., all valid parameterizations of exponential family distributions, marginal polytopes [winright *et al.* 2003].

Difficulty: H_q is intractable for general q

“solution”: approximate H_q
and/or,
relax or tighten \mathcal{Q}

Eric Xing

12



Bethe Free Energy/LBP

- But we do not optimize $q(\mathbf{X})$ explicitly, focus on the set of beliefs
 - *e.g.*, $b = \{b_{i,j} = \tau(x_i, x_j), b_i = \tau(x_i)\}$
 - Relax the optimization problem
 - approximate objective: $H_{\text{Bethe}} = H(b_{i,j}, b_i)$
 - relaxed feasible set: $\mathcal{M}_o = \{ \tau \geq 0 \mid \sum_{x_i} \tau(x_i) = 1, \sum_{x_i} \tau(x_i, x_j) = \tau(x_j) \}$
- $$b^* = \arg \min_{b \in \mathcal{M}_o} \left\{ \langle E \rangle_b - F(b) \right\}$$
- The loopy BP algorithm:
 - a fixed point iteration procedure that tries to solve b^*

Eric Xing

13



Mean field methods

- Optimize $q(\mathbf{X}_H)$ in the space of tractable families
 - *i.e.*, subgraph of G_p over which exact computation of H_q is feasible
- Tightening the optimization space
 - exact objective: H_q
 - tightened feasible set: $Q \rightarrow \mathcal{T} \quad (\mathcal{T} \subseteq Q)$

$$q^* = \arg \min_{q \in \mathcal{T}} \langle E \rangle_q - H_q$$

Eric Xing

14



Mean Field Approximation

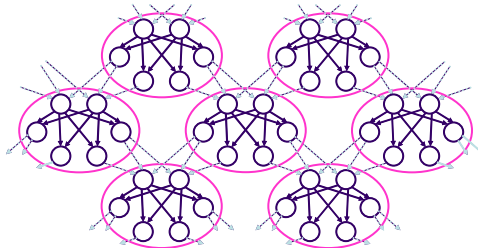
Cluster-based approx. to the Gibbs free energy

(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, $\{C_1, \dots, C_k\}$, of all variables

- Let
$$q(\mathbf{X}) = \prod_i q_i(\mathbf{X}_{C_i}),$$

- Mean-field free energy

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

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

- 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 $q_i(x_{C_i})$'s.
 - Variational calculus ...
 - Do inference in each $q_i(x_{C_i})$ using any tractable algorithm

Eric Xing

17

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:

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

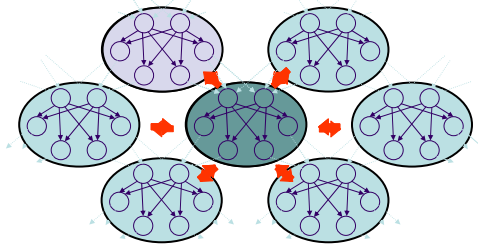
GMF algorithm: Iterate over each q_i

Eric Xing

18

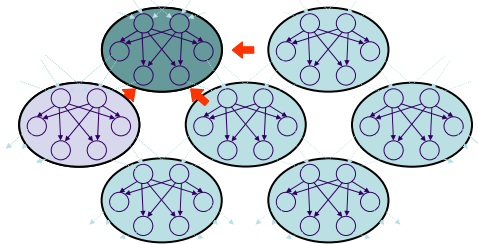
A generalized mean field algorithm

[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

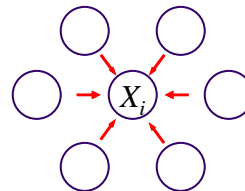


- Approximate $p(\mathbf{X})$ by fully factorized $q(\mathbf{X}) = \prod_i q_i(X_i)$
- For Boltzmann distribution $p(\mathbf{X}) = \exp\{\sum_{i < j} q_{ij} X_i X_j + \sum_i q_{i0} X_i\} / Z$:

mean field equation:

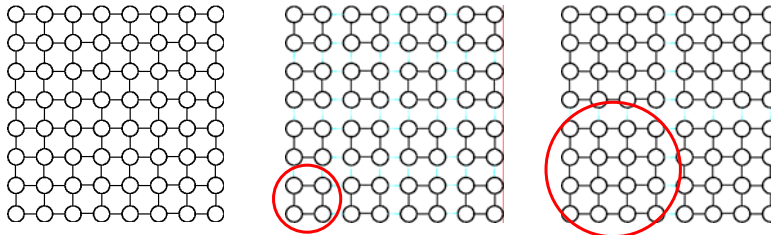
$$q_i(X_i) = \exp\left\{ \theta_{i0} X_i + \sum_{j \in \mathcal{N}_i} \theta_{ij} X_i \langle X_j \rangle_{q_j} + A_i \right\}$$

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



- $\langle X_j \rangle_{q_j}$ 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

Generalized MF approximation to Ising models



Cluster marginal of a square block C_k :

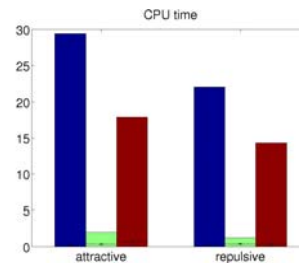
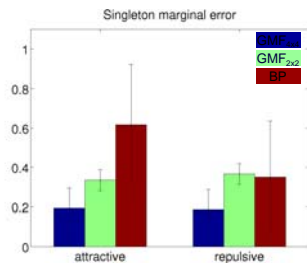
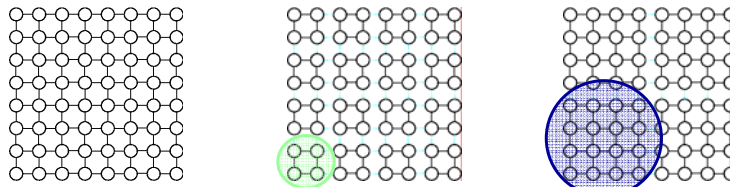
$$q(X_{C_k}) \propto \exp \left\{ \sum_{i,j \in C_k} \theta_{ij} X_i X_j + \sum_{i \in C_k} \theta_{i0} X_i + \sum_{\substack{i \in C_k, j \in MB_k \\ k \in MBC_k}} \theta_{ij} X_i \langle X_j \rangle_{q(X_{C_k})} \right\}$$

Virtually a reparameterized Ising model of small size.

Eric Xing

23

GMF approximation to Ising models

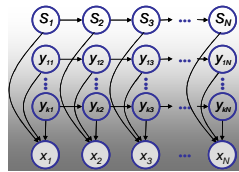


Attractive coupling: positively weighted
Repulsive coupling: negatively weighted

Eric Xing

24

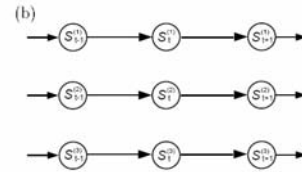
Automatic Variational Inference



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?

Eric Xing

25

Cluster-based MF (e.g., GMF)

- a general, iterative message passing algorithm
- clustering completely defines approximation
 - preserves dependencies
 - flexible performance/cost trade-off
 - clustering automatable
- recovers model-specific structured VI algorithms, including:
 - fHMM, LDA
 - variational Bayesian learning algorithms
- easily provides new structured VI approximations to complex models

Eric Xing

26

Example 1: Bayesian Gaussian Model



- Likelihood function

$$p(D|\mu, \tau) = \left(\frac{\tau}{2\pi}\right)^{N/2} \exp\left\{-\frac{\tau}{2} \sum_{n=1}^N (x_n - \mu)^2\right\}$$

mean

precision (inverse variance)

- Conjugate priors

$$p(\mu|\mu_0, \lambda_0) = \mathcal{N}(\mu|\mu_0, \lambda_0^{-1})$$

$$p(\tau|a_0, b_0) = \mathcal{G}(\tau|a_0, b_0)$$

- Factorized variational distribution

$$q(\mu, \tau) = q(\mu)q(\tau)$$

Eric Xing

27

Variational Posterior Distribution



$$q(\mu) = \mathcal{N}(\mu|\mu_N, \lambda_N^{-1})$$

$$q(\tau) = \mathcal{G}(\tau|a_N, b_N)$$

where

$$\mu_N = \frac{\lambda_0 \mu_0 + \langle \tau \rangle N \bar{x}}{\lambda_0 + N \langle \tau \rangle}$$

$$\lambda_N = \lambda_0 + N \langle \tau \rangle$$

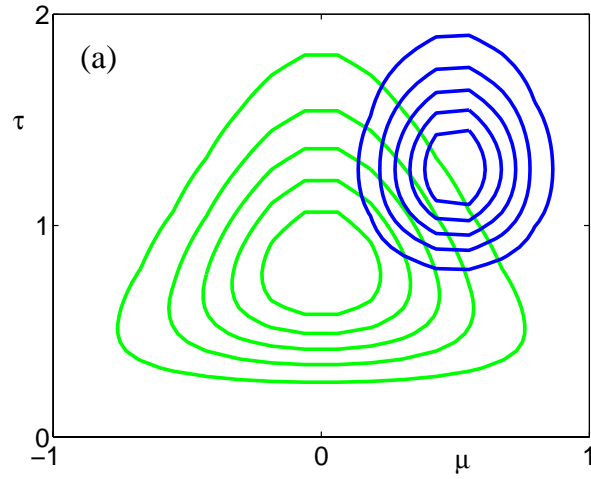
$$a_N = a_0 + \frac{N}{2}$$

$$b_N = b_0 + \frac{1}{2} \left\langle \sum_n (x_n - \mu)^2 \right\rangle_\mu$$

Eric Xing

28

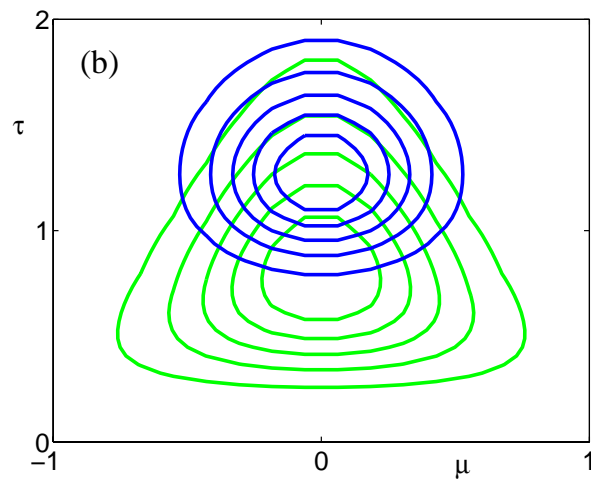
Initial Configuration



Eric Xing

29

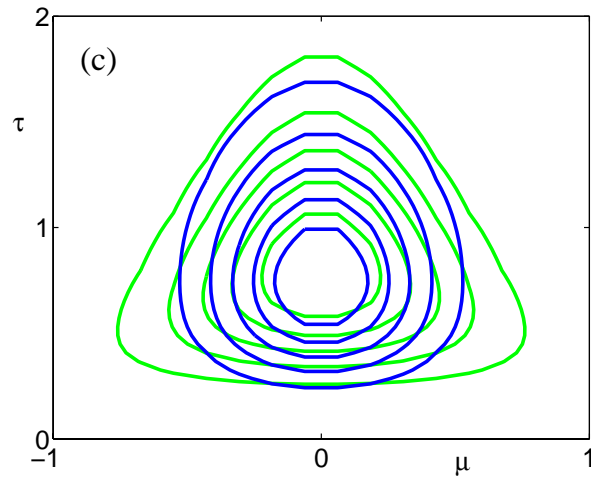
After Updating $q(\mu)$



Eric Xing

30

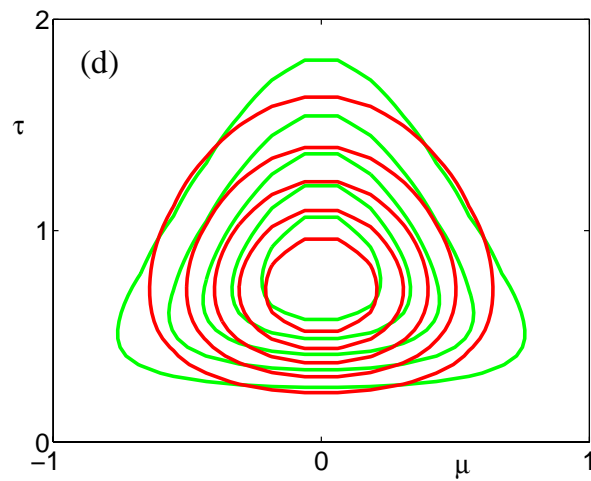
After Updating $q(\mu)$



Eric Xing

31

Converged Solution



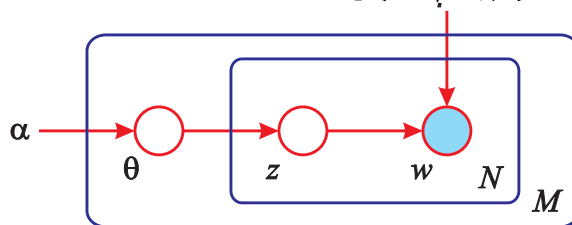
Eric Xing

32

Example 2: Latent Dirichlet Allocation



- Blei, Jordan and Ng (2003)
- Generative model of documents (but broadly applicable e.g. collaborative filtering, image retrieval, bioinformatics)
- Generative model:
 - choose $\theta \sim \text{Dir}(\alpha)$
 - choose topic $z_n \sim \text{Mult}(\theta)$
 - choose word $w_n \sim p(w_n | z_n, \beta)$



Eric Xing

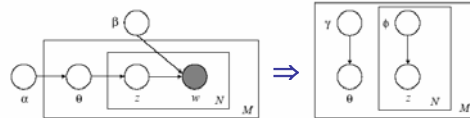
33

Latent Dirichlet Allocation



- Variational approximation

$$\begin{aligned}
 q(\theta, z) &= q_\theta(\theta)q_z(z) \\
 &= \text{Dir}(\theta | \gamma = f(\alpha, \langle z \rangle)) \times \\
 &\quad \text{Multi}(z | \phi = f(\beta_w, \langle \ln \theta \rangle))
 \end{aligned}$$



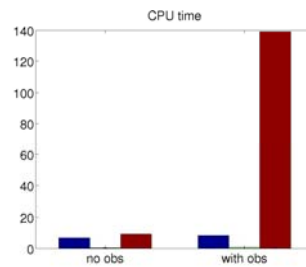
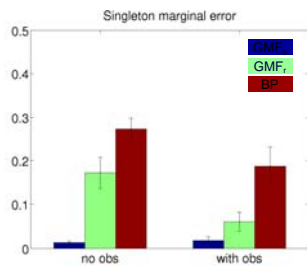
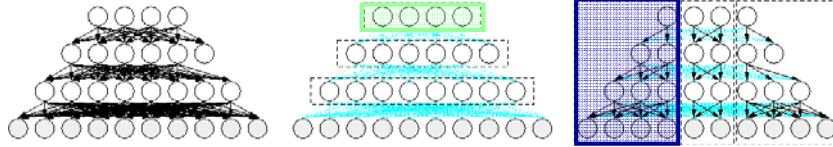
$$\begin{aligned}
 \phi_{ni} &\propto \beta_{iw_n} \exp\{E_q[\log(\theta_i) | \gamma]\} \\
 \gamma_i &= \alpha_i + \sum_{n=1}^N \phi_{ni}
 \end{aligned}$$

- Data set:
 - 15,000 documents
 - 90,000 terms
 - 2.1 million words
- Model:
 - 100 factors
 - 9 million parameters
- MCMC could be totally infeasible for this problem

Eric Xing

34

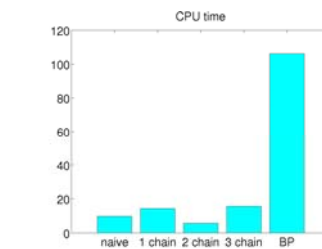
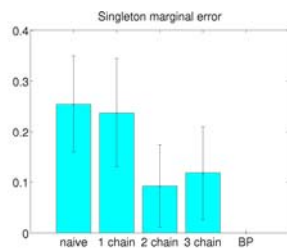
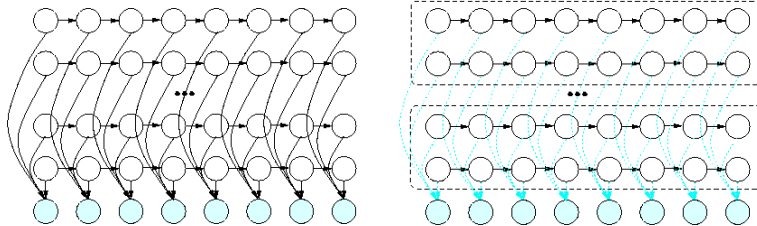
Example 3: Sigmoid belief network



Eric Xing

35

Example 4: Factorial HMM



Eric Xing

36