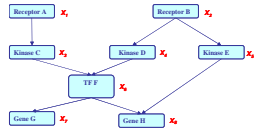


# Approximate Inference: Mean Field Methods

## Probabilistic Graphical Models (10-708)

Lecture 17, Nov 12, 2007



Eric Xing

Reading: KF-Chap. 12

- Questions????

- Kalman Filters

- Complex models

- LBP-Bethe Minimization

Handwritten notes and diagrams illustrating Kalman filters and graphical models.

Top right:  $N(x|\mu, \Sigma)$

Center:  $X \sim N(\mu, \Sigma)$  and  $Y \sim N(Ax, \phi)$

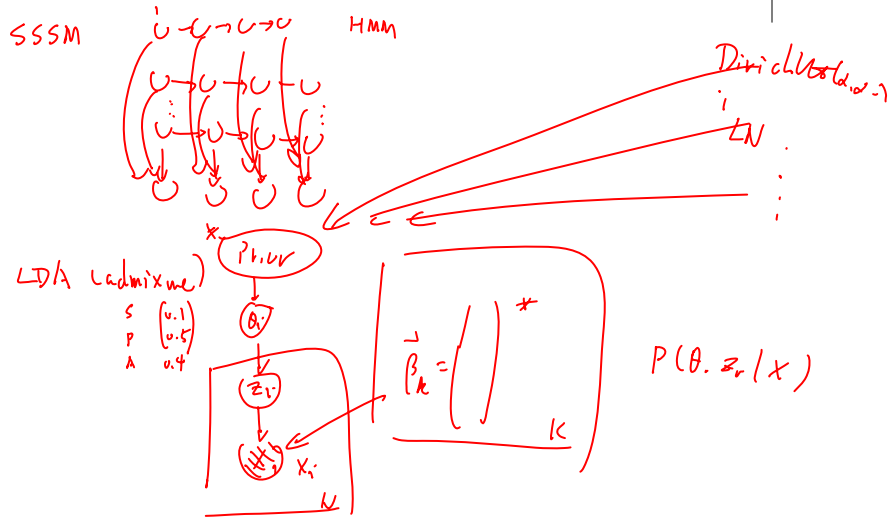
Center right: Covariance matrix  $\begin{bmatrix} \Sigma_{XX} & \Sigma_{XY} \\ \Sigma_{YX} & \Sigma_{YY} \end{bmatrix}$

Bottom left: Graphical model showing nodes  $X_t, X_{t+1}, Y_t, Y_{t+1}$  with directed edges.

Bottom right: Joint probability distributions  $P(X_{t+1} | y_{1:t+1})$ ,  $P(X_{t+1} | y_{1:t})$ , and  $P(X_{t+1} | X_{1:t})$ .

Bottom center: Matrix  $\begin{bmatrix} -1 & \\ & 1 \end{bmatrix}$  and nodes  $X_t, X_{t+1}, Y_t$ .

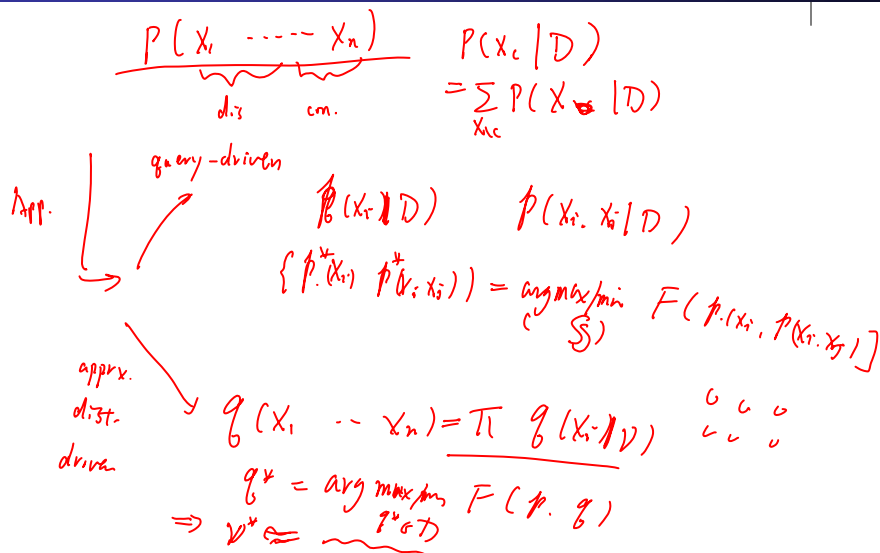
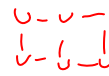
# Approximate Inference



Eric Xing

3

# Approx. Inf.



Eric Xing

4

## Variational Methods



- For a distribution  $p(\mathbf{X}/\theta)$  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

Eric Xing

5

## Exponential Family



- Exponential representation of graphical models:

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

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

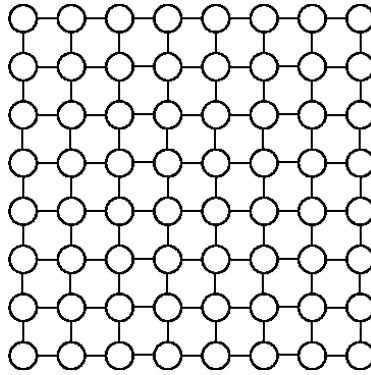
$$E(\mathbf{X}) = - \sum_{\alpha} \theta_{\alpha} \phi_{\alpha}(\mathbf{X}_{D_{\alpha}}) \text{ is referred to as the } \textit{energy} \text{ of state } \mathbf{x}$$

$$\Rightarrow \quad \underline{p(\mathbf{X} | \theta)} = \exp \{ -E(\mathbf{X}) - A(\theta) \} \\ = \exp \{ -E(\mathbf{X}_H, \mathbf{x}_E) - A(\theta, \mathbf{x}_E) \}$$

Eric Xing

6

## Example: the Boltzmann distribution on atomic lattice

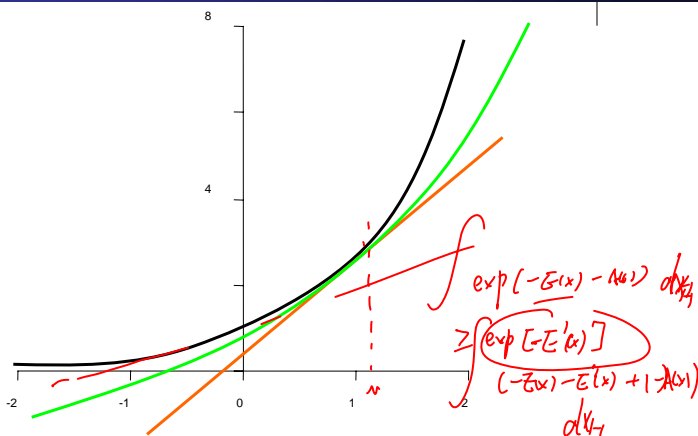


$$p(X) = \frac{1}{Z} \exp \left\{ \sum_{i < j} \theta_{ij} X_i X_j + \sum_i \theta_{i0} X_i \right\}$$

Eric Xing

7

## 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

8

## Lower bounding likelihood



Representing  $q(\mathbf{X}_H)$  by  $\exp\{-E'(\mathbf{X}_H)\}$ :  $f = \arg\min F(q, g)$

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

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

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

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

## 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

# 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”):

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

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

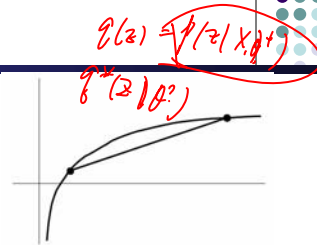
Gibbs Free Energy  $G(q)$ ;  
minimized when  $q(Z) = p(Z)$

# KL and Log Likelihood



- Jensen’s inequality

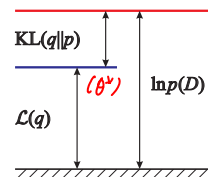
$$\begin{aligned} \mathcal{L}(\theta; x) &= \log p(x | \theta) \\ &= \log \sum_z p(x, z | \theta) \\ &= \log \sum_z q(z | x) \frac{p(x, z | \theta)}{q(z | x)} \\ &\geq \sum_z q(z | x) \log \frac{p(x, z | \theta)}{q(z | x)} \end{aligned}$$



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

- KL and Lower bound of likelihood

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



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

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

# A variational representation of probability distributions



Goal:  
To approximate  $p(x)$

$$q(x) = \prod q(x_i)$$

$$q = \arg \max_{q \in Q} \{ -\langle E \rangle_q + H_q \}$$

$$= \arg \min_{q \in Q} \{ \langle E \rangle_q - H_q \}$$

$\sum_x -q(x) \log q(x)$   
 $\downarrow$   
 $\sum_i \sum_j q(x_i) \log q(x_j)$

where  $Q$  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  $Q$

# Bethe Free Energy/LBP



- But we do not optimize  $q(\mathbf{X})$  explicitly, focus on the set of beliefs

- e.g.,  $b = \{b_{i,j} = \tau(x_i, x_j), b_i = \tau(x_i)\}$

$q(x)$   
global consistency  
 $p(x_i, x_j) = \sum_{x_{-ij}} q(x)$

- Relax the optimization problem

- approximate objective:  $H_{\text{Bethe}} = H(b_{i,j}, b_i)$

- relaxed feasible set:  $\mathcal{M}_o = \{ \tau \geq 0 \mid \sum_{x_i} \tau(x_i) = 1, \sum_{x_i} \tau(x_i, x_j) = \tau(x_j) \}$

$$b^* = \arg \min_{b \in \mathcal{M}_o} \{ \langle E \rangle_b - F(b) \}$$

- The loopy BP algorithm:
  - a fixed point iteration procedure that tries to solve  $b^*$

## 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:  $Q \rightarrow \mathcal{T} \quad (\mathcal{T} \subseteq Q)$

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

## Mean Field Approximation





# Cluster-based approx. to the Gibbs free energy

(Wiegerinck 2001, Xing et al 03,04)



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

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

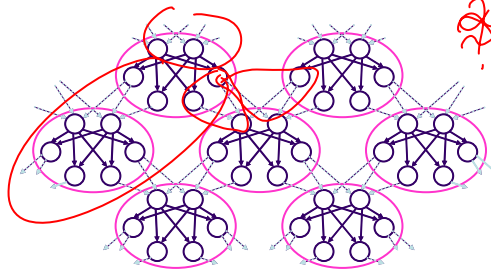
$$q(x) = \prod q_c(x_c)$$

$$p(x)$$

$$q_c(x_c)$$

$$= \exp(\xi_c)$$

$$\xi_c ?$$



$$\uparrow q_c = p(x)$$

$$\downarrow q_c = q(x_c)$$

Eric Xing

17

# 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

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

e.g.,  $G_{MF} = \sum_{i < j} \sum_{x_i, x_j} q(x_i) q(x_j) \psi(x_i, x_j) + \sum_i \sum_{x_i} q(x_i) \psi(x_i) + \sum_i \sum_{x_i} q(x_i) \ln q(x_i)$  (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  $q_i(x_c)$ 's.
  - Variational calculus ...
  - Do inference in each  $q_i(x_c)$  using any tractable algorithm

Eric Xing

18

Let  $\phi \in \exp\{\sum \theta_\alpha \phi_\alpha(x)\}$

Let:  $g(x) = \exp(E(x))$   
 $= \exp\{\sum \theta_\alpha \phi_\alpha(x)\}$

$$\begin{aligned} \text{Obj: } & -\langle E_\alpha(x) \rangle_g + H_g(x) \\ & = \langle \sum \theta_\alpha \phi_\alpha(x) \rangle_g - \langle \log g(x) \rangle_g \\ & = \sum \theta_\alpha \langle \phi_\alpha(x) \rangle_g - \sum_c \langle E(x_c) \rangle_g \\ & = \sum_{i \in I_\alpha} \langle \phi_\alpha(x_i) \rangle_{g_i} - \sum_c \langle E(x_c) \rangle_{g_c(x_c)} \end{aligned}$$

$I_c$ : the set of cliques that are either in or intersect with  $C$ .

$I_\alpha$ : the set of clusters that intersect with  $\alpha$ .

$x_{i,c}$ : variables other than those in  $C$

To optimize Obj w.r.t  $E_c(x_c)$   
 the relevant terms here led to a reduced objective

$$\Rightarrow \sum_{\alpha \in I_c} \theta_\alpha \langle \phi_\alpha(x_\alpha) \rangle_{g_i} - \langle E(x_c) \rangle_{g_c(x_c)}$$

To perform constrained optimization over  $E(x_c)$ ,  
 we write the follow Lagrangian:

$$\begin{aligned} \mathcal{J} = & \int dx_{i,c} \left( \sum_{\alpha \in I_c} \theta_\alpha \langle \phi_\alpha(x_\alpha) \rangle_{g_i} - E(x_c) \exp E(x_c) \right) \\ & - \lambda_c \left( \int \exp E(x_c) dx_c - 1 \right) \end{aligned}$$



blue: cliques  
 green: clusters

$$g(x_c) = \exp E(x_c)$$

$$\mathcal{J} : \int \left[ \prod_{c \in \mathcal{M}(x_c)} g(x_c) \left( \sum_{\alpha \in I_c} \theta_\alpha \phi_\alpha - \exp E(x_c) E(x_c) \right) dx_c - \lambda_c \left( \int \exp E(x_c) dx_c - 1 \right) \right]$$

$$\frac{d\mathcal{J}}{dE_c} : \int dx_c \left( \prod_{c \in \mathcal{M}(x_c)} g(x_c) \left( \sum_{\alpha \in I_c} \theta_\alpha \phi_\alpha - \exp E(x_c) - \exp E(x_c) E(x_c) \right) - \lambda_c \right) = 0$$

$$= \exp E(x_c) \int dx_c \left( \prod_{c \in \mathcal{M}(x_c)} g(x_c) \sum_{\alpha \in I_c} \theta_\alpha \phi_\alpha - 1 - E(x_c) \right) = \lambda_c \exp E(x_c)$$

= 0

$$\Rightarrow \int \left( \prod_{c \in \mathcal{M}(x_c)} g(x_c) \sum_{\alpha \in I_c} \theta_\alpha \phi_\alpha \right) dx_c = E(x_c) + C$$

$$\Rightarrow E(x_c) = \sum_{\alpha \in I_c'} \theta_\alpha \phi_\alpha + \sum_{\alpha \in I_c''} \theta_\alpha \phi_\alpha \cdot g - C$$





$$E(x_i) = \sum_{\alpha \in I_c} \theta_\alpha \phi_\alpha + \sum_{\alpha \in I_c} \theta_\alpha \langle \phi_\alpha \rangle_{q_i} + C$$

$$q(x_i) = \exp(E(x_i)) = \exp\left(\sum \theta_\alpha \phi_\alpha + \sum \theta_\alpha \langle \phi_\alpha \rangle_{q_i} + C\right)$$

$$p(x_c | x_{\alpha \in M(c)}) = \exp\left(\sum_{\alpha \in I_c} \theta_\alpha \phi_\alpha + C\right)$$



$$p(x) = \exp\left(\sum \theta_\alpha \phi_\alpha\right)$$

$$p(x_c | x_{\alpha \in M(c)})$$

## The Generalized Mean Field theorem



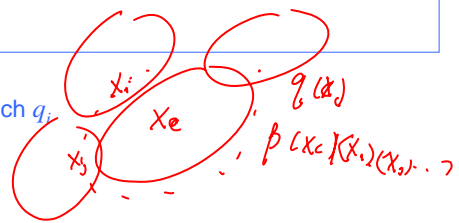
$$p(x)$$

$$q(x) = \prod q(x_i)$$

**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:

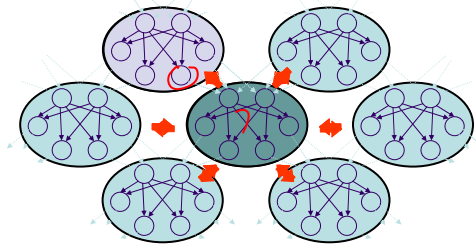
$$q_i^*(\mathbf{X}_{H,C_i}) = p(\mathbf{X}_{H,C_i} | \mathbf{x}_{E,C_i}, \langle \mathbf{X}_{H,MB_i} \rangle_{q_{j \neq i}})$$

GMF algorithm: Iterate over each  $q_i$



# A generalized mean field algorithm

[xing et al. UAI 2003]



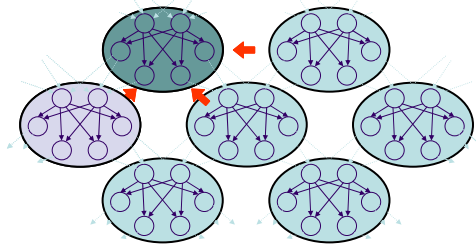
$$p(x_i) = \frac{1}{Z} \sum_{x_{-i}} P(x_i)$$

$$q(x_i) = p(x_i | \langle x_{-i} \rangle)$$

$$P(x_i)$$

# 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

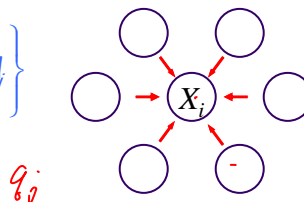


- Approximate  $p(\mathbf{X})$  by fully factorized  $q(\mathbf{X}) = \prod_i q_i(X_i)$
- For Boltzmann distribution  $p(\mathbf{X}) = \exp\{\sum_{i < j} q_{ij} X_i X_j + \sum_i q_{i0} X_i\} / Z$ :

mean field equation:

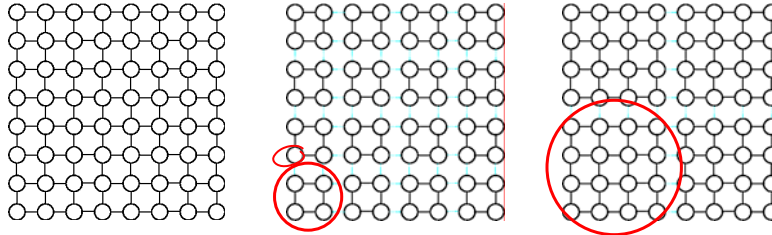
$$q_i(X_i) = \exp\left\{\theta_{i0} X_i + \sum_{j \in \mathcal{N}_i} \theta_{ij} X_i \langle X_j \rangle_{q_j} + A_i\right\}$$

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



- $\langle X_j \rangle_{q_j}$  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  $C_k$ :

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

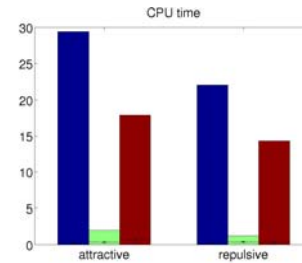
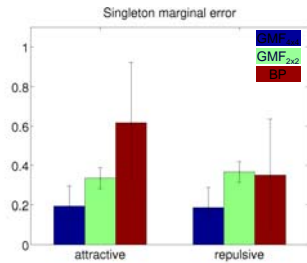
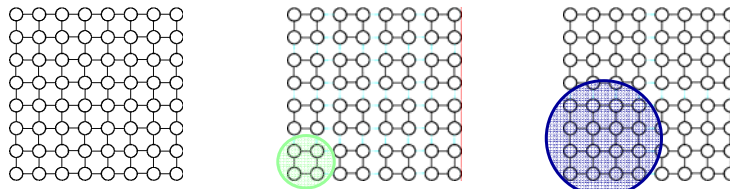
*Handwritten notes:*  
 $\int \mathcal{P}(x_i) x_i \cdot d x_i$   
 $\approx \int x_i \cdot q(x_i) d x_i$

Virtually a reparameterized Ising model of small size.

Eric Xing

27

# GMF approximation to Ising models

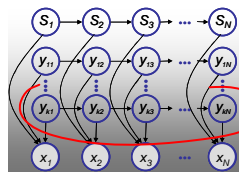


Attractive coupling: positively weighted  
 Repulsive coupling: negatively weighted

Eric Xing

28

# 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?

$$q(x) = q(x_1) q(x_2) \dots q(x_N)$$

$$q(x) = q(x_1 | \dots)$$

$$q(x_i)$$

Eric Xing

29

# 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

30

## 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

31

## Variational Posterior Distribution



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

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

where

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

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

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

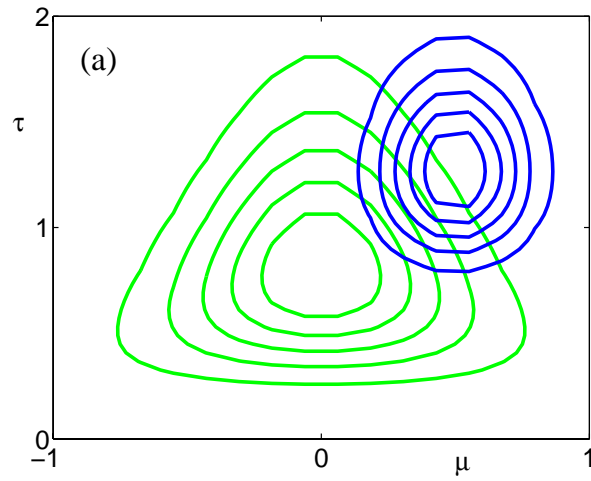
$$b_N = b_0 + \frac{1}{2} \left\langle \sum_n (x_n - \mu)^2 \right\rangle_\mu$$

Eric Xing

32



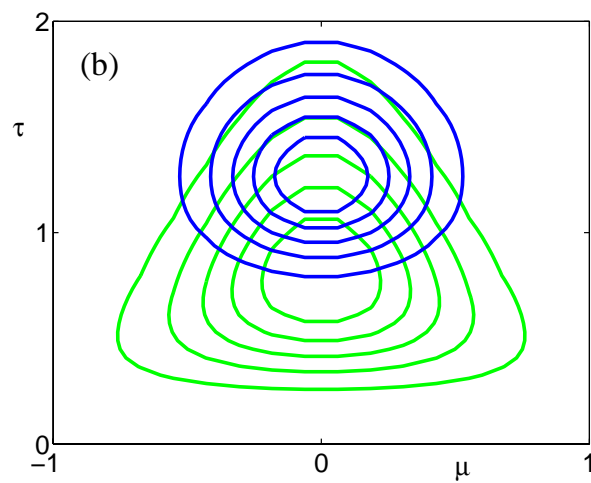
## Initial Configuration



Eric Xing

33

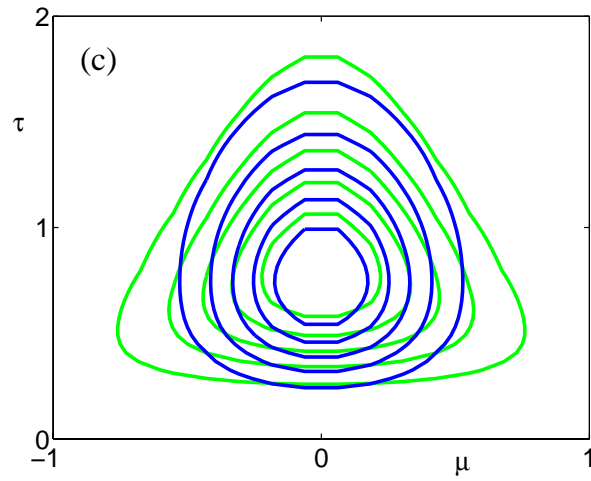
## After Updating $q(\mu)$



Eric Xing

34

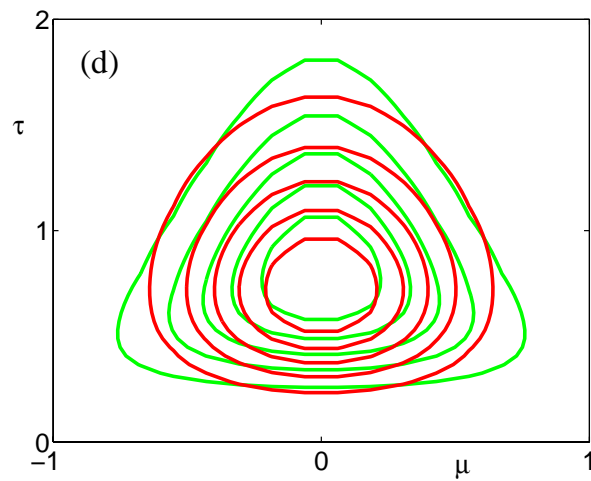
## After Updating $q(\mu)$



Eric Xing

35

## Converged Solution



Eric Xing

36

## Example 2: Latent Dirichlet Allocation

- Blei, Jordan and Ng (2003)
- Generative model of documents (but broadly applicable e.g. collaborative filtering, image retrieval, bioinformatics)

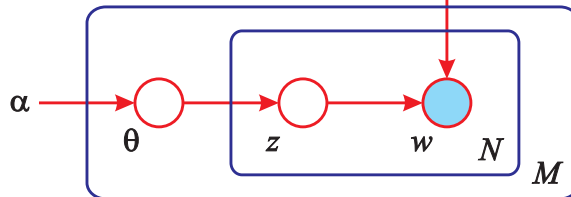
- Generative model:

- choose
- choose topic
- choose word

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

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

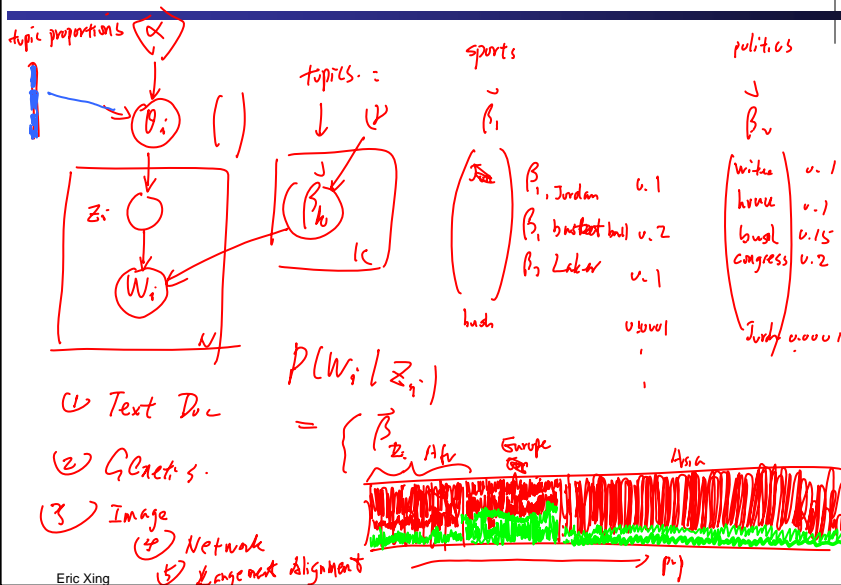
$$w_n \sim p(w_n | z_n, \beta)$$



Eric Xing


37

LDA → Admixture vs. Mixture



Eric Xing

38



$\{\beta\}$

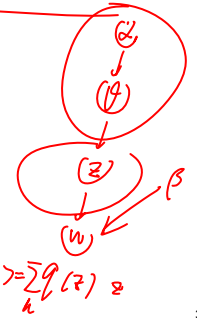
$$P(W_{n,i}, z_i | \theta_n) = P(W_{n,i} | z_{n,i}) \cdot P(z_{n,i} | \theta_n) \cdot P(\theta_n | \alpha) \cdot \phi$$

$$P(w) = \prod_{i=1}^K \int \sum_{z_i} \frac{1}{\Gamma(\alpha)} P(W_{n,i} | \beta_{n,k})^{\alpha} \delta(z_{n,i}, k) \frac{1}{\Gamma(\alpha)} \theta_{n,k}^{\alpha} \delta(z_{n,i}, k)$$

$$\times \prod_{k=1}^K \theta_{n,k}^{\alpha-1} \cdot C \cdot d\theta_n$$



---

$P(\theta, z | w) = \frac{P(\theta, z, w)}{P(w)}$   
 $q(\theta, z) = \frac{q(\theta) q(z)}{P(w)}$   
 ?  $q(\theta) \propto \frac{P(\theta | z, \alpha)}{P(z | \theta) P(\theta | \alpha)} = \frac{\prod_{k=1}^K \delta(z, k)^{\alpha-1}}{\prod_{k=1}^K \theta^{\alpha-1}}$   
 ?  $q(z) = P(z | \theta) P(\theta | \alpha)$



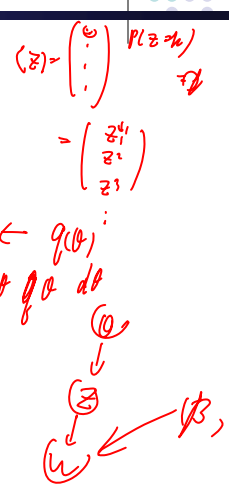
$\langle z \rangle = \sum_k q(z) z$

Eric Xing 39



$E(x_i) = \theta_{\alpha}(\phi)$   
 $p(x) = \exp(\theta \phi(x))$

suppose that  $q$  learn  $q(z)$  in iter.  $t$   
 $\rightarrow$  compute  $\langle z \rangle$  w.r.t  $q(z)$   
 $\rightarrow q(\theta) = \prod_{k=1}^K \theta^{\alpha-1} \delta(z, k)$   
 $\propto \prod_{k=1}^K \theta^{\alpha-1} \delta(z, k)$



$\langle z \rangle = \begin{pmatrix} z_1 \\ \vdots \\ z_t \end{pmatrix}$   
 $= \begin{pmatrix} z_1^t \\ z_2^t \\ \vdots \\ z_t^t \end{pmatrix}$

$q(z)$   
 $P(z | \theta, w, \beta)$   
 $\propto P(z | \theta) P(w | z, \beta)$   
 $\propto \prod_{k=1}^K \theta_{n,k}^{\alpha} \cdot \prod_{s=1}^M \beta_{n,k}^{\alpha} w_{n,i}^s = (\exp \ln \theta_{n,k})^{\alpha}$   
 $q(z) \propto \exp \sum_k \langle \ln \theta_{n,k} \rangle z_k^{\alpha} \exp \sum_{s=1}^M \ln \beta_{n,k} \cdot z_k^{\alpha} w_{n,i}^s$

Eric Xing 40

# Latent Dirichlet Allocation

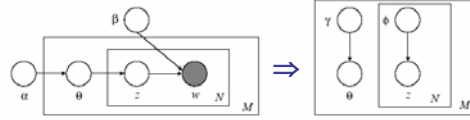


- Variational approximation

$$q(\theta, z) = q_\theta(\theta)q_z(z)$$

$$\text{Dir}(\theta | \gamma = f(\alpha, \langle z \rangle)) \times$$

$$\text{Multi}(z | \phi = f(\beta_w, \langle \ln \theta \rangle))$$



$$\phi_{ni} \propto \beta_{i w_n} \exp\{E_q[\log(\theta_i) | \gamma]\}$$

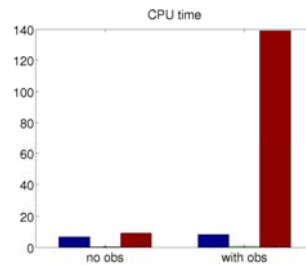
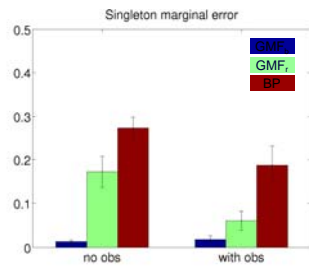
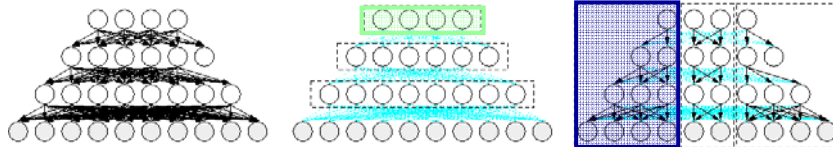
$$\gamma_i = \alpha_i + \sum_{n=1}^N \phi_{ni}$$

- 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

41

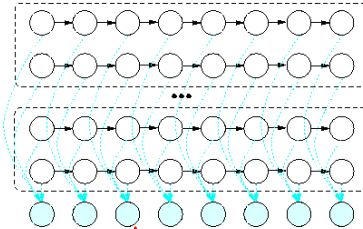
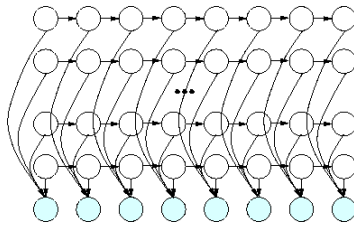
# Example 3: Sigmoid belief network



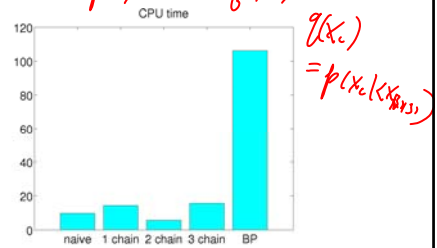
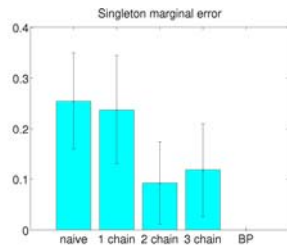
Eric Xing

42

# Example 4: Factorial HMM



$$p(x) \approx \prod q(x_i)$$



$$q(x_i) = p(x_i | x_{\setminus i})$$