

Probabilistic Graphical Models

10-708

Variational Inference

Eric Xing

Lecture 18, Nov 14, 2005

Reading: KF-Chap. 9



Variational Methods

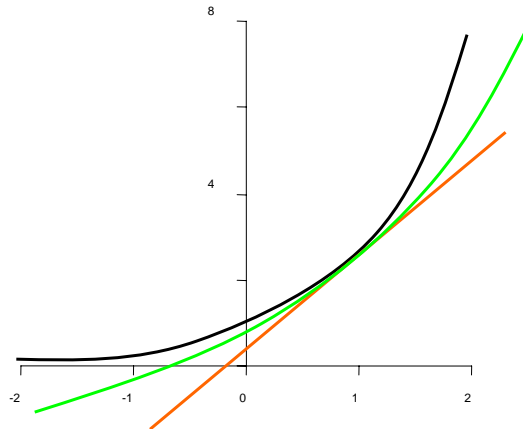
- For a distribution $p(X/\theta)$ associated with a complex graph, computing the marginal (or conditional) probability of arbitrary random variable(s) is intractable
- Variational methods
 - formulating probabilistic inference as an optimization problem:

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

\mathbf{f} : a (tractable) probability distribution
or, solutions to certain probabilistic queries



Lower bounds of exponential functions



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

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

Exponential Family



- Exponential representation of graphical models:

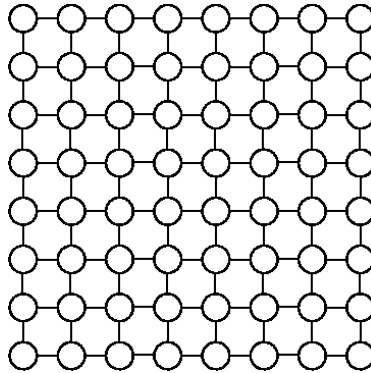
$$p(\mathbf{X} | \theta) = \exp \left\{ \sum_{\alpha} \theta_{\alpha} \phi_{\alpha}(\mathbf{X}_{D_{\alpha}}) - A(\theta) \right\}$$

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

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

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

Example: the Boltzmann distribution on atomic lattice



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

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]

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:

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

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”):

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

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

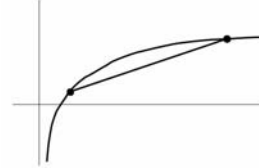
minimized when $q(Z) = p(Z)$



KL and Log Likelihood

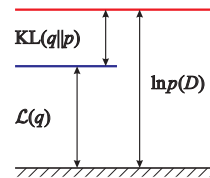
- Jensen's inequality

$$\begin{aligned} \ell(\theta; \mathbf{x}) &= \log p(\mathbf{x} | \theta) \\ &= \log \sum_z p(\mathbf{x}, z | \theta) \\ &= \log \sum_z q(z | \mathbf{x}) \frac{p(\mathbf{x}, z | \theta)}{q(z | \mathbf{x})} \\ &\geq \sum_z q(z | \mathbf{x}) \log \frac{p(\mathbf{x}, z | \theta)}{q(z | \mathbf{x})} \end{aligned} \Rightarrow \ell(\theta; \mathbf{x}) \geq \langle \ell_c(\theta; \mathbf{x}, z) \rangle_q + H_q = \mathcal{L}(q)$$



- KL and Lower bound of likelihood

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



- Setting $q()=q(z|\mathbf{x})$ closes the gap (c.f. EM)

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}



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

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



Belief Propagation

- Do not optimize $q(\mathbf{X}_H)$ explicitly, but focus on the set of beliefs
 - *e.g.*, $\mathbf{b} = \{b_{i,j} = \tau(x_i, x_j), b_i = \tau(x_i)\}$
- Relax the optimization problem
 - approximate objective: $H_{\text{beta}} = H(\mathbf{b}_{i,j}, b_i)$
 - 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^*



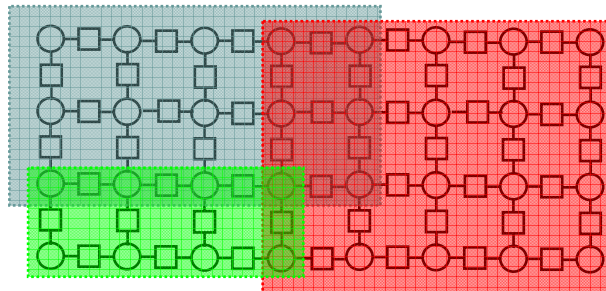
Loopy Belief Propagation

Region-based Approximations to the Gibbs Free Energy (Kikuchi, 1951)



Exact: $\mathcal{G}[q(X)]$ (*intractable*)

Regions: $\mathcal{G}[\{b_r(X_r)\}]$



Bethe Approximation to Gibbs Free Energy



- Bethe free energy

$$\mathcal{G}_{Bethe} = \sum_a \sum_{\mathbf{x}_a} b_a(\mathbf{x}_a) \ln f_a(\mathbf{x}_a) + \sum_i (1-d_i) \sum_{\mathbf{x}_i} b_i(\mathbf{x}_i) \ln b_i(\mathbf{x}_i)$$

- Equal to the exact Gibbs free energy when the factor graph is a tree because in that case,

$$b(\mathbf{x}) = \prod_a b_a(\mathbf{x}_a) \prod_i b_i(\mathbf{x}_i)^{1-d_i}$$

- But otherwise, it **may or may not** give a lower bound of the likelihood
- Optimize each $b(\mathbf{x}_a)$'s.
 - For discrete belief, constrained opt. with *Lagrangean* multiplier
 - For continuous belief, not yet a general formula
 - Not always converge

Minimizing the Bethe Free Energy



$$L = \mathcal{G}_{Bethe} + \sum_i \gamma_i \left\{ \sum_{\mathbf{x}_i} b_i(\mathbf{x}_i) - 1 \right\} + \sum_a \sum_{i \in N(a)} \sum_{\mathbf{x}_i} \lambda_{ai}(\mathbf{x}_i) \left\{ \sum_{\mathbf{x}_a \setminus \mathbf{x}_i} b_a(\mathbf{x}_a) - b_i(\mathbf{x}_i) \right\}$$

$$\frac{\partial L}{\partial b_i(\mathbf{x}_i)} = 0 \quad \Longrightarrow \quad b_i(\mathbf{x}_i) \propto \exp\left(\frac{1}{d_i-1} \sum_{a \in N(i)} \lambda_{ai}(\mathbf{x}_i)\right)$$

$$\frac{\partial L}{\partial b_a(\mathbf{x}_a)} = 0 \quad \Longrightarrow \quad b_a(\mathbf{x}_a) \propto \exp\left(-E_a(\mathbf{x}_a) + \sum_{i \in N(a)} \lambda_{ai}(\mathbf{x}_i)\right)$$

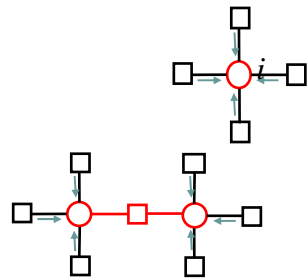
Bethe = BP



- Identify

$$\lambda_{ai}(x_i) = \ln \prod_{b \in N(i) \neq a} m_{b \rightarrow i}(x_i)$$

- to obtain BP equations:



$$b_i(x_i) \propto \prod_{a \in N(i)} m_{a \rightarrow i}(x_i)$$

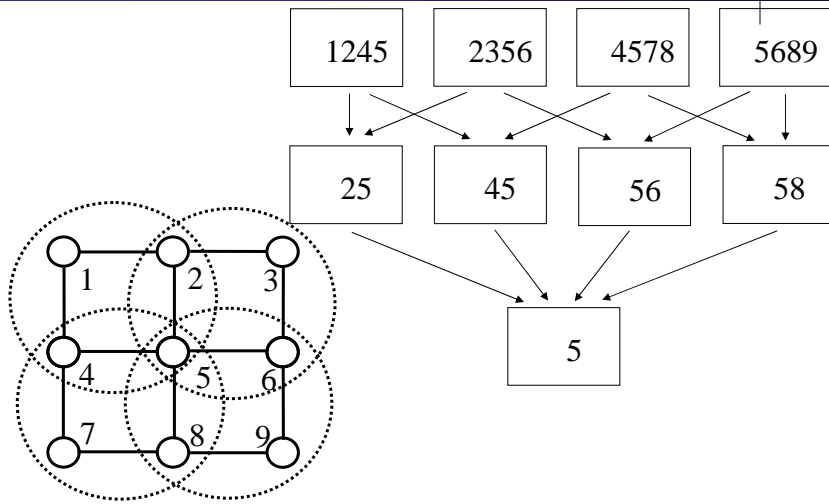
$$b_a(X_a) \propto f_a(X_a) \prod_{i \in N(a)} \prod_{b \in N(i) \setminus a} m_{b \rightarrow i}(x_i)$$

Generalized Belief Propagation

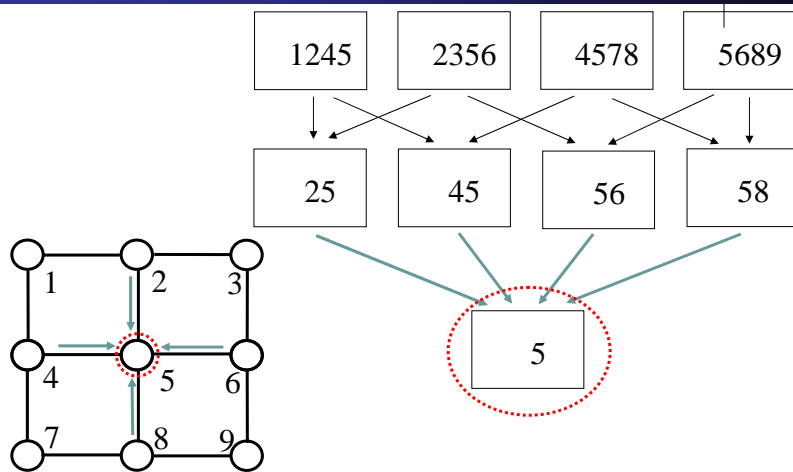


- Belief in a region is the product of:
 - Local information (factors in region)
 - Messages from parent regions
 - Messages into descendant regions from parents who are not descendants.
- Message-update rules obtained by enforcing marginalization constraints.
- Kikuchi free energy

Generalized Belief Propagation

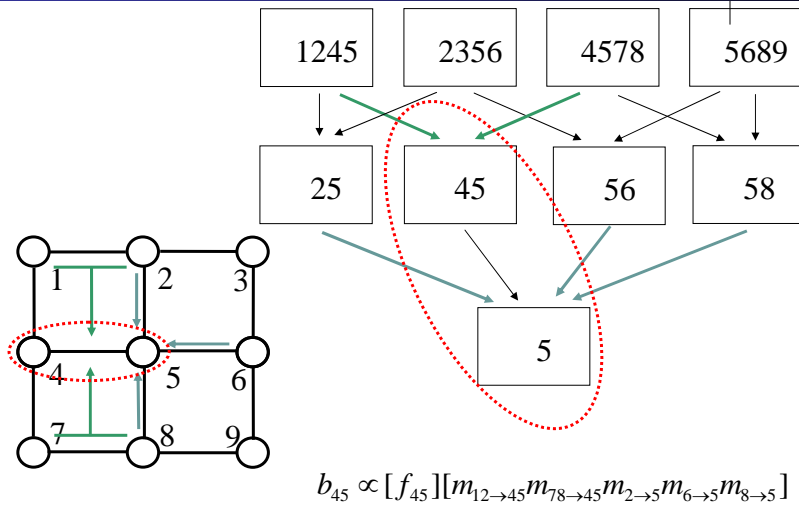


Generalized Belief Propagation

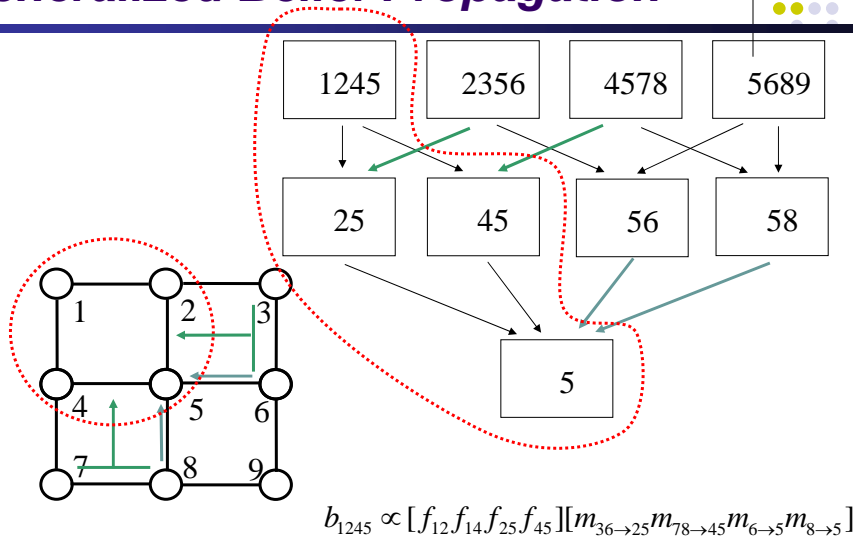


$$b_5 \propto m_{2 \rightarrow 5} m_{4 \rightarrow 5} m_{6 \rightarrow 5} m_{8 \rightarrow 5}$$

Generalized Belief Propagation



Generalized Belief Propagation





Mean Field Approximation

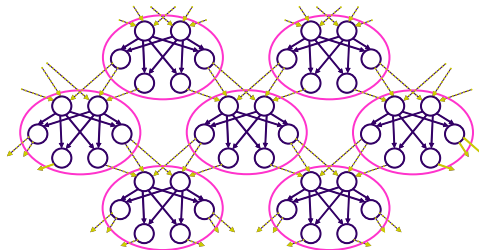
Cluster-based approx. to the Gibbs free energy

(Wiegerinck 2001,
Xing et al 03,04)



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

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



Mean field approx. to Gibbs free energy



- Given a disjoint clustering, $\{C_1, \dots, C_I\}$, of all variables

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

- Mean-field free energy

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

e.g.,
$$\mathcal{G}_{\text{MF}} = \sum_{i < j} \sum_{\mathbf{x}_i, \mathbf{x}_j} q(\mathbf{x}_i) q(\mathbf{x}_j) \psi(\mathbf{x}_i, \mathbf{x}_j) + \sum_i \sum_{\mathbf{x}_i} q(\mathbf{x}_i) \psi(\mathbf{x}_i) + \sum_i \sum_{\mathbf{x}_i} q(\mathbf{x}_i) \ln q(\mathbf{x}_i) \quad (\text{naive mean field})$$

- 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 $q_i(\mathbf{x}_{C_i})$'s.
 - Variational calculus ...
 - Do inference in each $q_i(\mathbf{x}_{C_i})$ 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:

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

GMF algorithm: Iterate over each q_i

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

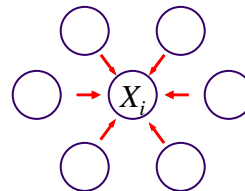


- Approximate $p(\mathbf{X})$ by fully factorized $q(\mathbf{X}) = \prod_i q_i(X_i)$
- For Boltzmann distribution $p(\mathbf{X}) = \exp\{\sum_{i < j} \theta_{ij} X_i X_j + \sum_i \theta_{i0} X_i\} / Z$:

mean field equation:

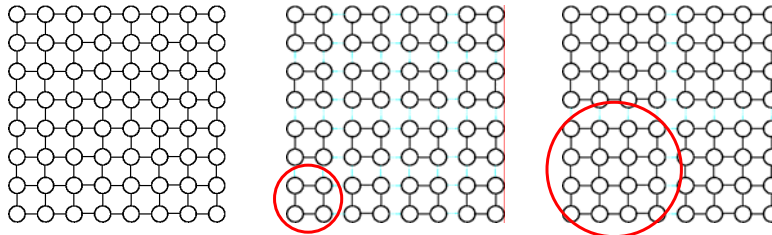
$$q_i(X_i) = \exp\left\{\theta_{i0} X_i + \sum_{j \in \mathcal{N}_i} \theta_{ij} X_i \langle X_j \rangle_{q_j} + A_i\right\}$$

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



- $\langle X_j \rangle_{q_j}$ 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

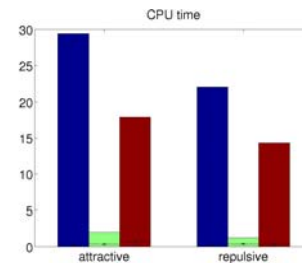
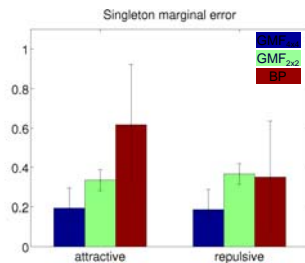
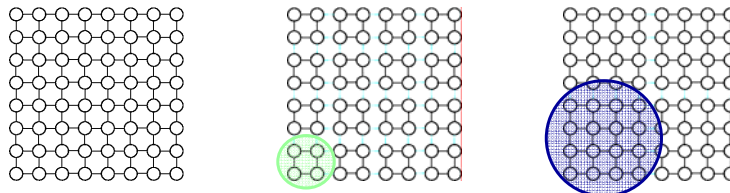


Cluster marginal of a square block C_k :

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

Virtually a reparameterized Ising model of small size.

GMF approximation to Ising models



Attractive coupling: positively weighted
 Repulsive coupling: negatively weighted

Automatic Variational Inference



- 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?

Cluster-based MF



- 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

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

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

Variational Posterior Distribution



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

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

where

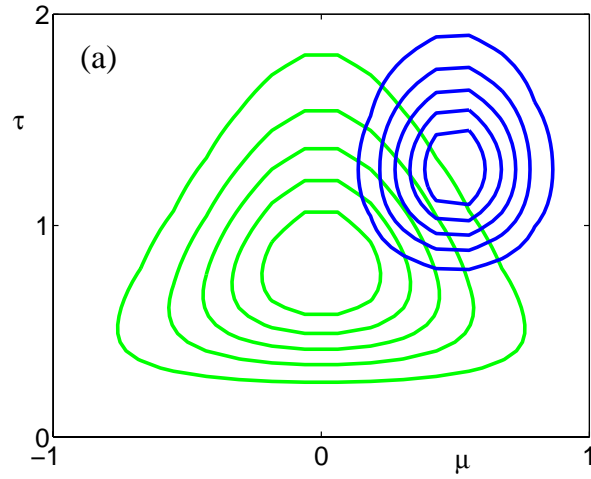
$$\mu_N = \frac{\lambda_0 \mu_0 + \langle \tau \rangle N \bar{x}}{\lambda_0 + N \langle \tau \rangle}$$

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

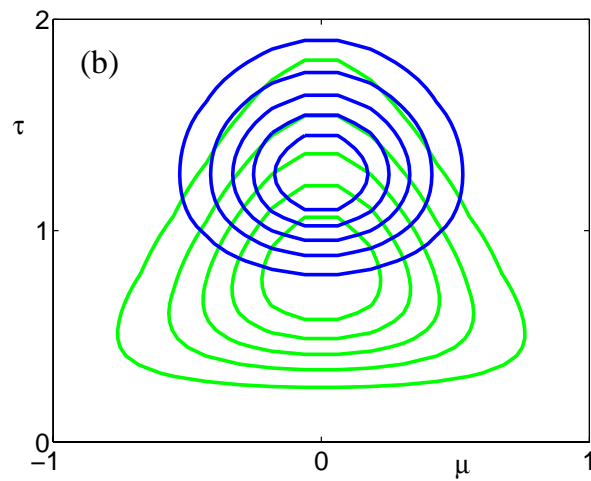
$$a_N = a_0 + \frac{N}{2}$$

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

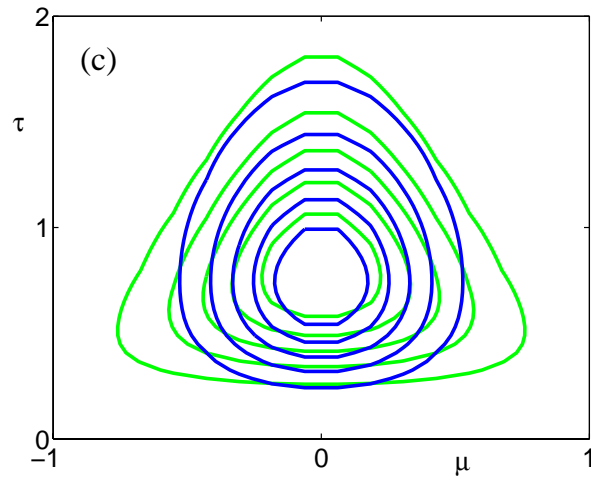
Initial Configuration



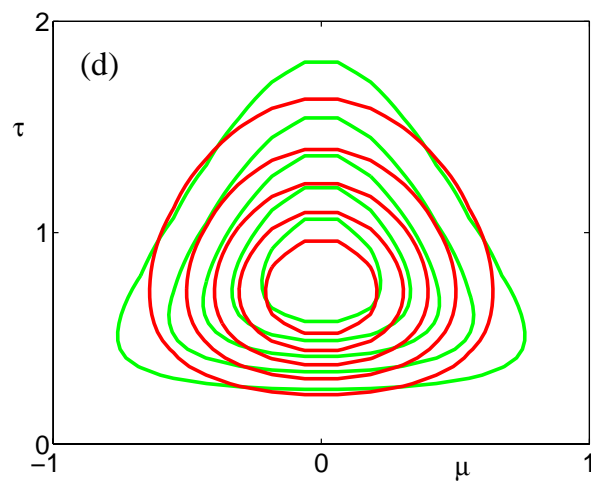
After Updating $q(\mu)$



After Updating $q(\tau)$



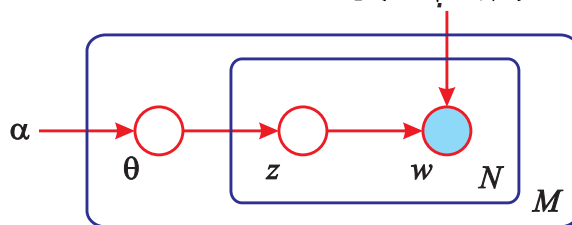
Converged Solution



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



Latent Dirichlet Allocation

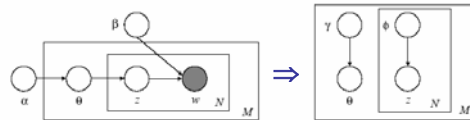


- Variational approximation

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

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

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

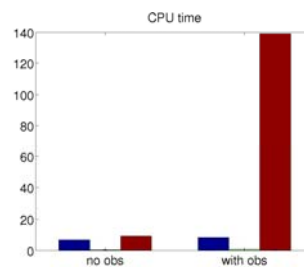
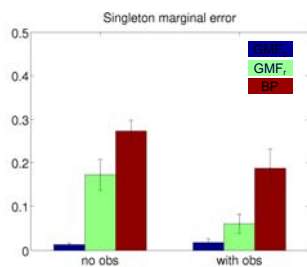
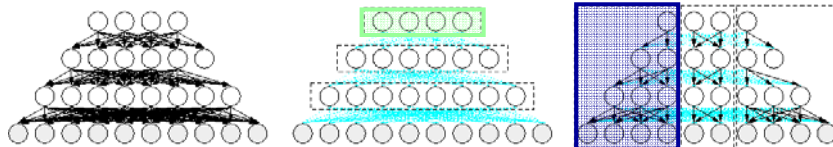


$$\phi_{ni} \propto \beta_{iw_n} \exp\{E_q[\log(\theta_i) | \gamma]\}$$

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

- 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

Example 3: Sigmoid belief network



Example 4: Factorial HMM

