

VARIATIONAL LOWER BOUND

⇒ Consider observed RV $X=x$ w/ marginal likelihood:

$$\begin{aligned} \log p(x) &= \log \int p(x, \theta) d\theta \\ &= \log \int q(\theta) \frac{p(x, \theta)}{q(\theta)} d\theta \\ &= \log \mathbb{E}_q \left[\frac{p(x, \theta)}{q(\theta)} \right] \\ &\geq \mathbb{E}_q \left[\log \frac{p(x, \theta)}{q(\theta)} \right] \quad (\text{Jensen's inequality}) \\ &= -\text{KL}(q \parallel p(x, \cdot)) \end{aligned}$$

TALK: Fri, 1-2pm, EPRZ 5395.
"PROB. REASONING IN COMPLEX SYSTEMS: ALG. & APP"

⇒ Variational optimization finds tightest lower bound:

$$\log p(x) \geq \max_{q \in \mathcal{Q}} -\text{KL}(q(\theta) \parallel p(x, \theta))$$

for some class of dist's \mathcal{Q}

- Bound is tight iff $q(\theta) = p(\theta|x)$

- KEY IDEA: VI turns statistical inference into optimization

⇒ Different VI methods amount to different \mathcal{Q}

- Mean field = \mathcal{Q}^{MF} marginal posterior independence
- Belief Propagation = \mathcal{Q}^{BP} "Tree-like" posterior
- Expectation Propagation = \mathcal{Q}^{EP} "Tree-like" in moments

⇒ ASIDE: Info-Theoretic interpretation suggests we should minimize

$$\min_q \text{KL}(p(\theta|x) \parallel q(\theta)) = \mathbb{E}_{p(\theta|x)} \left[\log \frac{p(\theta|x)}{q(\theta)} \right]$$

Can't compute this

Don't have this

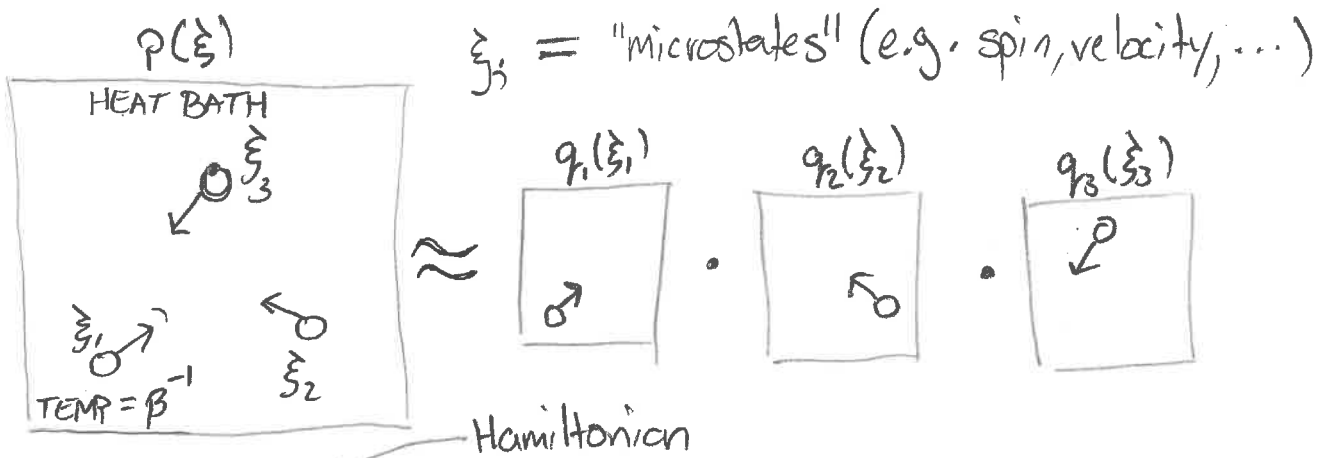
Today's Focus

MEAN FIELD VARIATIONAL

⇒ "Mean field" assumes d -dimensional θ , marginally independent.

$$q(\theta) = \prod_{i=1}^d q_i(\theta_i)$$

⇒ Originates from many-body problem in stat. mechanics:



$$P(\xi) = \underbrace{\frac{1}{\sum_{\xi} e^{-\beta H(\xi)}}}_{\text{GIBBS' DIST'N}} \approx \prod_i \frac{1}{\sum_{\xi_i} e^{-\beta h_i(\xi_i)}} \triangleq \prod_i q_i(\xi_i)$$

BOLTZMANN DIST'N

⇒ Mean field variational lower bound:

$$\log P(x) \geq \max_q \mathbb{E}_q [\log P(\theta, x)] + \sum_i H(q_i) \triangleq \mathcal{L}(q)$$

- Sometimes called "Evidence Lower Bound" (ELBO)

⇒ How do we optimize $\mathcal{L}(q)$? Coordinate ascent.

- Take derivative wrt each $q_i(\theta_i)$, set to 0, solve:

$$q_i^*(\theta_i) \propto \exp \left\{ \mathbb{E}_{q_{\gamma_i}} \left[\log P(\theta_i | \theta_{\gamma_i}, x) \right] \right\} \quad (*)$$

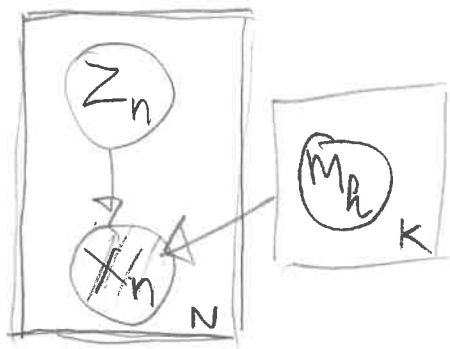
where $\gamma_i = \{1, \dots, d\} \setminus i$ KNOWN AS "COMPLETE CONDITIONAL"

COORDINATE ASCENT VI:

⇒ Monotonically increases $L(q)$

⇒ Fixed-points of $(*) \Leftrightarrow$ Local Optima of $L(q)$

EX. GAUSSIAN MIXTURE MODEL



$$m_h \sim N(0, \sigma^2) \text{ for } h=1, \dots, K$$

$$z_n \sim \text{CAT}(\frac{1}{K}, \dots, \frac{1}{K}) \text{ for } n=1, \dots, N$$

$$x_n | m_{z_n} \sim N(m_{z_n}, 1) \dots$$

- Variational distributions:

$$q_n(z_n | \pi_n) = \text{CAT}(z_n | \pi_n), \quad q_h(m_h | \mu_h, s_h) = N(m_h | \mu_h, s_h)$$

where $\{\pi_n, \mu_h, s_h\}$ are variational parameters

① Update Cluster assignments:

$$q_n(z_n | \pi_n) \propto \exp \left\{ \underbrace{\log p(z_n)}_{\text{PRIOR} = \frac{1}{K}} + \underbrace{\mathbb{E}_{q(m|\mu,s)} [\log p(x_n | z_n, m)]}_{\text{EXPECTED LOG-LIKELIHOOD}} \right\} \quad (**)$$

- Compute expected log-likelihood:

$$\mathbb{E}_{q(m)} [\log p(x_n | z_n, m)] = \mathbb{E} \left[\sum_h \mathbb{I}(z_n=h) \log N(x_n | m_h, 1) \right]$$

$$= \sum_h \mathbb{I}(z_n=h) \mathbb{E} \left[-\frac{1}{2} \log 2\pi - \frac{1}{2} x_n^2 + x_n m_h - \frac{1}{2} m_h^2 \right]$$

$$= \sum_h \mathbb{I}(z_n=h) \left(\mathbb{E}_{q(m_h)} [m_h] x_n - \frac{1}{2} \mathbb{E}_{q(m_h)} [m_h^2] \right) + \text{const.}$$

- Plug back into (**)

$$q_n(z_n=h) = \pi_n(h) \propto \exp \left\{ \mu_h x_n - \frac{1}{2} (s_h + \mu_h^2) \right\}$$

② Update component means $q(m_h | \mu_h, s_h)$ (we will skip that)

③

CONDITIONALLY CONJUGATE MODELS

⇒ In the GMM example our update was:

$$q_n(z_n) \propto p(z_n) \exp\left\{\mathbb{E}_m[\log p(x_n | z_n, m)]\right\} \propto \text{CAT}(\cdot)$$

Which is same family as prior $p(z_n)$

- We call this conditional conjugacy

⇒ In general this holds when:

1) Complete conditionals are expfam

$$2) q(\theta) \propto p(\theta) \exp\left\{\mathbb{E}_z[\log p(x, z | \theta)]\right\} \\ = p(\theta) f(x | \theta)$$

Is conjugate ⇒ Lies in same family as $p(\theta)$

⇒ Conditional conjugacy ⇒ Closed-form updates

VARIATIONAL OPTIMIZATION

⇒ How did we get fixed-point condition?

$$q_i(\theta_i) \propto \exp\left\{\mathbb{E}[\log p(\theta_i | \theta_{-i}, x)]\right\}$$

⇒ Form Lagrangian of variational problem:

$$\max_q \mathcal{L}(q) \triangleq \mathbb{E}_q[\log p(x, \theta)] + \sum_i H(q_i | \theta_i)$$

$$\text{s.t. } \int q_i(\theta_i) = 1 \quad \forall i=1, \dots, N$$

Lagrangian:

$$\mathcal{L}(q) + \sum_i \lambda_i (1 - \int q_i(\theta_i) d\theta_i) \triangleq J(q, \lambda)$$

⇒ Calculus of variations defines derivative of functional

$$J(f) = \int L(f(x)) dx$$

AS:

$$\frac{\partial J}{\partial f(x)} = \frac{\partial L}{\partial f}$$

NOTE: This is a simplification of a functional and the Euler-Lagrange equations.

⇒ Take functional derivative wrt $q_i(\theta_i)$:

$$\frac{\partial}{\partial q_i(\theta_i)} \mathcal{L}(q) + \frac{\partial}{\partial q_i(\theta_i)} \lambda_i (1 - \int q_i(\theta_i) d\theta_i)$$

$$= \frac{\partial}{\partial q_i(\theta_i)} \mathbb{E}_{q_i} \left[\mathbb{E}_{q_{-i}} [\log p(\theta_i | \theta_{-i}, x)] \right]$$

$$+ \frac{\partial}{\partial q_i(\theta_i)} H(q_i(\theta_i)) - \lambda_i$$

$$= \mathbb{E}_{q_{-i}} [\log p(\theta_i | \theta_{-i}, x)] - \log q_i(\theta_i) - 1 - \lambda_i = 0 \Rightarrow$$

$$q_i^*(\theta_i) = \exp \left\{ \mathbb{E}_{q_{-i}} [\log p(\theta_i | \theta_{-i}, x)] \right\} \div \mathcal{Z}(\lambda_i)$$

BETHE VARIATIONAL PROBLEM:

⇒ KEY IDEA: Optimize over distⁿs $q \in \mathcal{Q}^{\text{BETHE}}$ consistent w/ tree-structured MRFs.

⇒ Any tree-structured distⁿ $q(\theta)$ can be written as:

$$q(\theta) = \prod_{S \in \mathcal{V}} q_S(\theta_S) \prod_{(S,T) \in \mathcal{E}} \frac{q_{ST}(\theta_S, \theta_T)}{q_S(\theta_S) q_T(\theta_T)} \quad (*)$$

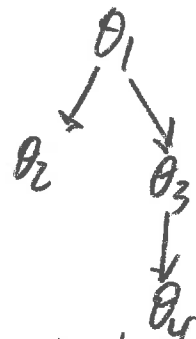
where $q_S(\theta_S)$ are marginals:

$$q_S(\theta_S) = \int q(\theta_S) d\theta_S$$

Ex 2

$$q(\theta) = \left(q_1(\theta_1) q_2(\theta_2) q_3(\theta_3) q_4(\theta_4) \right) \left(\frac{q_{12}(\theta_1, \theta_2)}{q_1(\theta_1) q_2(\theta_2)} \right) \left(\frac{q_{13}(\theta_1, \theta_3)}{q_1(\theta_1) q_3(\theta_3)} \right) \left(\frac{q_{34}(\theta_3, \theta_4)}{q_3(\theta_3) q_4(\theta_4)} \right)$$

$$= q_1(\theta_1) q_{2|1}(\theta_2|\theta_1) q_{3|1}(\theta_3|\theta_1) q_{4|3}(\theta_4|\theta_3)$$



⇒ Distⁿs of the form (*) necessarily satisfy Local Marginal Consistency:

$$q_S(\theta_S) = \int q_{ST}(\theta_S, \theta_T) d\theta_T \quad \forall T \in \Gamma(s) \quad \leftarrow \begin{array}{l} \text{Neighbors} \\ \text{of } S \end{array}$$

⇒ Bethe Variational Problem:

$$\max_q \mathbb{E}_q[\log p(\theta, x)] + H(q)^{\text{BETHE}}$$

s.t.

$$\int q_{ST}(\theta_S, \theta_T) d\theta_T = q_S(\theta_S) \quad \forall T \in \Gamma(s) \quad (\text{Local Consistency})$$

$$\int q_S(\theta_S) d\theta_S = 1 \quad \forall S \in \mathcal{V} \quad (\text{Normalization})$$