

PGM - Variational Inference

①

[Pres → Detail #1]

Detail #2: The objective is to minimize

$$KL(q \| p(\cdot|y)) = \mathbb{E}_q \left[\log \frac{q(x)}{p(x|y)} \right] = \frac{p(x,y)}{p(y)}$$

$$= \underbrace{\mathbb{E}_q [\log q(x)] - \mathbb{E}_q [\log p(x,y)]}_{\text{minimizing KL is the same as minimizing this expression, or equivalently maximizing}} + \underbrace{\log p(y)}_{\text{constant (wrt } q)}$$

minimizing KL is the same as minimizing this expression, or equivalently maximizing

$$ELBO(q) \stackrel{\text{def}}{=} \underbrace{\mathbb{E}_q [\log p(x,y)]}_{\text{encourages } q \text{ to concentrate at the mode of } p(x,y)} - \underbrace{\mathbb{E}_q [\log q(x)]}_{\text{entropy of } q, \text{ encourages } q \text{ to "spread out"}}$$

encourages q to concentrate at the mode of $p(x,y)$

= entropy of q , encourages q to "spread out"

ELBO = evidence lower bound

$$\log p(y) = \underbrace{KL(q \| p(\cdot|y))}_{\geq 0} + ELBO(q) \Rightarrow \log p(y) \geq ELBO(q)$$

$p(y)$ is often referred to as model evidence (or marginal likelihood).

[Pres → Detail #4, CAVI]

(2)

Consider updating the k th component $q_k(x_k)$ with the other components $q_{-k}(x_{-k}) \stackrel{\text{def}}{=} \prod_{j \neq k} q_j(x_j)$ being fixed.

$$\begin{aligned} \text{ELBO}(q_k; q_{-k}) &= \mathbb{E}_{q_k} [\mathbb{E}_{q_{-k}} [\log p(x_k, x_{-k}, y)]] \\ &\quad - \mathbb{E}_{q_k} [\log q_k] + \text{const} \end{aligned}$$

Define $q_k^*(x_k) \propto \exp(\mathbb{E}_{q_{-k}} [\log p(x_k, x_{-k}, y)])$
 $\propto \exp(\mathbb{E}_{q_{-k}} [\log p(x_k | x_{-k}, y)])$

Then $\text{ELBO}(q_k; q_{-k}) = -\text{KL}(q_k, q_k^*) + \text{const}$

which is maximized by setting $q_k^{\text{new}} = q_k^*$

[Pres \rightarrow "What does", reveal

$$\prod_{k=1}^K \theta_{d/k}^{x_k + c_{d/k}}]$$

(3)

We get

$$\log p(\theta_d | \beta, z, w)$$

$$= \sum_{k=1}^K (\alpha_k + c_{d,k} - 1) \log \theta_{d,k} + \text{const}$$

$$\Rightarrow \mathbb{E}_{q(\cdot | \theta_d)} [\log p(\theta_d | \beta, z, w)]$$

$$= \sum_{k=1}^K \{ \alpha_k + \mathbb{E}_q[c_{d,k}] - 1 \} \log \theta_{d,k} + \text{const}$$

Consequently, $q^*(\theta_d) \propto \exp(\cdot) = \prod_{k=1}^K \theta_{d,k}^{\alpha_k + \mathbb{E}_q[c_{d,k}] - 1}$

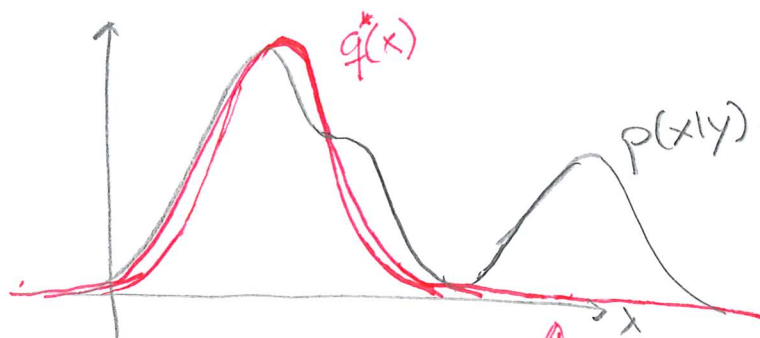
$$\propto \text{Dir}(\theta_d | \alpha + \mathbb{E}_q[c_d])$$

[Res \rightarrow Limitations]

④
"Mode seeking" behavior of

$$KL(q \| p(\cdot|y)) = \mathbb{E}_q \left[\log \frac{q(x)}{p(x|y)} \right]$$

↑ We "average the discrepancy" wrt the approximation itself.



↑ We don't care as much about the discrepancy here, where q is low

This typically leads to underestimation of the posterior variance, which VI is infamous for.