

$$P_\theta(x) = \int P_\theta(z) P_\theta(x|z) dz$$

$$\nabla_\theta P_\theta(x) = \int \nabla_\theta P_\theta(z) P_\theta(x|z) dz + P_\theta(z) \nabla_\theta P_\theta(x|z) dz.$$

$$\nabla_\theta \ln P_\theta(x) = \frac{\nabla_\theta P_\theta(x)}{P_\theta(x)} = \int P_\theta(z|x) \nabla_\theta (\ln P_\theta(z) + \frac{P_\theta(z)}{P_\theta(x)} \frac{P_\theta(x|z)}{P_\theta(x|z)} \nabla_\theta P_\theta(x|z) dz$$

$$= \int P_\theta(z|x) (\nabla_\theta \ln P_\theta(z) + \nabla_\theta \ln P_\theta(x|z)) dz.$$

$$= \mathbb{E}_{P_\theta(z|x)} \nabla_\theta (\ln P_\theta(x|z)) - \nabla_\theta U_\theta(z) + \mathbb{E}_{P_\theta(z)} \nabla_\theta U_\theta(z).$$

$$\mathbb{E}_{P_\theta(x)} \nabla_\theta \ln P_\theta(x) = \mathbb{E}_{P_\theta(z)} \nabla_\theta U_\theta(z) + \mathbb{E}_{z \sim P_\theta(z|x)} \nabla_\theta \ln P_\theta(x|z) - \nabla_\theta U_\theta(z) = -\nabla_\theta U_\theta(z) + \mathbb{E}_{P_\theta(z)} \nabla_\theta U_\theta(z)$$

$$\begin{cases} P_\theta(z) \sim q_\psi(z) & KL(q_\psi(z) || P_\theta(z)) \\ q_\psi(z|x) \sim P_\theta(z|x) & KL(q_\psi(z|x) || P_\theta(z|x)). \end{cases}$$

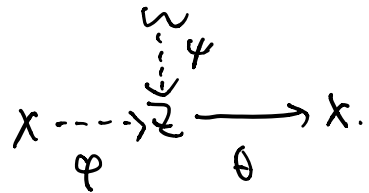
$$\approx \mathbb{E}_{q_\psi(w)} \nabla_\theta U_\theta(G(w)) + \mathbb{E}_{P_\theta(x)} \mathbb{E}_{q_\psi(z|x)} \nabla_\theta (\ln P_\theta(x|z) - \nabla_\theta U_\theta(z))$$

$$= \nabla_\theta \mathbb{E}_{P_\theta(x)} \mathbb{E}_{q_\psi(w)} U_\theta(G(w)) + \mathbb{E}_{q_\psi(z|x)} \ln P_\theta(x|z) - U_\theta(z)$$

$$= \nabla_\theta \mathcal{L}_{EVAE}$$

Part II:

$$\begin{aligned} & KL(q_\psi(z) || P_\theta(z)) \\ &= \mathbb{E} \ln q_\psi(z) + \mathbb{E} -U_\theta(z) - \ln Z(\theta) \end{aligned}$$



$$q_\psi(z) = \mathbb{E}_{q_\psi(w)} \delta(z - G(w))$$

$$q_\psi(z) \rightarrow P_\theta(z).$$

$$P_\theta(z) = \frac{1}{Z(\theta)} e^{-U_\theta(z)}$$

$$\nabla_\theta \ln P_\theta(z) = -\nabla_\theta U_\theta(z) - \frac{1}{Z(\theta)} \nabla_\theta Z(\theta)$$

$$q_\psi(z) \quad q_\psi(z)$$

$$= -H[q_\psi(z)] - \mathbb{E}_{q_\psi(z)} U_\theta(z) \quad \text{By Lemma 1.}$$

$$= -I(z, w) - \mathbb{E}_{q_\psi(z)} U_\theta(z) \quad \text{By Lemma 2.}$$

$$\sim \lambda_1 \mathbb{E}_{q_\psi(w)} \|w - E(G(w))\|^2 - \mathbb{E}_{q_\psi(z)} U_\theta(z).$$

$$\begin{aligned} & KL(q_\phi(z|x) \| p_\theta(z|x)) \\ &= H[q_\phi(z|x)] - \mathbb{E}_{q_\phi(z|x)} \ln \frac{p_\theta(x|z) p_\theta(z)}{p_\theta(x)} \\ &\sim H[q_\phi(z|x)] - \mathbb{E}_{q_\phi(z|x)} \ln p_\theta(x|z) - U_\theta(z) - \log \underbrace{p_\theta(z)}_{\text{Const}} \end{aligned}$$

两项加起来的为:

$$\min_{\psi, \phi} \lambda_1 \mathbb{E}_{q_\psi(w)} \|w - E(G(w))\|^2 - \mathbb{E}_{q_\psi(w)} U_\theta(G(z))$$

$$\lambda_2 (H[q_\phi(z|x)] + \mathbb{E}_{q_\phi(z|x)} U_\theta(z) - \ln p_\theta(x|z))$$

$$\max_{\theta} \mathbb{E}_{p(x)} \left[\mathbb{E}_{q_\psi(w)} U_\theta(G(w)) + \mathbb{E}_{q_\phi(z|x)} \ln p_\theta(x|z) - U_\theta(z) - \lambda_3 \| \nabla_z U_\theta(z) \|^2 \right]$$

在最后加上 $\lambda_3 \| \nabla_z U_\theta(z) \|^2$ 作为梯度惩罚项。

Lemma 1. $I(z, w) = H[q_\psi(z)] + \text{Const.}$

$$I(z, w) = \int q_\psi(z, w) [\ln q_\psi(z, w) - \ln q_\psi(w) q_\psi(z)] dz dw$$

$$= \mathbb{E}_{q_{\psi}(z, w)} \ln q_{\psi}(z|w) - \mathbb{E}_{q_{\psi}(z, w)} \ln q_{\psi}(z)$$

$$= \mathbb{E}_{q_{\psi}(w)} - H[q_{\psi}(z|w)] - \mathbb{E}_{q_{\psi}(z)} \ln q_{\psi}(z)$$

$$q_{\psi}(z|w) \approx \lim_{\delta \rightarrow 0} \mathcal{N}(G(w), \delta^2)$$

$$H[\mathcal{N}(G(w), \delta^2)] = n \log \delta$$

$$H[q_{\psi}(z|w)] = H[\lim_{\delta \rightarrow 0} \mathcal{N}(G(w), \delta^2)] = n \lim_{\delta \rightarrow 0} \log \delta = -\infty$$

$$\text{Then } I(z, w) = H(z) + \text{Const.} \quad \square$$

$$\text{Lemma 2. } I(z, w) = -\lambda_1 \mathbb{E}_{q_{\psi}(w)} \|w - \mathbb{E}[G(w)]\|^2$$

$$\text{Let } p_{\lambda}(w|z) = \mathcal{N}(\mathbb{E}(z), \sigma_{\lambda}^2).$$

$$I(z, w) = \iint q_{\psi}(z, w) \ln \frac{q_{\psi}(w|z)}{p_{\lambda}(w|z)} + \ln \frac{p_{\lambda}(w|z)}{q_{\psi}(w)} dw dz$$

$$= \mathbb{E}_{q_{\psi}(w)} \mathbb{E}_{q_{\psi}(w|z)} \ln \frac{q_{\psi}(w|z)}{p_{\lambda}(w|z)} + \mathbb{E}_{q_{\psi}(w)} \mathbb{E}_{q_{\psi}(z|w)} \ln \frac{p_{\lambda}(w|z)}{q_{\psi}(w)}$$

$$= \mathbb{E}_{q_{\psi}(w)} \text{KL}(q_{\psi}(w|z) \| p_{\lambda}(w|z)) + \mathbb{E}_{q_{\psi}(w)} \mathbb{E}_{q_{\psi}(z|w)} \ln \frac{p_{\lambda}(w|z)}{q_{\psi}(w)}$$

$$\geq \mathbb{E}_{q_{\psi}(w)} \mathbb{E}_{q_{\psi}(z|w)} \ln p_{\lambda}(w|z) - \mathbb{E}_{q_{\psi}(w)} \mathbb{E}_{q_{\psi}(z|w)} \ln q_{\psi}(w)$$

$$\begin{aligned}
 &\sim E_{q_{\psi}(w)} E_{q_{\psi}(z|w)} \sum_i^n \left(-\frac{\|w_i - E_i(z)\|^2}{2\sigma_i^2} - \frac{1}{2} \log \sigma_i^2 \right) - H[q_{\psi}(w)] \\
 &\sim E_{q_{\psi}(w)} \sum_i^n \left(-\frac{1}{2\sigma_i^2} \|w - E \circ G(z)\|^2 - \log \sigma_i \right) \square \quad \begin{matrix} \uparrow \\ \text{const.} \end{matrix}
 \end{aligned}$$