

Problem 3.

(1) Let

$$q_\phi(z|x, y) \sim \mathcal{N}(\mu_1, \Sigma_1) \quad p_\theta(z|y) \sim \mathcal{N}(\mu_2, \Sigma_2).$$

We use part of the results in (2)

$$p(x) = \frac{1}{(2\pi)^{\frac{n}{2}} \det(\Sigma)^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right).$$

$$D_{KL}(q_\phi \| p_\theta) = \frac{1}{2} \left[\log\left(\frac{\det \Sigma_2}{\det \Sigma_1}\right) - n + \text{tr}(\Sigma_2^{-1} \Sigma_1) + (\mu_2 - \mu_1)^T \Sigma_2^{-1} (\mu_2 - \mu_1) \right]$$

$$\mathbb{E}_{q_\phi} [\log p_\theta] = \mathbb{E}_{q_\phi} \left[-\frac{n}{2} \log(2\pi) - \frac{1}{2} \log(\det(\Sigma_2)) - \frac{1}{2} (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) \right].$$

$$= \mathbb{E}_{q_\phi} \left[-\frac{n}{2} \log(2\pi) - \frac{1}{2} \log(\det(\Sigma_2)) - \frac{1}{2} (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) \right].$$

$$= \mathbb{E}_{q_\phi} \left[-\frac{n}{2} \log(2\pi) - \frac{1}{2} \log(\det(\Sigma_2)) - \frac{1}{2} \text{tr}(\Sigma_2^{-1} (x x^T - 2x \mu_2^T + \mu_2 \mu_2^T)) \right].$$

$$= -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log(\det(\Sigma_2)) - \frac{1}{2} \text{tr}(\Sigma_2^{-1} (\Sigma_1 + \mu_1 \mu_1^T - 2\mu_1 \mu_2^T + \mu_2 \mu_2^T))$$

$$= -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log(\det(\Sigma_2)) - \frac{1}{2} \text{tr}(\Sigma_2^{-1} \Sigma_1) + \frac{1}{2} (\mu_1 - \mu_2)^T \Sigma_2^{-1} (\mu_1 - \mu_2).$$

$$\mathbb{E}_{q_\phi} - D_{KL} = -\frac{n}{2} \log(2\pi) + \frac{n}{2} - \log(\det(\Sigma_2)) + \frac{1}{2} \log(\det(\Sigma_1)) - \text{tr}(\Sigma_2^{-1} \Sigma_1)$$

$$\log p_\theta = -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log(\det(\Sigma_2)) - \frac{1}{2} (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2).$$

$$\Rightarrow \log p_\theta \geq \mathbb{E}_{q_\phi} [\log p_\theta] - D_{KL}(q_\phi \| p_\theta).$$

2. We directly cite John Duchi's lecture notes.

Assume P_1, P_2 Gaussian vectors,

$$P_1 \sim \mathcal{N}(\mu_1, \Sigma_1) \quad P_2 \sim \mathcal{N}(\mu_2, \Sigma_2)$$

$$D_{KL}(P_1 \| P_2) = \mathbb{E}_{P_1} [\log P_1 - \log P_2]$$

$$= \frac{1}{2} \mathbb{E}_{P_1} \left[-\log(\det \Sigma_1) - (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) + \log(\det \Sigma_2) + (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) \right]$$

This is from the density

$$p(x) = \frac{1}{(2\pi)^{\frac{n}{2}} \det(\Sigma)^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right)$$

Then

$$\begin{aligned} D_{KL}(P_1 \| P_2) &= \frac{1}{2} \log \left(\frac{\det \Sigma_2}{\det \Sigma_1} \right) + \frac{1}{2} \mathbb{E}_{P_1} \left[-\text{tr}(\Sigma_1^{-1} (x - \mu_1)(x - \mu_1)^T) + \text{tr}(\Sigma_2^{-1} (x - \mu_2)(x - \mu_2)^T) \right] \\ &= \frac{1}{2} \log \left(\frac{\det \Sigma_2}{\det \Sigma_1} \right) + \frac{1}{2} \mathbb{E}_{P_1} \left[-\text{tr}(\Sigma_1^{-1} \Sigma_1) + \text{tr}(\Sigma_2^{-1} (x x^T - 2x \mu_2^T + \mu_2 \mu_2^T)) \right] \\ &= \frac{1}{2} \log \left(\frac{\det \Sigma_2}{\det \Sigma_1} \right) - \frac{1}{2} n + \frac{1}{2} \text{tr} \left[\Sigma_2^{-1} (\Sigma_1 + \mu_1 \mu_1^T - 2\mu_2 \mu_1^T + \mu_2 \mu_2^T) \right] \\ &= \frac{1}{2} \left\{ \log \left(\frac{\det \Sigma_2}{\det \Sigma_1} \right) - n + \text{tr}(\Sigma_2^{-1} \Sigma_1) + \text{tr}(\mu_1^T \Sigma_2^{-1} \mu_1 - 2\mu_1^T \Sigma_2^{-1} \mu_2 + \mu_2^T \Sigma_2^{-1} \mu_2) \right\} \\ &= \frac{1}{2} \left[\log \left(\frac{\det \Sigma_2}{\det \Sigma_1} \right) - n + \text{tr}(\Sigma_2^{-1} \Sigma_1) + (\mu_2 - \mu_1)^T \Sigma_2^{-1} (\mu_2 - \mu_1) \right] \end{aligned}$$

Now take $\Sigma_1 = \text{diag}(\sigma_1^2, \dots, \sigma_J^2)$ $\mu_1 = (\mu_1, \dots, \mu_J)$
 $\Sigma_2 = I$ $\mu_2 = (0, \dots, 0)$

$$D_{KL}(P_1 \| P_2)$$

$$= \frac{1}{2} \left(\log \frac{1}{\prod_{i=1}^J \sigma_i^2} - J + \text{tr}(\Sigma_1) + \mu_1^T \mu_1 \right)$$

$$= \frac{1}{2} \left[-\sum_{i=1}^J \log \sigma_i^2 - \sum_{i=1}^J 1 + \sum_{i=1}^J \sigma_i^2 + \sum_{i=1}^J \mu_i^2 \right]$$

$$= -\frac{1}{2} \sum_{j=1}^J (1 + \log(\sigma_j^2) - \mu_j^2 - \sigma_j^2)$$