

# Homework 3

Deep Learning 2025 Spring

Due on 2025/3/24

## 1 True or False

**Problem 1.** If we optimize  $q_\theta$  w.r.t. a multi-modal distribution  $p$  using KL-divergence  $\text{KL}(q_\theta \| p)$ , we will get a distribution that uniformly covers all the modes.

**Problem 2.** The reparameterization trick applied in VAE helps passing gradient back to the encoder.

**Problem 3.** It is easy to compute the exact posterior  $p(\mathbf{z}|\mathbf{x})$  using VAE.

**Problem 4.** In  $\beta$ -VAE, large  $\beta$  enforces latent variables to be correlated with each other.

## 2 Q&A

**Problem 5.** (EM Algorithm) In statistics, expectation–maximization (EM) algorithm is an iterative method to find (local) maximum likelihood or maximum a posteriori (MAP) estimates of parameters in statistical models, where the model depends on unobserved latent variables.<sup>1</sup> Consider a latent variable model with parameter  $\theta$

$$p_\theta(\mathbf{x}, \mathbf{z}) = p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})$$

and we want to find the MLE of  $\theta$ , i.e.,

$$\hat{\theta}_{\text{MLE}} = \arg \max_{\theta} \log p_\theta(\mathbf{x}) = \arg \max_{\theta} \log \sum_{\mathbf{z}} p_\theta(\mathbf{x}, \mathbf{z})$$

The E-step (expectation) of EM algorithm is given by

$$Q(\theta|\theta^{(t)}) = \mathbb{E}_{\mathbf{z} \sim p_{\theta^{(t)}}(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}, \mathbf{z})]$$

and the M-step (maximization) is given by

$$\theta^{(t+1)} = \arg \max_{\theta} Q(\theta|\theta^{(t)}).$$

Prove that the following optimization process is equivalent to EM algorithm. Define  $F(\theta, q) = \mathbb{E}_{\mathbf{z} \sim q} [\log p_\theta(\mathbf{x}, \mathbf{z})] + H(q)$ , where  $H(\cdot)$  is Shannon entropy.

---

<sup>1</sup>[https://en.wikipedia.org/wiki/Expectation-maximization\\_algorithm](https://en.wikipedia.org/wiki/Expectation-maximization_algorithm)

(E-step)

$$q^{(t)} = \arg \max_q F(\theta^{(t)}, q)$$

(M-step)

$$\theta^{(t)} = \arg \max_{\theta} F(\theta, q^{(t)})$$

**Problem 6.** (KL-Divergence)

1. (Gaussian) Prove that the KL-divergence between two  $d$ -dimensional Gaussian distributions  $\mathcal{N}_0(\mu_0, \Sigma_0)$  and  $\mathcal{N}_1(\mu_1, \Sigma_1)$  has the following form:

$$\text{KL}(\mathcal{N}_0 \parallel \mathcal{N}_1) = \frac{1}{2} \left\{ \text{tr}(\Sigma_1^{-1} \Sigma_0) + (\mu_1 - \mu_0)^T \Sigma_1^{-1} (\mu_1 - \mu_0) - d + \log \frac{|\Sigma_1|}{|\Sigma_0|} \right\}.$$

2. (Convexity) Let  $\lambda \in [0, 1] \subset \mathbb{R}$ .  $p_1, p_2, q_1$  and  $q_2$  are discrete distributions over alphabet  $\mathcal{Y} = \{1, 2, \dots, n\}$  with nonzero probabilities. Prove

$$\text{KL}(\lambda p_1 + (1 - \lambda) p_2 \parallel \lambda q_1 + (1 - \lambda) q_2) \leq \lambda \text{KL}(p_1 \parallel q_1) + (1 - \lambda) \text{KL}(p_2 \parallel q_2).$$

3. (Inclusive/Exclusive) We can recognize the difference of inclusive and exclusive KL via a simple example. Consider the target distribution

$$p(x) = \frac{1}{3} \mathcal{N}(-3, 1) + \frac{2}{3} \mathcal{N}(3, 1)$$

which is a multi-modal Gaussian mixture. We model the variational distribution  $q(x)$  as a Gaussian distribution with mean  $\mu$  and variance  $\sigma^2$ , where  $\mu$  and  $\sigma$  are unknown parameters. Write a program to find the optimal  $\mu$  and  $\sigma$  w.r.t. inclusive and exclusive KL respectively. You need to submit a figure demonstrating the original distribution and two derived variational distributions.

4. (Variational Inference) While we used *reverse KL-divergence*  $\text{KL}(q_{\psi}(z|x) \parallel p(z|x))$  to conduct variational inference (i.e., optimize the first term  $q_{\psi}(z|x)$  with parameter  $\psi$  to approximate the second term  $p(z|x)$ ) in lecture, Bob proposes to use the *forward KL-divergence*  $\text{KL}(p(z|x) \parallel q_{\psi}(z|x))$ . In this case, what would be the objective for  $q_{\psi}(z|x)$ ? What are the pros and cons if we use this objective for  $q_{\psi}(z|x)$  in VAE?

**Hint:** The objective should be in the form of expectation.

**Problem 7.** (GM-VAE) In standard VAEs, the prior of the latent variables is assumed to be an isotropic Gaussian. In this problem, we use a mixture of Gaussian distributions as the prior to allow more complicated latent representations, named Gaussian Mixture Variational Auto-Encoder (GM-VAE).

Consider a latent variable model  $p_{\mu, \sigma, \theta}(x, w, z) = p(z) p_{\mu, \sigma}(w|z) p_{\theta}(x|w)$ , where an observable sample  $x$  is gener-

ated from latent variable  $w$  and  $z$ :

$$\begin{aligned}
 \mathbf{z} &\sim \text{Categorical}(\boldsymbol{\pi}), \mathbb{P}(\mathbf{z} = k) = \pi_k \text{ for } 1 \leq k \leq K \text{ and } \sum_{k=1}^K \pi_k = 1 \\
 \mathbf{w}|\mathbf{z} &\sim \prod_{k=1}^K \mathcal{N}(\mu_k, \sigma_k^2 I)^{\mathbb{I}(\mathbf{z}=k)} \\
 \mathbf{x}|\mathbf{w} &\sim \mathcal{N}(\mu_\theta(\mathbf{w}), \sigma_\theta^2(\mathbf{w}))
 \end{aligned}$$

where  $\boldsymbol{\mu} = [\mu_1, \dots, \mu_K]$ ,  $\boldsymbol{\sigma} = [\sigma_1, \dots, \sigma_K]$ , and  $\theta$  are trainable parameters. The prior distribution over  $\mathbf{z}$  is uniform over alphabet  $\{1, \dots, K\}$ . Define a variational model  $q_{\psi, \phi}(\mathbf{w}, \mathbf{z}|\mathbf{x}) = q_\psi(\mathbf{w}|\mathbf{x})q_\phi(\mathbf{z}|\mathbf{w}, \mathbf{x})$ , where  $\psi$  and  $\phi$  are trainable parameters.

1. Derive ELBO for  $\log p_{\boldsymbol{\mu}, \boldsymbol{\sigma}, \theta}(\mathbf{x})$ . Your answer should include 3 terms containing  $p(\mathbf{z})$ ,  $p_{\boldsymbol{\mu}, \boldsymbol{\sigma}}(\mathbf{w}|\mathbf{z})$  and  $p_\theta(\mathbf{x}|\mathbf{w})$  respectively.
2. Design a training procedure for GM-VAE.