

STA2311 (FALL 2023) - PRACTICE PROBLEMS FOR CLASS 5  
(VARIATIONAL INFERENCE)

1. Let  $p(x, y)$  and  $q(x, y)$  be two bivariate mass functions. Write  $p_1(x) = \sum_y p(x, y)$  and  $p_2^x(y) = p(y | x)$ , and write  $q_1(x)$  and  $q_2^x(y)$  similarly. Prove that

$$\text{KL}(p \parallel q) = \text{KL}(p_1 \parallel q_1) + \mathbb{E}_X[\text{KL}(p_2^X \parallel q_2^X)],$$

where  $X \sim p_1$ .

2. Show that the KL-divergence does not always satisfy the triangle inequality; that is, there exist distributions  $p, q, r$  such that

$$\text{KL}(p \parallel r) \not\leq \text{KL}(p \parallel q) + \text{KL}(q \parallel r)$$

3. Let  $p_1(\mathbf{x}) = \mathcal{N}_d(\mathbf{x} | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$  and  $p_2(\mathbf{x}) = \mathcal{N}_d(\mathbf{x} | \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$ . Compute  $\text{KL}(p_1 \parallel p_2)$ .

4. Consider the *differential entropy*  $H[\cdot]$  defined on the space of density functions.

(a) Show that the differential entropy is translation invariant in the sense that if  $X \sim f$  and  $X + c \sim f_c$ , then  $H[f] = H[f_c]$  for all  $c \in \mathbb{R}$ .

(b) Show that among all continuous univariate distributions  $f$  with mean  $\mu$  and variance  $\sigma^2$ , the  $\mathcal{N}(\mu, \sigma^2)$  distribution is the one that maximizes  $H[f]$ .

5. Suppose we approximate a  $d$ -dimensional distribution  $p(\mathbf{z})$  by a mean-field variational family  $q(\mathbf{z}) = \prod_{i=1}^d q_i(z_i)$ . Show that minimizing  $\text{KL}(p \parallel q)$  with respect to one factor  $q_i(z_i)$ , keeping all other factors fixed, leads to the optimal solution

$$q_i^*(z_i) = \int p(\mathbf{z}) dz_{-i}.$$

6. Consider linear regression: we have independent observations  $Y_1, \dots, Y_n$  and covariates  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^p$  with  $Y_i | \mathbf{x}_i \sim \mathcal{N}(\boldsymbol{\beta}^\top \mathbf{x}_i, \sigma^2)$  for some  $\boldsymbol{\beta} \in \mathbb{R}^p$ ; we assume that  $\sigma^2 > 0$  is known. We adopt a Bayesian model and impose a  $\mathcal{N}_p(\mathbf{0}, \alpha^{-1} \mathbf{I})$  prior on  $\boldsymbol{\beta}$  and a  $\text{Gamma}(a_0, b_0)$  prior on  $\alpha$ . Approximate the posterior  $p(\boldsymbol{\beta}, \alpha | \mathbf{y})$  by deriving a mean-field variational approximation of the form  $q(\boldsymbol{\beta}, \alpha) = q(\boldsymbol{\beta}) \cdot q(\alpha)$ .

7. Let  $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Gamma}(\theta, 1)$  for some  $\theta > 0$ , and let  $Z_i | X_i = x_i \stackrel{\text{indep}}{\sim} \text{Gamma}(x_i, \theta)$ . We are once again Bayesians and adopt a  $\text{Gamma}(a, b)$  prior on  $\theta$ . Approximate the posterior  $p(\theta, \mathbf{z} | \mathbf{x})$  by deriving a mean-field variational approximation of the form  $q(\theta, \mathbf{z}) = q(\theta) \cdot q(\mathbf{z})$ .