

## PROBABILISTIC INFERENCE &amp; LEARNING

## Exercise Sheet #11

## EM &amp; Variational Inference

1. **EM for MAP:** Suppose we wish to use the EM algorithm to maximize the *posterior* distribution over parameters  $p(\theta \mid x)$  for a model containing latent variables  $z$ , where  $x$  is the observed data set (rather than the *likelihood*  $p(x \mid \theta)$ , as done in the lectures). Show that the E step remains the same as in the maximum likelihood case, whereas in the M step the quantity to be maximized in  $\theta$  is given (up to constants) by  $\mathcal{L}(q, \theta) + \log p(\theta)$ , where the ELBO (negative variational free energy)  $\mathcal{L}(q, \theta)$  is defined as in the lecture, and  $q(z) = p(z \mid x, \theta)$ . 20 points
2. **Softer  $k$ -means:** Consider a special case of a Gaussian mixture model in which the covariance matrices  $\Sigma_k$  of the components are all constrained to have a common value  $\Sigma$ . Derive the EM equations for maximizing the likelihood function under such a model. 20 points
3. **Gibbs' inequality:** Jensen's inequality states that, for a convex function  $f : \mathbb{R} \rightarrow \mathbb{R}$  and a probability density  $p(x)$ ,

$$\mathbb{E}_p(f(x)) \geq f(\mathbb{E}_p(x)). \quad (1)$$

Use Jensen's inequality to show that the Kullback-Leibler divergence  $D_{\text{KL}}(p \parallel q)$  between two arbitrary<sup>1</sup> probability distributions  $p, q$  is nonnegative:

$$D_{\text{KL}}(p \parallel q) := \int p(x) \log \left( \frac{p(x)}{q(x)} \right) dx \geq 0 \quad (2)$$

This result is known as *Gibbs' inequality* and was quoted frequently in the lecture. Hint: Consider the convex function  $f(u) = \log(1/u)$  of the variable  $u(x) = p(x)/q(x) > 0$ . 20 points

4. **Free energy for Gaussians:** In the lecture, a sketch was shown to argue that an approximation  $q$  to a distribution  $p$  found by minimizing  $D_{\text{KL}}(p \parallel q)$  tends to be "too wide", while an approximation found by minimizing  $D_{\text{KL}}(q \parallel p)$  tends to be "too narrow". The following argument supports this insight using the simple case of Gaussian distributions:

- (a) Show that the Kullback-Leibler divergence between two scalar, centered Gaussian distributions is given by 10 points

$$D_{\text{KL}}(\mathcal{N}(x; 0, \sigma_q^2) \parallel \mathcal{N}(x; 0, \sigma_p^2)) = \frac{1}{2} \left( \log \left( \frac{\sigma_p^2}{\sigma_q^2} \right) - 1 + \frac{\sigma_q^2}{\sigma_p^2} \right). \quad (3)$$

- (b) Consider the two-dimensional Gaussian distribution  $p$  and a spherical approximation  $q$ , each given by (with  $\sigma_1 \neq \sigma_2$ )

$$p(x_1, x_2) = \mathcal{N} \left( \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}; \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix} \right), \quad p(x_1, x_2) = \mathcal{N} \left( \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}; \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_q^2 & 0 \\ 0 & \sigma_q^2 \end{bmatrix} \right), \quad (4)$$

Find the values of  $\sigma_q$  that minimize  $D_{\text{KL}}(p \parallel q)$  and  $D_{\text{KL}}(q \parallel p)$ , respectively. 30 points

<sup>1</sup>You may assume that  $p(x)/q(x) > 0$  exists for all  $x$ .