

Exercises for the tutorials: Jupyter notebook at <https://github.com/vsimkus/pmr2022-vaе>.

The other exercises are for self-study and exam preparation. All material is examinable unless otherwise mentioned.

Exercise 1. Mean field variational inference I

Let $\mathcal{L}_{\mathbf{x}}(q)$ be the evidence lower bound for the marginal $p(\mathbf{x})$ of a joint pdf/pmf $p(\mathbf{x}, \mathbf{y})$,

$$\mathcal{L}_{\mathbf{x}}(q) = \mathbb{E}_{q(\mathbf{y}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{y})}{q(\mathbf{y}|\mathbf{x})} \right]. \quad (1)$$

Mean field variational inference assumes that the variational distribution $q(\mathbf{y}|\mathbf{x})$ fully factorises, i.e.

$$q(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^d q_i(y_i|\mathbf{x}), \quad (2)$$

when \mathbf{y} is d -dimensional. An approach to learning the q_i for each dimension is to update one at a time while keeping the others fixed. We here derive the corresponding update equations.

- (a) Show that the evidence lower bound $\mathcal{L}_{\mathbf{x}}(q)$ can be written as

$$\mathcal{L}_{\mathbf{x}}(q) = \mathbb{E}_{q_1(y_1|\mathbf{x})} \mathbb{E}_{q(\mathbf{y}_{\setminus 1}|\mathbf{x})} [\log p(\mathbf{x}, \mathbf{y})] - \sum_{i=1}^d \mathbb{E}_{q_i(y_i|\mathbf{x})} [\log q_i(y_i|\mathbf{x})] \quad (3)$$

where $q(\mathbf{y}_{\setminus 1}|\mathbf{x}) = \prod_{i=2}^d q_i(y_i|\mathbf{x})$ is the variational distribution without $q_1(y_1|\mathbf{x})$.

- (b) Assume that we would like to update $q_1(y_1|\mathbf{x})$ and that the variational marginals of the other dimensions are kept fixed. Show that

$$\operatorname{argmax}_{q_1(y_1|\mathbf{x})} \mathcal{L}_{\mathbf{x}}(q) = \operatorname{argmin}_{q_1(y_1|\mathbf{x})} \operatorname{KL}(q_1(y_1|\mathbf{x}) || \bar{p}(y_1|\mathbf{x})) \quad (4)$$

with

$$\log \bar{p}(y_1|\mathbf{x}) = \mathbb{E}_{q(\mathbf{y}_{\setminus 1}|\mathbf{x})} [\log p(\mathbf{x}, \mathbf{y})] + \text{const}, \quad (5)$$

where const refers to terms not depending on y_1 . That is, $\bar{p}(y_1|\mathbf{x}) = \frac{1}{Z} \exp \left[\mathbb{E}_{q(\mathbf{y}_{\setminus 1}|\mathbf{x})} [\log p(\mathbf{x}, \mathbf{y})] \right]$, where Z is the normalising constant. Note that variables y_2, \dots, y_d are marginalised out due to the expectation with respect to $q(\mathbf{y}_{\setminus 1}|\mathbf{x})$.

- (c) Conclude that given $q_i(y_i|\mathbf{x})$, $i = 2, \dots, d$, the optimal $q_1(y_1|\mathbf{x})$ equals $\bar{p}(y_1|\mathbf{x})$.

This then leads to an iterative updating scheme where we cycle through the different dimensions, each time updating the corresponding marginal variational distribution according to:

$$q_i(y_i|\mathbf{x}) = \bar{p}(y_i|\mathbf{x}), \quad \bar{p}(y_i|\mathbf{x}) = \frac{1}{Z} \exp \left[\mathbb{E}_{q(\mathbf{y}_{\setminus i}|\mathbf{x})} [\log p(\mathbf{x}, \mathbf{y})] \right] \quad (6)$$

where $q(\mathbf{y}_{\setminus i}|\mathbf{x}) = \prod_{j \neq i} q_j(y_j|\mathbf{x})$ is the product of all marginals without marginal $q_i(y_i|\mathbf{x})$.

Exercise 2. Mean field variational inference II

Assume random variables y_1, y_2, x are generated according to the following process

$$y_1 \sim \mathcal{N}(y_1; 0, 1) \qquad y_2 \sim \mathcal{N}(y_2; 0, 1) \qquad (7)$$

$$n \sim \mathcal{N}(n; 0, 1) \qquad x = y_1 + y_2 + n \qquad (8)$$

where y_1, y_2, n are statistically independent.

- (a) y_1, y_2, x are jointly Gaussian. Determine their mean and their covariance matrix.
- (b) The conditional $p(y_1, y_2|x)$ is Gaussian with mean \mathbf{m} and covariance \mathbf{C} ,

$$\mathbf{m} = \frac{x}{3} \begin{pmatrix} 1 \\ 1 \end{pmatrix} \qquad \mathbf{C} = \frac{1}{3} \begin{pmatrix} 2 & -1 \\ -1 & 2 \end{pmatrix} \qquad (9)$$

Since x is the sum of three random variables that have the same distribution, it makes intuitive sense that the mean assigns $1/3$ of the observed value of x to y_1 and y_2 . Moreover, y_1 and y_2 are negatively corrected since an increase in y_1 must be compensated with a decrease in y_2 .

Let us now approximate the posterior $p(y_1, y_2|x)$ with mean field variational inference. Determine the optimal variational distribution using the method and results from Exercise 1. You may use that

$$p(y_1, y_2, x) = \mathcal{N}((y_1, y_2, x); \mathbf{0}, \Sigma) \qquad \Sigma = \begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \\ 0 & 1 & 3 \end{pmatrix} \qquad \Sigma^{-1} = \begin{pmatrix} 2 & 1 & -1 \\ 1 & 2 & -1 \\ -1 & -1 & 1 \end{pmatrix} \qquad (10)$$

Exercise 3. Variational posterior approximation I

We have seen that maximising the evidence lower bound (ELBO) with respect to the variational distribution q minimises the Kullback-Leibler divergence to the true posterior p . We here assume that q and p are probability density functions so that the Kullback-Leibler divergence between them is defined as

$$\text{KL}(q||p) = \int q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x})} d\mathbf{x} = \mathbb{E}_q \left[\log \frac{q(\mathbf{x})}{p(\mathbf{x})} \right]. \qquad (11)$$

- (a) You can here assume that \mathbf{x} is one-dimensional so that p and q are univariate densities. Consider the case where p is a bimodal density but the variational densities q are unimodal. Sketch a figure that shows p and a variational distribution q that has been learned by minimising $\text{KL}(q||p)$. Explain qualitatively why the sketched q minimises $\text{KL}(q||p)$.
- (b) Assume that the true posterior $p(\mathbf{x}) = p(x_1, x_2)$ factorises into two Gaussians of mean zero and variances σ_1^2 and σ_2^2 ,

$$p(x_1, x_2) = \frac{1}{\sqrt{2\pi\sigma_1^2}} \exp \left[-\frac{x_1^2}{2\sigma_1^2} \right] \frac{1}{\sqrt{2\pi\sigma_2^2}} \exp \left[-\frac{x_2^2}{2\sigma_2^2} \right]. \qquad (12)$$

Assume further that the variational density $q(x_1, x_2; \lambda^2)$ is parametrised as

$$q(x_1, x_2; \lambda^2) = \frac{1}{2\pi\lambda^2} \exp \left[-\frac{x_1^2 + x_2^2}{2\lambda^2} \right] \qquad (13)$$

where λ^2 is the variational parameter that is learned by minimising $\text{KL}(q||p)$. If σ_2^2 is much larger than σ_1^2 , do you expect λ^2 to be closer to σ_2^2 or to σ_1^2 ? Provide an explanation.

Exercise 4. *Variational posterior approximation II*

We have seen that maximising the evidence lower bound (ELBO) with respect to the variational distribution minimises the Kullback-Leibler divergence to the true posterior. We here investigate the nature of the approximation if the family of variational distributions does not include the true posterior.

- (a) Assume that the true posterior for $\mathbf{x} = (x_1, x_2)$ is given by

$$p(\mathbf{x}) = \mathcal{N}(x_1; \sigma_1^2) \mathcal{N}(x_2; \sigma_2^2) \quad (14)$$

and that our variational distribution $q(\mathbf{x}; \lambda^2)$ is

$$q(\mathbf{x}; \lambda^2) = \mathcal{N}(x_1; \lambda^2) \mathcal{N}(x_2; \lambda^2), \quad (15)$$

where $\lambda > 0$ is the variational parameter. Provide an equation for $J(\lambda) = \text{KL}(q(\mathbf{x}; \lambda^2) \| p(\mathbf{x}))$, where you can omit additive terms that do not depend on λ .

- (b) Determine the value of λ that minimises $J(\lambda) = \text{KL}(q(\mathbf{x}; \lambda^2) \| p(\mathbf{x}))$. Interpret the result and relate it to properties of the Kullback-Leibler divergence.