

# Variational Inference

Xingdong Zuo\*

Let  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  be a set of training inputs and let  $\mathbf{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_N\}$  be a set of training outputs. The objective is to find the parameters  $\omega$  of a parameterized function  $\mathbf{y} = f_\omega(\mathbf{x})$ . In Bayesian modelling, we place a prior distribution  $p(\omega)$  over the parameters, i.e. initial belief of how likely the parameters have generated the observations. We define a likelihood distribution  $p(\mathbf{y}|\mathbf{x}, \omega)$  to indicate how likely  $\mathbf{y}$  is generated by the input  $\mathbf{x}$  given the parameters  $\omega$ . The posterior distribution over the parameters is defined by

$$p(\omega|\mathbf{X}, \mathbf{Y}) = \frac{p(\mathbf{Y}|\mathbf{X}, \omega)p(\omega)}{p(\mathbf{Y}|\mathbf{X})}$$

Then the predictive distribution of a new test input  $\mathbf{x}^*$  can be obtained by integrating model parameters

$$p(\mathbf{y}^*|\mathbf{x}^*, \mathbf{X}, \mathbf{Y}) = \int p(\mathbf{y}^*|\mathbf{x}^*, \omega)p(\omega|\mathbf{X}, \mathbf{Y}) d\omega$$

However, the denominator of the posterior  $p(\mathbf{Y}|\mathbf{X}) = \int p(\mathbf{Y}|\mathbf{X}, \omega)p(\omega) d\omega$  is intractable to obtain analytically. Because true posterior is intractable, in VI, we use an approximating variational distribution  $q_\phi(\omega|\mathbf{X}, \mathbf{Y}) \approx p_\theta(\omega|\mathbf{X}, \mathbf{Y})$ . The optimum is obtained by minimizing the KL divergence between two distributions, leading to a fundamental theorem in VI. Intuitively, minimizing KL divergence is equivalent to maximizing the evidence lower bound (ELBO) w.r.t. the variational parameters  $\phi$ .

**Theorem 0.1.**

$$\begin{aligned} & \min_{\phi} D_{\text{KL}}(q_\phi(\omega|\mathbf{D})||p_\theta(\omega|\mathbf{D})) \\ \iff & \max_{\phi} \left\{ \mathbb{E}_{q_\phi(\omega|\mathbf{D})} [\log p_\theta(\mathbf{D}|\omega)] - D_{\text{KL}}(q_\phi(\omega|\mathbf{D})||p_\theta(\omega)) \right\} \end{aligned}$$

**First proof**

*Proof.* By the definition of KL divergence, we can obtain

$$\begin{aligned} & D_{\text{KL}}(q_\phi(\omega|\mathbf{D})||p_\theta(\omega|\mathbf{D})) \\ = & \mathbb{E}_{q_\phi(\omega|\mathbf{D})} [\log q_\phi(\omega|\mathbf{D})] + \mathbb{E}_{q_\phi(\omega|\mathbf{D})} [\log p_\theta(\mathbf{D})] - \mathbb{E}_{q_\phi(\omega|\mathbf{D})} [\log p_\theta(\mathbf{D}, \omega)]. \end{aligned}$$

Note that  $\log p_\theta(\mathbf{D}) = \mathbb{E}_{q_\phi(\omega|\mathbf{D})} [\log p_\theta(\mathbf{D})]$  and by rearranging the equality, we have

$$\begin{aligned} & \log p_\theta(\mathbf{D}) \\ = & \mathbb{E}_{q_\phi(\omega|\mathbf{D})} [\log p_\theta(\mathbf{D}, \omega)] - \mathbb{E}_{q_\phi(\omega|\mathbf{D})} [\log q_\phi(\omega|\mathbf{D})] + D_{\text{KL}}(q_\phi(\omega|\mathbf{D})||p_\theta(\omega|\mathbf{D})). \end{aligned}$$

Applying the non-negativity property of the KL divergence i.e.  $D_{\text{KL}}(q_\phi(\omega|\mathbf{D})||p_\theta(\omega|\mathbf{D})) \geq 0$  we have the following

$$\begin{aligned} \log p_\theta(\mathbf{D}) & \geq \mathbb{E}_{q_\phi(\omega|\mathbf{D})} [\log p_\theta(\mathbf{D}, \omega)] - \mathbb{E}_{q_\phi(\omega|\mathbf{D})} [\log q_\phi(\omega|\mathbf{D})] \\ & = \mathbb{E}_{q_\phi(\omega|\mathbf{D})} [\log p_\theta(\mathbf{D}|\omega)] - D_{\text{KL}}(q_\phi(\omega|\mathbf{D})||p_\theta(\omega)) \\ & = \mathcal{L} \end{aligned}$$

□

ELBO consists of reconstruction term, which tends to find variational parameters that likely generates the observation, and regularization term, which penalizes the model if it deviates too much from the prior distribution. ELBO is a trade-off between them.

---

\*December 18, 2018

### Second proof

We can also prove the theorem by using Jensen's inequality. Let  $X$  be a random variable and for a convex function  $f$ , the following inequality holds

$$f(\mathbb{E}[X]) \leq \mathbb{E}[f(X)]. \quad (1)$$

The inequality is reversed if  $f$  is a concave function

$$f(\mathbb{E}[X]) \geq \mathbb{E}[f(X)]. \quad (2)$$

*Proof.* By applying importance sampling of approximating variational distribution  $q_\phi(\omega|\mathbf{D})$  to the marginal likelihood  $\log p_\theta(\mathbf{D})$ , we have

$$\begin{aligned} \log p_\theta(\mathbf{D}) &= \log \int p_\theta(\mathbf{D}|\omega) p_\theta(\omega) \frac{q_\phi(\omega|\mathbf{D})}{q_\phi(\omega|\mathbf{D})} d\omega \\ &= \log \mathbb{E}_{q_\phi(\omega|\mathbf{D})} \left[ \frac{p_\theta(\mathbf{D}|\omega) p_\theta(\omega)}{q_\phi(\omega|\mathbf{D})} \right] \\ &\geq \mathbb{E}_{q_\phi(\omega|\mathbf{D})} \left[ \log \frac{p_\theta(\mathbf{D}|\omega) p_\theta(\omega)}{q_\phi(\omega|\mathbf{D})} \right] \\ &= \mathbb{E}_{q_\phi(\omega|\mathbf{D})} [\log p_\theta(\mathbf{D}|\omega)] - D_{\text{KL}}(q_\phi(\omega|\mathbf{D}) \| p_\theta(\omega)) \\ &= \mathcal{L}. \end{aligned}$$

□

### Third proof

We can also have a direct proof by expanding the definition of KL divergence.

*Proof.*

$$\begin{aligned} &\min_{\phi} D_{\text{KL}}(q_\phi(\omega|\mathbf{D}) \| p_\theta(\omega|\mathbf{D})) \\ \iff &\min_{\phi} \int q_\phi(\omega|\mathbf{D}) \log \frac{q_\phi(\omega|\mathbf{D})}{p_\theta(\mathbf{D}|\omega) p_\theta(\omega)} d\omega \\ \iff &\min_{\phi} \left\{ -\mathbb{E}_{q_\phi(\omega|\mathbf{D})} [\log p_\theta(\mathbf{D}|\omega)] + D_{\text{KL}}(q_\phi(\omega|\mathbf{D}) \| p_\theta(\omega)) \right\} \\ \iff &\max_{\phi} \left\{ \mathbb{E}_{q_\phi(\omega|\mathbf{D})} [\log p_\theta(\mathbf{D}|\omega)] - D_{\text{KL}}(q_\phi(\omega|\mathbf{D}) \| p_\theta(\omega)) \right\}. \end{aligned}$$

□