

# Notes On Variational Inference

---

Ludwig Winkler

December 4, 2017

Variational Bayesian Inference is a family of techniques for approximating intractable integrals arising in Bayesian inference and machine learning. They are typically used in complex statistical models consisting of observed variables as well as unknown parameters and latent variables, with various sorts of relationships among the three types of random variables.

## 1 PROBLEM SETTING

Let's assume that we have a set of observations  $x$  which we know where somehow produced by a process. We can call our set of observations  $\mathbf{x}$  our data set.

But the process that generated the observations  $\mathbf{x}$  also has a set of latent, random variables  $\mathbf{z}$  which we cannot observe. An intuitive example would be a classification task where we have an observation  $\mathbf{x}$  and a label  $\mathbf{z}$ . A standard classifier would construct a decision boundary from which we could obtain the probability  $p(\mathbf{z}|\mathbf{x})$ . An analogy would be a neural network classifier with a final softmax layer which outputs a probability of a class, dependent on the input.

Variational Inference (VI) trains a generative model on the data which constructs a joint probability  $p(\mathbf{x}, \mathbf{z})$  from which we can conveniently infer much more information than just the class-conditional probability. By building the entire generative process  $p(\mathbf{x}, \mathbf{z})$  and not just the class-dependent probability  $p(\mathbf{z}|\mathbf{x})$  we have the entire generative process at our disposal. Take for example studying for an exam: It is usually way better to actually

understand the entire topic intrinsically than just learning by heart a mapping from some question to some answer.

Furthermore we can assume that the true model  $p(\mathbf{x}, \mathbf{z})$  has a set of parameters  $\theta$  which parameterize the joint probability density function (PDF) and which we don't know. It is important to note that it is helpful to assume that joint PDF  $p(\mathbf{x}, \mathbf{z})$  consists of an arbitrary combination of parameterized PDFs. This combination of parameterized PDFs could be as simple as two normal distributions with independent parameters or a complicated combination of complex PDFs.

Variational Inference aims to find suitable parameters  $\theta$  which explain the generation of the observations  $\mathbf{x}$  with the latent  $\mathbf{z}$ . Using the analogy of images and labels, VI tries to find parameters which create great pictures for any label. Since we only receive the images, more generally  $\mathbf{x}$ , we want to model the generative process that produced them from labels, more generally  $\mathbf{z}$ , which we can't observe. Once we modeled the 'forward' process from  $\mathbf{z}$  to  $\mathbf{x}$  we can also go 'backward' and infer  $\mathbf{z}$ , the label, from  $\mathbf{x}$ , the image. This is possible because we are working with a generative model.

VI therefore aims to maximize the evidence of the observations  $\mathbf{x}$  by finding the right parameter  $\theta_{\max}$ . If we were to find the right parameters  $\theta_{\max}$  we could generate valid observations  $\mathbf{x}$  for any  $\mathbf{z}$ .

$$\theta_{\max} = \operatorname{argmax}_{\theta} \log p_{\theta}(\mathbf{x}) \quad (1.1)$$

Another nice thing to know, given some observations  $\mathbf{x}$ , for example images, would be to know the latent factors  $\mathbf{z}$ , for example the labels, that generated them. So we would like to know

$$p_{\theta_{\max}}(\mathbf{z}|\mathbf{x}) = \frac{p_{\theta_{\max}}(\mathbf{x}, \mathbf{z})}{\int p_{\theta_{\max}}(\mathbf{x}, \mathbf{z}) d\mathbf{z}} = \frac{p_{\theta_{\max}}(\mathbf{x}, \mathbf{z})}{p_{\theta_{\max}}(\mathbf{x})} \quad (1.2)$$

The problem is that the denominator  $\int p_{\theta_{\max}}(\mathbf{x}, \mathbf{z}) d\mathbf{z}$  is usually an intractable integral. It might work for simpler models but for models which try to infer something on high-dimensional data, calculating the integral becomes quickly very difficult. Due to the usually intractable denominator it is difficult to correctly estimate  $p_{\theta_{\max}}(\mathbf{z}|\mathbf{x})$ .

Because we condition the distribution  $p_{\theta_{\max}}(\mathbf{z}|\mathbf{x})$  on  $\mathbf{x}$  we have to calculate the evidence  $\int p_{\theta_{\max}}(\mathbf{x}, \mathbf{z}) d\mathbf{z}$ . In a nutshell, the trick of VI is to try to estimate an 'unconditional' distribution  $q_{\Psi}(\mathbf{z})$  which doesn't require an intractable integral like  $\int p_{\theta_{\max}}(\mathbf{x}, \mathbf{z}) d\mathbf{z}$  does.

## 2 ELBO

The Kullback-Leibler-Divergence  $\text{KL}(q(\mathbf{x}) || p(\mathbf{x}))$  measures the distance between two probability distributions  $q(\mathbf{x})$  and  $p(\mathbf{x})$ . In variational inference it is used as a criterion with which we minimize the difference between  $q_\Psi(\mathbf{z})$  and  $p_\theta(\mathbf{z}|\mathbf{x})$ .

$$\min_{\Psi} \text{KL}(q_\Psi(\mathbf{z}) || p_\theta(\mathbf{z}|\mathbf{x})) = \mathbb{E}_{q_\Psi(\mathbf{z})} \left[ \log \frac{q_\Psi(\mathbf{z})}{p_\theta(\mathbf{z}|\mathbf{x})} \right] \quad (2.1)$$

$$= \mathbb{E}_{q_\Psi(\mathbf{z})} [\log q_\Psi(\mathbf{z})] - \mathbb{E}_{q_\Psi(\mathbf{z})} [\log p_\theta(\mathbf{z}|\mathbf{x})] \quad (2.2)$$

$$= \mathbb{E}_{q_\Psi(\mathbf{z})} [\log q_\Psi(\mathbf{z})] - \mathbb{E}_{q_\Psi(\mathbf{z})} \left[ \log \frac{p_\theta(\mathbf{x}, \mathbf{z})}{p_\theta(\mathbf{x})} \right] \quad (2.3)$$

$$= \mathbb{E}_{q_\Psi(\mathbf{z})} [\log q_\Psi(\mathbf{z})] - \mathbb{E}_{q_\Psi(\mathbf{z})} [\log p_\theta(\mathbf{x}, \mathbf{z})] + \mathbb{E}_{q_\Psi(\mathbf{z})} [\log p_\theta(\mathbf{x})] \quad (2.4)$$

$$= \underbrace{\mathbb{E}_{q_\Psi(\mathbf{z})} [\log q_\Psi(\mathbf{z})] - \mathbb{E}_{q_\Psi(\mathbf{z})} [\log p_\theta(\mathbf{x}, \mathbf{z})]}_{-\text{ELBO}[q_\Psi(\mathbf{z})]} + \log p_\theta(\mathbf{x}) \quad (2.5)$$

$$= \mathbb{E}_{q_\Psi(\mathbf{z})} \left[ \log \frac{q_\Psi(\mathbf{z})}{p_\theta(\mathbf{x}, \mathbf{z})} \right] + \log p_\theta(\mathbf{x}) \quad (2.6)$$

$$= \text{KL}(q_\Psi(\mathbf{z}) || p_\theta(\mathbf{x}, \mathbf{z})) + \log p_\theta(\mathbf{x}) \quad (2.7)$$

We can intuitively see from 2.5 and 2.7 that we can reformulate the minimization problem as a minimization of the KL-Divergence between the variational distribution  $q_\Psi(\mathbf{z})$  and  $p_\theta(\mathbf{x}, \mathbf{z})$ . The term  $\log p_\theta(\mathbf{x})$  is a constant with respect to our variational distribution  $q_\Psi(\mathbf{z})$ .

The name 'Evidence Lower Bound' (ELBO) is derived from a property of the Kullback-Leibler-Divergence, namely that  $\text{KL}(q(\mathbf{x}) || p(\mathbf{x})) \geq 0$  for any  $q(\mathbf{x})$  and  $p(\mathbf{x})$ . For variational inference this implicates the following identity:

$$\text{KL}(q_\Psi(\mathbf{z}) || p_\theta(\mathbf{z}|\mathbf{x})) \geq 0 \quad (2.8)$$

$$\text{KL}(q_\Psi(\mathbf{z}) || p_\theta(\mathbf{x}, \mathbf{z})) + \log p_\theta(\mathbf{x}) \geq 0 \quad (2.9)$$

$$-\text{ELBO} + \log p_\theta(\mathbf{x}) \geq 0 \quad (2.10)$$

$$\log p_\theta(\mathbf{x}) \geq \text{ELBO} \quad (2.11)$$

In order to minimize our original objective in 2.1 we have to maximize the ELBO in 2.5. With the property of the KL-Divergence, namely  $\text{KL}(q(\mathbf{x}) || p(\mathbf{x})) \geq 0$ , we can see that the ELBO is bounded from above by the log-probability of  $p_\theta(\mathbf{x})$ . So we can only maximize the ELBO up to the log-probability of  $p_\theta(\mathbf{x})$ .

In order to minimize our original objective in 2.1 we will try to approximate the joint distribution  $p_\theta(\mathbf{x}, \mathbf{z})$  with a simpler distribution  $q_\Psi(\mathbf{z})$ . The KL-Divergence is always  $\geq 0$  so we will be always left with the constant  $\log p_\theta(\mathbf{x})$  which we will not be able to optimize. In a nutshell by approximating the complicated conditional distribution  $p_\theta(\mathbf{z}|\mathbf{x})$  with a simpler  $q_\Psi(\mathbf{z})$  we don't have to deal with the intractable integral but we also obtain a term in our optimization problem which we cannot reduce.

To gain more intuition about the ELBO we can go two steps back and look at the following:

$$\text{ELBO}(q_\Psi(\mathbf{z})) = \mathbb{E}_{q_\Psi(\mathbf{z})} [\log p_\theta(\mathbf{x}, \mathbf{z})] - \mathbb{E}_{q_\Psi(\mathbf{z})} [\log q_\Psi(\mathbf{z})] \quad (2.12)$$

$$= \mathbb{E}_{q_\Psi(\mathbf{z})} \left[ \log \frac{p_\theta(\mathbf{x}, \mathbf{z}) p_\theta(\mathbf{z})}{p_\theta(\mathbf{z})} \right] - \mathbb{E}_{q_\Psi(\mathbf{z})} [\log q_\Psi(\mathbf{z})] \quad (2.13)$$

$$= \mathbb{E}_{q_\Psi(\mathbf{z})} [\log p_\theta(\mathbf{x}|\mathbf{z})] + \mathbb{E}_{q_\Psi(\mathbf{z})} [p_\theta(\mathbf{z})] - \mathbb{E}_{q_\Psi(\mathbf{z})} [\log q_\Psi(\mathbf{z})] \quad (2.14)$$

$$= \mathbb{E}_{q_\Psi(\mathbf{z})} [\log p_\theta(\mathbf{x}|\mathbf{z})] - \text{KL}(q_\Psi(\mathbf{z}) \parallel p_\theta(\mathbf{z})) \quad (2.15)$$

Keeping in mind that we have to maximize the ELBO term in 2.5 can furthermore see in 2.15 that  $q_\Psi(\mathbf{z})$  will be balanced between putting weight to the likelihood as well as the prior. The expected likelihood emphasizes  $q_\Psi(\mathbf{z})$  putting its probability on configurations of  $\mathbf{z}$  that explain the observed data  $\mathbf{x}$ . If too much emphasis is paid to the expected likelihood the KL-Divergence between the prior  $p_\theta(\mathbf{z})$  and  $q_\Psi(\mathbf{z})$  will grow. The ELBO term therefore is encouraged to find a balanced approximation of the prior  $p_\theta(\mathbf{z})$  and the likelihood  $p_\theta(\mathbf{x}|\mathbf{z})$ . This is corroborated by the fact that we try to approximate the entire joint probability  $p_\theta(\mathbf{x}, \mathbf{z})$  with  $q_\Psi(\mathbf{z})$ .

### 3 STOCHASTIC GRADIENT OPTIMIZATION

Similarly to how neural networks are nowadays trained with stochastic gradient descent we can train variational inference algorithms with stochastic gradient optimization [1]. Recall from 2.5 that in order to minimize our original objective we need to maximize the Evidence Lower Bound.

$$\nabla_{\Psi} [-\text{ELBO}(q_{\Psi}(\mathbf{z}))] = \nabla_{\Psi} [\text{KL}(q_{\Psi}(\mathbf{z}) \parallel p_{\theta}(\mathbf{x}, \mathbf{z}))] \quad (3.1)$$

$$= \nabla_{\Psi} \left[ \int_{\mathbf{z}} q_{\Psi}(\mathbf{z}) \log \frac{q_{\Psi}(\mathbf{z})}{p_{\theta}(\mathbf{x}, \mathbf{z})} \right] \quad (3.2)$$

$$= \int_{\mathbf{z}} \nabla_{\Psi} [q_{\Psi}(\mathbf{z})] \log \frac{q_{\Psi}(\mathbf{z})}{p_{\theta}(\mathbf{x}, \mathbf{z})} + \int_{\mathbf{z}} q_{\Psi}(\mathbf{z}) \nabla_{\Psi} \left[ \log \frac{q_{\Psi}(\mathbf{z})}{p_{\theta}(\mathbf{x}, \mathbf{z})} \right] \quad (3.3)$$

$$= \int_{\mathbf{z}} \nabla_{\Psi} [q_{\Psi}(\mathbf{z})] \log \frac{q_{\Psi}(\mathbf{z})}{p_{\theta}(\mathbf{x}, \mathbf{z})} + \int_{\mathbf{z}} q_{\Psi}(\mathbf{z}) \left( \nabla_{\Psi} [\log q_{\Psi}(\mathbf{z})] - \underbrace{\nabla_{\Psi} [p_{\theta}(\mathbf{x}, \mathbf{z})]}_{=0} \right) \quad (3.4)$$

$$= \int_{\mathbf{z}} \nabla_{\Psi} [q_{\Psi}(\mathbf{z})] \log \frac{q_{\Psi}(\mathbf{z})}{p_{\theta}(\mathbf{x}, \mathbf{z})} + \underbrace{\int_{\mathbf{z}} q_{\Psi}(\mathbf{z}) \nabla_{\Psi} [\log q_{\Psi}(\mathbf{z})]}_{=0} \quad (3.5)$$

$$= \int_{\mathbf{z}} q_{\Psi}(\mathbf{z}) \nabla_{\Psi} [\log q_{\Psi}(\mathbf{z})] \log \frac{q_{\Psi}(\mathbf{z})}{p_{\theta}(\mathbf{x}, \mathbf{z})} \quad (3.6)$$

$$\approx \frac{1}{N} \sum_{\mathbf{z}_j} \nabla_{\Psi} [q_{\Psi}(\mathbf{z}_j)] \left( \log \frac{q_{\Psi}(\mathbf{z}_j)}{p_{\theta}(\mathbf{x}, \mathbf{z}_j)} + K \right) \quad (3.7)$$

This gives us a training algorithm in which we don't have to use the entire set, but with which we can train batch by batch just like stochastic gradient descent in neural networks.

The identity  $\nabla \log p(x) = \frac{\nabla p(x)}{p(x)}$  was used for

$$\int_{\mathbf{z}} q_{\Psi}(\mathbf{z}) \nabla_{\Psi} [\log q_{\Psi}(\mathbf{z})] = \int_{\mathbf{z}} \nabla_{\Psi} [q_{\Psi}(\mathbf{z})] \quad (3.8)$$

$$= \nabla_{\Psi} \left[ \int_{\mathbf{z}} q_{\Psi}(\mathbf{z}) \right] \quad (3.9)$$

$$= \nabla_{\Psi} [1] \quad (3.10)$$

$$= 0 \quad (3.11)$$

## 4 EXAMPLE

This is an example which is taken from the tutorials on the probabilistic programming framework 'Pyro' by Uber AI Labs. In this example we want to infer the probability of a coin to land on its head (1) or tail (0). For that we have a data set of observations which consists of 60% heads and 40% tails.

It is straight forward to see from the distributin of observed heads and tails that the coin is biased towards heads but for more complex problems and observations the inference might not be as obvious.

In order to make variational inference work properly we need to ingredients: data and a model that structurally captures the generative process of the data. Once the data and the model are defined, variational inference will optimize the generative model to find the parameters of the model which best explain the data. Once the approximate solution for the latent parameters is found we can calculate the mean and variance of our latent factors.

The model for the coin problem is shown in Figure 4.1. The values of  $\alpha$  and  $\beta$  are parameters for a Beta distribution which has the convenient property of having support of  $x \in [0, 1]$ . This means that when we sample from a Beta distribution we will always obtain a value between zero and one, which we can use as a probability. The Beta distribution is often used for priors since it encodes the belief how likely certain values for the latent random variables are.

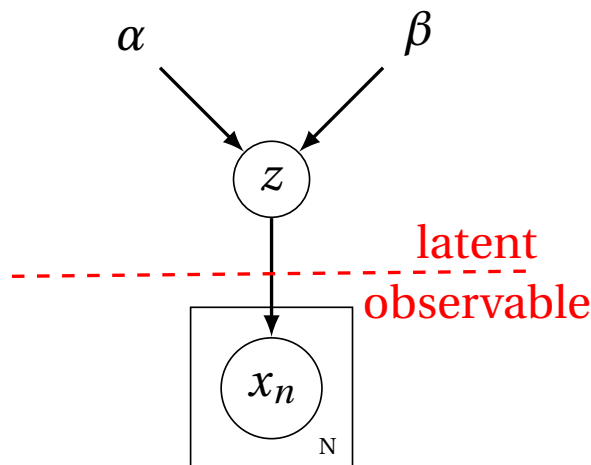


Figure 4.1: The probability  $z$  of the coin to give heads is given by the Beta distribution  $z \sim p_{Beta}(\cdot | \alpha, \beta)$ . Once the  $z$  is drawn, a number of  $N$  samples  $x_n$  are drawn according to a Bernoulli distribution,  $x_n \sim p_{Bernoulli}(\cdot | z)$ . Variational inference aims to find the parameters  $\alpha$  and  $\beta$ .

The Beta distribution is in turn used as input for the Bernoulli distribution which tells us how probable each of the two observations ('head' or 'tail') is for the latent random variable  $x$ . If the prior distribution is very concentrated around the value 0.9 for example, the likelihood distribution will sample heads with a probability of 90%. Due to our very concentrated prior, the samples will have a high probability of being heads.

$$p(x, z; \alpha, \beta) = p_{\text{Bernoulli}}(x | z) p_{\text{Beta}}(z; \alpha, \beta) \quad (4.1)$$

Equation 4.1 formulates the generative model into mathematical form. The probability  $p(x, z; \alpha, \beta)$  is the joint probability of the random variable  $x$  which we can observe and the latent random variable  $z$ . The parameters  $\alpha$  and  $\beta$  define the prior probability  $p(z; \alpha, \beta)$ . The prior probability in turn defines the likelihood of the observations.

With the generative process defined, variational inference tries to find a minimum of the Kullback-Leibler divergence between the joint distribution  $p(x, z; \alpha, \beta)$  and a 'variational' distribution  $q_{\Psi}(z)$ .

## REFERENCES

- [1] D. Wingate and T. Weber, “Automated variational inference in probabilistic programming,” *arXiv preprint arXiv:1301.1299*, 2013.