

Solutions to Problem Set 8

Data Compression With Deep Probabilistic Models

Prof. Robert Bamler, University of Tuebingen

Course material available at <https://robamler.github.io/teaching/compress21/>

Problem 8.1: Understanding the ELBO

In the lecture, we introduced the evidence lower bound, or ELBO,

$$\text{ELBO}(\theta, \phi) = \mathbb{E}_{\mathbf{z} \sim Q_\phi(\mathbf{z}|\mathbf{x}=\mathbf{x})} [\log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x})] \quad (1)$$

where $p_\theta(\mathbf{x}, \mathbf{z}) = p_\theta(\mathbf{z}) p_\theta(\mathbf{x}|\mathbf{z})$ is the joint probability *density* of the generative model P_θ , which has model parameters θ , and $q_\phi(\mathbf{z}|\mathbf{x})$ is the probability density of the variational distribution (or “approximate posterior”) Q_ϕ , which has amortized variational parameters ϕ .

We arrived at the ELBO by considering the negative expected net bit rate of a modified variant of bits-back coding that uses Q_ϕ as a stand-in for the true posterior distribution. Based on the fact that the regular bits-back coding algorithm is optimal, we argued that our modification to the algorithm cannot reduce the net bit rate. This led us to the conclusion that the ELBO is indeed a lower bound on the evidence, i.e.,

$$\text{ELBO}(\theta, \phi) \leq \log p_\theta(\mathbf{x}) \quad \forall \theta, \phi. \quad (2)$$

In this problem, you will derive important equivalent formulations of the ELBO that will allow you to interpret what happens when we maximize the ELBO over θ and ϕ . In doing so, you will also prove the important Eq. 2 in a more direct way.

(a) Show by simple regrouping of the terms in Eq. 1 that

$$\text{ELBO}(\theta, \phi) = \mathbb{E}_{\mathbf{z} \sim Q_\phi(\mathbf{z}|\mathbf{x}=\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] - D_{\text{KL}}[Q_\phi(\mathbf{Z}|\mathbf{X}=\mathbf{x}) \parallel P_\theta(\mathbf{Z})]. \quad (3)$$

What would the encoder and decoder networks learn if, instead of maximizing the ELBO, one would optimize only the first or only the second term on the right-hand side of Eq. 3, respectively?

Solution: Eq. 3 follows directly by inserting $p_\theta(\mathbf{x}, \mathbf{z}) = p_\theta(\mathbf{z}) p_\theta(\mathbf{x}|\mathbf{z})$ into Eq. 1.

Optimizing only the first term on the right-hand side of Eq. 3 and ignoring the KL-term would amount to *maximum likelihood estimation*. The maximization over ϕ would learn a variational distribution that is peaked at $\mathbf{z}^* := \arg \max_{\mathbf{z}} p_\theta(\mathbf{x}|\mathbf{z})$ for the given data point \mathbf{x} , and that is as sharp as possible (since any deviation from the peak \mathbf{z}^* would make the expectation value smaller). If the variational family (i.e., the class of variational distributions that can be parameterized by the encoder model) admits this, $Q_\phi(\mathbf{Z}|\mathbf{X}=\mathbf{x})$ would therefore collapse to a δ -peak at \mathbf{z}^* and

thus the expectation over Q_ϕ would simply evaluate to $\log p_\theta(\mathbf{x}|\mathbf{z}^*)$. Therefore, the (concurrently executed) maximization over θ would fit a likelihood model such that $p_\theta(\mathbf{x}|\mathbf{z}^*)$ is maximized for the given data point \mathbf{x} . Thus, ignoring the KL-term on the right-hand side of Eq. 3 would treat \mathbf{z} as *point-estimated* model parameters, i.e., we would ignore their posterior uncertainty. In this sense, we would be treating \mathbf{z} on the same level as the model parameters θ .

Optimizing only the second term on the right-hand side of Eq. 3 would mean minimizing the KL-divergence from the prior to the approximate posterior. The KL-divergence takes its minimum value of zero if the two distributions are equal. Thus, such a training objective would completely ignore the data and just train an encoder network that always predicts the prior distribution. In this sense, the KL-term in Eq. 3 can be regarded as a regularizer. ■

(b) Show again by simple regrouping of the terms in Eq. 1 that

$$\text{ELBO}(\theta, \phi) = \log p_\theta(\mathbf{x}) - D_{\text{KL}}[Q_\phi(\mathbf{Z}|\mathbf{X} = \mathbf{x}) || P_\theta(\mathbf{Z}|\mathbf{X} = \mathbf{x})]. \quad (4)$$

Notice that this confirms Eq. 2 since we know that the Kullback-Leibler divergence D_{KL} is always nonnegative.

What is the name of the first term on the right-hand side of Eq. 4, and why would we want to maximize it?

In many applications of variational inference, the generative model $P(\mathbf{X}, \mathbf{Z})$ is fixed (i.e., there are no free model parameters θ). In these applications, maximizing the ELBO just amounts to minimizing the KL-term on the right-hand side of Eq. 4 over the variational parameters ϕ . What does this minimization achieve, and why is the method called “variational inference”?

Solution: Eq. 4 follows directly by inserting $p_\theta(\mathbf{x}, \mathbf{z}) = p_\theta(\mathbf{x})p_\theta(\mathbf{z}|\mathbf{x})$ into Eq. 1 and noticing that $\log p_\theta(\mathbf{x})$ can be pulled out of the expectation since it does not depend on \mathbf{z} .

The term $\log p_\theta(\mathbf{x})$ is called the *evidence*. In a compression setup, we’d want to maximize the evidence because the negative evidence is the theoretical lower bound on the expected bit rate. Also beyond compression, one often seeks to maximize the evidence because this means finding the model parameters θ that best describe the observed data \mathbf{x} .

If the generative model is fixed (i.e., there are not free model parameters θ) then we only have to maximize the ELBO over the variational parameters ϕ . This is equivalent to minimizing the KL-divergence from the true posterior $P(\mathbf{Z}|\mathbf{X} = \mathbf{x})$ to the approximate posterior $Q_\phi(\mathbf{Z}|\mathbf{X} = \mathbf{x})$. In other words, it amounts to searching among all distributions that can be expressed as $Q_\phi(\mathbf{Z}|\mathbf{X} = \mathbf{x})$ (this set of distributions is called the *variational family*) for the member that is closest (in KL-distance) to the true posterior. Thus, maximizing the ELBO is an approximate form of Bayesian inference, which explains the name “variational inference”. ■

Problem 8.2: Black Box Variational Inference

In this problem, we discuss the actual task of *maximizing* the ELBO in Eq. 1.

The most efficient way to maximize the ELBO is via the so-called coordinate ascent variational inference (CAVI) algorithm [Blei et al., 2017]. Roughly speaking, this algorithm can be derived by solving the equation $\nabla_{\phi_i} \text{ELBO}(\theta, \phi) = 0$ for one coordinate ϕ_i at a time, by writing out the expectation \mathbb{E} on the right-hand side of Eq. 1 as an explicit integral over \mathbf{z} , taking the derivative, and solving the resulting integrals analytically. While this CAVI algorithm is extremely fast (and should therefore be preferred whenever possible!), its application is limited because the resulting integrals admit an analytic solution only for very special models (e.g., so-called conditional conjugate models).

Mainstream adoption of variational inference only occurred after the invention of so-called black box variational inference (BBVI), which generalizes the method to arbitrary model architectures. In this problem, you will prove the validity of two different approaches to BBVI.

- (a) Let's first understand the problem: the ELBO in Eq. 1 is an expectation of a known quantity $\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})$ under a known distribution Q_{ϕ} . This seems very similar to the typical situation in supervised learning, where we usually have to minimize some loss function that is the expectation of some known expression over samples from the training set.

So why couldn't we just use the same techniques that we know from supervised learning and maximize the ELBO with regular stochastic gradient descent¹? In other words, why can't we just do the following:

- draw some sample $\mathbf{z}_{\text{sample}} \sim Q_{\phi}(\mathbf{Z}|\mathbf{X} = \mathbf{x})$;
- evaluate the gradients of $\log p_{\theta}(\mathbf{x}, \mathbf{z}_{\text{sample}}) - \log q_{\phi}(\mathbf{z}_{\text{sample}}|\mathbf{x})$ with respect to θ and ϕ ;
- use these gradients as an estimate of the gradient of $\text{ELBO}(\theta, \phi)$, and update θ and ϕ with a small gradient step.

Hint: Focus on the optimization over ϕ and look at all places where ϕ appears in the ELBO.

Solution: The gradient step for ϕ in stochastic gradient descent has to be constructed from an unbiased gradient estimate \hat{g} , i.e., an estimate that satisfies $\mathbb{E}_{Q_{\phi}(\mathbf{z}|\mathbf{x}=\mathbf{x})}[\hat{g}] = \nabla_{\phi} \text{ELBO}(\theta, \phi)$. However, the above estimate does not satisfy this requirement because it only takes the gradient of the term inside the expectation in Eq. 1. This neglects the fact that the distribution $Q_{\phi}(\mathbf{Z}|\mathbf{X} = \mathbf{x})$ over which the expectation is taken depends on ϕ itself. This dependency also contributes to the gradient.

¹More precisely, stochastic gradient *ascent* since we want to *maximize*, but that's not the issue here.

More formally, we can write out the ELBO as follows:

$$\begin{aligned} \text{ELBO}(\theta, \phi) &= \mathbb{E}_{\mathbf{z} \sim Q_\phi(\mathbf{z}|\mathbf{x}=\mathbf{x})} [\log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x})] \\ &= \int q_\phi(\mathbf{z}|\mathbf{x}) (\log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x})) d\mathbf{z}. \end{aligned}$$

When we then take the gradient $\nabla_\phi \text{ELBO}(\theta, \phi)$, we mustn't forget the contribution from the first factor in the integral, $q_\phi(\mathbf{z}|\mathbf{x})$.

This complication doesn't arise when stochastic gradient descent is used in the usual supervised learning setup, because the expectation there is over a *fixed* training set, which does not depend on the model parameters over which one optimizes. ■

In the following parts, we will consider two possible solutions to the problem from part (a). We will limit the discussion to the differentiation with respect to ϕ , since differentiation with respect to θ does not pose a problem.

- (b) The simplest form of BBVI uses so-called reparameterization gradients [Kingma and Welling, 2014]. Assume, for example, that the variational distribution is a normal distribution,

$$q_\phi(\mathbf{z}|\mathbf{x}) = \prod_{i=1}^K \mathcal{N}(z_i; \mu_i(\mathbf{x}, \phi), \sigma_i(\mathbf{x}, \phi)^2) \quad (5)$$

where the means $\mu_i(\mathbf{x}, \phi)$ and standard deviations $\sigma_i(\mathbf{x}, \phi)$ together comprise the output $g_\phi(\mathbf{x})$ of the encoder network.

Convince yourself that, for such a variational distribution, the expectation of any function $t(\mathbf{z})$ can be written as follows,

$$\mathbb{E}_{\mathbf{z} \sim Q_\phi(\mathbf{z}|\mathbf{x}=\mathbf{x})} [t(\mathbf{z})] = \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(0, I)} [t(\boldsymbol{\mu}(\mathbf{x}, \phi) + \boldsymbol{\sigma}(\mathbf{x}, \phi) \odot \boldsymbol{\epsilon})]. \quad (6)$$

Here, $\boldsymbol{\mu} = (\mu_1, \dots, \mu_k)$ and $\boldsymbol{\sigma} = (\sigma_1, \dots, \sigma_k)$ are the concatenations into vectors of the means and standard deviations, respectively, $\mathcal{N}(0, I)$ is a k -dimensional standard Normal distribution (i.e., with zero mean and unit variance), and \odot denotes elementwise multiplication.

Now use Eq. 6 to fix the problem from part (a), i.e., to come up with an *unbiased* estimate of $\nabla_\phi \text{ELBO}(\theta, \phi)$.

Solution: Eq. 6 just expresses the normal distributed random variables z_i with mean μ_i and standard deviation σ_i as a scaled and shifted variant of a standard-normal distributed random variables ϵ_i (you can formally prove the equivalence either by noting that a scaled and shifted Gaussian is still a Gaussian, and then calculating the mean and standard deviation of $\mu_i + \sigma_i \epsilon_i$, or by comparing the probability density functions).

Using Eq. 1, we can express the gradient of the ELBO as an expectation under a distribution that is independent of the variational parameters ϕ :

$$\begin{aligned}\nabla_{\phi} \text{ELBO}(\theta, \phi) &= \nabla_{\phi} \mathbb{E}_{\mathbf{z} \sim Q_{\phi}(\mathbf{z}|\mathbf{X}=\mathbf{x})} [\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})] \\ &= \nabla_{\phi} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, I)} \left[\log p_{\theta}(\mathbf{x}, \boldsymbol{\mu}(\mathbf{x}, \phi) + \boldsymbol{\sigma}(\mathbf{x}, \phi) \odot \epsilon) \right. \\ &\quad \left. - \log q_{\phi}(\boldsymbol{\mu}(\mathbf{x}, \phi) + \boldsymbol{\sigma}(\mathbf{x}, \phi) \odot \epsilon | \mathbf{x}) \right] \\ &= \mathbb{E}_{\epsilon \sim \mathcal{N}(0, I)} \left[\nabla_{\phi} \left(\log p_{\theta}(\mathbf{x}, \boldsymbol{\mu}(\mathbf{x}, \phi) + \boldsymbol{\sigma}(\mathbf{x}, \phi) \odot \epsilon) \right. \right. \\ &\quad \left. \left. - \log q_{\phi}(\boldsymbol{\mu}(\mathbf{x}, \phi) + \boldsymbol{\sigma}(\mathbf{x}, \phi) \odot \epsilon | \mathbf{x}) \right) \right]\end{aligned}$$

where, in the last step, we pulled the gradient inside the expectation, which is now allowed since the distribution $\mathcal{N}(0, I)$ over which we take the expectation no longer depends on ϕ . ■

- (c) While the approach from part (b) can be generalized to *some* variational distributions other than the normal distribution, it does not work on arbitrary variational distributions, in particular not on variational distributions over discrete \mathbf{z} . For such variational distributions, an alternative and more general approach called score function gradient estimates (or the “REINFORCE method”) can be used [Ranganath et al., 2014].

Similar to the approach in part (a), one first draws some sample $\mathbf{z}_{\text{sample}} \sim Q_{\phi}(\mathbf{Z}|\mathbf{X} = \mathbf{x})$. However, in the next step, one does *not* simply evaluate $\nabla_{\phi}(\log p_{\theta}(\mathbf{x}, \mathbf{z}_{\text{sample}}) - \log q_{\phi}(\mathbf{z}_{\text{sample}}|\mathbf{x}))$. Instead, one evaluates

$$\hat{g} := \hat{g}^{(1)} + \hat{g}^{(2)} \tag{7}$$

where

$$\begin{aligned}\hat{g}^{(1)} &:= \left(\nabla_{\phi} \log q_{\phi}(\mathbf{z}_{\text{sample}}|\mathbf{x}) \right) \left(\log p_{\theta}(\mathbf{x}, \mathbf{z}_{\text{sample}}) - \log q_{\phi}(\mathbf{z}_{\text{sample}}|\mathbf{x}) \right); \\ \hat{g}^{(2)} &:= -\nabla_{\phi} \log q_{\phi}(\mathbf{z}_{\text{sample}}|\mathbf{x}).\end{aligned} \tag{8}$$

Show that \hat{g} is an unbiased gradient estimate of the ELBO, i.e., that

$$\mathbb{E}_{\mathbf{z}_{\text{sample}} \sim Q_{\phi}(\mathbf{z}|\mathbf{X}=\mathbf{x})} [\hat{g}] = \nabla_{\phi} \text{ELBO}(\theta, \phi). \tag{9}$$

Thus, \hat{g} can be used to optimize the ELBO with stochastic gradient descent.

Hint: write out the expectation \mathbb{E} in the definition of the ELBO (Eq. 1) as an integral, pull the gradient operation ∇_{ϕ} into the integral, and apply the product rule of differential calculus. Then compare the result to the left-hand side of Eq. 9.

Solution:

$$\begin{aligned}
\nabla_\phi \text{ELBO}(\theta, \phi) &= \mathbb{E}_{\mathbf{z} \sim Q_\phi(\mathbf{z}|\mathbf{x}=\mathbf{x})} [\log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x})] \\
&= \nabla_\phi \left(\int q_\phi(\mathbf{z}|\mathbf{x}) (\log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x})) d\mathbf{z} \right) \\
&= \int \left((\nabla_\phi q_\phi(\mathbf{z}|\mathbf{x})) (\log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x})) \right. \\
&\quad \left. + q_\phi(\mathbf{z}|\mathbf{x}) \nabla_\phi (\log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x})) \right) d\mathbf{z} \\
&= \int \left((\nabla_\phi q_\phi(\mathbf{z}|\mathbf{x})) (\log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x})) \right. \\
&\quad \left. - q_\phi(\mathbf{z}|\mathbf{x}) \nabla_\phi \log q_\phi(\mathbf{z}|\mathbf{x}) \right) d\mathbf{z} \\
&= \int (\nabla_\phi q_\phi(\mathbf{z}|\mathbf{x})) (\log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x})) d\mathbf{z} \\
&\quad + \mathbb{E}_{\mathbf{z} \sim Q_\phi(\mathbf{z}|\mathbf{x}=\mathbf{x})} [\hat{g}^{(2)}]
\end{aligned}$$

Our goal is to express $\nabla_\phi \text{ELBO}(\theta, \phi)$ as an expectation under Q_ϕ so that we can estimate it by sampling from Q_ϕ (since evaluating the integral over \mathbf{z} would be infeasible in practice). The second term on the right-hand side of the above equation is already an expectation under Q_ϕ . To get the first term into this form as well, we have to multiply the integrand with $1 = \frac{q_\phi(\mathbf{z}|\mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})}$. We can then simplify by using the relation $\frac{\nabla_\phi q_\phi(\mathbf{z}|\mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})} = \nabla_\phi \log q_\phi(\mathbf{z}|\mathbf{x})$. Thus,

$$\begin{aligned}
&\int (\nabla_\phi q_\phi(\mathbf{z}|\mathbf{x})) (\log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x})) d\mathbf{z} \\
&= \int q_\phi(\mathbf{z}|\mathbf{x}) \frac{\nabla_\phi q_\phi(\mathbf{z}|\mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})} (\log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x})) d\mathbf{z} \\
&= \int q_\phi(\mathbf{z}|\mathbf{x}) (\nabla_\phi \log q_\phi(\mathbf{z}|\mathbf{x})) (\log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x})) d\mathbf{z} \\
&= \mathbb{E}_{\mathbf{z} \sim Q_\phi(\mathbf{z}|\mathbf{x}=\mathbf{x})} \left[(\nabla_\phi \log q_\phi(\mathbf{z}|\mathbf{x})) (\log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x})) \right] \\
&= \mathbb{E}_{\mathbf{z} \sim Q_\phi(\mathbf{z}|\mathbf{x}=\mathbf{x})} [\hat{g}^{(1)}].
\end{aligned}$$

Combining the last two equations proves Eq. 9. ■

- (d) It turns out that the score-function gradients from Eqs. 7-8 can be simplified: we don't even need $\hat{g}^{(2)}$. Show that

$$\mathbb{E}_{\mathbf{z}_{\text{sample}} \sim Q_\phi(\mathbf{z}|\mathbf{x}=\mathbf{x})} [\hat{g}^{(2)}] = 0. \tag{10}$$

Hint: Write out the expectation again as an integral, apply the chain rule of differentiation and then pull the gradient operation out of the integral and use the fact that the density q_ϕ is normalized.

Solution:

$$\begin{aligned}\mathbb{E}_{\mathbf{z} \sim Q_\phi(\mathbf{z}|\mathbf{x}=\mathbf{x})}[\hat{g}^{(2)}] &= \mathbb{E}_{\mathbf{z} \sim Q_\phi(\mathbf{z}|\mathbf{x}=\mathbf{x})}[-\nabla_\phi \log q_\phi(\mathbf{z}|\mathbf{x})] \\ &= -\int q_\phi(\mathbf{z}|\mathbf{x}) \frac{\nabla_\phi q_\phi(\mathbf{z}|\mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})} d\mathbf{z} \\ &= -\int \nabla_\phi q_\phi(\mathbf{z}|\mathbf{x}) d\mathbf{z} \\ &= -\nabla_\phi \left(\int q_\phi(\mathbf{z}|\mathbf{x}) d\mathbf{z} \right) \\ &= -\nabla_\phi(1) = 0\end{aligned}$$

Note: Such contributions to a gradient estimate whose expectation value is zero may still be useful because they may (if constructed well) reduce the variance of the gradient estimate, which speeds up convergence of stochastic gradient optimization because it allows using larger learning rates. Terms with this property are called “control variates”, and there is still a lot of ongoing research about finding good control variates (e.g., in automatic ways). ■

References

- [Blei et al., 2017] Blei, D. M., Kucukelbir, A., and McAuliffe, J. D. (2017). Variational inference: A review for statisticians. *Journal of the American statistical Association*, 112(518).
- [Kingma and Welling, 2014] Kingma, D. P. and Welling, M. (2014). Auto-encoding variational Bayes. In *International Conference on Learning Representations*.
- [Ranganath et al., 2014] Ranganath, R., Gerrish, S., and Blei, D. (2014). Black box variational inference. In *Artificial intelligence and statistics*, pages 814–822.