

Chapter 14

Variational Inference

Reading

1. Sections 1-2 of “Variational Inference: A Review for Statisticians” by Blei et al. (2017).
2. Sections 1-5 of “Auto-Encoding Variational Bayes” by Kingma and Welling (2013)
3. Chapter 2 of Durk Kingma’s thesis: <https://pure.uva.nl/ws/files/17891313/Thesis.pdf>.
4. Bishop Chapter 11.5-11.6
5. Bishop Chapter 10-10.3
6. Lots of great intuition at <http://ruishu.io/2018/03/14/vae/>

We have been primarily concerned with models for classification and regression as yet in this course. The task there is to match the target (a class identity or a real-valued outcome). We now change tracks to consider generative modeling, these are models that are trained to synthesize new data. Effectively, the task here is not *match* a target datum, but given a training dataset of images/text, create a model that outputs similar images/text at test time. We will first take a look at variational methods and generative modeling using these methods in this chapter and do implicit generative models such as Generative Adversarial Networks in the next chapter.

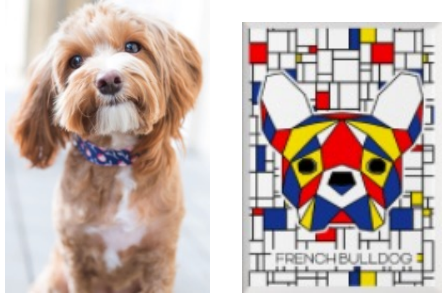
14.1 The model

Imagine how you would draw the image of a dog x on paper. First, you would decide in your mind, its breed, its age, the color of its fur etc. Let us call these quantities “latent factors”. Latent factors can also include things that are not specific to the dog, e.g., the background of your painting (grass, house, beach

etc.), the weather on that day (cloudy, sunny etc.), the viewpoint (zoomed in/far away). We will denote all such quantities by

$z :=$ latent factors.

Having decided upon all these factors, you realize your painting x . The painting x is not unique given latent factors z , e.g., two people can start off with the same latent factors and draw two totally different pictures.



We therefore model the generative process as obtaining samples from a probability distribution

$$p(x|z).$$

Given a latent factor z and an image x , the quantity $p(x|z)$ denotes the likelihood of the sample. Given the painting image x , we do not know what the latent factors are. For instance, it is not easy to say whether the following image is that of a cat or a dog.



In other words, the latent factors of data x are not known to us if we do not take part in the generative process. Nature is in charge of generating the data and our goal here is to guess the parameters of this generative model to be able to synthesize new samples that look as if Nature generated them.

There can be lots of latent factors z . So let us control this complexity and assume that we know a prior over the latent factors

$$\text{prior } p(z)$$

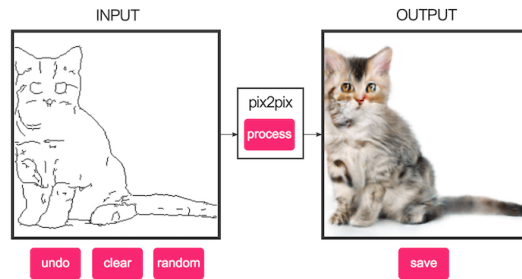
that models our belief of how likely a factor “dog with color blue” is in Nature.

Let us imagine Nature’s generative model as running in two steps

1. First, sample a latent factor z from some distribution, and then
2. sample a datum $x \sim p(x|z)$.

The central point to appreciate is that we know neither Nature's distribution for sampling latents z nor its generative model $p(x|z)$. We will need to fit both these quantities using a training dataset of images/text.

The purpose of doing so can be many-fold, e.g., we may want to generate new data to amplify the size of our training set, given a part of the input image (say due to occlusions, or image corruption) we may want to complete the rest of it.



Most such applications require the knowledge of the latent factors that generated the data. Therefore, formally, we are interested in computing the posterior distribution

$$\text{posterior } p(z|x)$$

given the prior distribution $p(z)$ and samples in a training dataset $D = \{x^i\}_{i=1}^n$. Notice that we do not need labels for this problem, effectively labels $y^i = x^i$ itself because our generative model should of course be very good at generating samples from the training data.

14.2 Some technical basics

14.2.1 Variational calculus

We will first take a brief look at what is called variational calculus.

A function is something takes in a variable as input and returns the value of the function as the output, e.g., $\mathbb{R} \ni f(x) = 5x^2$ for $x \in \mathbb{R}$. Similarly, a

51 *functional* is an object that takes in a *function* as an input and returns a real
52 number as the output. An example of this is entropy

$$\mathbb{R} \ni H[p] = - \int p(x) \log p(x) \, dx$$

53 which takes in a probability density p as the input and returns a real num-
54 ber. Entropy is therefore a *functional*. Just like standard calculus where we
55 take derivatives/minimize over functions, we can also take derivatives of the
56 functional.

57 The functional derivative $\frac{\delta H[p]}{\delta p}(x)$ is defined in a funny way as

$$\int \frac{\delta H[p]}{\delta p}(x) \varphi(x) \, dx = \lim_{\epsilon \rightarrow 0} \frac{H[p + \epsilon \varphi] - H[p]}{\epsilon}$$

58 for any arbitrary function φ . Essentially, you perturb the argument to the
59 functional p by some epsilon and see how much the functional changes. The
60 change in the functional is measured using the test function φ by integrating
61 its changes $\frac{\delta H(p)}{\delta p}(x)$ at each point x in the domain. There may be certain
62 conditions that the perturbation φ needs to satisfy, e.g., since $p + \epsilon \varphi$ should
63 also be probability density, the functional derivative above should only consider
64 test functions φ such that

$$\forall \epsilon \int (p(x) + \epsilon \varphi(x)) \, dx = 1 \Rightarrow \int \varphi(x) \, dx = 0.$$

65 The KL-divergence between two probability densities,

$$\text{KL}(p \parallel q) = \int p(x) \log \frac{p(x)}{q(x)} \, dx,$$

66 is another such functional; it has two arguments p and q .

67 Variational optimization is concerned with minimizing functionals. For
68 instance, if we have a problem

$$w^* = \underset{w \in \mathbb{R}^p}{\operatorname{argmin}} \ell(w)$$

69 in standard optimization, a variational optimization problem with KL-divergence
70 as the loss given a fixed density p looks like

$$q^* = \underset{q \in \mathcal{Q}}{\operatorname{argmin}} \text{KL}(q \parallel p). \quad (14.1)$$

71 The variable of optimization is the probability density q and we will denote the
72 domain of the variable by \mathcal{Q} . Since we want q to be a legitimate probability
73 density, we should choose

$$\mathcal{Q} \subseteq \mathcal{P}(\mathcal{X})$$

74 where $\mathcal{P}(\mathcal{X})$ denotes the set of all probability densities on some domain \mathcal{X} .

75 **Picking the domain and objective in variational optimization** Picking a
76 good domain \mathcal{Q} to minimize over is important. It is similar to the notion of
77 the a hypothesis class in machine learning. If \mathcal{Q} is too big, it is difficult to
78 solve the optimization problem but we obtain a better value to $\text{KL}(q \parallel p)$. If \mathcal{Q}

is too small, the optimization problem may be easy but we may not match the desired distribution p very well. Imagine if p is a mixture of two Gaussians and we pick \mathcal{Q} to be a family of uni-modal Gaussian distributions. Since the KL-divergence is zero if and only if the two distributions are equal, we are never going to be able to minimize it completely. On the other hand, if we pick \mathcal{Q} to be the family of distributions with 2 or more Gaussian modes, then we can perfectly match p . Essentially, the crux of variational inference boils down to picking a good family of distributions \mathcal{Q} that makes solving (14.1) easy.

What functional should we use to measure the distance between q and p ?

The KL-divergence is popular and easy to use in practice but there are many others. For example, when we studied the Gibbs distribution we briefly talked about something called “Wassserstein metric”: if one imagines a mountain of dirt given by distribution q and another mountain of dirt p , the Wassserstein distance $W_2(q, p)$ is the amount of work done in transporting the dirt from q to p ; it is also called the “earth mover’s distance”. The Wassserstein metric is as legitimate a distance between two distributions as the Kullback-Leibler divergence.

14.2.2 Laplace approximation

Laplace approximation is a very useful trick that is similar to variational inference. Here is how it works: suppose we have to estimate approximately an expectation of our random variable $\varphi(w)$

$$\mathbb{E}_{w \sim e^{-n\ell(w)}} [\varphi(w)] = \int e^{-nf(w)} \varphi(w) dw$$

for some large value of n . The above integral takes many values, some have small $\ell(w)$ and some have large $\ell(w)$. The values of w where $\ell(w)$ is small are the ones that have the highest $e^{-n\ell(w)}$, especially as $n \rightarrow \infty$, and therefore the ones that we should pay most attention while approximating the integral. For large n , Laplace approximation replaces the above integral by simply

$$\begin{aligned} \int e^{-n\ell(w)} \varphi(w) dw &\approx \int \varphi(w) e^{-n(\ell(w^*) + \frac{1}{2}(w-w^*)^\top \nabla^2 \ell(w^*)(w-w^*))} dw \\ &= e^{-n\ell(w^*)} \int \varphi(w) e^{-\frac{n}{2}(w-w^*)^\top \nabla^2 \ell(w^*)(w-w^*)} dw \end{aligned} \quad (14.2)$$

where $w^* = \operatorname{argmin} \ell(w)$ is the global minimum of $\ell(w)$. The integral is now with respect to a Gaussian distribution and can be done more easily.

How does a variational approximation differ from the Laplace approximation? Let us look at an example.

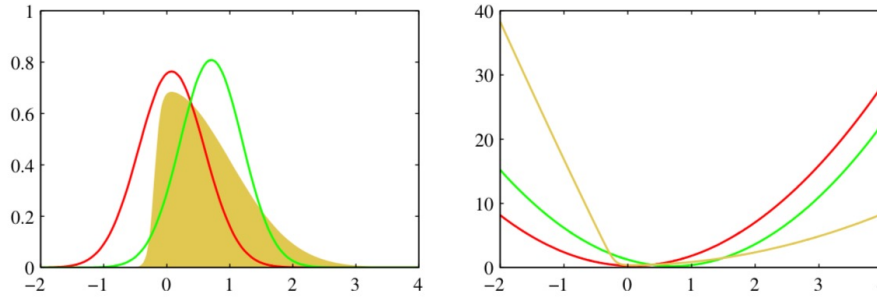


Figure 10.1 Illustration of the variational approximation for the example considered earlier in Figure 4.14. The left-hand plot shows the original distribution (yellow) along with the Laplace (red) and variational (green) approximations, and the right-hand plot shows the negative logarithms of the corresponding curves.

14.2.3 Digging deeper into KL-divergence

Let us take an example to understand KL-divergence better.

Figure 14.1 compares two forms of KL-divergence. The green contours represent equi-probability lines (1,2,3 standard deviations) for a two-dimensional correlated Gaussian $p(z_1, z_2)$. Red contours represent similar equi-probability lines for the variational approximation of this distribution using an uncorrelated Gaussian distribution

$$q(z) = q_1(z_1)q_2(z_2)$$

where both q_1, q_2 are one-dimensional Gaussians. The variational family $q \in \mathcal{Q}$ thus consists of factored uncorrelated Gaussians and we are trying to find the best member of this family that approximates the *correlated* true distribution $p(z)$.

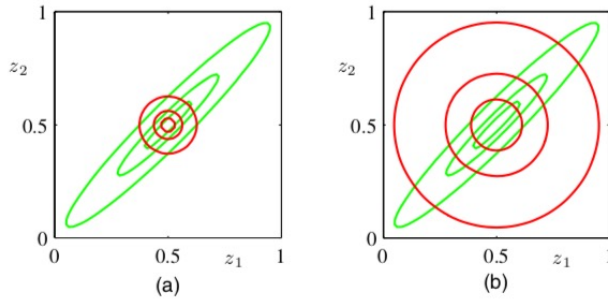


Figure 14.1: Comparison between the variational approximation of a correlated Gaussian using forward and reverse KL divergence and a factored Gaussian family.

Left panel (a) in Figure 14.1 shows the result using the forward KL-divergence minimization

$$q^* = \text{KL}(q \parallel p).$$

while the right panel (b) shows the result for the reverse KL-divergence minimization

$$q^* = \text{KL}(p \parallel q).$$

We see that both these forms capture the mean of the true distribution $p(z)$ correctly. The variance of the two approximations is quite different depending upon which form we employ.

🔗 Use the expression of the KL-divergence to convince yourself why the forward KL under-estimates the variance while the reverse KL over-estimates the variance in Figure 14.1.

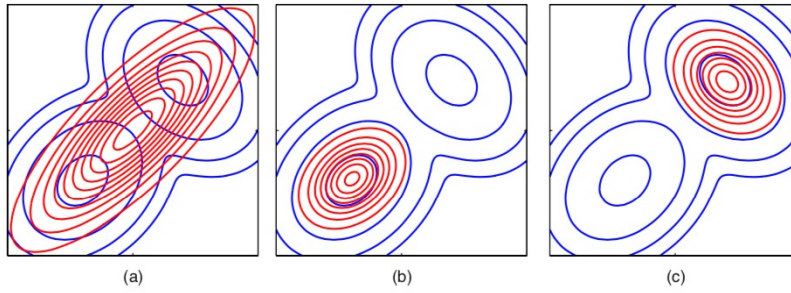


Figure 14.2: Approximating a multi-modal distribution using a uni-modal variational family.

We next consider the case when a multi-modal probability distribution $p(z)$ is approximated using a unimodal Gaussian distribution. Both these examples are very often seen in practice, the distribution of true data/latent factors is often correlated and multi-modal. We have seen one instance of this: the distribution of weights of a deep network in the Gibbs distribution is multi-modal because of multiple global minima.

In panel (a) of Figure 14.2, the reverse KL divergence $\text{KL}(p \parallel q)$ is used to obtain the red contours of q^* . In contrast the forward KL divergence $\text{KL}(p \parallel q)$ is used to obtain q^* in panels (b) and (c). Note that the distribution p is bi-modal and the variational problem is no longer convex in this case; depending upon the initial condition using q , one may get different solutions shown in panels (b) and (c).

KL-divergence is not the only distance used in variational inference and there are many many other ones. You should think of these different ways to measure distances between probability distributions in variational inference as different surrogate losses; which one we use is highly problem dependent although the forward KL-divergence $\text{KL}(q \parallel p)$ is the most common.

14.3 Evidence Lower Bound (ELBO)

We now go back to the generative model.

We will formalize our goal in generative modeling as computing

Nature's true posterior distribution of latent factors

$$p(z|x).$$

We have access to a training dataset $D = \{(x^i)\}_{i=1}^n$. We do not know (i) what form Nature's posterior distribution takes, e.g., Gaussian, multi-modal distribution etc. and (ii) we do not know the true latent factors z that Nature uses. So we are going to approximate the true posterior using some variational family of our choice

$$\mathcal{Q} \ni q^*(z|x) \approx p(z|x).$$

This is the basic idea of variational inference: to approximate a complex distribution $p(z|x)$ using a member of from a simpler family of our choosing \mathcal{Q} . In practice, this variational family \mathcal{Q} will be parameterized by a deep network.

148 With this background, the mathematical process of executing the above
149 program is quite simple. We will simply minimize the KL-divergence

$$q^*(z|x) = \operatorname{argmin}_{q \in \mathcal{Q}} \frac{1}{n} \sum_{i=1}^n \text{KL}(q(z|x^i) || p(z|x^i)). \quad (14.3)$$

150 We next rewrite this KL-divergence above in a special form.

$$\begin{aligned} 0 &\leq \text{KL}(q(z|x^i) || p(z|x^i)) \\ &= \mathbb{E}_{z \sim q(z|x^i)} \left[\log \frac{q(z|x^i)}{p(z|x^i)} \right] \\ &= - \mathbb{E}_{z \sim q(z|x^i)} [\log p(z|x^i)] + \mathbb{E}_{z \sim q(z|x^i)} [\log q(z|x^i)] \\ &= - \mathbb{E}_{z \sim q(z|x^i)} [\log p(z, x^i) - \log p(x^i)] + \mathbb{E}_{z \sim q(z|x^i)} [\log q(z|x^i)] \\ &= \log p(x^i) - \mathbb{E}_{z \sim q(z|x^i)} [\log p(z, x^i)] + \mathbb{E}_{z \sim q(z|x^i)} [\log q(z|x^i)] \\ \Rightarrow \log p(x^i) &\geq \mathbb{E}_{z \sim q(z|x^i)} [\log p(z, x^i)] - \mathbb{E}_{z \sim q(z|x^i)} [\log q(z|x^i)] \end{aligned}$$

151 This is quite interesting. The left-hand side of this inequality is the log-
152 likelihood of the data under Nature's distribution, i.e., it is fixed and inde-
153 pendent of what we do. The left-hand side is also called the evidence. The
154 right hand-side

$$\text{ELBO}(q, x^i) := \mathbb{E}_{z \sim q(z|x^i)} [\log p(z, x^i)] - \mathbb{E}_{z \sim q(z|x^i)} [\log q(z|x^i)]. \quad (14.4)$$

155 is a lower bound on the evidence and therefore called the Evidence Lower
156 Bound (ELBO).

Next comes a key step: a good generative model should be such that

the evidence of the training data, i.e., the log-likelihood of this data under Nature's distribution, should be large under the model. We therefore want to maximize the ELBO on our training data

$$q^*(z|x) = \operatorname{argmax}_{q \in \mathcal{Q}} \frac{1}{n} \sum_{i=1}^n \text{ELBO}(q, x^i). \quad (14.5)$$

to find the posterior distribution of the latent factors $q^*(z)$. Maximizing ELBO is equivalent to minimizing the KL-divergence $\text{KL}(q(z|x^i) || p(z|x^i))$.

We will again solve the optimization problem in (14.5) using stochastic gradient descent. Before we study how to do that, let us consider what model we have developed so far. The solution to this problem

$$q^*(z|x) \approx p(z|x)$$

approximates Nature's posterior distribution. If we maximize ELBO well, given an input x , samples $z \sim q^*(z|x)$ are likely to be the latent factors that Nature could have chosen while rendering this image. But we still do not know how to synthesize an image x for these latent factors. We now rewrite ELBO in a different form to understand this.

$$\begin{aligned} \text{ELBO}(q, x^i) &= \mathbb{E}_{z \sim q(z|x^i)} [\log p(z, x^i)] - \mathbb{E}_{z \sim q(z|x^i)} [\log q(z|x^i)] \\ &= \mathbb{E}_{z \sim q(z|x^i)} [\log p(x^i|z) + \log p(z)] - \mathbb{E}_{z \sim q(z|x^i)} [\log q(z|x^i)] \\ &= \mathbb{E}_{z \sim q(z|x^i)} [\log p(x^i|z)] - \text{KL}(q(z|x^i) || p(z)). \end{aligned}$$

This form of ELBO

$$\text{ELBO}(q, x^i) = \mathbb{E}_{z \sim q(z|x^i)} [\log p(x^i|z)] - \text{KL}(q(z|x^i) || p(z)) \quad (14.6)$$

is very interesting. The first term is Nature's log-likelihood of datum x^i given the latent factor z sampled from *our* candidate posterior q . The second term is the discrepancy between our variational approximation of the posterior $q^*(z|x^i) \approx p(z|x^i)$ and Nature's true marginal distribution over latent factors $p(z)$. This alternative form of ELBO is conceptually very similar to what we do in standard classification, e.g.,

$$\operatorname{argmin}_w \left\{ \ell(w) + \frac{\alpha}{2} \|w\|^2 \right\}.$$

We would like our $q(z|x^i)$ to be close to Nature's prior distribution $p(z)$ but at the same time be such that samples from $q(z|x^i)$ have a high log-likelihood $p(x^i|z)$ of synthesizing images in the training set. The KL-term is therefore a regularizer for the first data-fitting term.

14.3.1 Parameterizing ELBO

What variational family \mathcal{Q} should we choose? Say we parametrized each distribution $q(z|x^i)$ by its mean and diagonal of the covariance.

$$\mathbb{R}^m \ni z \sim q(z|x^i) = N(\mu(x^i), \sigma^2(x^i)I) \in \mathcal{Q}(x^i)$$

179 where $\mu(x^i), \sigma^2(x^i) \in \mathbb{R}^m$. The ELBO in (14.6) is totally independent for
 180 each x^i in the training dataset, so all $i \in \{1, \dots, n\}$ we can solve for

$$\mu^*(x^i), \sigma^2(x^i) = \operatorname{argmax}_{\mu, \sigma^2} \operatorname{ELBO}(N(\mu(x^i), \sigma^2(x^i)I), x^i).$$

181 But this is not a good idea: the parameters μ, σ^2 are distinct for each input x^i
 182 and effectively they are being trained using a dataset of only input image x^i .

Amortized variational inference is a clever trick that ties together the variational families $Q(x^i)$. We will be using a deep network with parameters $u \in \mathbb{R}^p$ that takes x^i as the input and gives $\mu(x^i; u), \sigma^2(x^i; u)$ as the outputs

$$\text{Encoder} : x^i \xrightarrow[\text{parameters } u]{} \mu(x^i; u), \sigma^2(x^i; u).$$

The variational family $Q(x^i)$ that we are considering is therefore the set of distributions expressed by this deep network with p parameters. The family $Q(x^i)$ is still distinct for each datum x^i but they are all tied together by the same weights u .

Encoder. We will call this deep network the encoder because it takes in an input x^i and encodes it into $\mu(x^i; u), \sigma^2(x^i; u)$ which parameterize the distribution of the latent factors.

183 **Decoder.** Observe that although we have now parameterized the distribution
 184 $q(z|x^i)$ using a deep network with weights u , we still do not know how to
 185 model the term $p(x^i|z)$. After all, this is Nature's log-likelihood.

186 We have a dataset $\{(x^i, z^i)\}_{i=1}^n$ that consists of the images x^i and their
 187 corresponding latents z^i sampled from our encoder. We are going to model
 188 Nature's rendering process $p(x|z)$ using a deep network. This is a program that
 189 we have done many times in the past, e.g., we model the targets in classification
 190 y^i as samples from the softmax distribution with images x^i as the input
 191 and train the weights using maximum-likelihood (as you may recall, this is
 192 equivalent to the cross-entropy loss).

193 We can repeat that program here: we are going to learn a deep network

$$\text{Decoder} : p_v(x^i|z) \approx p(x^i|z).$$

194 with parameters $v \in \mathbb{R}^p$ that models Nature's likelihood $p(x^i|z)$.

195 **Different possible decoders for MNIST** Depending upon the type of data
 196 x^i , we will code up the deep network in different ways. For instance, if each
 197 pixel of $x^i \in \mathbb{R}^{28 \times 28}$ is grayscale $[0, 255]$ like it is in MNIST, the output of
 198 the decoder is a multinomial with size $28 \times 28 \times 256$.

199 If we take the training dataset as binarized MNIST (if pixel jk is less
 200 than 128 set it to 0, else set it to 1), then the output of the decoder has size
 201 $28 \times 28 \times 2$ and we can fit this using a logistic distribution at each pixel

$$p_v(x^i|z) = \prod_{j,k=1}^{28} \underbrace{p_v(x_{jk}^i|z)}_{\text{logistic distribution for pixel } x_{jk}^i \in \{0,1\}}$$

▲ The distribution of labels y^i in classification was one-hot vectors, so the softmax layer created a multinomial distribution on the classes.

The log-likelihood term in (14.6) will then correspond to the logistic loss as discussed in the Homework.

Using a mean-field prior $p(z)$. We do not know what the prior distribution $p(z)$ in (14.6) is. We will choose a simple prior

$$p(z) = \prod_{j=1}^m p_j(z_j) \quad (14.7)$$

where $p_i(z_i)$ is the distribution of the i^{th} latent factor z_i . Such distributions are called mean-field priors (where the distribution of a vector $z \in \mathbb{R}^m$ is modeled as independent distributions on its components). We will further choose each distribution

$$p_j(z_j) = N(0, 1)$$

to be a zero-mean standard Gaussian distribution. This is a Gaussian mean-field prior. Just like the choice of a regularizer is critical in machine learning for obtaining good generalization, the choice of a prior is critical in variational inference for synthesizing good images from the generative model.

14.4 Gradient of the ELBO

We now have all the ingredients in place for training a variational generative model. Let us summarize our setup.

1. Encoder parameters u are weights of a deep network that takes in x^i as input and outputs parameters $\mu(x^i), \sigma^2(x^i)$ of the latent distribution. We have tacitly assumed the latent posterior $p(z|x^i)$ to be a Gaussian here; if you have a problem where you wish to have a different latent, e.g., all the latent genes that could have caused a particular cancer, then you want to output the parameters of that distribution from the encoder.
2. The decoder models the likelihood $p_v(x^i|z)$ using parameters v .
3. The prior $p(z)$ will be a mean-field Gaussian distribution. The prior has no parameters in our case, although you may see research papers where the prior also has its own parameters. A popular choice is to use

$$\text{ELBO}_\beta(q, x^i) = \mathbb{E}_{z \sim q(z|x^i)} [\log p(x^i|z)] - \beta^{-1} \text{KL}(q(z|x^i) || p(z))$$

in place of the standard ELBO. The hyper-parameter $\beta > 0$ gives more control over the strength of the prior; this is of course akin to picking the weight-decay coefficient.

⚠ The concept of variational inference and ELBO are much more general than generative models or the encoder-decoder structure that we have developed. Go through the assigned reading material to learn more.

The ELBO when rewritten in terms of the encoder and decoder pa-

rameters looks as follows.

$$\text{ELBO}(u, v; x^i) = \mathbb{E}_{z \sim q_u(z|x^i)} [\log p_v(x^i|z)] - \text{KL}(q_u(z|x^i) || p(z)). \quad (14.8)$$

Our goal is to fit the weights u, v using

$$u^*, v^* = \underset{u, v \in \mathbb{R}^p}{\text{argmax}} \frac{1}{n} \sum_{i=1}^n \text{ELBO}(u, v; x^i). \quad (14.9)$$

The number of parameters of the encoder and decoder can be different but for clarity we imagine them to be the same.

(14.9) is an optimization problem and in this section, we will see how to compute the gradient of the objective so that we can solve the problem using SGD.

14.4.1 The Reparameterization Trick

Focus on the gradient with respect to u of the first term of ELBO

$$\nabla_u \mathbb{E}_{z \sim q(z|x^i)} [\varphi(z)].$$

We have written $\log p_v(x^i|z) = \varphi(z)$ to keep the notation clear; we do not care about the exact form of the integrand in this section.

If we draw a computational graph for the forward propagation of this term, it looks as follows

$$u, x^i \rightarrow \text{sample } z \text{ from } q_u(z|x^i) \rightarrow \varphi(z).$$

The intermediate sampling step is troublesome, we do not really know how to use the chain rule of calculus across sampling, i.e., given

$$\overline{\varphi(z)} := \frac{d}{d_u} \varphi(z)$$

we need to compute $\overline{u} = d\ell/d_u u$. We only know how to apply the chain rule for *deterministic operations* of the form

$$u, x^i \rightarrow z = \text{some deterministic function } g(u, x^i) \rightarrow \varphi(z),$$

in which case we use the standard backprop across the function g .

The Reparameterization Trick enables us to obtain backpropagation gradients across sampling operations via a creative use of the Laplace approximation of the distribution $q_u(z|x^i)$.

We known from the Laplace approximation that we can compute an expectation over z using a Gaussian centered at the global maximum of the

246 distribution $q_u(z|x^i)$ with variance equal to the inverse Hessian at that maxi-
 247 mum. Motivated by this, the Reparameterization Trick *rewrites* the random
 248 variable z as

$$z = \mu(x^i; u) + \sigma(x^i; u) \odot \epsilon$$

249 where

$$\epsilon \sim N(0, I_{m \times m})$$

250 is a sample from a standard multi-variate Gaussian distribution and the notation
 251 \odot denotes element-wise product. Effectively, we imagine that the encoder
 252 outputs

$$\begin{aligned} \mu(x^i; u) &= \underset{z}{\operatorname{argmax}} q_u(z|x^i) \\ \sigma^2(x^i; u) &= \operatorname{diag} \left([\nabla_z^2 q_u(z|x^i)]^{-1} \right). \end{aligned}$$

253 Just like the integral in (14.2) was performed over the Gaussian, the integral
 254 over z can be rewritten as an integral over ϵ

$$\begin{aligned} \nabla_u \mathbb{E}_{z \sim q_u(z|x^i)} [\varphi(z)] &= \nabla_u \mathbb{E}_{\epsilon \sim N(0, I)} [\varphi(\mu(x^i; u) + \sigma(x^i; u) \odot \epsilon)] \\ &= \mathbb{E}_{\epsilon \sim N(0, I)} [\nabla_u \varphi(\mu(x^i; u) + \sigma(x^i; u) \odot \epsilon)] \\ &\approx \frac{1}{N} \sum_{j=1}^N \nabla_u \varphi(\mu(x^i; u) + \sigma(x^i; u) \odot \epsilon^j), \text{ where } \epsilon^j \sim N(0, I). \end{aligned}$$

255 We can take the gradient operator inside the expectation in this case because ϵ
 256 no longer depends on the weights u . The term $\nabla_u \varphi(\mu(x^i; u) + \sigma(x^i; u) \odot \epsilon^j)$
 257 is a deterministic operation given a sample z^j and can be computed using
 258 standard backpropagation.

259 14.4.2 Score-function estimator of the gradient

260 Let us look at an alternative way to compute the same gradient.

$$\begin{aligned} \nabla_u \mathbb{E}_{z \sim q_u(z|x^i)} [\varphi(z)] &= \nabla_u \int \varphi(z) q_u(z|x^i) \mathrm{d}z \\ &= \int \varphi(z) \nabla_u q_u(z|x^i) \mathrm{d}z \\ &= \int \varphi(z) \frac{\nabla_u q_u(z|x^i)}{q_u(z|x^i)} q_u(z|x^i) \mathrm{d}z \\ &= \int \varphi(z) \nabla_u \log q_u(z|x^i) q_u(z|x^i) \mathrm{d}z \\ &= \mathbb{E}_{z \sim q_u(z|x^i)} [\varphi(z) \nabla_u \log q_u(z|x^i)] \\ &\approx \frac{1}{N} \sum_{j=1}^N \varphi(z^j) \nabla_u \log q_u(z^j|x^i), \text{ with } z^j \sim q_u(z|x^i). \end{aligned} \tag{14.10}$$

261 The term

$$\frac{\nabla_u q_u(z|x^i)}{q_u(z|x^i)} = \nabla_u \log q_u(z|x^i) \tag{14.11}$$

262 is called the score function of a probability distribution q_u . The above cal-
 263 culation is quite beautiful: calculating the gradient of the expectation of any

264 quantity $\varphi(z)$ is equal to the expectation of the same quantity weighted by the
 265 score function

$$\nabla_u \mathbb{E}_{z \sim q_u} [\varphi(z)] = \mathbb{E}_{z \sim q_u} [\varphi(z) \nabla_u \log q_u].$$

266 Due to this trick, we can compute the gradient using N samples

$$z^j \sim p_u(z|x^i) \quad (14.12)$$

267 from the encoder; this is easy if, say, the encoder outputs the mean and standard-
 268 deviation of the distribution of the latents. Given z^j , the gradient

$$\nabla_u \log q_u(z^j|x^i)$$

269 is just the standard back-propagation gradient of the quantity $\log q_u(z^j|x^i)$
 270 with respect to weights u of the deep network and can be computed using
 271 autograd.

The key difference between the Reparameterization Trick and the score-function estimator is that in the latter, we do not need to make sure that the gradient $d\ell/dz^j$ can be back-propagated across the sampling operation. The score-function estimator directly computes the gradient of the entire expectation by a weighted average across the samples.

Having two different ways of computing the same gradient may seem redundant but they both are suited to very different applications. The Reparameterization Trick is not accurate in cases when the distribution $q_u(z|x^i)$ is multi-modal because we have only one mean $\mu(x^i)$ around which the samples are drawn. The score-function trick does not have this problem because so long as iid samples are drawn in (14.12) (using any method, e.g., importance sampling) we obtain true estimate of the gradient. The problem in score-function estimator lies in that the denominator $q_u(z|x^i)$ in (14.11) can take very small values if the particular sample z is unlikely. The summation (14.10) is a combination of many N , some very large in magnitude and some very small; the variance of score-function estimate of the gradient in (14.10) can therefore be quite large in most problems.

Typically, the Reparameterization Trick is commonly used in generative models while both the Reparameterization Trick and the score-function estimator are used widely in Reinforcement Learning.

272 14.4.3 Gradient of the remaining terms in ELBO

273 The gradient with respect to weights v of the decoder of the first term in ELBO

$$\nabla_v \mathbb{E}_{z \sim q_u(z|x^i)} [\log p_v(x^i|z)]$$

274 is simply the standard backpropagation gradient (the sampling distribution of
 275 the encoder does not depend on the weights of the decoder).

276 Let us focus on the second term

$$\text{KL} \left(q_u(z|x^i) \parallel \prod_{j=1}^m p_j(z_j) \right). \quad (14.13)$$

where $p_j(z_j) = N(0, 1)$ are terms of the mean-field prior. The gradient of this term with respect to weights of the decoder is zero

$$\nabla_v \text{KL} \left(q_u(z|x^i) \parallel \prod_{j=1}^m p_j(z_j) \right) = 0.$$

Following the reasoning in the Reparameterization Trick, we are positing that $q_u(z|x^i)$ is a Gaussian distribution:

$$q_u(z|x^i) = N(\mu(x^i; u), \sigma^2(x^i; u)I).$$

Notice that $\sigma^2(x^i; u) \in \mathbb{R}^m$ is the diagonal of the covariance and therefore the individual marginals $q_u(z_j|x^i)$ and $q_u(z_{j'}|x^i)$ for two indices j, j' are independent. We can therefore write

$$q_u(z|x^i) = \prod_{j=1}^m N(\mu_j(x^i; u), \sigma_j^2(x^i; u)). \quad (14.14)$$

The KL-divergence of a univariate Gaussian $N(\mu_1, \sigma_1^2)$ with respect to the standard Gaussian is

$$\text{KL}(N(\mu, \sigma^2) \parallel N(0, 1)) = \log \frac{1}{\sigma} + \frac{\sigma^2 + \mu^2}{2} - \frac{1}{2}. \quad (14.15)$$

The general formula is

$$\text{KL}(N(\mu_1, \sigma_1^2) \parallel N(\mu_2, \sigma_2^2)) = \log \frac{\sigma_2}{\sigma_1} + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2}.$$

Due to (14.14), the KL-divergence in (14.13) is a sum of the KL-divergences of the individual univariate Gaussians

$$\text{KL}(q_u(z|x^i) \parallel p(z)) = -\frac{1}{2} \sum_{j=1}^m (\log \sigma_j^2(x^i; u) - \sigma_j^2(x^i; u) + \mu_j^2(x^i; u) + 1). \quad (14.16)$$

The right-hand side of this equation is only a function of u and its gradient can be calculated using standard back-propagation.

This completes our development of ELBO. Using the gradient calculated in this section, we can use SGD to maximize the objective in (14.5) and train a generative model.

14.5 Some comments

Although the mathematics of ELBO seems complicated, it is quite easy to implement generative models using variational inference in practice. You did for a simple MNIST problem in the homework/recitation but if the encoder and decoder are convolutional and deconvolutional architectures respectively, we can get very sophisticated generative models.

🔗 Prove that

$$\begin{aligned} \text{KL} \left(\prod_{j=1}^m q_j(z_j) \parallel \prod_{j=1}^m p_j(z_j) \right) \\ = \sum_{j=1}^m \text{KL}(q_j(z_j) \parallel p_j(z_j)). \end{aligned}$$

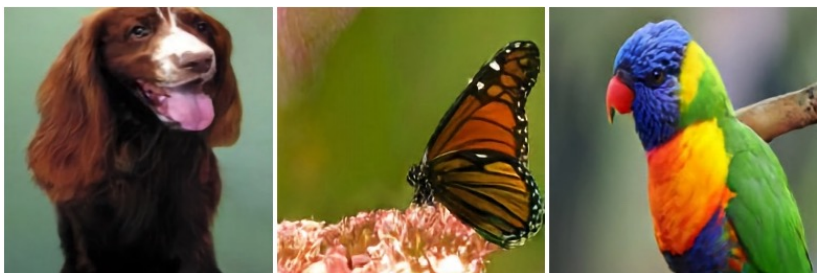


Figure 14.3: Samples from a state-of-the-art VAE trained on ImageNet ([Razavi et al., 2019](#))

300 Variational inference and information-theoretic methods are a rich (and
301 old) area of research and there are many modifications/innovations to ELBO,
302 e.g., read [Alemi et al. \(2018\)](#) for some simple yet deep modifications.

303

Bibliography

- 304 Alemi, A., Poole, B., Fischer, I., Dillon, J., Saurous, R. A., and Murphy, K.
305 (2018). Fixing a broken elbo. In *International Conference on Machine*
306 *Learning*, pages 159–168. PMLR.
- 307 Blei, D. M., Kucukelbir, A., and McAuliffe, J. D. (2017). Variational inference:
308 A review for statisticians. *Journal of the American Statistical Association*,
309 112(518):859–877.
- 310 Kingma, D. P. and Welling, M. (2013). Auto-encoding variational bayes. *arXiv*
311 *preprint arXiv:1312.6114*.
- 312 Razavi, A., van den Oord, A., and Vinyals, O. (2019). Generating diverse
313 high-fidelity images with vq-vae-2. In *Advances in Neural Information*
314 *Processing Systems*, pages 14866–14876.