

# The Reparamaterisation Trick

Greg Feldmann

May 26, 2020

## Abstract

The reparamatisation trick is used to enable variational autoencoders (VAEs) to be trained with gradient descent.

## 1 Quick Overview of Variational Autoencoders

Let  $\mathbf{z}$  be a vector of latent variables,  $\mathbf{x}$  be a row vector from a dataset  $X$ ,  $q_\phi(\mathbf{z}|\mathbf{x})$  be the encoder and  $p_\theta(\mathbf{x}|\mathbf{z})$  be the decoder. The loss function used is referred to as the evidence lower bound (ELBO). ELBO is defined as follows

$$\mathcal{L}_{\theta,\phi}(\mathbf{x}) = \log(p_\theta(\mathbf{x})) - D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}|\mathbf{x})) \quad (1)$$

where  $D_{KL}$  is the Kullback-Leibler divergence.

## 2 Optimising ELBO with SGD

To optimise our VAE, we need to be able to take gradients of the expected value of ELBO with respect to the network weights  $\theta$  and  $\phi$ . This is easy enough for  $\theta$ :

$$\nabla_\theta \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\mathcal{L}_{\theta,\phi}] = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\nabla_\theta \mathcal{L}_{\theta,\phi}] \quad (2)$$

Once the grad operator is inside the expectation, all we have to do is approximate the expectation with a sample.

Things are trickier with  $\phi$ . As the expectation assumes that  $p(\mathbf{z}|\mathbf{x}) = q_\phi(\mathbf{z}|\mathbf{x})$ , we cannot simply move  $\nabla_\phi$  in and out of the expectation arbitrarily.

$$\nabla_\phi \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\mathcal{L}_{\theta,\phi}] \neq \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\nabla_\phi \mathcal{L}_{\theta,\phi}] \quad (3)$$

That the left and right hand sides are not equal can be seen as follows:

$$\nabla_\phi \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\mathcal{L}_{\theta,\phi}] = \nabla_\phi \int q_\phi(\mathbf{z}|\mathbf{x}) \mathcal{L}_{\theta,\phi} d\mathbf{x} \quad (4)$$

$$= \int \nabla_\phi (q_\phi(\mathbf{z}|\mathbf{x}) \mathcal{L}_{\theta,\phi}) d\mathbf{x} \quad (5)$$

$$= \int [q_\phi(\mathbf{z}|\mathbf{x})(\nabla_\phi \mathcal{L}_{\theta,\phi}) + \mathcal{L}_{\theta,\phi}(\nabla_\phi q_\phi(\mathbf{z}|\mathbf{x}))] d\mathbf{x} \quad (6)$$

$$= \int q_\phi(\mathbf{z}|\mathbf{x})(\nabla_\phi \mathcal{L}_{\theta,\phi}) d\mathbf{x} + \int \mathcal{L}_{\theta,\phi}(\nabla_\phi q_\phi(\mathbf{z}|\mathbf{x})) d\mathbf{x} \quad (7)$$

$$= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\nabla_\phi \mathcal{L}_{\theta,\phi}] + \int \mathcal{L}_{\theta,\phi}(\nabla_\phi q_\phi(\mathbf{z}|\mathbf{x})) d\mathbf{x} \quad (8)$$

This isn't in a form that we can easily handle. In particular,  $\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})$ , meaning  $\mathbf{z}$  is stochastic rather than deterministic. This is where the reparameterisation trick comes into play. We can make  $\mathbf{z}$  deterministic by making it a function of a stochastic variable,  $\epsilon$ :