

# Assignment 3. VAE's and Adversarial Attacks

University of Amsterdam – Deep Learning Course

December 2, 2024

**The deadline for this assignment is December 15<sup>th</sup> at 23:59.**

This assignment consists of two parts. The first part will be about Deep Generative Models. Modelling distributions in high dimensional spaces is difficult. Simple distributions such as multivariate Gaussians or mixture models are not powerful enough to model complicated high-dimensional distributions. The question is: How can we design complicated distributions over high-dimensional data, such as images or audio? In concise notation, how can we model a distribution  $p(\mathbf{x}) = p(x_1, x_2, \dots, x_M)$ , where  $M$  is the number of dimensions of the input data  $\mathbf{x}$ ? The solution: Deep Generative Models.

Deep generative models come in many flavors, but all share a common goal: to model the probability distribution of the data. Examples of well-known generative models are Variational Autoencoders (VAEs) [Kingma and Welling, 2014], Generative Adversarial Networks (GANs) [Goodfellow et al., 2014], Adversarial Autoencoder Networks (AAEs) [Makhzani et al., 2015], and Normalizing Flows (NF) [Rezende and Mohamed, 2015]. In this assignment, we will focus on VAEs. The assignment guides you through the theory of VAEs with questions along the way, and finally you will implement one yourself in PyTorch.

In the second part of this assignment, we will dive into adversarial attacks and what they can teach us about neural networks. You will both implement attacks on a ResNet, and defend against attacks.

This assignment contains 50 points: 32 on VAEs, and 18 on Adversarial Attacks.

**Note:** for this assignment you are **not** allowed to use the `torch.distributions` package. You are, however, allowed to use standard, stochastic PyTorch functions like `torch.randn` and `torch.multinomial`, and all other PyTorch functionalities (especially from `torch.nn`). Moreover, try to stay as close as you can to the template files provided as part of the assignment.

# 1 Variational Autoencoders

(Total: 35 points)

## 1.1 Prerequisites

VAEs leverage the flexibility of neural networks (NN) to learn and specify a latent variable model. Before we get into the details of VAEs, we first provide a brief recap of latent variable models and the KL divergence.

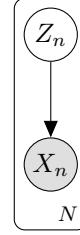
### 1.1.1 Latent Variable Models

A latent variable model is a statistical model that contains both observed and unobserved (i.e. latent) variables. Mathematically, they are connected to a distribution  $p(\mathbf{x})$  over  $\mathbf{x}$  in the following way:  $p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$ . This integral is typically too expensive to evaluate. However, in this assignment, you will learn an approximate solution via VAEs.

Assume a dataset  $\mathcal{D} = \{\mathbf{x}_n\}_{n=1}^N$ , where  $\mathbf{x}_n \in \{0, 1, \dots, k-1\}^M$ . For example,  $\mathbf{x}_n$  could be the pixel values for an image, in which each pixel can take values 0 through  $k-1$  (for example,  $k = 256$ ). A simple latent variable model for this data is shown in Figure 1, which we can also summarize with the following generative story:

$$\mathbf{z}_n \sim \mathcal{N}(0, \mathbf{I}_D), \quad (1)$$

$$\mathbf{x}_n \sim p_X(f_\theta(\mathbf{z}_n)), \quad (2)$$



**Figure 1.** Graphical model of VAE.  $N$  denotes the dataset size. The gray variable is observed.

where  $f_\theta$  is some function – parameterized by  $\theta$  – that maps  $\mathbf{z}_n$  to the parameters of a distribution over  $\mathbf{x}_n$ . For example, if  $p_X$  would be a Gaussian distribution we will use  $f_\theta : \mathbb{R}^D \rightarrow (\mathbb{R}^M, \mathbb{R}_+^M)$  for the mean and diagonal elements of the covariance matrix, or if  $p_X$  is a product of Bernoulli distributions, we have  $f_\theta : \mathbb{R}^D \rightarrow [0, 1]^M$ . Here,  $D$  denotes the dimensionality of the latent space. Likewise, if pixels can take on  $k$  discrete values,  $p_X$  could be a product of Categorical distributions, so that  $f_\theta : \mathbb{R}^D \rightarrow (p_1, \dots, p_k)^M$ , where  $p_1, \dots, p_k$  are event probabilities of the pixel belonging to value  $k$ , where  $p_i \geq 0$  and  $\sum_{i=1}^k p_i = 1$ . Note that our dataset  $\mathcal{D}$  does not contain  $\mathbf{z}_n$ , hence  $\mathbf{z}_n$  is a latent (or unobserved) variable in our statistical model. In the case of a VAE, a (deep) NN is used for  $f_\theta(\cdot)$ .

#### Food for thought

How does the VAE relate to a standard autoencoder (see e.g. [Tutorial 9](#))?

1. Are they different in terms of their main purpose? How so?
2. A VAE is generative. Can the same be said of a standard autoencoder?
3. Can a VAE be used in place of a standard autoencoder for its purpose you mentioned above?

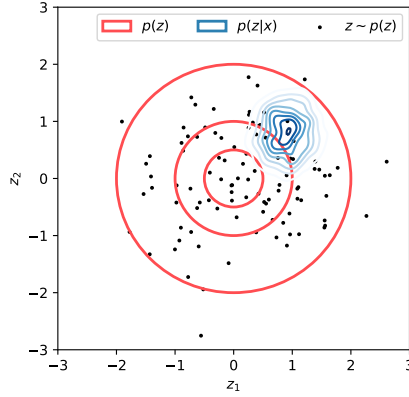
### 1.1.2 KL Divergence

Before covering VAEs, we will need to learn about one more concept that will help us later: the Kullback-Leibler divergence (KL divergence). It measures how different one probability distribution is from another:

$$D_{\text{KL}}(p||q) = \mathbb{E}_{p(x)} \left[ \log \frac{p(X)}{q(X)} \right] = \sum_x p(x) \left[ \log \frac{p(x)}{q(x)} \right], \quad (3)$$

where  $q$  and  $p$  are probability distributions in the space of some random variable  $X$ .<sup>1</sup> In Appendix A, we provide more intuition on how the KL divergence achieves this effect of measuring the difference between two distributions. Note that, while such an intuition can be useful in general, it is not strictly necessary to make this assignment.

<sup>1</sup>The KL divergence of a continuous RV is obtained by replacing the summation with an integral.



**Figure 2.** Plot of 2-dimensional latent space and contours of prior and posterior distributions. The red contour shows the prior  $p(z)$  which is a Gaussian distribution with zero mean and standard deviation of one. The black points represent samples from the prior  $p(z)$ . The blue contour shows the posterior distribution  $p(z|x)$  for an arbitrary  $x$ , which is a complex distribution and here, for example, peaked around  $(1, 1)$ .

## 1.2 Decoder: The Generative Part of the VAE

In Section 1.1.1, we described a general graphical model which also applies to VAEs. In this section, we will define a more specific generative model that we will use throughout this assignment. This will later be referred to as the decoding part (or decoder) of a VAE. For this assignment we will assume the pixels of our images  $\mathbf{x}_n$  in the dataset  $\mathcal{D}$  are Categorical( $p$ ) distributed.

$$p(\mathbf{z}_n) = \mathcal{N}(0, \mathbf{I}_D) \quad (4)$$

$$p(\mathbf{x}_n | \mathbf{z}_n) = \prod_{m=1}^M \text{Cat}(\mathbf{x}_n^{(m)} | f_\theta(\mathbf{z}_n)_m) \quad (5)$$

where  $\mathbf{x}_n^{(m)}$  is the  $m$ -th pixel of the  $n$ -th image in  $\mathcal{D}$ ,  $f_\theta : \mathbb{R}^D \rightarrow (p_1, \dots, p_k)^M$  is a neural network parameterized by  $\theta$  that outputs the probabilities of the Categorical distributions for each pixel in  $\mathbf{x}_n$ . In other words,  $\mathbf{p}_m = (p_{m1}, \dots, p_{mk})$  are event probabilities of the  $m$ -th pixel having intensities  $1, \dots, k$  respectively. Hence, for all  $m$ , we have  $p_{mi} \geq 0$  and  $\sum_{i=1}^k p_{mi} = 1$ .

### Question 1.1 (2 points)

Given a decoder  $f_\theta$ , describe the steps needed to sample an image.

Now that we have defined the model, we can write out an expression for the log probability of the data  $\mathcal{D}$  under this model:

$$\begin{aligned} \log p(\mathcal{D}) &= \sum_{n=1}^N \log p(\mathbf{x}_n) \\ &= \sum_{n=1}^N \log \int p(\mathbf{x}_n | \mathbf{z}_n) p(\mathbf{z}_n) d\mathbf{z}_n \\ &= \sum_{n=1}^N \log \mathbb{E}_{p(\mathbf{z}_n)} [p(\mathbf{x}_n | \mathbf{z}_n)] \end{aligned} \quad (6)$$

This quantity is useful for obtaining the parameters of the decoder  $\theta$  by Maximum Likelihood Estimation. However, evaluating  $\log p(\mathbf{x}_n) = \log \mathbb{E}_{p(\mathbf{z}_n)} [p(\mathbf{x}_n|\mathbf{z}_n)]$  involves a very expensive integral. Equation 6 hints at a method for approximating it, namely Monte-Carlo Integration. The log-likelihood can be approximated by drawing samples  $\mathbf{z}_n^{(l)}$  from  $p(\mathbf{z}_n)$ :

$$\log p(\mathbf{x}_n) = \log \mathbb{E}_{p(\mathbf{z}_n)} [p(\mathbf{x}_n|\mathbf{z}_n)] \quad (7)$$

$$\approx \log \frac{1}{L} \sum_{l=1}^L p(\mathbf{x}_n|\mathbf{z}_n^{(l)}), \quad \mathbf{z}_n^{(l)} \sim p(\mathbf{z}_n), \quad (8)$$

where ‘ $\sim$ ’ means ‘sampled from’. As the number of samples  $L$  tends to infinity, the gap between the approximation and the actual expectation becomes tight. This estimator has the nice property of being unbiased when approximating expectations, though not necessarily log-expectations. Nonetheless, it can be used to approximate  $\log p(\mathbf{x}_n)$  with a sufficiently large number of samples.

#### Question 1.2 (3 points)

Although Monte-Carlo Integration with samples from  $p(\mathbf{z}_n)$  can be used to approximate  $\log p(\mathbf{x}_n)$ , it is not used for training VAE type of models, because it is inefficient. In a few sentences, describe why it is inefficient and how this efficiency scales with the dimensionality of  $\mathbf{z}$ . (Hint: you may use Figure 2 in your explanation.)

### 1.3 The Encoder: $q_\phi(\mathbf{z}_n|\mathbf{x}_n)$ - Efficiently evaluating the integral

In Section 1.2, we have developed the intuition why we need the posterior  $p(\mathbf{z}_n|\mathbf{x}_n)$ . Unfortunately, the true posterior  $p(\mathbf{z}_n|\mathbf{x}_n)$  is as difficult to compute as  $p(\mathbf{x}_n)$  itself. To solve this problem, instead of modeling the true posterior  $p(\mathbf{z}_n|\mathbf{x}_n)$ , we can learn an approximate posterior distribution, which we refer to as the variational distribution. This variational distribution  $q(\mathbf{z}_n|\mathbf{x}_n)$  is used to approximate the (very expensive) posterior  $p(\mathbf{z}_n|\mathbf{x}_n)$ .

Now we have all the tools to derive an efficient bound on the log-likelihood  $\log p(\mathcal{D})$ . We start from Equation 6 where the log-likelihood objective is written, but for simplicity in notation we write the log-likelihood  $\log p(\mathbf{x}_n)$  only for a single datapoint.

$$\begin{aligned} \log p(\mathbf{x}_n) &= \log \mathbb{E}_{p(\mathbf{z}_n)} [p(\mathbf{x}_n|\mathbf{z}_n)] \\ &= \log \mathbb{E}_{p(\mathbf{z}_n)} \left[ \frac{q(\mathbf{z}_n|\mathbf{x}_n)}{q(\mathbf{z}_n|\mathbf{x}_n)} p(\mathbf{x}_n|\mathbf{z}_n) \right] \quad (\text{multiply by } q(\mathbf{z}_n|\mathbf{x}_n)/q(\mathbf{z}_n|\mathbf{x}_n)) \\ &= \log \mathbb{E}_{q(\mathbf{z}_n|\mathbf{x}_n)} \left[ \frac{p(\mathbf{z}_n)}{q(\mathbf{z}_n|\mathbf{x}_n)} p(\mathbf{x}_n|\mathbf{z}_n) \right] \quad (\text{switch expectation distribution}) \\ &\geq \mathbb{E}_{q(\mathbf{z}_n|\mathbf{x}_n)} \log \left[ \frac{p(\mathbf{z}_n)}{q(\mathbf{z}_n|\mathbf{x}_n)} p(\mathbf{x}_n|\mathbf{z}_n) \right] \quad (\text{Jensen's inequality}) \\ &= \mathbb{E}_{q(\mathbf{z}_n|\mathbf{x}_n)} [\log p(\mathbf{x}_n|\mathbf{z}_n)] + \mathbb{E}_{q(\mathbf{z}_n|\mathbf{x}_n)} \log \left[ \frac{p(\mathbf{z}_n)}{q(\mathbf{z}_n|\mathbf{x}_n)} \right] \quad (\text{re-arranging}) \\ &= \underbrace{\mathbb{E}_{q(\mathbf{z}_n|\mathbf{x}_n)} [\log p(\mathbf{x}_n|\mathbf{z}_n)] - KL(q(\mathbf{z}_n|\mathbf{x}_n)||p(\mathbf{z}_n))}_{\text{Evidence Lower Bound (ELBO)}} \quad (\text{writing 2nd term as KL}) \end{aligned} \quad (9)$$

This is awesome! We have derived a bound on  $\log p(\mathbf{x}_n)$ , exactly the thing we want to optimize, where all terms on the right hand side are computable. Let's put together what

we have derived again in a single line:

$$\log p(\mathbf{x}_n) \geq \mathbb{E}_{q(\mathbf{z}_n|\mathbf{x}_n)} [\log p(\mathbf{x}_n|\mathbf{z}_n)] - KL(q(\mathbf{z}_n|\mathbf{x}_n)||p(\mathbf{z}_n)).$$

The right side of the equation is referred to as the *evidence lowerbound* (ELBO) on the log-probability of the data.

This leaves us with the question: **How close is the ELBO to  $\log p(\mathbf{x}_n)$ ?** With an alternate derivation<sup>2</sup>, we can find the answer. It turns out the gap between  $\log p(\mathbf{x}_n)$  and the ELBO is exactly  $KL(q(\mathbf{z}_n|\mathbf{x}_n)||p(\mathbf{z}_n|\mathbf{x}_n))$  such that:

$$\text{ELBO}(\mathbf{q}) = \log p(\mathbf{x}_n) - KL(q(\mathbf{z}_n|\mathbf{x}_n)||p(\mathbf{z}_n|\mathbf{x}_n)) = \mathbb{E}_{q(\mathbf{z}_n|\mathbf{x}_n)} [\log p(\mathbf{x}_n|\mathbf{z}_n)] - KL(q(\mathbf{z}_n|\mathbf{x}_n)||p(\mathbf{z}_n)) \quad (10)$$

Now, let's optimize the ELBO. For this, we define our loss as the mean negative lower bound over samples:

$$\mathcal{L}(\theta, \phi) = -\frac{1}{N} \sum_{n=1}^N \mathbb{E}_{q_\phi(\mathbf{z}_n|\mathbf{x}_n)} [\log p_\theta(\mathbf{x}_n|\mathbf{z}_n)] - D_{KL}(q_\phi(\mathbf{z}_n|\mathbf{x}_n)||p_\theta(\mathbf{z}_n)) \quad (11)$$

Note, that we make an explicit distinction between the generative parameters  $\theta$  of the decoder, and the variational parameters  $\phi$  of the encoder.

#### Question 1.3 (2 points)

Explain how you can see from Equation 10 that the right hand side has to be a *lower bound* on the log-probability  $\log p(\mathbf{x}_n)$ ?

#### Question 1.4 (2 points)

Describe what happens to the ELBO as the variational distribution  $q(\mathbf{z}_n|\mathbf{x}_n)$  approaches the true posterior  $p(\mathbf{z}_n|\mathbf{x}_n)$ ?

#### Food for thought

Consider the third line of Equation 9, before we apply Jensen's inequality. We change the expectation distribution after multiplying and dividing by the variational distribution.

Intuitively, this process corresponds to sampling from a different distribution ( $q(\mathbf{z}_n|\mathbf{x}_n)$  instead of  $p(\mathbf{z}_n)$ ) with the correction weight  $\frac{p(\mathbf{z}_n)}{q(\mathbf{z}_n|\mathbf{x}_n)}$  offsetting the change in distribution we made. Assuming that  $q(\mathbf{z}_n|\mathbf{x}_n)$  is a good approximation to the true  $p(\mathbf{z}_n|\mathbf{x}_n)$ , why would Monte-Carlo integration described in Section 1.2 be more sample efficient after changing the expectation distribution?

(Hint: Consider again Figure 2 and how the samples would be distributed under the variational distribution.)

## 1.4 Defining the optimization objective

The loss in Equation 11:

$$\mathcal{L}(\theta, \phi) = -\frac{1}{N} \sum_{n=1}^N \mathbb{E}_{q_\phi(\mathbf{z}_n|\mathbf{x}_n)} [\log p_\theta(\mathbf{x}_n|\mathbf{z}_n)] - D_{KL}(q_\phi(\mathbf{z}_n|\mathbf{x}_n)||p_\theta(\mathbf{z}_n))$$

<sup>2</sup>This derivation is not done here, but can be found in for instance Bishop sec 9.4.

can be rewritten in terms of per-sample losses:

$$\mathcal{L} = \frac{1}{N} \sum_{n=1}^N (\mathcal{L}_n^{\text{recon}} + \mathcal{L}_n^{\text{reg}}),$$

where

$$\begin{aligned}\mathcal{L}_n^{\text{recon}} &= -\mathbb{E}_{q_\phi(\mathbf{z}_n|\mathbf{x}_n)}[\log p_\theta(\mathbf{x}_n|\mathbf{z}_n)] \\ \mathcal{L}_n^{\text{reg}} &= D_{\text{KL}}(q_\phi(\mathbf{z}_n|\mathbf{x}_n) || p_\theta(\mathbf{z}_n))\end{aligned}$$

can be seen as a reconstruction loss term and a regularization term, respectively.

#### Question 1.5 (2 points)

Explain shortly why the names reconstruction and regularization are appropriate for these two losses.

(Hint: Suppose we use just one sample to approximate the expectation  $\mathbb{E}_{q_\phi(\mathbf{z}_n|\mathbf{x}_n)}[p_\theta(\mathbf{x}_n|\mathbf{z}_n)]$  – as is common practice in VAEs.)

First, we write down the **reconstruction term**:

$$\begin{aligned}\mathcal{L}_n^{\text{recon}} &= -\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x}_n)}[\log p_\theta(\mathbf{x}_n|Z)] \\ &= -\frac{1}{L} \sum_{l=1}^L \log p_\theta(\mathbf{x}_n|\mathbf{z}_n^{(l)}), \quad \mathbf{z}_n^{(l)} \sim q(\mathbf{z}_n|\mathbf{x}_n)\end{aligned}$$

here we used Monte-Carlo integration to approximate the expectation

$$= -\frac{1}{L} \sum_{l=1}^L \sum_{m=1}^M \log \text{Cat}(\mathbf{x}_n^{(m)} | f_\theta(\mathbf{z}_n^{(l)}))$$

Remember that  $f_\theta(\cdot)$  denotes our decoder. Now let  $\mathbf{p}_{nl}^{(m)} = f_\theta(\mathbf{z}_n^{(l)})_m$ , then

$$= -\frac{1}{L} \sum_{l=1}^L \sum_{m=1}^M \sum_{k=1}^K \mathbf{x}_{nk}^{(m)} \log p_{nlk}^{(m)}.$$

where  $\mathbf{x}_{nk}^{(m)} = 1$  if the  $m$ -th pixel has the value  $k$ , and zero otherwise. In other words, the equation above represents the common cross-entropy loss term. When setting  $L = 1$  (i.e. only one sample for  $\mathbf{z}_n$ ), we obtain:

$$= -\sum_{m=1}^M \sum_{k=1}^K \mathbf{x}_{nk}^{(m)} \log p_{nk}^{(m)}$$

where  $\mathbf{p}_n^{(m)} = f_\theta(\mathbf{z}_n)_m$  and  $\mathbf{z}_n \sim q(\mathbf{z}_n|\mathbf{x}_n)$ . Thus, we can use the cross-entropy loss with respect to the original input  $\mathbf{x}_n$  to optimize  $\mathcal{L}_n^{\text{recon}}$

To **compute the regularization term**, we must specify the posterior distribution of the latent variable  $p(\mathbf{z})$  as well the encoder distribution  $q_\phi(\mathbf{z}|\mathbf{x})$ . Similarly to having the flexibility to model the output (i.e. image) distribution  $p(\mathbf{x}|\mathbf{z})$ , we also have some freedom

here. We usually set the prior to be a normal distribution with a zero mean and unit variance:  $p = \mathcal{N}(\mathbf{0}, \mathbf{I}_D)$  and the encoder to be:

$$q_\phi(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}|\mu_\phi(\mathbf{x}), \text{diag}(\sigma_\phi(\mathbf{x}))), \quad (12)$$

where  $\mu_\phi : \mathbb{R}^M \rightarrow \mathbb{R}^D$  maps an input image to the mean of the multivariate normal over  $\mathbf{z}_n$  and  $\sigma_\phi : \mathbb{R}^M \rightarrow \mathbb{R}_+^D$  maps the input image to the diagonal of the covariance matrix of that same distribution. For this case, we can actually find a closed-form solution of the KL divergence:

$$\begin{aligned} \mathcal{L}_n^{\text{reg}} &= D_{\text{KL}}(q_\phi(\mathbf{z}_n|\mathbf{x}_n) || p_\theta(\mathbf{z}_n)) \\ &= D_{\text{KL}}(\mathcal{N}(\mathbf{z}|\mu_\phi(\mathbf{x}), \text{diag}(\sigma_\phi(\mathbf{x}))) || \mathcal{N}(\mathbf{0}, \mathbf{I}_D)) \end{aligned}$$

Using the fact that both probability distributions factorize and that the KL-divergence of two factorizable distributions is a sum of KL terms, we can rewrite this to

$$\begin{aligned} &= \sum_{d=1}^D D_{\text{KL}}(\mathcal{N}(z_{nd}|\mu_\phi(\mathbf{x}_n)_d, \sigma_\phi(\mathbf{x}_n)_d) || \mathcal{N}(z_{nd}|0, 1)) \\ &= \sum_{d=1}^D \frac{\sigma_\phi(\mathbf{x}_n)_d^2 + \mu_\phi(\mathbf{x}_n)_d^2 - 1 - \log \sigma_\phi(\mathbf{x}_n)_d^2}{2} \\ &= \frac{1}{2} \sum_{d=1}^D \sigma_{nd}^2 + \mu_{nd}^2 - 1 - \log \sigma_{nd}^2 \\ &= \frac{1}{2} \sum_{d=1}^D \exp(2 \log \sigma_{nd}) + \mu_{nd}^2 - 1 - 2 \log \sigma_{nd}. \end{aligned}$$

For simplicity, we skipped most of the steps in the derivation, but you can find more details [here](#) if you are interested (it is not essential for understanding the VAE).

Plugging everything together, our final loss is:

$$\mathcal{L}_n(\theta, \phi) = - \underbrace{\sum_{m=1}^M \sum_{k=1}^K \mathbf{x}_{nk}^{(m)} \log p_{nk}^{(m)}}_{\mathcal{L}_n^{\text{recon}}} + \underbrace{\sum_{d=1}^D \frac{1}{2} (\exp(2 \log \sigma_{nd}) + \mu_{nd}^2 - 1 - 2 \log \sigma_{nd})}_{\mathcal{L}_n^{\text{reg}}} \quad (13)$$

$$\mathcal{L}(\theta, \phi) = \frac{1}{N} \sum_{n=1}^N \mathcal{L}_n(\theta, \phi). \quad (14)$$

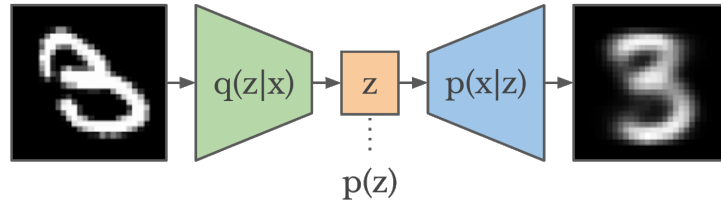
#### Question 1.6 (3 points)

The above derivation of closed form expression for the regularization term requires Gaussian prior and variational distributions. Assume that we want to model the prior  $p(\mathbf{z})$  with a more complex distribution — it is likely the closed form expression would not exist.

Keeping in mind that  $D_{\text{KL}}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z})) = \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} [\log \frac{q(\mathbf{z}|\mathbf{x})}{p(\mathbf{z})}]$ , propose an alternative way of estimating the regularization term for a given sample  $\mathbf{x}_n$ .

## 1.5 The Reparametrization Trick

Although we have written down (the terms of) an objective above, we still cannot simply minimize this by taking gradients with regard to  $\theta$  and  $\phi$ . This is due to the fact that we



**Figure 3.** A VAE architecture on MNIST. The encoder distribution  $q(z|x)$  maps the input image into latent space. This latent space should follow a unit Gaussian prior  $p(z)$ . A sample from  $q(z|x)$  is used as input to the decoder  $p(x|z)$  to reconstruct the image.

sample from  $q_\phi(z_n|x_n)$  to approximate the  $\mathbb{E}_{q_\phi(z|x_n)}[\log p_\theta(x_n|Z)]$  term. Yet, we need to pass the derivative through these samples if we want to compute the gradient of the encoder parameters, i.e.,  $\nabla_\phi \mathcal{L}(\theta, \phi)$ . Our posterior approximation  $q_\phi(z_n|x_n)$  is parameterized by  $\phi$ . If we want to train  $q_\phi(z_n|x_n)$  to maximize the lower bound, and therefore approximate the posterior, we need to have the gradient of the lower-bound with respect to  $\phi$ .

#### Question 1.7 (2 points)

Passing the derivative through samples can be done using the *reparameterization trick* — the process of sampling  $z$  directly from  $\mathcal{N}(\mu(x), \Sigma(x))$  is commonly replaced by calculating  $z = \Sigma(x)\epsilon + \mu(x)$ , where  $\epsilon \sim \mathcal{N}(0, 1)$ . In a few sentences, explain why the act of direct way of sampling usually prevents us from computing  $\nabla_\phi \mathcal{L}$ , and how the reparameterization trick solves this problem.

## 1.6 Putting things together: Building a VAE

Given everything we have discussed so far, we now have an objective (the evidence lower bound or ELBO) and a way to backpropagate to both  $\theta$  and  $\phi$  (i.e., the reparameterization trick). Thus, we can now implement a VAE in PyTorch to train on MNIST images. We will model the encoder  $q(z|x)$  and decoder  $p(x|z)$  by a deep neural network each, and train them to maximize the data likelihood. Figure 3 provides an overview of the components you need to consider in your implementation of the VAE.

In the code directory `part1`, you can find the templates to use for implementing the VAE. For your implementation of the training loop, you will work with **PyTorch Lightning** - a framework that simplifies your code needed to train, evaluate, and test a model in PyTorch. You do not need to be familiar with PyTorch Lightning to the lowest level, but a high-level understanding as from the introduction in **Tutorial 5** is recommended for implementing the template.

You also need to implement additional functions in `utils.py`, and the encoder and decoder in the files `cnn_encoder_decoder.py`. We specified a recommended architecture to start with, but you are allowed to experiment with your own ideas for the models. For the sake of the assignment, it is sufficient to use the recommended architecture to achieve full points. Use the provided unit tests to ensure the correctness of your implementation. Details on the files can be found in the README of part 1.

As a loss objective and test metric, we will use the bits per dimension score (bpd). Bpd is motivated from an information theory perspective and describes how many bits we would need to encode a particular example in our modeled distribution. You can see it as how many bits we would need to store this image on our computer or send it over a network, if we have given our model. The less bits we need, the more likely the example is in our distribution. Hence, we can use bpd as loss metric to minimize. When we test for the bits per dimension on our test dataset, we can judge whether our model generalizes to new samples of the dataset and didn't in fact memorize the training dataset. In order



to calculate the bits per dimension score, we can rely on the negative log-likelihood we got from the ELBO, and change the log base (as bits are binary while NLL is usually exponential):

$$\text{bpd} = \text{nll} \cdot \log_2(e) \cdot \left( \prod_{i=1}^K d_i \right)^{-1}$$

where  $d_1, \dots, d_K$  are the dimensions of the input **excluding any batch dimension**. For images, this would be the height, width and channel number. In other words, for an image of size  $28 \times 28$  with one channel, we have  $d_1 = 28, d_2 = 28, d_3 = 1$  (the order does not matter). We average over those dimensions in order to have a metric that is comparable across different image resolutions. The nll represents the negative log-likelihood loss  $\mathcal{L}$  from Equation 11 for a single data point. You should implement this function in `utils.py`.

#### Question 1.8 (12 points)

Build a Variational Autoencoder in the provided templates, and train it on the MNIST dataset. Both the encoder and decoder should be implemented as a CNN. For the architecture, you can use the same as used in [Tutorial 9](#) about Autoencoders. Note that you have to adjust the output shape of the decoder to output  $1 \times 28 \times 28$  for MNIST. You can do this by adjusting the output padding of the first transposed convolution in the decoder. Use a latent space size of `z_dim=20`. Read the provided README to become familiar with the code template.

In your submission, plot the estimated bit per dimension score of the lower bound on the training and validation set as training progresses, and the final test score. You are allowed to take screenshots of a TensorBoard plot if the axes values are clear.

Note: using the default hyperparameters is sufficient to obtain full points. As a reference, the training loss should start at around 4 bpd, reach below 2.0 after 2 epochs, and end around 0.52 after 80 epochs.

#### Question 1.9 (3 points)

Plot 64 samples ( $8 \times 8$  grid) from your model at three points throughout training (before training, after training 10 epochs, and after training 80 epochs). You should observe an improvement in the quality of samples. Describe shortly the quality and/or issues of the generated images.

#### Question 1.10 (4 points)

Train a VAE with a 2-dimensional latent space (`z_dim=2` in the code). Use this VAE to plot the data manifold as is done in Figure 4b of [\[Kingma and Welling, 2014\]](#) and was discussed in Lecture 6. This is achieved by taking a two dimensional grid of points in  $Z$ -space, and plotting  $f_\theta(Z) = \mu|Z$ . Use the percent point function (ppf, or the inverse CDF) to cover the part of  $Z$ -space that has significant density. Implement it in the function `visualize_manifold` in `utils.py`, and use a grid size of 20. Are you recognizing any patterns of the positions of the different digits?

## 2 Adversarial Attacks

(Total: 18 points)

Deep learning models are used to perform tasks by processing the input. Picture classification. The network outputs the probability that the input belongs to a specific class. With this in mind, a desirable property is that, given two perceptually similar inputs, the output probability also does not change significantly. To put this into more precise terms, assume an input  $\mathbf{x}$ , and a small radius  $\epsilon \in \mathbb{R}_{>0}$ . We expect that for any vector  $\mathbf{r}$  with  $\|\mathbf{r}\| < \epsilon$ , then the new input  $\mathbf{x} + \mathbf{r}$  does not affect the class predicted by the model. This is often called the smoothness bias of neural networks. This, however, does not hold for deep networks. There exist very small changes in the input that do not effectively change its content but that can drastically affect the prediction of the model. This concept was explored in depth by Szegedy [2013], where the concept of **adversarial attack** was introduced (and several very pedagogical experiments were carried out).

**Definition.** An adversarial attack is an imperceptible perturbation performed to an input that causes a network to make a mistake.

Such attacks are not just a quirk of neural networks. They are a dangerous effect of having a very complex system that we do not entirely understand trusted to perform critical tasks. This could be modifying the prompt to an LLM-assisted agent to compel it to give you private information, bypass the fraud checks of a bank, or stop a face-recognition system from identifying your face.

The fact that adversarial attacks are possible also highlights that, given inputs that to us look perfectly reasonable, there could be something in them which throws off the performance of a model, even when we have not been attacked by a malicious agent. Something in the data could have changed, and the model picks up on it, while we do not.

### 2.1 Fast Gradient Sign Method (FGSM)

In this assignment we will be attacking an old ResNet with a very straightforward attack. The FGSM uses the gradients with respect to the input to make a barely noticeable perturbation to the input, while significantly changing the output [Goodfellow et al., 2015]. The formula for the perturbation is given by:

$$\boldsymbol{\eta} = \epsilon \cdot \text{sign}(\nabla_{\mathbf{x}} J(\boldsymbol{\theta}, \mathbf{x}, y))$$

### Question 2.1 (6 points)

In `adversarial_attack.py`, implement the `fgsm_attack` function, and implement the FGSM part in `test_attack`. For a pretrained ResNet18, compute the accuracy on the test set with and without the FGSM attack. You can use the `train.py` for this, report the accuracies and the configuration you used. In addition, answer the following questions:

- (1pt) Why would adding a random perturbation with size  $\epsilon$  not have a similar effect as a FGSM perturbation of the same size?
- (1pt) Say we split the training data into two sets A and B, and train models A and B on those datasets, respectively. Then, we have two models trained for the same task but using different subsets of the same dataset. Now for an instance  $x$  in the test set, a perturbation built using the gradients from model A, will likely have a similar effect on model B, even though the models don't share weights or the exact training data. What is likely to be the cause of this phenomenon?
- (1pt) What is the effect on using data augmentation compared to using no data augmentation, how would you explain this effect?

Note: you can leave hyperparameters as they are, but feel free to experiment to get a better idea what is happening, try the `visualise` flag for example. (What do you think would happen if we increase  $\epsilon$  ?) Or you can try using a pretrained resnet, instead of training from scratch.

It is also possible to defend a model against a certain adversarial attack, this can be done in different ways. In the original papers, the authors proposed an "adversarial objective function". Here the loss takes adversarial examples into account:

$$\tilde{J}(\theta, x, y) = \alpha J(\theta, x, y) + (1 - \alpha) J(\theta, x + \epsilon \text{sign}(\nabla_x J(\theta, x, y))).$$

### Question 2.2 (4 points)

Also train a ResNet18 using this modified loss, this can be done by implementing the `fgsm_loss` function, and calling the necessary configuration in `train.py`. Report the accuracies with and without attacks, of the models with and without the defense.

Hint: You should see improvements for the defense, but expect it to still perform significantly worse on attacked examples than on normal examples. Hint: Try it using pretraining and augmentation, in which case does the defense work significantly?

- (1pt) Describe the tradeoff between defending against the attack or not, why does this tradeoff occur?

## 2.2 Projected Gradient Descent (PGD)

FGSM is a lightweight, simple attack. Using a similar principle, we can achieve even less detectable perturbations. In this part we will be implementing PGD, which iteratively follows the gradient w.r.t. the data, but clips change at each timestep to make sure the perturbed image stays "close" to the original. This is done by iteratively taking a FGSM step with step size  $\alpha$ , and then keeping the perturbed image within the ball of size  $\epsilon$  around the original image.

This time the defense against the attack will have a similar idea, but different implementation. Instead of integrating an adversarial defense into the loss function,

we will add adversarial examples to the batch during training, and train with the original loss.

**Question 2.3** (8 points)

Implement the *pgd\_attack* function in *adversarial\_attack.py*, and the PGD section in *test\_attack*. Moreover, for the defense implement adding adversarial examples to the batch in *train* in *utils.py*.

- (1pt) In this implementation, how do "using an adversarial loss" and "adding adversarial examples to the batch" compare? Are they equivalent or different? Provide reasoning for your answer, including under what conditions they might align or diverge.
- (1pt) Describe a tradeoff between using FGSM and PGD. For each method, identify one advantage it has over the other.

## References

- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In Advances in neural information processing systems, pages 2672–2680, 2014. 1
- Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples, 2015. URL <https://arxiv.org/abs/1412.6572>. 10
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. International Conference on Learning Representations (ICLR), 2014. 1, 9
- Alireza Makhzani, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey. Adversarial autoencoders. arXiv preprint arXiv:1511.05644, 2015. 1
- Danilo Jimenez Rezende and Shakir Mohamed. Variational inference with normalizing flows. In Proceedings of the 32Nd International Conference on International Conference on Machine Learning - Volume 37, ICML’15, pages 1530–1538. JMLR.org, 2015. URL <http://dl.acm.org/citation.cfm?id=3045118.3045281>. 1
- C Szegedy. Intriguing properties of neural networks. arXiv preprint arXiv:1312.6199, 2013. 10

## A Information Theoretic View on the KL Divergence

In order to get a better intuition on how the KL divergence measures how different two probability distributions are, it is useful to treat it from an information theoretic perspective. When we observe an event (e.g., a random variable takes on a specific value), we can think of it as receiving a certain amount of information. Formally, the information content of an event  $x$  is

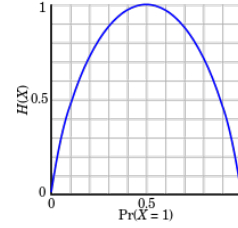
$$h(x) = \log \left( \frac{1}{p(x)} \right) = -\log p(x). \quad (15)$$

Notice that the lower the probability of that event, the more information we will have received after observing that event.

The **entropy** of a discrete random variable  $X$  with distribution  $p$  is the *average* amount of information we effectively receive if we draw a single sample from the distribution  $p$ :

$$H(p) = \mathbb{E}_{p(x)} [-\log p(X)] = -\sum_x p(x) \log p(x). \quad (16)$$

The larger the uncertainty ('spread') of a random variable, the higher its entropy. For instance, a uniform distribution has a high entropy, whereas any delta-function has a low entropy. Similarly, as shown in Figure 4, Bernoulli( $p$ ) has a high entropy for  $p = 0.5$ , and a low entropy for  $p = 0$  and  $p = 1$ .



**Figure 4.** The entropy of a Bernoulli distribution for different values of its parameter  $p$ . Source: [Wikipedia](#).

Now, suppose we model the true distribution  $p$  with another distribution  $q$ , and we use  $q$  to construct a 'coding scheme' <sup>3</sup> to transmit information about samples drawn from  $p$ . Then, the **cross-entropy** is the average amount of information required (e.g. in bits) to specify the value of a sample drawn from  $p$  if we use the coding scheme based on  $q$ :

$$H(p, q) = \mathbb{E}_{p(x)} [-\log q(X)] = -\sum_x p(x) \log q(x). \quad (17)$$

Notice that  $H(p, p) = H(p)$  is a lower bound on  $H(p, q)$ , i.e., the coding scheme based on  $p$  is optimal if we want to specify a value of a sample drawn from  $p$ . Using these definitions, we can derive an expression of the KL divergence  $D_{\text{KL}}(p||q)$  in terms of entropy  $H(p)$  and cross-entropy  $H(p, q)$ :

$$\begin{aligned} D_{\text{KL}}(p||q) &= \sum_x p(x) \left[ \log \frac{p(x)}{q(x)} \right] \\ &= \sum_x p(x) \log p(x) - \sum_x p(x) \log q(x) \\ &= \sum_x p(x) \log p(x) - \sum_x p(x) \log q(x) \\ &= H(p, q) - H(p). \end{aligned}$$

### Food for thought

We have derived that  $D_{\text{KL}}(p||q) = H(p, q) - H(p)$ . What does this expression tell you about the meaning of the KL divergence from an information theoretic point of view? How does this relate to the idea that the KL divergence measures how different one probability distribution is from the other?

<sup>3</sup>Please have a look at [this](#) video for more details.