



# Investigating the Robustness of Disentangled Variational Autoencoders to Adversarial Examples

Sten Sootla<sup>1</sup>

Computational Statistics and Machine Learning

Supervised by Dr. Cristina Calnegru and Dr. John Shawe-Taylor

Submission date: 10.08.2018

<sup>1</sup>**Disclaimer:** This report is submitted as part requirement for the MSc degree in Computational Statistics and Machine Learning at UCL. It is substantially the result of my own work except where explicitly indicated in the text. The report may be freely copied and distributed provided the source is explicitly acknowledged.

## **Abstract**

Summarise your report concisely.

# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Related work and background</b>	<b>5</b>
2.1	Probabilistic learning and inference . . . . .	5
2.1.1	Discriminative approach . . . . .	5
2.1.2	Generative approach . . . . .	6
2.2	Variational autoencoders . . . . .	9
2.2.1	General framework . . . . .	9
2.2.2	Specific example . . . . .	11
2.2.3	Incorporating discrete variables . . . . .	12
2.3	Disentangled representations . . . . .	14
2.4	Adversarial examples . . . . .	17
<b>3</b>	<b>Methods</b>	<b>20</b>
<b>4</b>	<b>Results</b>	<b>21</b>
<b>5</b>	<b>Conclusion</b>	<b>25</b>
<b>6</b>	<b>Appendix</b>	<b>26</b>

# Chapter 1

## Introduction

In recent years, deep neural networks have shown impressive results in various supervised pattern recognition tasks, including image classification [42, 31], object detection [24], and machine translation [2]. Capitalizing on the resurgence of deep neural models in the supervised learning domain, unsupervised methods have also gone through a renaissance, as they have moved from modelling relatively simple toy datasets to being able to capture exceedingly complex real-world distributions [7, 55]. One of the most prominent generative models behind these successes is the variational autoencoder (VAE) [40], which merges conventional probabilistic methods with deep learning, using neural networks to parameterize both the generative and approximate posterior distributions in the the evidence lower bound.

The successes of these models, which mainly consist of stacks of layers, each of which is composed of a linear mapping followed by a nonlinearity, can be attributed to the fact that they can theoretically approximate arbitrarily complex functions [15]. However, as the appropriate functions are discovered automatically via a black-box optimization algorithm (e.g. gradient descent or an evolutionary method), the models can have counter-intuitive properties [58].

First, the resulting decision boundaries of the trained networks are often non-interpretable and unintuitive. Indeed, one of the most well-known and intriguing failure modes of most machine learning algorithms is the existence of adversarial examples - datapoints that are perceptually very close to examples from the original training distribution, but which fool the model with high confidence [58]. The fact that state-of-the-art deep learning models are subject to such adversarial perturbations challenges the narrative that they are unreasonably effective in generalizing outside the distribution they are trained on [4]. More practically, as the models under discussion are being pushed into production for use in semi-autonomous vehicles [9], drones [52], surveillance cameras [18] and malware detection pipelines [29], the existence of adversarial examples poses a security threat with possible real-world consequences.

Second, the internal representations learnt by deep neural networks are infamously complex and hard to analyse, giving them a reputation of being completely black-box learning machines. Indeed, Szegedy et. al. [58] have shown that directions along the coordinate axes are often no more

semantically meaningful than other directions in the latent space, while Morcos et. al. [50] discovered that even if there exist individual units that are interpretable, they are no more important to the task at hand than neurons that don't exhibit clear correlation with specific semantic concepts. In other words, the models learn representations where units are *entangled* - each individual semantically meaningful factor of variation in the input is explained by a combination of units, as opposed to a single neuron.

It has been argued [5] that such entangled representations are suboptimal for generalization, and incorporating inductive biases which would steer the representations towards *disentanglement* would be beneficial. Indeed, it seems intuitively plausible that generalizing to novel configuration of concepts that was not encountered in the training distribution is easiest when relevant factors of variation in the input data are aligned with individual neurons in the latent space. Validating the intuition, such representations have recently shown promising results in zero-shot transfer learning [33] and semi-supervised learning [16].

This thesis investigates the possible relationship between the two intriguing properties of neural networks outlined above: the counter-intuitive decision boundaries on one hand, and the apparent overly complex internal representations on the other. It is reasonable to assume the existence of such a connection, since the leading conjecture explaining the presence of adversarial examples posits that they come about due to excessive linearity of today's pattern recognizers [27]. However, models whose internal representations are disentangled are arguably less linear, since they are implicitly guided to capture the main generative factors of the data, which in datasets with enough complexity are highly non-linear with respect to inputs in the pixel space.

Moreover, one way to construct adversarial examples for encoder-decoder type models is to perturb the source image in such a way that its internal representation would be close to the representation of some target image, while simultaneously meddling with the pixels of the source image as little as possible [41, 59]. The success of this methodology implies that variational autoencoders do not preserve the locality of inputs in their latent space. It has been argued, however, that models whose internal representation is disentangled encode nearby points in pixel space to similar representations in latent space as well [32]. If true, we would expect such models to be significantly more robust to adversarial examples of this type.

To test the preceding hypothesis, we empirically evaluated the robustness of state-of-the-art disentanglement-inducing VAEs to adversarial examples <sup>1</sup>, and compared the results to regular, entangled VAEs. The analyses were carried out for the MNIST handwritten images dataset, as well as for CelebA, a large-scale dataset of celebrity faces. Surprisingly, it was found that disentangled generative models are actually *less* robust to attacks on both datasets, contradicting our initial rationale. This stands in stark contrast to previous work by Alemi et. al. [1], which performed similar analysis for supervised models, but arrived at the opposite conclusion, high-

---

<sup>1</sup>In the context of generative models, an adversarial example is a datapoint which is very similar to the original, untampered image, but whose reconstruction differs markedly from the original's. As an example, an adversarial image could look like a regular depiction of the digit two, but which is reconstructed as the digit four by the generative model.

lighting a fundamental difference between generative and discriminative models when it comes to adversarial examples.

The thesis is organized as follows. In chapter 1, a sufficiently in-depth overview of variational autoencoders, disentanglement, and adversarial examples is given. Chapter 2 outlines the experimental methodology, and provides exact details of the analyses done in this thesis for reproducibility. Chapter 3 gives the results of the experiments.

# Chapter 2

## Related work and background

This chapter introduces the preliminary topics that are integral to understanding the methods and results of this thesis. In addition, it provides an overview of the most important previous works that are closely related to ours. However, the reader is assumed to be comfortable with the basics of discriminative methods in machine learning. The first section of this chapter familiarizes the reader with a probabilistic approach to pattern recognition. The second section builds on the first, introducing the variational autoencoder. The third section focusses on how to regularize generative models such that the learnt representations would be disentangled. The last section introduces the general problem of adversarial examples, and discusses their applicability to generative models.

### 2.1 Probabilistic learning and inference

To fully appreciate variational autoencoders and the problems they are best suited to tackle, it is paramount to first be aware of the context that they operate in. To that end, a brief introduction to the probabilistic view of pattern recognition follows.

#### 2.1.1 Discriminative approach

Usually, in the case of conventional supervised learning, one is concerned about capturing the probability distribution of some predefined labels, conditioned on a set of inputs. For example, in classification, we try to model the categorical distribution of possible classes, given an image (or some set of higher-level features derived from raw pixels) belonging to one of the categories. Relatedly, we could be interested in continuous conditional distributions, as one might be when trying to estimate the future earnings of a company based on its current key performance indicators.

As a concrete example, due to Bishop [8], of how viewing machine learning through a probabilistic lens can yield useful insights into well-known methodologies, let us see how the ubiquitous sum-of-squares error function comes about naturally from a principled probabilistic approach. Suppose we have  $N$  possibly high-dimensional training inputs  $\mathbf{X} = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)}\}$ , and

their corresponding one-dimensional target values  $\mathbf{t} = \{t^{(1)}, t^{(2)}, \dots, t^{(n)}\}$ . Our uncertainty about the targets can be modelled as a Gaussian with a fixed precision  $\beta$ , whose mean is represented by a parameterized function  $y(\mathbf{x}, \boldsymbol{\theta})$  (e.g. a neural network):

$$p(t^{(i)}|\mathbf{x}^{(i)}, \boldsymbol{\theta}, \beta) = \mathcal{N}(t^{(i)}|y(\mathbf{x}^{(i)}, \boldsymbol{\theta}), \beta^{-1}).$$

If we further assume that the given targets were drawn independently from this normal distribution, we can estimate the unknown parameters  $\boldsymbol{\theta}$  and  $\beta$  by maximizing the log-likelihood of the targets:

$$\begin{aligned} \arg \max_{\boldsymbol{\theta}} \log p(\mathbf{t}; \boldsymbol{\theta}) &= \\ \arg \max_{\boldsymbol{\theta}} \log \prod_{n=1}^N \mathcal{N}(t^{(n)}|y(\mathbf{x}^{(n)}, \boldsymbol{\theta}), \beta^{-1}) &= \\ \arg \max_{\boldsymbol{\theta}} \sum_{n=1}^N \log \mathcal{N}(t^{(n)}|y(\mathbf{x}^{(n)}, \boldsymbol{\theta}), \beta^{-1}) &= \\ \arg \max_{\boldsymbol{\theta}} \sum_{n=1}^N \log \left( \frac{\beta}{2\pi} \right)^{\frac{1}{2}} e^{-\frac{\beta}{2}(t^{(n)} - y(\mathbf{x}^{(n)}, \boldsymbol{\theta}))^2} &= \\ \arg \max_{\boldsymbol{\theta}} \sum_{n=1}^N \frac{1}{2} (\log \beta - \log 2\pi) - \frac{\beta}{2} (t^{(n)} - y(\mathbf{x}^{(n)}, \boldsymbol{\theta}))^2 &= \\ \arg \min_{\boldsymbol{\theta}} \sum_{n=1}^N (t^{(n)} - y(\mathbf{x}^{(n)}, \boldsymbol{\theta}))^2, & \end{aligned} \tag{2.1}$$

where the last equation follows from the fact that the first term inside the summation does not depend on the parameter  $\boldsymbol{\theta}$ , and that scaling the parameter-dependent term by a negative constant is equivalent to dropping the constant and turning the maximization problem into a minimization. With that, we have convinced ourselves that even in the usual discriminative setting, we are in fact making implicit use of a probabilistic approach to machine learning.

### 2.1.2 Generative approach

So far we have dealt with modelling the conditional distribution of the targets, given the inputs. Sometimes, however, we might be after the marginal distribution of the input data itself. Having such an object in hand, we are able to determine if an observed example has an unusually low probability of occurrence for anomaly detection [22], fill out missing dimensions of incomplete datapoints [23], or sample new, realistic datapoints from the learnt density function [45].

Capturing the marginal distribution of the inputs is often somewhat more involved than it was in the conditional case discussed previously. This is largely due the fact that high-dimensional, multi-modal datapoints have complex densities whose modelling with well-known distributions would yield suboptimal results. In such situations, additional structure is often incorporated into



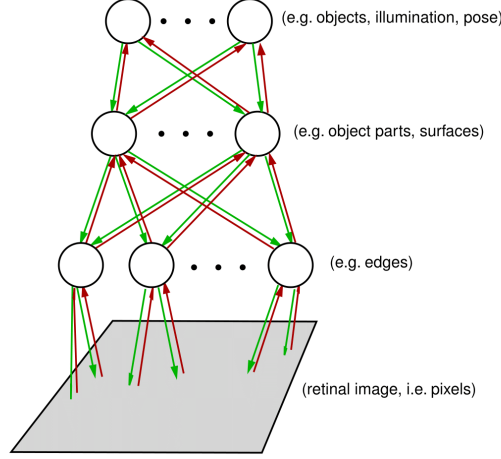


Figure 2.1: Illustration of a probabilistic model with latent variables. Green arrows represent the generative process, in which higher-level concepts guide the development of lower-level ones, eventually leading to the generation of each pixel in an image. Red arrows depict inference, in which case the goal is to infer the underlying causes of an obtained input.

the probabilistic model in the form of lower-dimensional latent variables.

The inclusion of such auxiliary variables is often motivated by the hierarchical nature of physical processes. For example, by thinking about how images could be constructed from a generative point of view, we might imagine that first it is decided which objects to lay onto the scene. Second, one might go into more detail and think about what individual parts do the objects consist of. Finally, the actual orientations of the edges might be decided upon, after which the actual image could be drawn out onto the canvas. Moreover, having imposed such hierarchical dependency structure on our model, we can also traverse the dependency graph bottom-up, reducing our uncertainty about the latents, given an observed datapoint. This inversion of the generative process is called inference, and plays a central role in this thesis. Figure 2.1 illustrates the generative and inferential processes just described.

One immediate drawback to using more complicated latent variable models is that evaluating the likelihood turns out to be intractable, as one has to formally consider the probability of the observed datapoint under every possible setting of the latents. More specifically, consider an observable random variable  $\mathbf{x}$ , an unobserved latent  $\mathbf{z}$ , and their joint distribution  $p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})$ , where the vector  $\boldsymbol{\theta}$  denotes the parameters of the generative model<sup>1</sup>. Then, using the sum rule of probability, we can express the marginal of  $\mathbf{x}$  as  $p_{\boldsymbol{\theta}}(\mathbf{x}) = \int p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z}) d\mathbf{z}$ .

<sup>1</sup>Note that we can focus on the marginal likelihood of individual datapoints, rather than the marginal likelihood of the whole dataset, since under the independence assumption, the former is composed of a sum of the latter (as shown in derivation 2.1):  $\log p_{\boldsymbol{\theta}}(x^{(1)}, x^{(2)}, \dots, x^{(n)}) = \sum_{n=1}^N \log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})$ .

Naturally, as we are unable to evaluate the likelihood function, explicit optimization of its monotonic transform, the log-likelihood, presents difficulties as well. Fortunately, there exists a workaround, which entails optimizing a lower-bound of the log-likelihood instead. By introducing a new distribution  $q_\phi(\mathbf{z}|\mathbf{x})$ , parameterized by  $\phi$ , into our framework, we can lower-bound the original likelihood as follows:

$$\begin{aligned}
& \log p_\theta(\mathbf{x}) \stackrel{\text{sum rule}}{=} \\
& \log \int p_\theta(\mathbf{x}, \mathbf{z}) d\mathbf{z} \stackrel{\text{multiply by identity}}{=} \\
& \log \int p_\theta(\mathbf{x}, \mathbf{z}) \frac{q_\phi(\mathbf{z}|\mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})} d\mathbf{z} \stackrel{\text{Jensen's inequality}}{\geq} \\
& \int q_\phi(\mathbf{z}|\mathbf{x}) \log \frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} d\mathbf{z} =: \mathcal{L}(\mathbf{x}; \theta, \phi).
\end{aligned} \tag{2.2}$$

To see that this lower bound could indeed make the optimization problem easier, it helps to write it in two different forms. The first form illustrates that the lower bound can be written as the expected log-joint under the approximate posterior, plus an additive entropy term that does not depend on the parameters of the generative model:

$$\begin{aligned}
& \mathcal{L}(\mathbf{x}; \theta, \phi) = \\
& \int q_\phi(\mathbf{z}|\mathbf{x}) \log p_\theta(\mathbf{x}, \mathbf{z}) d\mathbf{z} - \int q_\phi(\mathbf{z}|\mathbf{x}) \log(\mathbf{z}|\mathbf{x}) d\mathbf{z} = \\
& \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}, \mathbf{z})] + \mathcal{H}(q_\phi(\mathbf{z}|\mathbf{x})).
\end{aligned} \tag{2.3}$$

The second form reveals that the free energy is nothing more than the original log-likelihood that we are lower-bounding, minus an additional penalty term which increases the gap between the bound and the log-likelihood in proportion to the distance between the approximate posterior and the real posterior, as measured by Kullback-Leibler divergence:

$$\begin{aligned}
& \mathcal{L}(\mathbf{x}; \theta, \phi) = \\
& \int q_\phi(\mathbf{z}|\mathbf{x}) \log \frac{p_\theta(\mathbf{z}|\mathbf{x}) p_\theta(\mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})} d\mathbf{z} = \\
& \int q_\phi(\mathbf{z}|\mathbf{x}) \log p_\theta(\mathbf{x}) + \int q_\phi(\mathbf{z}|\mathbf{x}) \log \frac{p_\theta(\mathbf{z}|\mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})} d\mathbf{z} = \\
& \log p_\theta(\mathbf{x}) - \mathcal{D}_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x}) || p_\theta(\mathbf{z}|\mathbf{x})).
\end{aligned} \tag{2.4}$$

Having laid out the two forms of the lower bound, one should notice that in equation 2.3, only the first term depends on the generative parameters  $\theta$ , while in equation 2.4, the parameters of the approximate posterior feature only in the divergence term. This observation suggests a coordinate-wise optimization scheme where we take turns optimizing  $\theta$  and  $\phi$ , leading us to the well-known Expectation-Maximization (EM) algorithm [23]. This approach is applicable when

we have analytically tractable forms for the expectation of the log-joint, as well as for the true posterior under the current generative parameters  $\theta$ .

What makes EM particularly appealing is that each full step of the algorithm either increases the log-likelihood, or the algorithm has converged to the latter's local maximum. From equation 2.4, we see that free energy is maximized with respect to  $\phi$  when the KL vanishes (as it is always non-negative), which is accomplished when the approximate posterior <sup>2</sup>is set equal to the true posterior. After that, the lower bound is equal to the log-likelihood (under the current best-guess parameters). Since the next step, maximization with respect to  $\theta$  can not decrease the free energy, we observe from equation 2.4 that the log-likelihood can also not decrease, again due to the non-negativity of KL divergences.

Some of the latent variable models that have tractable expected log-joints and posteriors, yielding EM an appropriate tool for learning their parameters, include factor analysis, mixture of Gaussians, hidden markov models, and linear gaussian state-space models. However, often we are interested in using more complicated graphical models, particularly ones where continuous priors and conditionals have non-linear interactions. Assuming that we have the required derivatives of the lower bound in hand, one solution would be to optimize the bound jointly using regular gradient-based methods. Keeping in mind that this approach is usually inferior to EM, as it loses the theoretical guarantees the latter provides, it is often the method of choice for handling models with intractabilities. A concrete example of such an approach is the variational autoencoder, the topic of the next section.

## 2.2 Variational autoencoders

### 2.2.1 General framework

As already alluded to in the previous section, the variational autoencoder is a model that optimizes the lower bound of the log-likelihood (equation 2.2) using a gradient-based optimization scheme. The form of the lower bound that it is explicitly working with is as follows:

---

<sup>2</sup>It is important to distinguish between the best-guess parameters that we are working with when carrying out the optimization scheme, and the maximum-likelihood parameters. Our objective is to find the latter, but we are always operating with an estimated set of parameters. When referring to the *true* posterior, we are not implying that we have access to the true, maximum-likelihood parameters, but the "true" form of the posterior, i.e.  $p(x|z)$  as opposed to  $q(x|z)$ .

$$\begin{aligned}
\mathcal{L}(\mathbf{x}; \boldsymbol{\theta}, \phi) &= \\
&\int q_\phi(\mathbf{z}|\mathbf{x}) \log \frac{p_\theta(\mathbf{x}|\mathbf{z})p_\theta(\mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} d\mathbf{z} = \\
&\int q_\phi(\mathbf{z}|\mathbf{x}) \log p_\theta(\mathbf{x}|\mathbf{z}) d\mathbf{z} + \int q_\phi(\mathbf{z}|\mathbf{x}) \log \frac{p_\theta(\mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} d\mathbf{z} = \\
&\int q_\phi(\mathbf{z}|\mathbf{x}) \log p_\theta(\mathbf{x}|\mathbf{z}) d\mathbf{z} - \mathcal{D}_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z})) = \\
&\mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] - \mathcal{D}_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z})) \simeq \\
&\frac{1}{L} \sum_{l=1}^L \log p_\theta(\mathbf{x}|\mathbf{z}^{(l)}) - \mathcal{D}_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z})), \quad \text{where } \mathbf{z}^{(l)} \sim q_\phi(\mathbf{z}^{(l)}|\mathbf{x}),
\end{aligned} \tag{2.5}$$

where the first term of the last line is a Monte-Carlo estimate of the expectation above. The KL divergence is left as-is, because as we will soon see, for common distributions, it can be calculated analytically. In general, however, this term can be estimated by a finite set of samples as well.

An immediate problem that presents itself is that sampling directly from the approximate posterior is not a continuous operation. This presents a problem, as we need the gradient of the Monte-Carlo estimate of the lower bound with respect to  $\phi$  for optimization. In principle, the necessary gradients could be estimated using the score function estimator, an ubiquitous method in the reinforcement learning literature [62]:

$$\begin{aligned}
&\nabla_\phi \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] = \\
&\int \nabla_\phi q_\phi(\mathbf{z}|\mathbf{x}) \log p_\theta(\mathbf{x}|\mathbf{z}) d\mathbf{z} \stackrel{\text{multiply by identity}}{=} \\
&\int \frac{q_\phi(\mathbf{z}|\mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})} \nabla_\phi q_\phi(\mathbf{z}|\mathbf{x}) \log p_\theta(\mathbf{x}|\mathbf{z}) d\mathbf{z} \stackrel{\text{derivative of log}}{=} \\
&\int q_\phi(\mathbf{z}|\mathbf{x}) \nabla_\phi \log q_\phi(\mathbf{z}|\mathbf{x}) \log p_\theta(\mathbf{x}|\mathbf{z}) d\mathbf{z} = \\
&\mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})} [\nabla_\phi \log q_\phi(\mathbf{z}|\mathbf{x}) \log p_\theta(\mathbf{x}|\mathbf{z})] \simeq \\
&\frac{1}{L} \sum_{l=1}^L \nabla_\phi \log q_\phi(\mathbf{z}|\mathbf{x}) \log p_\theta(\mathbf{x}|\mathbf{z}), \quad \text{where } \mathbf{z}^{(l)} \sim q_\phi(\mathbf{z}^{(l)}|\mathbf{x}).
\end{aligned} \tag{2.6}$$

Unfortunately, due to its generality, the preceding estimator has a very high variance, rendering its use impractical in the context of generative modelling. **Elaborate on the variance of the score function gradient estimator!**

A better approach to estimating the variance, introduced concurrently by Kingma and Welling [40], as well as by Rezende et. al. [54], is to reparameterize the lower bound in equation 2.5 such that the necessary stochasticity that is injected into the model would not depend on any of its parameters. As a simple illustrative example of such an approach, suppose our approximate posterior is taken to be a univariate Gaussian:  $z \sim p(z|x) = \mathcal{N}(\mu, \sigma^2)$ . If we sampled  $z$  naively,

we would not be able to calculate the required derivatives of this random variable, neither  $\frac{\partial z}{\partial \mu}$  nor  $\frac{\partial z}{\partial \sigma}$ . However, by reparameterizing the variable as  $z = \mu + \sigma\epsilon$ , where  $\epsilon \sim \mathcal{N}(0, 1)$ , we decouple the variables of our model from the stochasticity, resulting in differentiability:  $\frac{\partial(\mu+\sigma\epsilon)}{\partial \mu} = 1$  and  $\frac{\partial(\mu+\sigma\epsilon)}{\partial \sigma} = \epsilon$ .

Let us denote the differentiable transformation of the random variable as  $\mathbf{z} = g_\phi(\epsilon, \mathbf{x})$ , where  $\epsilon \sim p(\epsilon)$  is the auxiliary noise variable that does not depend on the model's parameters. The lower bound of a mini-batch of  $M$  samples, randomly drawn from the full dataset, is then given by equation 2.7. In practice, Kingma and Welling [40] found that the number of samples drawn from the approximate posterior can be taken to be 1 (i.e.  $L = 1$ ), provided that the mini-batch is large enough. As the resulting bound is jointly differentiable with respect to  $\theta$  and  $\phi$ , it can be optimized with any of the popular stochastic optimization methods, such as SGD, Adam [39], or RMSProp [60].

$$\begin{aligned} \frac{1}{M} \sum_{i=1}^M \mathcal{L}(\mathbf{x}^{(i)}; \theta, \phi) = \\ \frac{1}{M} \sum_{i=1}^M \left\{ \mathbb{E}_{\epsilon \sim p(\epsilon)} [\log p_\theta(\mathbf{x}^{(i)} | g_\phi(\epsilon, \mathbf{x}^{(i)}))] - \mathcal{D}_{\text{KL}}(q_\phi(\mathbf{z} | \mathbf{x}^{(i)}) || p_\theta(\mathbf{z})) \right\} \simeq \\ \frac{1}{ML} \sum_{i=1}^M \left\{ \sum_{l=1}^L \log p_\theta(\mathbf{x}^{(i)} | \mathbf{z}^{(i,l)}) - \mathcal{D}_{\text{KL}}(q_\phi(\mathbf{z} | \mathbf{x}^{(i)}) || p_\theta(\mathbf{z})) \right\}, \\ \text{where } \mathbf{z}^{(i,l)} \sim g_\phi(\epsilon, \mathbf{x}), \epsilon \sim p(\epsilon). \end{aligned} \quad (2.7)$$

With that, we have laid out all the fundamental components that make up a variational autoencoder. The next subsection covers some of the specifics that are still left to be filled in, particularly the forms and parameterizations of the distributions in play.

### 2.2.2 Specific example

The name "variational autoencoder" is usually reserved for a particular realization of the preceding framework, one where both the approximate posterior  $q_\phi(\mathbf{x} | \mathbf{z})$ , as well as the conditional in the generative part of the model,  $p_\theta(\mathbf{z} | \mathbf{x})$ , are parameterized by artificial neural networks. These distributions themselves are usually taken to be multivariate normals with diagonal covariances. The prior does not generally have learnable parameters, as it is taken to be an isotropic Gaussian:  $p_\theta(\mathbf{x}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$ . Figure 2.2 gives a schematic of this setup.

A nice property that follows from having Gaussian distributions is the analytic tractability of the KL divergence featuring in the lower bound. **Derive the corresponding closed-form equation.**

Looking at the general outline of the model in figure 2.2, it becomes apparent why it is called an autoencoder. An input  $\mathbf{x}$  is fed through a neural network, which outputs the mean and diagonal covariance of the approximate posterior. Making use of the reparameterization trick, a sample is then drawn from this distribution. As the dimension of the latent variable is usually much

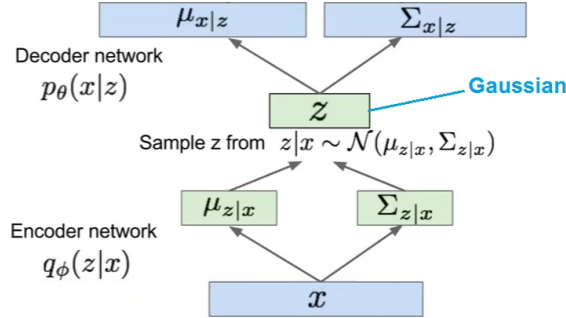


Figure 2.2: A schematic illustration of a variational autoencoder.

lower than that of the input, it is akin to the bottleneck layer of a regular autoencoder, and the whole underlying subnetwork (shown in green) is dubbed an encoder. The drawn vector  $z$  is subsequently fed through a second neural network, which outputs the estimated parameters of  $p_{\theta}(x|z)$ . As the estimated mean is obtained by upsampling the latent, and the conditional is maximized if it is equal to the original input, it bears resemblance to the decoding part of a conventional autoencoder, justifying the use of the name "decoder" for the underlying neural network.

### 2.2.3 Incorporating discrete variables

In the previous subsections, we focussed on variational autoencoders with continuous variables, namely ones with Gaussian distributions. Sometimes, however, it would be appropriate for the posterior to include some discrete variables as well. For example, when modelling handwritten digits, some of dimensions of the posterior could correspond to thickness, angle, size, and digit class. While the first three factor of variation are reasonably modelled by continuous distributions, a discrete one is more natural for the last. Unfortunately, the reparameterization trick that allowed joint optimization of the parameters in the Gaussian case is not applicable for discrete variables.

To get around taking direct samples from a discrete distribution, one might employ the Gumbel-Max trick [47]. Suppose our one-hot encoded random variable of interest,  $Y$ , has a discrete distribution with  $K$  choices, parameterized by unnormalized probabilities  $(\alpha_1, \alpha_2, \dots, \alpha_K)$ , where  $\alpha_i \in (0, \infty)$ . A sample  $y$  can then be drawn from this distribution using the following scheme:

1. Draw  $K$  independent samples from a uniform distribution whose supported lies on the open unit interval:  $U_i \sim \text{Uniform}(0, 1)$ , where  $i = 1, 2, \dots, K$ .
2. For each of the drawn samples, evaluate  $x_i = \log \alpha_i + G_i$ , where  $G_i = -\log(-\log U_i)$  is a sample from the Gumbel distribution (hence the name of the trick).
3. Set  $z_{\arg \max_i \{x_i\}}$  to 1, and all other dimensions of  $z$  to 0.

This looks promising, as the uniform distribution that we sample from does not depend on any of the parameters of the encoder. However, it is still of no use for our purposes, as we have now

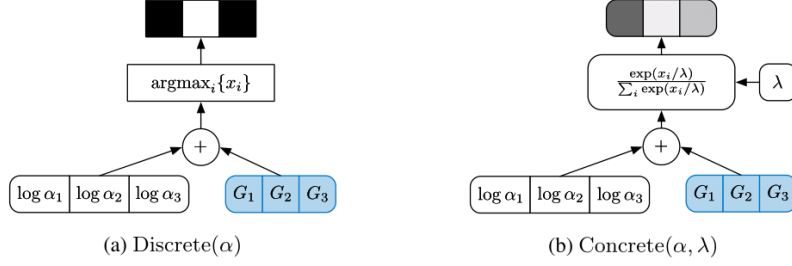


Figure 2.3: Concrete

introduced a non-continuous  $\text{arg max}$  operation. A solution, proposed concurrently by Maddison et. al. [48], and Jang et. al. [36], is to approximate it by a softmax function with a tunable temperature parameter  $\lambda \in (0, \infty)$ , as shown in Equation 2.8. Note that in this case, our sample, drawn from what the authors call a Concrete distribution, is no longer discrete, but rather a continuous relaxation of it that approximates a one-hot vector as  $\lambda \rightarrow 0$ . The difference between the original Gumbel-Max trick, and its Concrete relaxation is illustrated in Figure 2.3.

$$z_i = \frac{\exp((\log \alpha_i + G_i)/\lambda)}{\sum_{j=1}^K \exp((\log \alpha_j + G_j)/\lambda)}. \quad (2.8)$$

Two additional considerations arise when Concrete variables are used in VAEs. First, since no gradients need to be computed at test time, the distributions can be discretized, allowing real categorical variables to be fed through the graph. Second, the KL term in the objective needs to be modified, as simply "plugging" the continuous samples into discrete probability mass functions does not result in a valid KL divergence. This is because the samples are not drawn from the same distribution as is in the numerator under the logarithm, as shown by Equation 2.9, where the Concrete posterior is denoted as  $q_{\alpha, \lambda}^C$ , while  $q_{\alpha}$  and  $p_{\alpha}$  stand for the approximate categorical posterior and prior, respectively:

$$\mathbb{E}_{\mathbf{z} \sim q_{\alpha, \lambda}^C(\mathbf{z}|\mathbf{x})} \left[ \log \frac{q_{\alpha}(\mathbf{z}|\mathbf{x})}{p_{\alpha}(\mathbf{z})} \right]. \quad (2.9)$$

The authors of the original paper [48] provide several options for modifying the KL term, while arguing for replacing the discrete masses in the KL with Concrete densities. However, unlike the Gaussian case, there is no closed form expression for the KL divergence between two Concrete distributions, so it is necessary to approximate it via sampling. However, for the types of models considered in this thesis, a variety of estimators have been found to exhibit too high of a variance, producing poor results as a consequence (Emilien Dupont, personal communication, August 1, 2018). Thus, we opted for replacing the original penalty with an analytic discrete KL term instead:

$$\mathbb{E}_{\mathbf{d} \sim q_{\alpha}(\mathbf{d}|\mathbf{x})} \left[ \log \frac{q_{\alpha}(\mathbf{d}|\mathbf{x})}{p_{\alpha}(\mathbf{d})} \right] = \sum_{i=1}^K q_{\alpha}(d_i|\mathbf{x}) \log \frac{q_{\alpha}(d_i|\mathbf{x})}{p_{\alpha}(d_i)}, \quad (2.10)$$

where the vector  $\mathbf{d}$  represents the one-hot categorical variable. While this strategy provides good results empirically, it is important to emphasize that the use of this KL term in the VAE objective does not necessarily result in a lower-bound of the marginal log-likelihood.

## 2.3 Disentangled representations

general intro and difference from invariance

Most earlier work on learning disentangled representations has tackled the problem with weakly supervised methods, which require some a-priori knowledge of the data generating factors [34, 53, 63]. DC-IGN [43] is perhaps the most well-known method among such approaches. It adopts a clever curriculum learning [6] strategy to encode the data generating factors into a small subset latent variables in a VAE, each of which is a-priori assigned a symbolic interpretation. During training, a random latent from this predefined group is chosen per each mini-batch of training data, which in turn is carefully constructed such that all training examples in it would vary in only the chosen factor. Finally, the bottleneck layer of the underlying VAE is regularized to indeed only capture all the variance in the current batch with a single latent.

Independently controllable factors

Generative models like attend, infer and repeat

A particularly interesting approach that requires no supervision is Shmidhuber’s predictability minimization [57]. They introduce additional submodules for each unit in an autoencoder’s bottleneck layer which try to predict the unit from the remaining ones. In turn, the model is concurrently optimized to minimize the predictability of the units. The scheme was motivated by the idea that such a minimax game would give rise to an internal representation where each dimension is statistically independent of the others. While this ‘competitive learning’ framework was indeed found to be able to learn such factorial codes on toy binary datasets, it has not been shown to achieve promising results on more complicated inputs.

The first scalable, purely unsupervised method that showed reasonable level of disentanglement on real-world datasets was InfoGAN, proposed by Chen et. al. [14]. They extended the powerful generative adversarial network (GAN) framework [26] with an additional objective, which would encourage the information between a small subset of the input noise variables and the generated output to be maximized. In practice, this was achieved by adding an additional auxiliary network into the system. The subnetwork, optimized concurrently with the original GAN objective, was fed with the generated images, and tasked to predict the values of a subset of the noise variables that gave rise to it. Contextually, InfoGAN can be seen as a close cousin of both DC-IGN and predictability minimization, borrowing the idea that only a subset of latents should be interpretable from the former, while having an additional subnetwork that predicts the latent activations is akin to the approach of the latter.



The work most relevant to this thesis is that of Higgins et. al. [32], which demonstrated that simply adding more weight to the KL term in the VAE objective function encourages the posterior to become disentangled. In particular, they multiplied the divergence measure by a constant  $\beta > 1$ , giving the model its name,  $\beta$ -VAE:

$$\mathcal{L}(\mathbf{x}; \boldsymbol{\theta}, \phi) = \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})] - \beta \mathcal{D}_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p_{\boldsymbol{\theta}}(\mathbf{z})). \quad (2.11)$$

The method superseded InfoGAN, proving to be superior in disentangling, as well as enjoying better training stability. In particular, the use of GANs as a backbone has proved to be prohibitive for InfoGAN, as they are notorious for their extremely unstable behaviour in training. Adding further complications, vanilla GANs do not optimize an explicit density function, making it hard to measure the training progress quantitatively [28]. While some recent works have addressed these problems with some success [37, 30, 7], Kim et. al. [38] found that plugging them into the InfoGAN framework gives inconsistent results. As  $\beta$ -VAE uses a principled likelihood-based approach in place of a minimax game whose solution lies at a saddle point in a high-dimensional space, its training dynamics are much wicdly.

In a closely related follow up work to  $\beta$ -VAE, Burgess et. al. [11] build upon it by proposing an enhanced optimization scheme based on some well-developed intuitive arguments. They analyse the  $\beta$ -VAE objective via the lens of the information bottleneck principle [61]. The principle posits that a good internal representation  $Z$  of a learning system is one which is maximally informative about the output  $Y$ , while also compressing the information in the input  $X$ . Mathematically, this intuition is expressed as maximizing a Lagrangian, where "informativeness" is naturally measured by mutual information, denoted as  $I(\cdot; \cdot)$ :

$$\max [I(Z; Y) - \beta I(X; Z)] \quad (2.12)$$

Taking this information-theoretic view of the variational autoencoder, we can conceptualize the probabilistic encoder parameterizing a diagonal Gaussian  $q(\mathbf{z}|\mathbf{x})$  as an independent set of noisy channels  $\{z_i\}$ , each transmitting information about the input  $\mathbf{x}$ . Accordingly, the KL divergence in Equation 2.11 can be thought of as the maximum capacity of said channels per data sample. When the posterior collapses to a unit Gaussian prior  $p(\mathbf{z})$ , the KL term vanishes, and no information about the input is let through the latent bottleneck. Armed with this intuitive framework, it is interesting to analyse how does the reconstruction loss of VAE influence the development of the latent channels.

Burgess et. al. [11] argue that trading off the two terms in the  $\beta$ -VAE objective encourages an encoder to be learned which would have the property that nearby points in data space are close together in the latent space as well. Indeed, suppose that two very different datapoints (as measured by the L2 norm) in pixel-space have similar representations in the latent space. As random noise that is injected into the bottleneck layer, it is likely that the two codes get "mixed up", resulting in the decoder reconstructing the wrong input and incurring a high loss. The situation is illustrated in Figure 2.4. To avoid such a scenario, we would like the encoding of each datapoint to

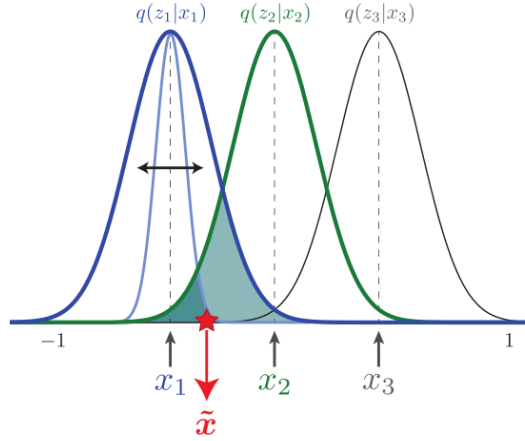


Figure 2.4: **Caption!!**

be as far from others’ encodings as possible, which is achieved by squeezing the posterior variances and dispersing their means across datapoints. This, however, is in turn discouraged by the KL penalty term, which pressures the channels to be overlapping for all datapoints. A reasonable compromise between these two opposing forces would be an arrangement where the latent space is such that when we do make the occasional mistake of decoding the wrong datapoint, we are at least outputting something that is not too far from the original.

A second fundamental observation Burgess et. al. [11] make is that when the latent channel capacities are constrained to be very low, the only information that gets through the bottleneck is one whose utilization by the decoder increases the value of the reconstruction term  $\mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})]$  the most. As an example, suppose we have a dataset of white shapes drawn on a black canvas, each of which has 5 factors of variation: translation along the vertical and horizontal axes, scale, orientation, and class (heart, square, or ellipse) [49]. Now, for a high value of  $\beta$ , only positional information of the sprite might get through the overly constrained latent bottleneck, as getting the position wrong incurs the most cost in terms of log-likelihood. Subsequently, by gradually increasing the capacities of the latent channels, the model starts to allocate additional latents to new factors of variation, like scale and rotation. Crucially, due to the data locality pressure described before, all the factors of variation are encoded into distinct latents, satisfying our working definition of disentanglement. Figure 2.5 illustrates this phenomenon.

The intuitive framework laid out above suggests a natural modification to the  $\beta$ -VAE objective, given by Equation 2.13. In the new formulation, the capacities of the channels, denoted as  $C$ , are made explicit, and increased gradually in the course of training. The new hyperparameter  $\gamma$  specifies the degree of pressure put on the latent channels to be close to their desired capacity. It is usually chosen big enough that the KL term would be relatively close to  $C$  at all times, as this makes the external control of the channel capacities rather convenient.

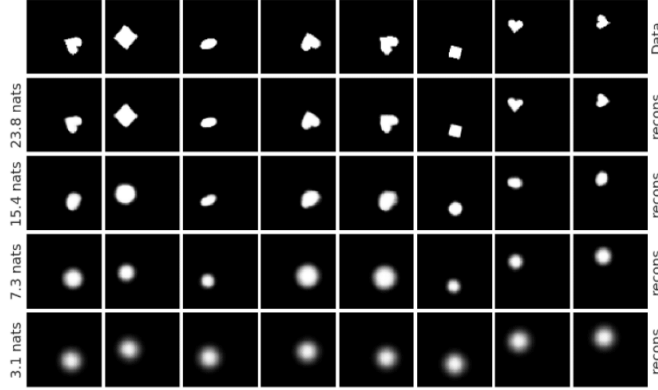


Figure 2.5: **Caption!!**

$$\mathcal{L}(\mathbf{x}; \boldsymbol{\theta}, \phi) = \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})] - \gamma |\mathcal{D}_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p_{\boldsymbol{\theta}}(\mathbf{z})) - C|. \quad (2.13)$$

While  $\beta$ -VAE with a controlled capacity increase showed good results disentangling continuous factors of variation, it struggled to capture discrete generative factors. For example, it was unable to completely disentangle the classes of the shapes in the sprites dataset. As a way to overcome this difficulty, Dupont introduced JointVAE [19], which augmented the original multivariate normal latent distribution of  $\beta$ -VAE with categorical Concrete random variables, leading to the following objective:

$$\begin{aligned} \mathcal{L}(\mathbf{x}; \boldsymbol{\theta}, \phi) = & \mathbb{E}_{\mathbf{z}, \mathbf{c} \sim q_{\phi}(\mathbf{z}, \mathbf{c}|\mathbf{x})} [\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z}, \mathbf{c})] - \\ & \gamma \{ |\mathcal{D}_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p_{\boldsymbol{\theta}}(\mathbf{z})) - C_z| + |\mathcal{D}_{\text{KL}}(q_{\phi}(\mathbf{c}|\mathbf{x}) || p_{\boldsymbol{\theta}}(\mathbf{c})) - C_c| \}, \end{aligned} \quad (2.14)$$

where  $\mathbf{z}$  and  $\mathbf{c}$  denote the continuous Gaussian and Concrete latent variables, respectively. The new model was shown to be able to outperform the original formulation of  $\beta$ -VAE on datasets where capturing discrete generative factors is paramount to achieving good reconstruction quality, as is the case for MNIST handwritten digit dataset [46]. However, complete disentanglement of the sprites dataset is still beyond the capabilities of current models, as the types of shapes are all quite similar in terms of the area they cover, making the possible increase in log-likelihood marginal.

## 2.4 Adversarial examples

The problem of adversarial examples in the context of deep neural networks was first articulated by Szegedy et. al. [58]. They showed that by perturbing the pixels of the input image in a by an imperceptibly small amount, it is possible to cause neural-network based classifiers to assign high likelihoods to arbitrary classes. It is important to stress that the pixels of the original image are not changed randomly, but are rather done so in a principled way. More specifically, the authors



Figure 2.6: **Caption!!**

use box-constrained L-BFGS [51] to solve the following optimization problem <sup>3</sup>:

$$\begin{aligned} & \text{minimize } c\|\mathbf{r}\|_2 + \text{loss}_f(\mathbf{x} + \mathbf{r}, l) \\ & \text{subject to } \mathbf{x} + \mathbf{r} \in [0, 1]^m, \end{aligned} \quad (2.15)$$

where  $\mathbf{r}$  is the  $m$ -dimensional perturbation vector that is added to the original image  $\mathbf{x}$ , such that the classifier  $f$  would output the class  $l$ .  $\text{loss}_f$  is some continuous measure of classification loss, commonly taken to be cross-entropy, while the coefficient  $c > 0$  balances the size of the perturbation and the loss for the generated adversarial image. Of course, the problem is only of interest when  $f(\mathbf{x}) \neq l$ . Figure 2.6 illustrates some adversarial examples  $\mathbf{x} + \mathbf{r}$  generated by this procedure.

In subsequent work, Goodfellow et. al. [27] argue that the existence of adversarial examples can be at least partially explained by their linear nature. Linear models have the property that their outputs change proportionally to the change in the intensity of any of the pixels. Thus, given an input with high enough dimensionality, one can make small perturbations to a large number of pixels, all of which add up to produce a large shift in the output. Goodfellow et. al. make the observation that this same phenomenon can manifest itself in modern deep neural networks, as they are often explicitly constructed to be rather linear for ease of optimization. Illustrative examples include the rectified linear units as general-purpose activation functions [25], LSTMs for time-series analysis [35], and residual networks for computer vision [31]. Moreover, adding further evidence in favour of the linearity hypothesis, they observed that RBF networks [10] - a method that classifies inputs based on their similarity, as measured by the exponentiated quadratic kernels, to a set of predefined centroids - are significantly more robust to adversarial examples.

One of the most surprising properties of adversarial examples is their generality and transferability. Early work on the problem found that they generalize across different partitions of the same dataset [58] and across models with different architectures [27]. More recently, Kurakin et. al. [44] discovered that they can also be extrapolated to the physical world, showing that models trained on the regular ImageNet [17] dataset are misclassifying adversarial examples even if they

<sup>3</sup>Note that this specific formulation constrains the images to have a single channel and either be normalized or grayscale. However, generalizing it to unnormalized RGB images is straightforward.

are printed on paper and viewed through a physical camera. Somewhat reassuringly, they have also been found to fool time-limited humans [21]. While most research on adversarial examples has focussed on image classification, models trained to solve other kinds of machine learning tasks are subject to them as well, including ones specializing in speech-to-text [13], neural machine translation [20], and deep reinforcement learning [3].

Pertinent to this thesis is the work by Kos. et. al. [41], who investigate adversarial examples in the context of generative models. It was indeed found that VAEs and VAE-GANs [45], a derivative of VAEs which enjoys better sample quality due to the fact that its reconstruction loss operates in the feature space of a GAN, are also suspect to them. To generate adversarial examples for such models, they took inspiration from Sabour et. al. [56], and perturbed the source image in such a way that its predicted posterior mean would be close to the posterior mean of some target image. Mathematically, this optimization problem is equivalent to the one in Equation 2.15, with the classification loss replaced by the L2 norm measuring the discrepancy of the latent representations. They did, however, use Adam instead of L-BFGS, as advocated by Carlini & Wagner [12].

Another closely related work is that of Alemi et. al. [1], in which they parameterize the information bottleneck in Equation 2.12 using neural networks. By subsequently deriving a variational lower bound to this objective and making use of the reparameterization trick, they end up with a model analogous to that of  $\beta$ -VAE, given that the output variable  $Y$  in Equation 2.12 is taken to be the reconstruction of the input. However, as their treatment is more broad, they are also able to generalize this approach to a classification setting, by taking  $Z$  to be the class label. Intuitively, this amounts to swapping the decoder of  $\beta$ -VAE with an appropriate classifier. To demonstrate the usefulness of the bottleneck principle, they measured the robustness of the resulting classifier to adversarial examples on the MNIST and ImageNet datasets. They found that the model becomes progressively harder to fool as the capacity of the bottleneck is increased.

## Chapter 3

# Methods

To measure the effect that disentangling has on the ability to generate adversarial examples for VAEs, we used  $\beta$ -VAE as the appropriate model. The analysis was carried out on the MNIST dataset, due to its several desirable properties.

First, it is one of the few datasets that have been extensively analysed both in the disentanglement, as well as in the adversarial example literature. Indeed, MNIST was one of the first dataset for which adversarial examples were constructed [58], and it has been shown that at least 3 factors of variation of handwritten digits can be disentangled on this dataset: class, angle, and thickness [14], [19]. Moreover, the two most closely related works to this thesis, that of Kos et. al. and Alemi et. al, concentrate on MNIST as well, making drawing comparisons easier.

Second, MNIST allows

Encoder	Decoder
$\mathbf{x} \in \mathbb{R}^{32 \times 32}$ grayscale image	$\mathbf{z} \in \mathbb{R}^{10}$ , and $\mathbf{c}$ is one-hot 10-dim
Conv. 32 ReLU, $4 \times 4$ kernels, 2 stride	Dense. 128 ReLU
Conv. 32 ReLU, $4 \times 4$ kernels, 2 stride	Dense. $64 \times 4 \times 4$ ReLU
Conv. 64 ReLU, $4 \times 4$ kernels, 2 stride	Upconv. 64 ReLU, $4 \times 4$ kernels, 2 stride
Dense. 128 ReLU	Upconv. 32 ReLU, $4 \times 4$ kernels, 2 stride
Dense. 128 linear	Upconv. 32 Sigmoid, $4 \times 4$ kernels, 2 stride

Table 3.1: VAE architecture on MNIST.

# Chapter 4

## Results

	AS ignore-target	AS target
<b>Disentangled</b>	36.86%	12.54%
<b>Entangled</b>	29.83%	9.49%
<b>Entangled wdec</b>	26.88%	7.85%

Table 4.1: Means

Source	0	1	2	3	4	5	6	7	8	9
<b>Disentangled</b>	8.89% (3.18%)	82.62% (33.76%)	20.92% (5.36%)	16.72% (6.20%)	33.73% (13.02%)	44.25% (11.24%)	38.49% (13.48%)	40.57% (11.65%)	18.23% (8.47%)	64.23% (18.97%)
<b>Entangled</b>	6.12% (1.93%)	63.75% (28.58%)	11.12% (3.06%)	12.11% (4.30%)	64.40% (13.46%)	37.79% (11.04%)	27.83% (9.55%)	39.70% (7.89%)	16.27% (6.43%)	19.19% (8.70%)
<b>Entangled wdec</b>	3.30% (0.97%)	59.78% (23.44%)	7.09% (1.41%)	10.83% (3.02%)	53.99% (10.94%)	34.74% (8.65%)	24.94% (6.10%)	18.45% (4.26%)	23.46% (7.61%)	32.28% (12.09%)

Table 4.2: Source marginals

Target	0	1	2	3	4	5	6	7	8	9
<b>Disentangled</b>	40.36% (1.78%)	25.91% (0.93%)	46.50% (25.91%)	41.81% (10.19%)	31.76% (14.23%)	38.98% (10.04%)	37.46% (7.80%)	27.77% (7.42%)	46.23% (32.19%)	31.87% (14.87%)
<b>Entangled</b>	29.39% (3.52%)	13.88% (0.24%)	35.42% (11.50%)	40.39% (18.68%)	24.71% (5.53%)	32.75% (7.38%)	22.94% (4.23%)	28.20% (4.54%)	43.49% (25.44%)	27.11% (13.87%)
<b>Entangled wdec</b>	28.36% (4.50%)	12.62% (1.06%)	37.50% (16.60%)	36.40% (12.30%)	19.14% (1.88%)	27.39% (4.44%)	20.63% (3.89%)	29.92% (12.57%)	33.21% (12.55%)	23.67% (8.72%)

Table 4.3: Target marginals

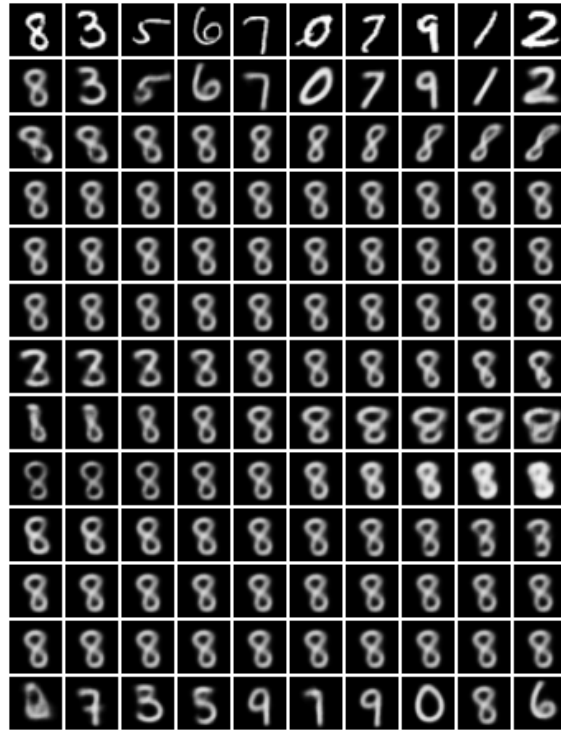


Figure 4.1: mnist disentangled traversal



Figure 4.2: MNIST entangled reconstruction



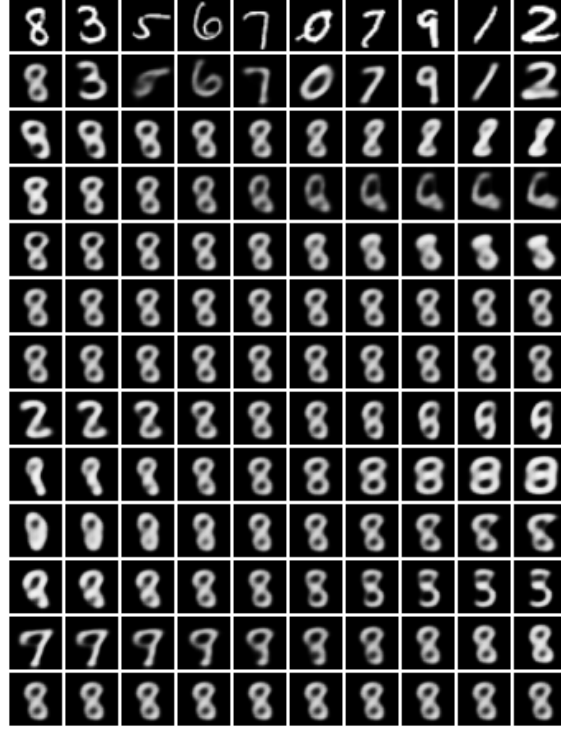


Figure 4.3: MNIST entangled weight decay reconstruction

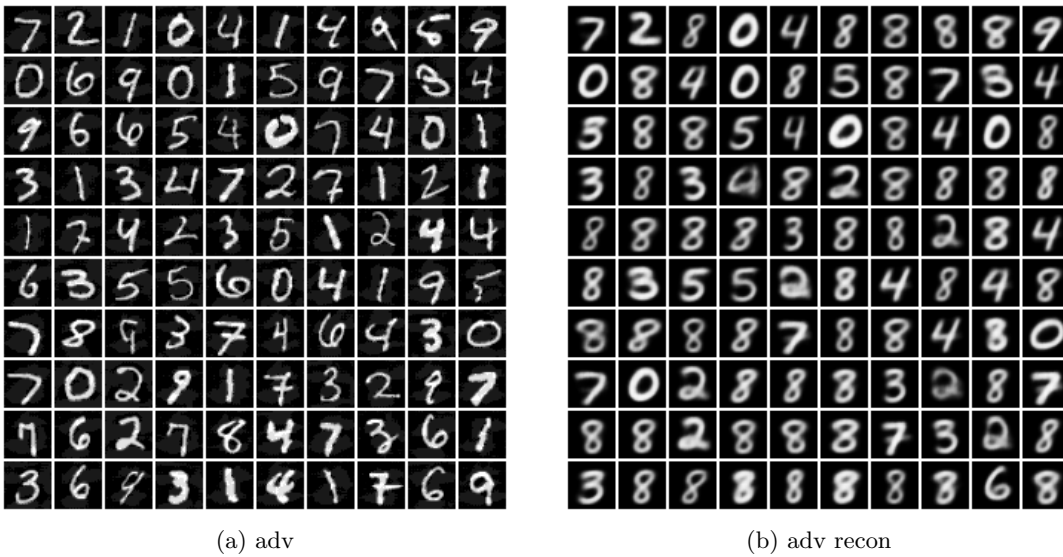
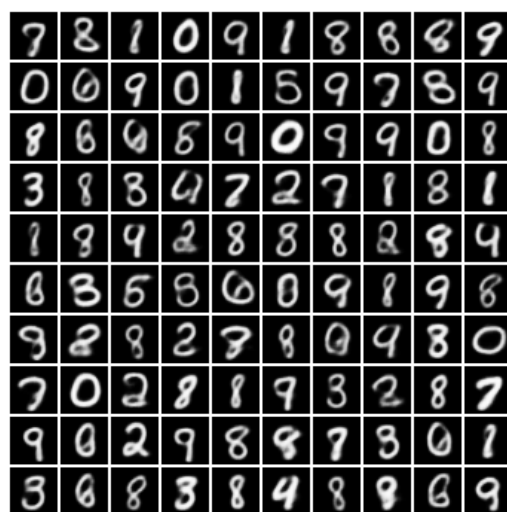


Figure 4.4: MNIST disentangled adv and adv recon



(a) adv



(b) adv recon

Figure 4.5: MNIST ent adv and adv recon

## Chapter 5

## Conclusion

# Chapter 6

## Appendix

Source	Targets									
	0	1	2	3	4	5	6	7	8	9
0	-	6.58% (0.00%)	11.12% (5.98%)	8.36% (3.34%)	9.21% (1.99%)	11.15% (6.44%)	8.13% (4.71%)	7.60% (0.31%)	8.74% (5.89%)	9.10% (0.00%)
1	99.48% (0.00%)	-	97.62% (93.64%)	95.63% (4.98%)	73.25% (42.63%)	95.49% (10.87%)	97.64% (26.85%)	23.43% (13.36%)	95.48% (72.58%)	65.52% (38.95%)
2	7.74% (0.66%)	24.84% (0.42%)	-	14.63% (4.33%)	27.98% (4.24%)	23.36% (1.58%)	6.64% (2.29%)	26.38% (12.29%)	26.85% (22.36%)	29.84% (0.11%)
3	2.90% (0.58%)	20.00% (1.14%)	16.61% (11.86%)	-	19.27% (0.80%)	17.33% (14.38%)	12.00% (1.43%)	19.86% (1.47%)	24.27% (19.49%)	18.22% (4.61%)
4	24.70% (0.85%)	44.48% (0.48%)	37.18% (25.28%)	31.27% (8.15%)	-	21.63% (4.27%)	21.31% (11.65%)	51.10% (6.48%)	34.71% (26.55%)	37.16% (33.50%)
5	46.66% (3.54%)	45.71% (0.65%)	67.40% (8.54%)	34.48% (26.92%)	31.63% (8.79%)	-	51.50% (12.26%)	42.93% (2.36%)	46.55% (30.24%)	31.36% (7.87%)
6	28.87% (7.41%)	34.20% (2.72%)	39.89% (33.11%)	39.65% (6.35%)	38.06% (24.21%)	33.49% (10.85%)	-	38.18% (2.74%)	50.00% (27.84%)	44.10% (6.11%)
7	47.31% (1.43%)	4.44% (0.23%)	40.25% (24.94%)	47.69% (2.08%)	38.66% (4.94%)	52.41% (6.63%)	41.66% (0.37%)	-	49.03% (27.36%)	43.71% (36.91%)
8	25.11% (0.79%)	13.13% (1.71%)	30.48% (24.20%)	24.06% (14.71%)	11.16% (5.69%)	18.67% (12.36%)	22.80% (6.02%)	10.87% (5.03%)	-	7.79% (5.73%)
9	80.48% (0.80%)	39.78% (1.01%)	77.94% (5.65%)	80.51% (20.82%)	36.64% (34.79%)	77.28% (22.97%)	75.44% (4.60%)	29.54% (22.74%)	80.48% (57.39%)	-

Table 6.1: Disentangled

Source	Targets									
	0	1	2	3	4	5	6	7	8	9
0	-	3.46% (0.00%)	8.93% (2.60%)	7.47% (1.89%)	6.85% (0.31%)	5.59% (1.60%)	5.74% (2.34%)	4.58% (0.42%)	8.82% (6.51%)	3.65% (1.67%)
1	58.13% (0.21%)	-	76.68% (52.17%)	69.11% (41.26%)	60.94% (24.38%)	66.59% (17.45%)	50.53% (19.64%)	64.42% (28.05%)	73.72% (55.29%)	53.65% (18.76%)
2	13.13% (1.89%)	4.53% (0.10%)	-	13.32% (9.73%)	9.25% (0.54%)	14.37% (0.36%)	11.17% (1.02%)	7.71% (1.04%)	16.65% (12.13%)	9.94% (0.74%)
3	8.26% (1.00%)	4.71% (0.11%)	6.93% (4.01%)	-	20.51% (0.11%)	14.46% (10.78%)	12.30% (0.68%)	4.75% (0.22%)	26.07% (19.89%)	11.00% (1.92%)
4	65.87% (3.26%)	42.55% (0.65%)	66.48% (9.88%)	89.27% (10.51%)	-	79.10% (4.01%)	26.20% (5.15%)	71.54% (3.79%)	79.87% (30.09%)	58.75% (53.80%)
5	37.52% (4.06%)	23.74% (0.12%)	49.51% (6.25%)	40.81% (30.51%)	38.44% (0.49%)	-	26.43% (6.82%)	33.75% (0.62%)	55.09% (42.21%)	34.78% (8.29%)
6	25.22% (17.32%)	6.80% (0.44%)	29.56% (8.57%)	41.88% (9.42%)	14.19% (6.22%)	27.86% (15.47%)	-	38.40% (0.33%)	42.39% (25.33%)	24.19% (2.82%)
7	33.82% (1.12%)	24.40% (0.21%)	30.49% (8.33%)	28.37% (12.51%)	50.22% (6.47%)	40.85% (1.09%)	55.27% (0.11%)	-	56.51% (9.71%)	37.35% (31.43%)
8	12.83% (1.03%)	7.80% (0.00%)	21.08% (9.91%)	36.89% (30.99%)	8.67% (0.57%)	21.05% (6.97%)	10.69% (2.21%)	16.72% (0.80%)	-	10.72% (5.36%)
9	9.74% (1.79%)	6.89% (0.54%)	29.15% (1.82%)	36.42% (21.34%)	13.30% (10.64%)	24.85% (8.65%)	8.11% (0.11%)	11.96% (5.60%)	32.30% (27.78%)	-

Table 6.2: Entangled

Source	Targets									
	0	1	2	3	4	5	6	7	8	9
0	-	2.85% (0.00%)	6.19% (1.47%)	3.70% (1.37%)	1.99% (0.00%)	4.59% (1.38%)	2.78% (1.71%)	2.43% (0.32%)	3.41% (2.13%)	1.78% (0.31%)
1	68.30% (1.36%)	-	73.49% (59.25%)	67.68% (25.63%)	39.62% (3.75%)	57.65% (10.70%)	38.69% (16.51%)	61.00% (47.46%)	76.51% (35.59%)	55.09% (10.69%)
2	9.64% (2.15%)	4.82% (0.21%)	-	6.37% (3.08%)	6.12% (0.91%)	8.84% (0.25%)	6.18% (0.93%)	7.40% (1.39%)	6.90% (2.94%)	7.53% (0.87%)
3	9.82% (1.75%)	5.70% (0.35%)	10.50% (6.96%)	-	14.30% (0.00%)	12.83% (8.60%)	14.10% (0.96%)	8.55% (2.66%)	9.54% (4.05%)	12.09% (1.86%)
4	55.38% (2.22%)	30.42% (2.12%)	58.08% (17.51%)	72.58% (4.84%)	-	70.14% (1.37%)	27.46% (4.53%)	56.47% (12.21%)	66.54% (10.78%)	48.83% (42.91%)
5	38.78% (12.70%)	21.45% (1.37%)	43.55% (3.36%)	34.33% (19.27%)	40.65% (0.27%)	-	32.88% (7.37%)	32.69% (5.75%)	34.68% (20.38%)	33.61% (7.41%)
6	17.76% (12.83%)	8.63% (1.55%)	35.99% (11.12%)	33.48% (3.71%)	16.85% (5.03%)	18.91% (5.44%)	-	30.75% (0.66%)	37.07% (12.98%)	25.00% (1.54%)
7	13.38% (1.70%)	9.21% (1.41%)	22.97% (12.09%)	19.82% (9.03%)	19.93% (0.68%)	18.98% (1.36%)	21.93% (0.00%)	-	26.29% (3.21%)	13.52% (8.83%)
8	19.32% (3.26%)	16.40% (1.59%)	37.96% (26.31%)	34.68% (20.08%)	19.24% (0.38%)	16.34% (4.40%)	21.80% (2.42%)	29.87% (6.03%)	-	15.54% (4.05%)
9	22.89% (2.49%)	14.12% (0.92%)	48.74% (11.30%)	54.99% (23.67%)	13.57% (5.88%)	38.27% (6.44%)	19.83% (0.61%)	40.14% (36.64%)	37.96% (20.88%)	-

Table 6.3: Entangled with weight decay

# Bibliography

- [1] Alexander A. Alemi, Ian Fischer, Joshua V. Dillon, and Kevin Murphy. Deep variational information bottleneck. *CoRR*, abs/1612.00410, 2016.
- [2] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *CoRR*, abs/1409.0473, 2014.
- [3] Vahid Behzadan and Arslan Munir. Vulnerability of deep reinforcement learning to policy induction attacks. *CoRR*, abs/1701.04143, 2017.
- [4] Yoshua Bengio. Learning deep architectures for ai. *Found. Trends Mach. Learn.*, 2(1):1–127, January 2009.
- [5] Yoshua Bengio, Aaron C. Courville, and Pascal Vincent. Unsupervised feature learning and deep learning: A review and new perspectives. *CoRR*, abs/1206.5538, 2012.
- [6] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. In *Proceedings of the 26th Annual International Conference on Machine Learning*, ICML '09, pages 41–48, New York, NY, USA, 2009. ACM.
- [7] David Berthelot, Tom Schumm, and Luke Metz. BEGAN: boundary equilibrium generative adversarial networks. *CoRR*, abs/1703.10717, 2017.
- [8] Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer-Verlag, Berlin, Heidelberg, 2006.
- [9] Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Bernhard Firner, Beat Flepp, Prasoon Goyal, Lawrence D. Jackel, Mathew Monfort, Urs Muller, Jiakai Zhang, Xin Zhang, Jake Zhao, and Karol Zieba. End to end learning for self-driving cars. *CoRR*, abs/1604.07316, 2016.
- [10] D.S. Broomhead and D. Lowe. Multivariable functional interpolation and adaptive networks. *Complex Systems*, 2:321–355, 1988.
- [11] Christopher P. Burgess, Irina Higgins, Arka Pal, Loïc Matthey, Nick Watters, Guillaume Desjardins, and Alexander Lerchner. Understanding disentangling in  $\beta$ -vae. *CoRR*, abs/1804.03599, 2018.
- [12] Nicholas Carlini and David A. Wagner. Towards evaluating the robustness of neural networks. *CoRR*, abs/1608.04644, 2016.

- [13] Nicholas Carlini and David A. Wagner. Audio adversarial examples: Targeted attacks on speech-to-text. *CoRR*, abs/1801.01944, 2018.
- [14] Xi Chen, Yan Duan, Rein Houthooft, John Schulman, Ilya Sutskever, and Pieter Abbeel. Info-gan: Interpretable representation learning by information maximizing generative adversarial nets. *CoRR*, abs/1606.03657, 2016.
- [15] G. Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of Control, Signals and Systems*, 2(4):303–314, Dec 1989.
- [16] Rodrigo de Bem, Arnab Ghosh, Thalaiyasingam Ajanthan, Ondrej Miksik, N. Siddharth, and Philip H. S. Torr. Dgpose: Disentangled semi-supervised deep generative models for human body analysis. *CoRR*, abs/1804.06364, 2018.
- [17] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR09*, 2009.
- [18] Pratik Dubal, Rohan Mahadev, Suraj Kothawade, Kunal Dargan, and Rishabh Iyer. Deployment of customized deep learning based video analytics on surveillance cameras. *CoRR*, abs/1805.10604, 2018.
- [19] Emilien Dupont. Learning disentangled joint continuous and discrete representations, 2018.
- [20] Javid Ebrahimi, Daniel Lowd, and Dejing Dou. On adversarial examples for character-level neural machine translation. 2018.
- [21] Gamaleldin F. Elsayed, Shreya Shankar, Brian Cheung, Nicolas Papernot, Alex Kurakin, Ian J. Goodfellow, and Jascha Sohl-Dickstein. Adversarial examples that fool both human and computer vision. *CoRR*, abs/1802.08195, 2018.
- [22] Yaxiang Fan, Gongjian Wen, Deren Li, ShaoHua Qiu, and Martin D. Levine. Video anomaly detection and localization via gaussian mixture fully convolutional variational autoencoder. *CoRR*, abs/1805.11223, 2018.
- [23] Zoubin Ghahramani and Michael I. Jordan. Supervised learning from incomplete data via an em approach. In J. D. Cowan, G. Tesauro, and J. Alspector, editors, *Advances in Neural Information Processing Systems 6*, pages 120–127. Morgan-Kaufmann, 1994.
- [24] Ross B. Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. *CoRR*, abs/1311.2524, 2013.
- [25] Xavier Glorot, Antoine Bordes, and Yoshua Bengio. Deep sparse rectifier neural networks. In Geoffrey Gordon, David Dunson, and Miroslav Dudk, editors, *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, volume 15 of *Proceedings of Machine Learning Research*, pages 315–323, Fort Lauderdale, FL, USA, 11–13 Apr 2011. PMLR.
- [26] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In Z. Ghahramani,

- M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 2672–2680. Curran Associates, Inc., 2014.
- [27] Ian Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *International Conference on Learning Representations*, 2015.
- [28] Ian J. Goodfellow. NIPS 2016 tutorial: Generative adversarial networks. *CoRR*, abs/1701.00160, 2017.
- [29] Kathrin Grosse, Nicolas Papernot, Praveen Manoharan, Michael Backes, and Patrick McDaniel. Adversarial examples for malware detection. In Simon N. Foley, Dieter Gollmann, and Einar Sneekkenes, editors, *Computer Security – ESORICS 2017*, pages 62–79, Cham, 2017. Springer International Publishing.
- [30] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. Improved training of wasserstein gans. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 5767–5777. Curran Associates, Inc., 2017.
- [31] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *CoRR*, abs/1512.03385, 2015.
- [32] Irina Higgins, Loïc Matthey, Xavier Glorot, Arka Pal, Benigno Uria, Charles Blundell, Shakir Mohamed, and Alexander Lerchner. Early visual concept learning with unsupervised deep learning. *CoRR*, abs/1606.05579, 2016.
- [33] Irina Higgins, Arka Pal, Andrei A. Rusu, Loic Matthey, Christopher P Burgess, Alexander Pritzel, Matthew Botvinick, Charles Blundell, and Alexander Lerchner. Darla: Improving zero-shot transfer in reinforcement learning, 2017.
- [34] Geoffrey E. Hinton, Alex Krizhevsky, and Sida D. Wang. Transforming auto-encoders. In *Proceedings of the 21th International Conference on Artificial Neural Networks - Volume Part I*, ICANN’11, pages 44–51, Berlin, Heidelberg, 2011. Springer-Verlag.
- [35] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Comput.*, 9(8):1735–1780, November 1997.
- [36] Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. *CoRR*, abs/1611.01144, 2016.
- [37] Hyunjik Kim and Andriy Mnih. Disentangling by factorising. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 2649–2658, Stockholmsmssan, Stockholm Sweden, 10–15 Jul 2018. PMLR.
- [38] Hyunjik Kim and Andriy Mnih. Disentangling by factorising. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 2649–2658, Stockholmsmssan, Stockholm Sweden, 10–15 Jul 2018. PMLR.



- [39] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2014.
- [40] Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. *CoRR*, abs/1312.6114, 2013.
- [41] Jernej Kos, Ian Fischer, and Dawn Song. Adversarial examples for generative models, 2017.
- [42] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 1097–1105. Curran Associates, Inc., 2012.
- [43] Tejas D. Kulkarni, Will Whitney, Pushmeet Kohli, and Joshua B. Tenenbaum. Deep convolutional inverse graphics network. *CoRR*, abs/1503.03167, 2015.
- [44] Alexey Kurakin, Ian J. Goodfellow, and Samy Bengio. Adversarial examples in the physical world. *CoRR*, abs/1607.02533, 2016.
- [45] Anders Boesen Lindbo Larsen, Søren Kaae Sønderby, and Ole Winther. Autoencoding beyond pixels using a learned similarity metric. *CoRR*, abs/1512.09300, 2015.
- [46] Yann LeCun and Corinna Cortes. MNIST handwritten digit database. 2010.
- [47] R. Duncan Luce. *Individual Choice Behavior: A Theoretical analysis*. Wiley, New York, NY, USA, 1959.
- [48] Chris J. Maddison, Andriy Mnih, and Yee Whye Teh. The concrete distribution: A continuous relaxation of discrete random variables. *CoRR*, abs/1611.00712, 2016.
- [49] Loic Matthey, Irina Higgins, Demis Hassabis, and Alexander Lerchner. dsprites: Disentanglement testing sprites dataset. <https://github.com/deepmind/dsprites-dataset/>, 2017.
- [50] Ari S. Morcos, David G.T. Barrett, Neil C. Rabinowitz, and Matthew Botvinick. On the importance of single directions for generalization. In *International Conference on Learning Representations*, 2018.
- [51] Jorge Nocedal and Stephen J. Wright. *Numerical Optimization*. Springer, New York, NY, USA, second edition, 2006.
- [52] Daniele Palossi, Antonio Loquercio, Francesco Conti, Eric Flamand, Davide Scaramuzza, and Luca Benini. Ultra low power deep-learning-powered autonomous nano drones. *CoRR*, abs/1805.01831, 2018.
- [53] Scott Reed, Kihyuk Sohn, Yuting Zhang, and Honglak Lee. Learning to disentangle factors of variation with manifold interaction. In Eric P. Xing and Tony Jebara, editors, *Proceedings of the 31st International Conference on Machine Learning*, volume 32 of *Proceedings of Machine Learning Research*, pages 1431–1439, Beijing, China, 22–24 Jun 2014. PMLR.

- [54] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In Eric P. Xing and Tony Jebara, editors, *Proceedings of the 31st International Conference on Machine Learning*, volume 32 of *Proceedings of Machine Learning Research*, pages 1278–1286, Beijing, China, 22–24 Jun 2014. PMLR.
- [55] Adam Roberts, Jesse Engel, Colin Raffel, Curtis Hawthorne, and Douglas Eck. A hierarchical latent vector model for learning long-term structure in music. *CoRR*, abs/1803.05428, 2018.
- [56] Sara Sabour, Yanshuai Cao, Fartash Faghri, and David J. Fleet. Adversarial manipulation of deep representations. *CoRR*, abs/1511.05122, 2015.
- [57] Jrgen Schmidhuber. Learning factorial codes by predictability minimization. *Neural Computation*, 4(6):863–879, 1992.
- [58] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *CoRR*, abs/1312.6199, 2013.
- [59] Pedro Tabacof, Julia Tavares, and Eduardo Valle. Adversarial images for variational autoencoders. *CoRR*, abs/1612.00155, 2016.
- [60] Tijmen Tieleman and Geoffrey Hinton. Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude. *COURSERA: Neural networks for machine learning*, 4(2):26–31, 2012.
- [61] Naftali Tishby, Fernando C. Pereira, and William Bialek. The information bottleneck method. pages 368–377, 1999.
- [62] Ronald J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8(3):229–256, May 1992.
- [63] Jimei Yang, Scott E. Reed, Ming-Hsuan Yang, and Honglak Lee. Weakly-supervised disentangling with recurrent transformations for 3d view synthesis. *CoRR*, abs/1601.00706, 2016.