

Imperial College of Science, Technology and Medicine
Department of Computing

On the Feasibility of Using Fully-Convolutional Variational Autoencoders to Advance Deep Symbolic Reinforcement Learning

D. G. Sherburn

Submitted in part fulfilment of the requirements for the degree of
Master of Engineering in Joint Mathematics and Computing
of Imperial College, June 2017

Acknowledgements

I would like to express (whatever feelings I have) to:

- My supervisor
- My second supervisor
- Other researchers
- My family and friends

Dedication

Dedication here.

‘Quote text here.’

Guy Quoted

Contents

Acknowledgements	i
1 Introduction	1
1.1 Motivation	1
1.2 Objectives	4
1.3 Contributions	5
2 Background	6
2.1 Loss functions	6
2.1.1 Euclidean Distance	7
2.1.2 Binary Cross-Entropy	7
2.2 Stella	8
2.3 Arcade Learning Environment	8
2.4 Keras	9
2.5 TensorBoard	9
3 Related Work	10
3.1 Autoencoders	10

3.1.1	Fully-Connected Autoencoders	11
3.1.2	Fully-Convolutional Autoencoders	13
3.2	Variational Autoencoders	15
3.2.1	A Probabilistic Perspective	15
3.2.2	Overcoming the Intractable Posterior	16
3.2.3	Finding a Suitable Loss Function: the ELBO	16
3.2.4	Writing the ELBO in Closed-Form	18
3.2.5	Implementing the Variational Autoencoder	20
3.2.6	Intuition Behind the Variational Autoencoder	22
3.3	Disentangled Representations	27
3.4	Unsupervised Learning of Generative Factors	27
3.4.1	InfoGAN	27
3.4.2	β -VAE	27
3.5	Improving Sampling from Generative Autoencoders with Markov Chains	31
3.6	Regularizing CNNs with Locally Constrained Decorrelations	33
3.7	Batch Normalisation	33
4	Implementation	34
4.1	Dimensionality Reduction	34
4.1.1	Pre-processing Pipeline	35
4.2	Qualitative Assessment Using GUIs	35
4.3	Ensuring Numerical Stability in the Latent Space	37
4.4	Activation Functions in the Latent Space	37

5 Methods	39
5.1 Single Latent Filter	39
5.1.1 Architecture	40
5.1.2 Derivation	40
5.2 Disentangling Latent Neurons	40
5.2.1 Architecture	41
5.2.2 Derivation	42
5.3 Disentangling Latent Filters Using Averages	42
5.3.1 Architecture	42
5.3.2 Derivation	43
5.4 Decoupling Latent Filters Using Weighted-Averages	43
5.4.1 Architecture	43
5.4.2 Derivation	43
5.5 Separating Colour Spaces	43
5.5.1 Architecture	44
5.5.2 Derivation	44
5.6 Orthogonal Convolutions	44
5.6.1 Architecture	44
5.6.2 Derivation	44
5.7 Winner Takes All	44
5.7.1 Architecture	45
5.7.2 Derivation	45

6	Results	46
6.1	The No Free Lunch Relationship Between Reconstruction Loss and KL Divergence	46
6.2	Using Batch Normalisation With Convolutional Variational Latent Spaces	46
6.3	Using Batch Normalisation With Convolutional Variational Latent Spaces	46
7	Conclusion	48
7.1	Summary of Thesis Achievements	48
7.2	Applications	48
7.3	Future Work	48
Bibliography		48

List of Tables

1.1	Low-dimensional symbolic representation	5
3.1	A simple fully-connected autoencoder with one hidden layer. After 15 epochs, the validation score was recorded to be 71.94.	11
3.2	A simple fully-convolutional autoencoder with 2D convolutions and max pooling, plus the corresponding deconvolutional layers. After 15 epochs, the validation score was recorded to be 64.89.	13

List of Figures

1.1 May 1997: Gary Kasparov makes his first move against IBM's Deep Blue. Deep Blue would later emerge the victor in the best of six games; the first time a reigning world chess champion is defeated by a computer. [26]	2
1.2 March 2016: Lee Sedol, one of the greatest modern Go players, plays his first move of game three against AlphaGo. AlphaGo won four of five games. This feat was considered by many to be a decade away. [22]	2
1.3 Overview of deep symbolic reinforcement learning system architecture. A: The neural back end maps high-dimensional raw input data to a compositionally structured symbolic representation. B: The compositionally structured symbolic representation. C: Reinforcement learning of mapping from symbolic representation to action with maximum expected reward over time. <i>Source: Garnelo et al.</i> [13].	4
1.4 A toy example of a raw high-dimensional input.	4
2.1	7
2.2	7
2.3	7
3.1 A black-box description of an autoencoder. The autoencoder learns the identity function, and in turn, the encoder and decoder learn suitable encoding and decoding algorithms respectively.	10

3.2	An example architecture of a fully-connected autoencoder. The latent space is constrained by having fewer neurons than the input and output layers.	11
3.3	A collection of images from the MNIST data set and their respective reconstructions using the fully-connected autoencoder specified in Table 3.1. The original MNIST images are in odd columns, and their reconstructions to their immediate right.	12
3.4	An example architecture of a fully-convolutional autoencoder. The latent space is constrained by reducing the number and/or size of the filters.	13
3.5	A collection of images from the MNIST data set and their respective reconstructions using the fully-convolutional autoencoder specified in Table 3.2. The original MNIST images are in odd columns, and their reconstructions to their immediate right.	14
3.6	A naïve implementation of the variational autoencoder. The input \mathbf{x} is mapped to intermediate layers taking the values of $\boldsymbol{\mu}$ and σ^2 . The latent variable \mathbf{z} is then sampled from the probabilistic encoder $\mathbf{z} \sim q_\phi(\mathbf{z} \mathbf{x})$. Finally \mathbf{z} is mapped back to the input dimension to give reconstruction $\tilde{\mathbf{x}}$	21
3.7	A viable implementation of the variational autoencoder. Sampling from the probabilistic encoder $q_\phi(\mathbf{z} \mathbf{x})$ is simulated by evaluating $\mathbf{z} = g_\phi(\mathbf{x}, \boldsymbol{\epsilon}) = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \boldsymbol{\epsilon}$	22
3.8	The encoder takes a data point and returns a normal distribution (orange); some samples of which are shown (blue). A sample is drawn from the normal distribution (red) and decoded. [12]	24
3.9	The prior distribution should approximate the standard multivariate Gaussian $\mathcal{N}(\mathbf{0}, \mathbf{I})$. Samples of the prior are shown (yellow); two of which are decoded (red) and approximate an isotropic Gaussian. [12]	26
3.10	The GAN over the high latent space distribution of all data points (in approx 3D) for the data set. Rows correspond to a data point and columns the value of the latent variable (varied from -1 to 1). Each section (a), (b), (c) and (d) consider a different latent variable. [5]	28

3.12 A comparison of InfoGAN, β -VAE ($\beta = 20$) and VAE on the 3D face data set. Different latent variables are varied for sections (a), (b) and (c). All models learnt lighting and elevation, but only InfoGAN and β -VAE managed to continuously vary the azimuth. [32]	30
3.13 The probabilistic encoder $q_\phi(\mathbf{z} \mathbf{x})$ maps a given data point \mathbf{x} to the unknown distribution $\hat{p}(\mathbf{z})$. The probabilistic decoder $p_\theta(\mathbf{x} \mathbf{z})$ is trained to map samples from $\hat{p}(\mathbf{z})$ back to $p(\mathbf{x})$, since its inputs are drawn from the probabilistic encoder $q_\phi(\mathbf{z} \mathbf{x})$. A sample from the prior $p(\mathbf{z})$ will not be mapped back by $p_\theta(\mathbf{x} \mathbf{z})$ to $p(\mathbf{x})$ exactly if $p(\mathbf{z}) \neq \hat{p}(\mathbf{z})$. [8]	32
3.14 Samples from a variational autoencoder trained on the CelebA data set after $t = 0, 1, 5$ and 10 steps of the generative procedure (3.45). The MCMC chain was initialised with a sample from the prior $\mathbf{z}_{t=0} \sim p(\mathbf{z})$, which often improves the quality of the samples. [8]	33
4.1 A collection of frames captured from Space Invaders emulated on Stella. Left column: an even frame. Middle column: the (odd) frame following. Right column: Combining the even and odd frames by taking the maximal value over each channel (RGB). Clearly the bullets visible in one frame fail to persist in the next. As mentioned, this is due to the limited number of sprites Atari 2600 can load in a single frame.	36
4.2 Pre-processed frames captured from Pong emulated on Stella. These frames were originally 84×84 , but are printed here as 168×168 to emphasise distortions. Left: The JPEG format distorts the ball, paddle and score sprites. Right: The PNG format displays the frame without such distortions.	37
5.1 Caption.	40
5.2 Caption.	41
5.3 Caption.	41
5.4 Caption.	42

5.5 Caption. 42

5.6 Caption. 44

Chapter 1

Introduction

1.1 Motivation

A long term goal of artificial intelligence (AI) is the development of artificial general intelligence (AGI). Since the field's inception in the 1950s, it has swung between hype and breakthroughs, followed by disappointment and reduced funding, known as AI winters [17]. During the first period of hype from the 50s to the early 70s, Marvin Minsky made the following prediction: [9]

“In from three to eight years we will have a machine with the general intelligence
of an average human being.” - Marvin Minsky, 1970

This prediction was clearly not realised, and the first AI winter would shortly follow.

Symbolic AI was developed during this winter, which encodes knowledge as human-readable rules and facts, making it easy to comprehend chains of actions and abstract relationships [23]. For instance, given the unary relations `red` and `strawberry`, and the binary relation `bigger`, we can say that `A` is the smallest red strawberry by writing

```
red(A)  strawberry(A)  ∀B bigger(B, A)
```

But given the unary relations `yellow` and `banana` we could also write that `A` is the third biggest yellow strawberry, or a red banana, and so on. We can see that the rules and facts in symbolic logic can be endlessly recombined and extended. This allows for the manipulation of high-level abstract concepts, which is key to AGI [13].

However, symbolic AI has a major philosophical problem: the facts and rules are only meaningful to the human writing them; their meaning is not intrinsic to the system itself. This is known as the *symbol grounding problem*.

Today we find ourselves in yet another period of hype and exciting breakthroughs not afflicted by the symbol grounding problem. Reinforcement learning (RL) has become a prominent area of research, with many considering it fundamental for AGI [15], as have deep neural networks. Recently, deep reinforcement learning (DRL) systems have achieved impressive feats, including mastering a wide range of Atari 2600 games to a superhuman level using only raw pixels and score as input, and the board game Go [21, 28].



Figure 1.1: May 1997: Gary Kasparov makes his first move against IBM’s Deep Blue. Deep Blue would later emerge the victor in the best of six games; the first time a reigning world chess champion is defeated by a computer. [26]



Figure 1.2: March 2016: Lee Sedol, one of the greatest modern Go players, plays his first move of game three against AlphaGo. AlphaGo won four of five games. This feat was considered by many to be a decade away. [22]

Though DRL systems are not afflicted by the same problems as symbolic AI, they have a number of drawbacks of their own. Namely, they are: [13]

1. **Slow to learn.** Neural networks require large data sets and are therefore slow to learn.

2. **Unable to transfer past experience.** They often fail to perform well on tasks very similar to those they have mastered.
3. **Unable to reason abstractly.** They fail to exploit statistical regularities in the data.
4. **Hard to reason about.** It's often difficult to extract a comprehensible chain of reasons for why a deep neural network operated in the way it did.

Deep symbolic reinforcement learning (DSRL) is a marrying of DRL and symbolic AI; a recent advance which overcomes the symbol grounding problem and the drawbacks associated with DRL [13]. That is, DSRL systems overcome the symbol grounding problem, and are:

1. **Fast to learn.** Large data sets are not necessary.
2. **Able to transfer past experience.** Symbolic AI lends itself to multiple processes associated with high-level reasoning, including transfer learning.
3. **Able to reason abstractly.** The agent is able to exploit statistical regularities in the training data by using high-level processes like planning or causal reasoning.
4. **Easy to reason about.** Since the front end uses symbolic AI, its knowledge is encoded as human-readable facts and rules, making the extraction of comprehensible chains of logic much easier.

An overview of DSRL is shown in Figure 1.3. The neural back end takes a high-dimensional input and outputs a symbolic representation. This symbolic representation is then fed to the symbolic front end, whose role is action selection. The agent then acts on the environment and obtains a reward and the sensory input of the next time step. As the neural back end learns how to represent the raw input data in a compositionally structured representation in an unsupervised manner, and the symbolic front end learns to select the action with maximum expected reward over time, the system as a whole learns end-to-end.

TODO: Finish description of DSRL

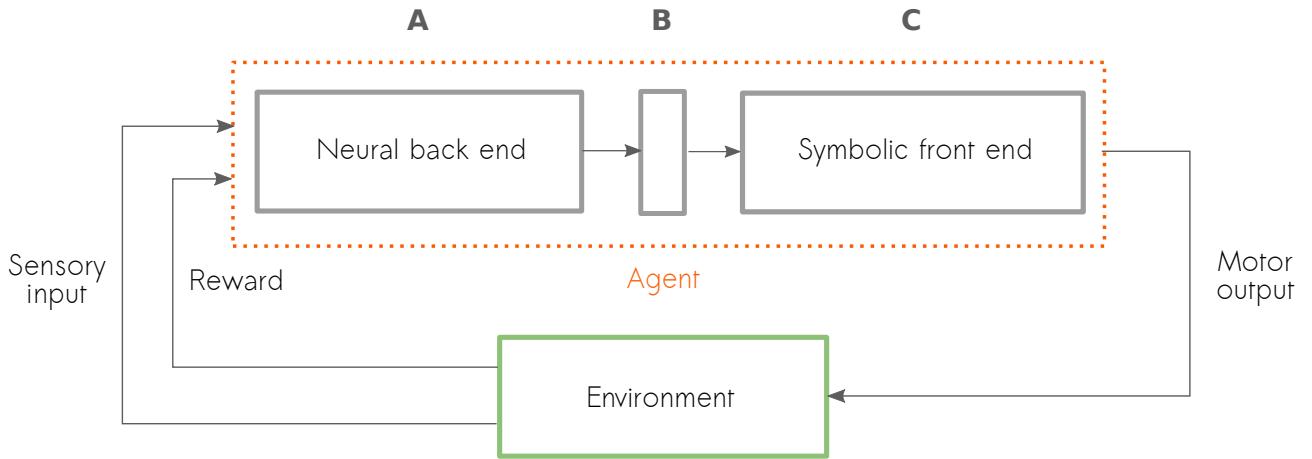


Figure 1.3: Overview of deep symbolic reinforcement learning system architecture. **A:** The neural back end maps high-dimensional raw input data to a compositionally structured symbolic representation. **B:** The compositionally structured symbolic representation. **C:** Reinforcement learning of mapping from symbolic representation to action with maximum expected reward over time. *Source: Garnelo et al. [13].*

1.2 Objectives

We'll use the image in Figure 1.4 as an example of a high-dimensional input to the neural back end. This world consists of only two shapes (circle and square) and four spaces occupied by at most one shape (top left, top right, bottom left and bottom right). The neural back end maps this raw high-dimensional input to a low-dimensional symbolic representation, shown in Table ??.

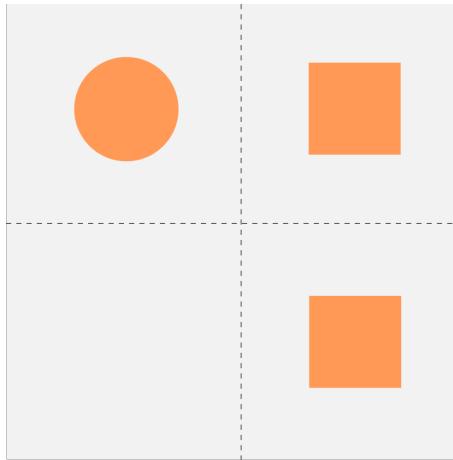


Figure 1.4: A toy example of a raw high-dimensional input.

How this is done will be explained in Chapter 2, but as for now we can just take it as fact that the current method doesn't scale. That is, for very simple scenes, as in Figure ..., This is done

Type	Location
1	[0, 0]
2	[0, 1]
2	[1, 1]

Table 1.1: Low-dimensional symbolic representation

by passing the high-dimensional input through a series of convolutional layers and extracting the activation spectra in the latent space. These spectra are then used to classify the

relies on the unsupervised extraction of disentangled features, allowing for transfer learning and high-level cognitive processes. However, the unsupervised extraction of features from a wide range of scenes is still a challenge in AI research [?]. Fortunately methods are getting better, and the first unsupervised scalable model β -VAE was developed recently.

TODO: Finish objectives

1.3 Contributions

Contributions here.

Chapter 2

Background

We will cover how deep symbolic reinforcement learning extracts symbolic representations from raw input data, which will motivate a discussion of loss functions and the introduction of autoencoders. Seeing the limitations of the current approach in extracting symbolic representations, we can appreciate the recent development of β -VAE, a variant of the variational autoencoder used to learn disentangled representations. Finally, we can conclude by mentioning less technical matters, such as libraries and hardware used.

2.1 Loss functions

The idea of image reconstruction plays a vital role throughout this project. Although it's possible to qualitatively compare the original to its reconstruction, it's important to be able to quantify the difference, which lends itself to automation. The loss function will quantify how similar two images are.

To compare loss functions, we'll use the MNIST data set. MNIST is a collection of 70,000 black-and-white images of handwritten digits, with 60,000 in the training set containing and 10,000 in the test set. These images will be represented as vectors without loss of generality.

2.1.1 Euclidean Distance

The Euclidean distance between two vectors \mathbf{x} and \mathbf{y} is defined by

$$\sqrt{\sum_i (x_i - y_i)^2}$$

where x_i and y_i are the i^{th} components of \mathbf{x} and \mathbf{y} respectively.

Euclidean distance is an intuitive measure of the distance between two points in space. Unfortunately, this doesn't also translate to visual similarity, as illustrated by Doersch et al. [11]. Figure 2.1 is a digit drawn from the MNIST dataset, and Figures 2.3 and 2.2 are attempted reconstructions. Of the reconstructions, Figure 2.2 looks most like the original, but Figure 2.3 is closer in space.

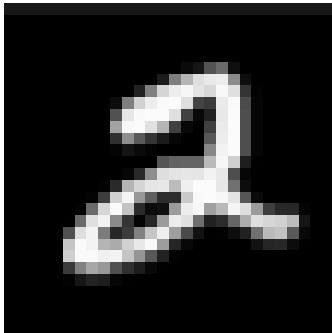


Figure 2.1

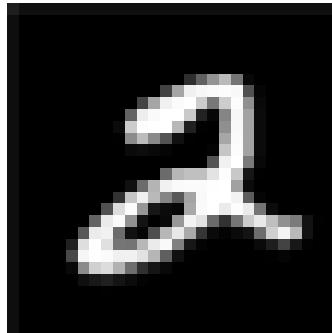


Figure 2.2

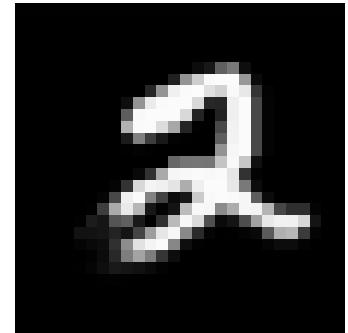


Figure 2.3

This leads to an alternative measure, binary cross-entropy, which gives a much better quantification of how visually similar two images are.

2.1.2 Binary Cross-Entropy

Consider a single black-and-white pixel with probability $p(0) = c$ of being 0 and $p(1) = 1 - c$ of being 1. Here $p(x)$ is a probability distribution over the possible pixel values $x \in \{0, 1\}$. Suppose a given model tries to learn the distribution described by $p(x)$, and says that the pixel has probability $q(0) = \hat{c}$ of being 0 and $q(1) = 1 - \hat{c}$ of being 1. The model is perfect if it learns

the true distribution, that is, if $q(x) = p(x)$ for $x \in \{0, 1\}$. We'd like to quantify how similar the distributions p and q are.

This is done by computing the binary cross-entropy between p and q , which is defined by

$$H(p, q) = -c \log \hat{c} - (1 - c) \log(1 - \hat{c})$$

To see how we may use this as a similarity measure among images, consider a 1×1 image. Normalising this image yields a pixel value in the interval $[0, 1]$, which may now be interpreted as a probability, corresponding to c above. In the normalised reconstructed image, the pixel value corresponds to \hat{c} . We simply compute the binary cross-entropy to measure the similarity of these two distributions, and in turn, the similarity of the images themselves! (Note: we could have also assigned the probabilities to $1 - y$ and $1 - \hat{y}$ by symmetry of binary cross-entropy).

For images larger than 1×1 , we may take the component-wise binary cross-entropy, then, for example, average the components. How the component-wise binary cross-entropies are suitably combined to give a single floating point number will vary from problem to problem.

2.2 Stella

Video games are becoming increasingly popular in the generation of data sets for machine learning purposes. One reason for this is that they do not inherit many operational drawbacks from real world data, such as noise, the stability of the camera or the observation of rare circumstances (the video game may just be queried for rare cases) [34, 24].

2.3 Arcade Learning Environment

[1]

2.4 Keras

Keras is a high-level deep learning library which interfaces with a Tensorflow or Theano backend.

It provides a simple API which contains [6]

2.5 TensorBoard

Chapter 3

Related Work

3.1 Autoencoders

An autoencoder is a neural network that learns a compression algorithm for its input data in an unsupervised manner [20]. This is achieved by placing constraints on a hidden layer, called the latent space, and setting the target values to the input values, effectively learning the identity function. Since the network is trying to reconstruct the original input from the constrained latent space, over time the latent space corresponds to a meaningful compression of the network's input.

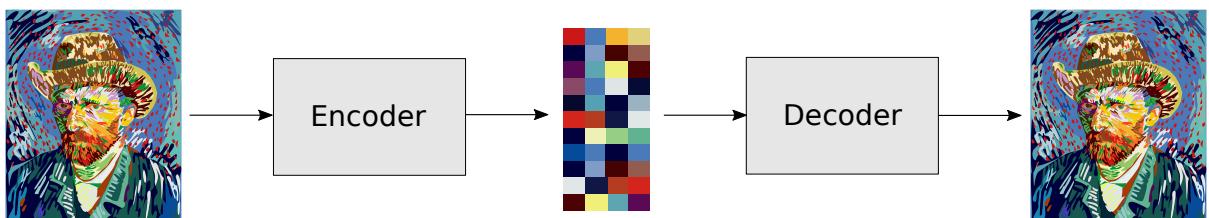


Figure 3.1: A black-box description of an autoencoder. The autoencoder learns the identity function, and in turn, the encoder and decoder learn suitable encoding and decoding algorithms respectively.

As before, we will use the MNIST data set to compare architectures. Unless specified in the example, the Adam optimiser is used with a learning rate of $1e - 4$, the batch size is 1 and the loss function is binary cross-entropy. Intermediate layers use the ReLU activation function, while the final layer uses sigmoid.

3.1.1 Fully-Connected Autoencoders

In dense feed-forward neural networks we may place a constraint on the latent space by reducing the number of neurons, as shown in Figure 3.2. Images must be flattened into vectors to be fed as input. Consequently, any spatial information is destroyed in dense feed-forward neural networks.

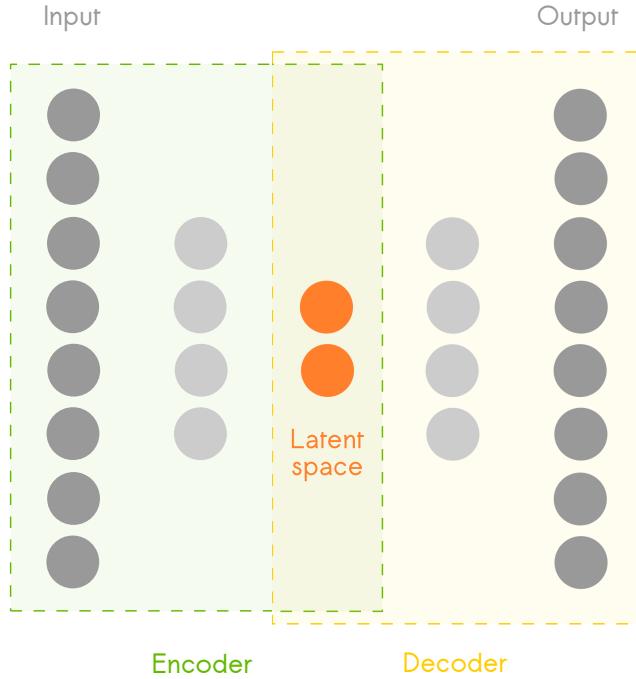


Figure 3.2: An example architecture of a fully-connected autoencoder. The latent space is constrained by having fewer neurons than the input and output layers.

An example architecture is given in Table 3.1, which was trained on MNIST. Despite the latent space being $\sim 4\%$ of the size of the input space, the network is capable of producing realistic reconstructions. For verification, a collection of samples from the dataset and their corresponding reconstructions are shown in Figure 3.3.

Layer Type	Output Shape
InputLayer	(1, 28, 28)
Flatten	(784,)
Dense	(32,)
Dense	(784,)
Reshape	(1, 28, 28)

Table 3.1: A simple fully-connected autoencoder with one hidden layer. After 15 epochs, the validation score was recorded to be 71.94.



Figure 3.3: A collection of images from the MNIST data set and their respective reconstructions using the fully-connected autoencoder specified in Table 3.1. The original MNIST images are in odd columns, and their reconstructions to their immediate right.

3.1.2 Fully-Convolutional Autoencoders

In fully-convolutional feed-forward neural networks, we may place a constraint on the latent space by reducing the number and/or size of the filters, as shown in Figure 3.4. To compare the fully-convolutional autoencoder to the fully-connected, we'll train the architecture in Table 3.2 on MNIST. As before, we'll compare the reconstructions to the originals, which can be found in Figure 3.5.

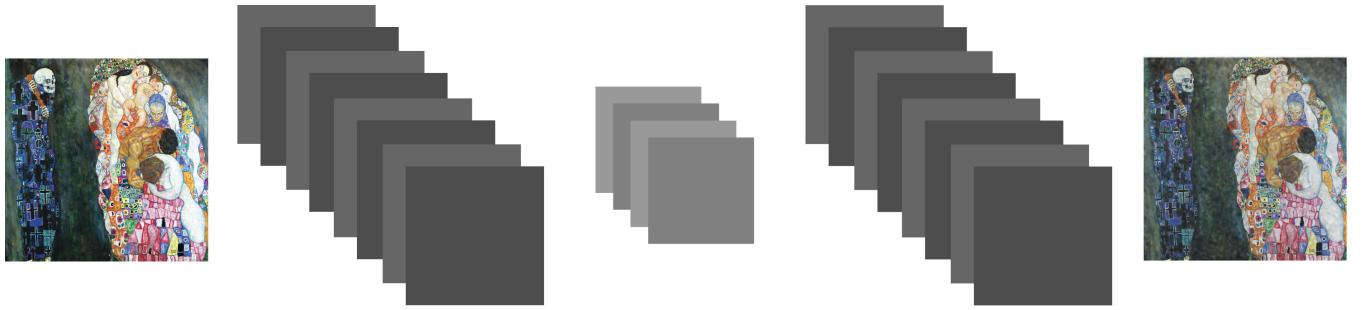


Figure 3.4: An example architecture of a fully-convolutional autoencoder. The latent space is constrained by reducing the number and/or size of the filters.

Layer Type	Output Shape
InputLayer	(1, 28, 28)
Conv2D	(32, 28, 28)
MaxPooling2D	(32, 14, 14)
Conv2D	(4, 14, 14)
MaxPooling2D	(4, 7, 7)
UpSampling2D	(4, 14, 14)
Conv2DTranspose	(32, 14, 14)
UpSampling2D	(32, 28, 28)
Conv2DTranspose	(1, 28, 28)

Table 3.2: A simple fully-convolutional autoencoder with 2D convolutions and max pooling, plus the corresponding deconvolutional layers. After 15 epochs, the validation score was recorded to be 64.89.

Convolutional layers have been shown to be effective in tasks with images as input [18, 33, 30]. This is because spatial information is preserved in convolutional layers, and the number of trainable parameters is far less in a convolutional layer than it is in a fully connected layer. Convolutional layers will be used from here on as we'll be using images as input.



Figure 3.5: A collection of images from the MNIST data set and their respective reconstructions using the fully-convolutional autoencoder specified in Table 3.2. The original MNIST images are in odd columns, and their reconstructions to their immediate right.

3.2 Variational Autoencoders

[16]

The variational autoencoder is central to this project, and we'll therefore dedicate a considerable amount of time exploring it. First we'll precisely define the problems the variational autoencoder solves, after which we may develop its loss function and detail its implementation. This section will conclude by developing an intuition of the theory covered by way of examples.

3.2.1 A Probabilistic Perspective

We'll begin by making necessary definitions and describe the input data as samples from a generative process.. This sets the context to detail the problems the variational autoencoder solves.

Let $X = \{\mathbf{x}^{(i)}\}_{i=1}^N$ be the data set of N independent and identically distributed samples of the variable \mathbf{x} . (X may be a data set of images, for instance). Let us assume that these samples are generated by a random process with parameters θ^* involving an unobserved latent variable \mathbf{z} in the following way:

Algorithm 1 Generate data set X

```

1: for  $i = 1 \rightarrow N$  do
2:    $\mathbf{z}^{(i)} \sim p_{\theta^*}(\mathbf{z})$                                 // Sample from true prior
3:    $\mathbf{x}^{(i)} \sim p_{\theta^*}(\mathbf{x}|\mathbf{z}^{(i)})$                 // Sample from true conditional
4:   Append  $\mathbf{x}^{(i)}$  to  $X$ 
5: end for

```

We only observe the data set X in this process. The parameters θ^* and latent variables $Z = \{\mathbf{z}^{(i)}\}_{i=1}^N$ are unknown to us. Let us assume that the prior $p_{\theta^*}(\mathbf{z})$ and conditional $p_{\theta^*}(\mathbf{x}|\mathbf{z})$ are parameterised by the distributions $p_{\theta}(\mathbf{z})$ and $p_{\theta}(\mathbf{x}|\mathbf{z})$ respectively. In this context, the variational autoencoder provides [16]:

1. ML or MAP estimation for the parameters θ

2. An approximation of the latent variable $\mathbf{z}^{(i)}$ given $\mathbf{x}^{(i)}$ and set of parameters $\boldsymbol{\theta}$
3. Approximate marginal inference of the variable \mathbf{x}

3.2.2 Overcoming the Intractable Posterior

The variational autoencoder solves the problems above by approximate inference of the latent variable \mathbf{z} . Exact inference is not possible, and to see this we may use Bayes' theorem to find an expression for the posterior:

$$p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}) = \frac{p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})p_{\boldsymbol{\theta}}(\mathbf{z})}{p_{\boldsymbol{\theta}}(\mathbf{x})} \quad (3.1)$$

The marginal likelihood

$$p_{\boldsymbol{\theta}}(\mathbf{x}) = \int p_{\boldsymbol{\theta}}(\mathbf{z})p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})d\mathbf{z} \quad (3.2)$$

involves an exponential-time integral over every combination of the latent variable \mathbf{z} , and is therefore computationally intractable [16]. Instead, we define an approximation $q_{\phi}(\mathbf{z}|\mathbf{x})$ to the intractable posterior. Since $q_{\phi}(\mathbf{z}|\mathbf{x})$ gives a distribution over the possible latent variables \mathbf{z} that generated the given data point \mathbf{x} , it is known as the probabilistic decoder. By the same reasoning, $p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})$ is known as the probabilistic encoder.

3.2.3 Finding a Suitable Loss Function: the ELBO

The variational autoencoder's ability to learn the generative parameters $\boldsymbol{\theta}^*$ relies on how closely $q_{\phi}(\mathbf{z}|\mathbf{x})$ approximates $p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})$. In the interest of training a model, this difference will be quantified. For this we use the KL divergence

$$D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})) = \mathbf{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log_e \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})} \right] \quad (3.3)$$

which measures how much information is lost when we represent $p_{\theta}(\mathbf{z}|\mathbf{x})$ with $q_{\phi}(\mathbf{z}|\mathbf{x})$ (measured in nats) [4]. Using the KL divergence, our problem now amounts to the optimisation problem [19]:

$$\boldsymbol{\theta}^*, \boldsymbol{\phi}^* = \arg \min_{\boldsymbol{\theta}^*, \boldsymbol{\phi}^*} D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}|\mathbf{x})) \quad (3.4)$$

To see how we can start to minimise the KL divergence, we'll start by rewriting it in a different form:

$$D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}|\mathbf{x})) = \mathbf{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log_e \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{p_{\theta}(\mathbf{z}|\mathbf{x})} \right] \quad (3.5)$$

$$= \mathbf{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log_e q_{\phi}(\mathbf{z}|\mathbf{x}) \right] - \mathbf{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log_e p_{\theta}(\mathbf{z}|\mathbf{x}) \right] \quad (3.6)$$

$$= \mathbf{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log_e q_{\phi}(\mathbf{z}|\mathbf{x}) \right] - \mathbf{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log_e \frac{p_{\theta}(\mathbf{z}, \mathbf{x})}{p_{\theta}(\mathbf{x})} \right] \quad (3.7)$$

$$= \mathbf{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log_e \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{p_{\theta}(\mathbf{z}, \mathbf{x})} \right] + \mathbf{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log_e p_{\theta}(\mathbf{x})] \quad (3.8)$$

$$= \mathbf{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log_e \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{p_{\theta}(\mathbf{z}, \mathbf{x})} \right] + \log_e p_{\theta}(\mathbf{x}) \quad (3.9)$$

Here we see that the KL divergence depends on the intractable marginal likelihood $p_{\theta}(\mathbf{x})$! There's no way we can minimise it if we can't write down $p_{\theta}(\mathbf{x})$. However, we can get around this: we'll minimise the KL divergence, but not directly. Instead, we try to find a quantity which we can maximise, and show that in turn this minimises the KL divergence. The trick is not obvious, but is simply done by finding a lower bound on the log marginal likelihood.

Using Jensen's inequality

$$f(\mathbf{E}[X]) \geq \mathbf{E}[f(X)] \quad (3.10)$$

we can write down a lower bound on the log marginal likelihood:

$$\log_e p_{\theta}(\mathbf{x}) = \log_e \int p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z} \quad (3.11)$$

$$= \log_e \int p_{\theta}(\mathbf{x}, \mathbf{z}) \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{q_{\phi}(\mathbf{z}|\mathbf{x})} d\mathbf{z} \quad (3.12)$$

$$= \log_e \mathbf{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] \quad (3.13)$$

$$\geq \mathbf{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log_e \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] \quad (3.14)$$

$$:= \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}) \quad (3.15)$$

Expression (3.14) is called the *ELBO* (short for expected lower bound) [2, 16].

How does the *ELBO* help us with minimising the KL divergence? First recall the alternative form of the KL divergence:

$$D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}|\mathbf{x})) = \mathbf{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log_e \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{p_{\theta}(\mathbf{z}, \mathbf{x})} \right] + \log_e p_{\theta}(\mathbf{x}) \quad (3.9)$$

Writing this in terms of the *ELBO* we have:

$$D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}|\mathbf{x})) = -\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}) + \log_e p_{\theta}(\mathbf{x}) \quad (3.16)$$

Since the KL divergence is the negative of the *ELBO* up to an additive constant (with respect to q), minimising the KL divergence is equivalent to maximising the *ELBO* [3]. Now we may make a revision to our original optimisaiton problem (3.4). Our problem is written as [19]:

$$\boldsymbol{\theta}^*, \boldsymbol{\phi}^* = \arg \max_{\boldsymbol{\theta}^*, \boldsymbol{\phi}^*} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}) \quad (3.17)$$

3.2.4 Writing the ELBO in Closed-Form

We've found that we can maximise the *ELBO* to minimise the KL divergence between the approximation $q_{\phi}(\mathbf{z}|\mathbf{x})$ and the intractable posterior $p_{\theta}(\mathbf{z}|\mathbf{x})$. Now we will seek to write the *ELBO* in closed-form, after which we will be able to implement the variational autoencoder.

To write the *ELBO* in closed-form, we'll start with a useful manipulation:

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}) = \mathbf{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})} \left[\log_e \frac{p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})}{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})} \right] \quad (3.18)$$

$$= -\mathbf{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})} \left[\log_e \frac{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})}{p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})} \right] \quad (3.19)$$

$$= -\mathbf{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})} \left[\log_e \frac{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})}{p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})p_{\boldsymbol{\theta}}(\mathbf{z})} \right] \quad (3.20)$$

$$= -\mathbf{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})} \left[\log_e \frac{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})}{p_{\boldsymbol{\theta}}(\mathbf{z})} \right] + \mathbf{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})} [\log_e p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})] \quad (3.21)$$

$$= -D_{KL}(q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})||p_{\boldsymbol{\theta}}(\mathbf{z})) + \mathbf{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})} [\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})] \quad (3.22)$$

Thus for a single data point $\mathbf{x}^{(i)}$, the *ELBO* becomes:

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}^{(i)}) = -D_{KL}(q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}^{(i)})||p_{\boldsymbol{\theta}}(\mathbf{z})) + \mathbf{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}^{(i)})} [\log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}|\mathbf{z})] \quad (3.23)$$

With some assumptions and a bit of work, we're finally in a position to write the *ELBO* in closed-form. The first term is the KL divergence between the probabilistic encoder and the prior. To make our lives simpler, we can choose the prior to be the isotropic multivariate Gaussian:

$$p_{\boldsymbol{\theta}}(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (3.24)$$

We will also assume that the posterior is a k -dimensional Gaussian with diagonal covariance. It follows that the approximate posterior should take the same form. That is,

$$q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}^{(i)}) = \mathcal{N}(\boldsymbol{\mu}^{(i)}, \boldsymbol{\sigma}^{2(i)}\mathbf{I}) \quad (3.25)$$

where $\boldsymbol{\mu}^{(i)}$ and $\boldsymbol{\sigma}^{(i)}$ are the mean and standard deviation (respectively) for data point $\mathbf{x}^{(i)}$.

With these two assumptions, and the KL divergence for k -dimensional Gaussians,

$$D_{KL}(\mathcal{N}_0||\mathcal{N}_1) = \frac{1}{2} \left[\text{tr}(\boldsymbol{\Sigma}_1^{-1}\boldsymbol{\Sigma}_0) + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}_1^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) - k + \ln \left(\frac{\det \boldsymbol{\Sigma}_1}{\det \boldsymbol{\Sigma}_0} \right) \right] \quad (3.26)$$

we may write the first term in closed-form:

$$D_{KL}(q_\phi(\mathbf{z}|\mathbf{x}^{(i)})||p_\theta(\mathbf{z})) = \frac{1}{2} \left[\sum_{j=1}^k \sigma_j^{2(i)} + \sum_{j=1}^k \mu_j^{2(i)} - k - \ln \prod_{j=1}^k \sigma_j^{2(i)} \right] \quad (3.27)$$

Now we turn our attention to the second term of the *ELBO*:

$$\mathbf{E}_{q_\phi(\mathbf{z}|\mathbf{x}^{(i)})} [\log p_\theta(\mathbf{x}^{(i)}|\mathbf{z})] \quad (3.28)$$

Luckily, maximising terms like this is encountered regularly in statistics, and is known as maximum likelihood estimation. This can be taken to be the reconstruction loss (defined it earlier) [19].

Therefore, using an unspecified reconstruction loss, the *ELBO* is:

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}) = \frac{1}{2} \left[\sum_{j=1}^k \sigma_j^{2(i)} + \sum_{j=1}^k \mu_j^{2(i)} - k - \ln \prod_{j=1}^k \sigma_j^{2(i)} \right] + \mathbf{E}_{q_\phi(\mathbf{z}|\mathbf{x}^{(i)})} [\log p_\theta(\mathbf{x}^{(i)}|\mathbf{z})] \quad (3.29)$$

3.2.5 Implementing the Variational Autoencoder

The variational autoencoder uses a neural network for the probabilistic encoder $q_\phi(\mathbf{z}|\mathbf{x})$ [16]. First a naïve implementation of the variational autoencoder will be proposed. As we shall see, this implementation needs one adjustment, known as the “reparameterisation trick”. This will lead to the final implementation of the variational autoencoder, and conclude our formal study.

The probabilistic encoder

$$q_\phi(\mathbf{z}|\mathbf{x}^{(i)}) = \mathcal{N}(\boldsymbol{\mu}^{(i)}, \boldsymbol{\sigma}^{2(i)} \mathbf{I}) \quad (3.25)$$

requires the mean $\boldsymbol{\mu}^{(i)}$ and standard deviation $\boldsymbol{\sigma}^{(i)}$ for a given data point $\mathbf{x}^{(i)}$. These two vectors will be represented by distinct layers in the latent space of a neural network, as shown in Figure (3.6) [19].

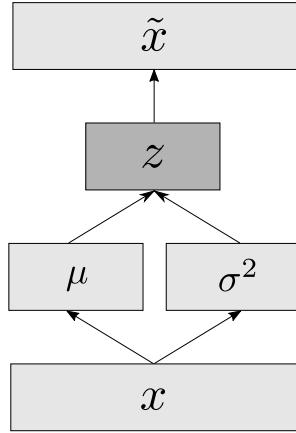


Figure 3.6: A naïve implementation of the variational autoencoder. The input \mathbf{x} is mapped to intermediate layers taking the values of $\boldsymbol{\mu}$ and σ^2 . The latent variable \mathbf{z} is then sampled from the probabilistic encoder $\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})$. Finally \mathbf{z} is mapped back to the input dimension to give reconstruction $\tilde{\mathbf{x}}$.

Unfortunately we cannot perform backpropagation, as it makes no sense to differentiate the the random variable \mathbf{z} wrt ϕ (\mathbf{z} is drawn from a distribution, and not a function of ϕ). To solve this, we introduce an auxillary variable ϵ and vector-valued function $g_\phi(\mathbf{x}, \epsilon)$ parameterised by ϕ [16]. The auxillary variable ϵ is governed by a parameterless distribution, which we will take to be the multivariate isotropic Gaussian.

The sampling step can now be written as

$$\mathbf{z} = g_\phi(\mathbf{x}, \epsilon) = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \epsilon \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (3.30)$$

where \odot represents the element-wise product. This is a reparameterisation of the random variable $\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})$ with a differentiable function $g_\phi(\mathbf{x}, \epsilon)$, and is suitable known as the “reparameterisation trick” [16]. The reparameterisation trick allows backpropagation to be used, and thus completes our method of solving optimisation problem (3.17).

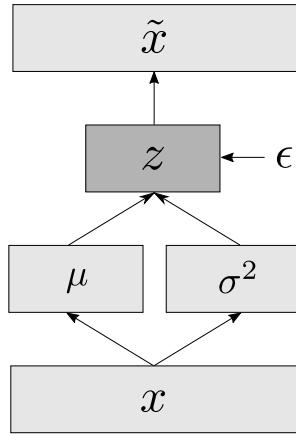


Figure 3.7: A viable implementation of the variational autoencoder. Sampling from the probabilistic encoder $q_\phi(\mathbf{z}|\mathbf{x})$ is simulated by evaluating $\mathbf{z} = g_\phi(\mathbf{x}, \epsilon) = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \epsilon$.

3.2.6 Intuition Behind the Variational Autoencoder

In this section, we'll develop the intuition behind the variational autoencoder. We've been considering the data set $X = \{\mathbf{x}^{(i)}\}_{i=1}^N$. Now we'll choose X to be the MNIST data set (the vector $\mathbf{x}^{(i)}$ now corresponds to a flattened MNIST image of length $28 \times 28 = 784$). For visualisation purposes, we'll also take the latent space dimension to be $k = 2$, that is, $\mathbf{z} = (z_1, z_2)$.

The probabilistic encoder

$$q_\phi(\mathbf{z}|\mathbf{x}^{(i)}) = \mathcal{N}(\boldsymbol{\mu}^{(i)}, \boldsymbol{\sigma}^{2(i)} \mathbf{I}) \quad (3.25)$$

gives a normal distribution over the values of the latent variable $\mathbf{z}^{(i)}$ given data point $\mathbf{x}^{(i)}$. A sample is drawn from this normal distribution and passed through to the decoder. The decoder then reconstructs the original image from the latent representation. Since we've chosen our latent space to be two-dimensional, it's possible to visualise this process, which is done in Figure (3.8). It's then possible to compare the reconstruction to the original (using an appropriate reconstruction loss function), and therefore to train the autoencoder as a whole.

Recall that we chose the prior to be the standard multivariate Gaussian

$$p_{\theta}(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (3.24)$$

We also derived the loss function to be

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}^{(i)}) = -D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})||p_{\theta}(\mathbf{z})) + \mathbf{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})} [\log p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z})] \quad (3.23)$$

which we want to maximise. Since

$$D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})||p_{\theta}(\mathbf{z})) \geq 0 \quad (3.31)$$

maximising the *ELBO* is equivalent to reducing the KL divergence between the probabilistic encoder $q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})$ and the prior $p_{\theta}(\mathbf{z})$. Therefore after training, we expect that the probabilistic encoder should map samples to the standard multivariate Gaussian. (In practice, $\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}^{(i)}) > 0 \quad \forall i$, so we therefore only expect it to be *approximately* mapped to the standard Gaussian). Decoding samples from the prior should correspond to meaningful reconstructions.

TODO: FINISH

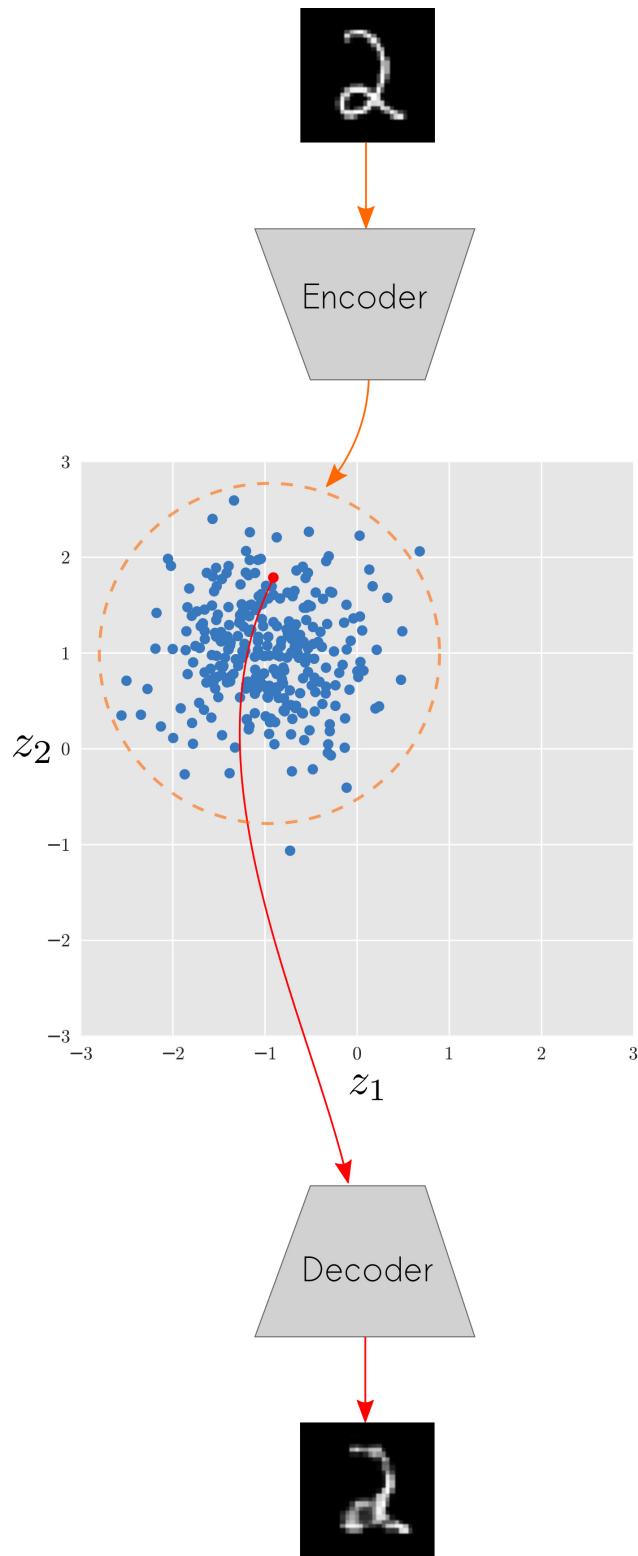


Figure 3.8: The encoder takes a data point and returns a normal distribution (orange); some samples of which are shown (blue). A sample is drawn from the normal distribution (red) and decoded. [12]

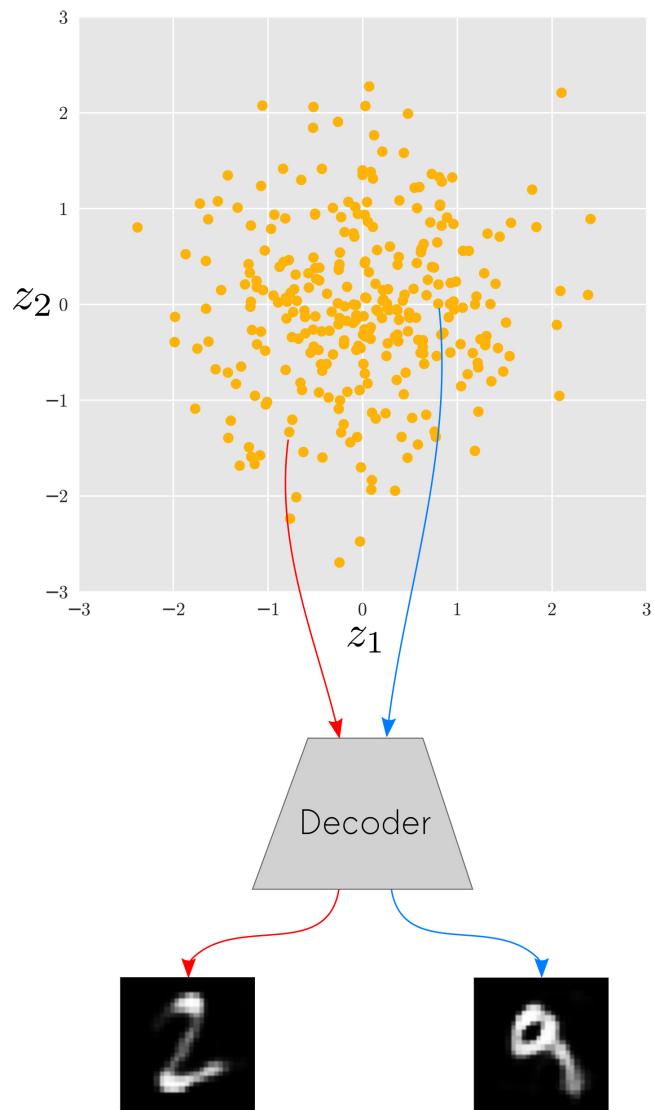


Figure 3.9: The prior distribution should approximate the standard multivariate Gaussian $\mathcal{N}(\mathbf{0}, \mathbf{I})$. Samples of the prior are shown (yellow); two of which are decoded (red and blue).

[12]

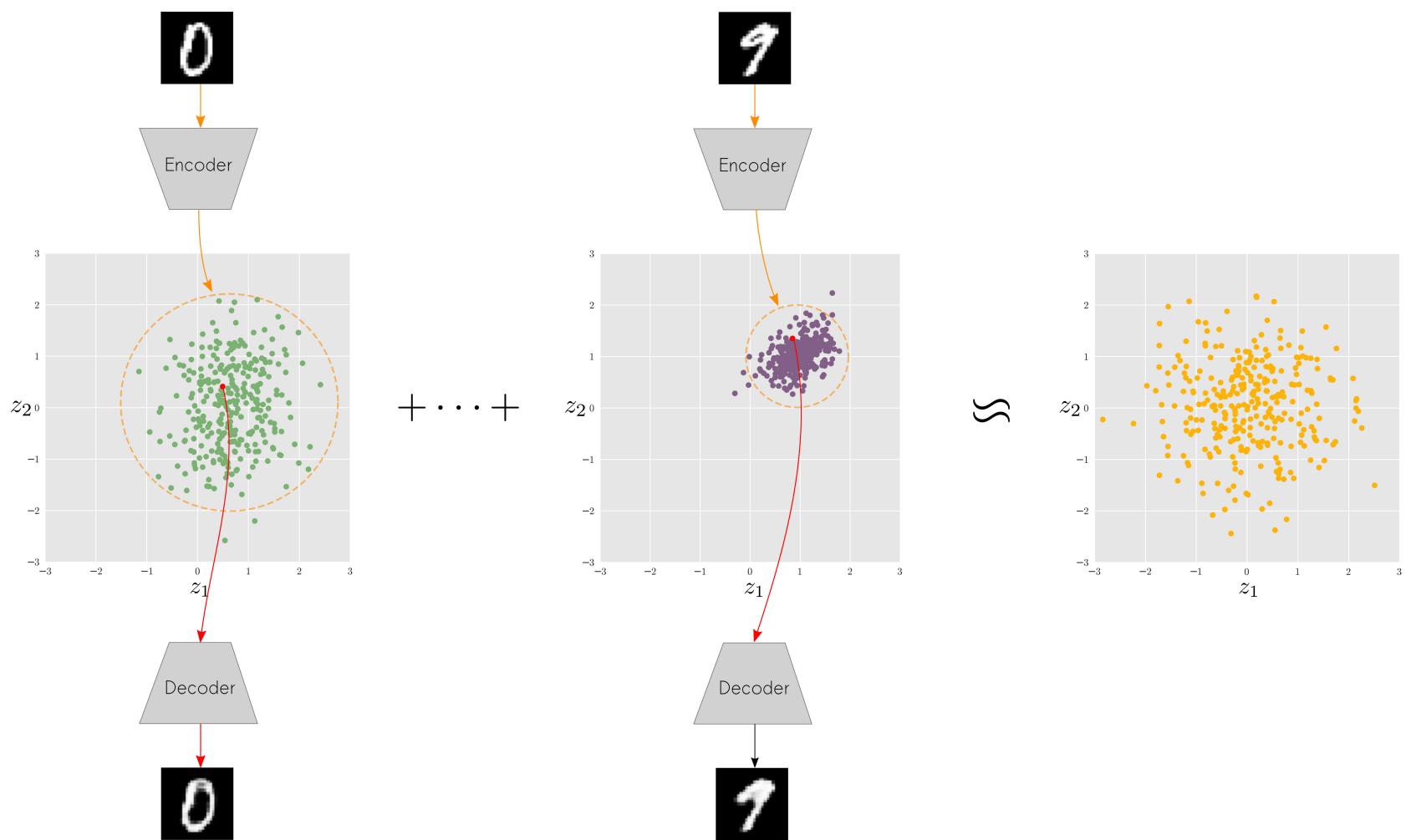


Figure 3.10: The sum over the latent space distributions of all data points $\mathbf{x}^{(i)}$ approximates the multivariate isotropic Gaussian. [12]

3.3 Disentangled Representations

[14, 32]

3.4 Unsupervised Learning of Generative Factors

Learning disentangled generative factors of a scene in an unsupervised manner is an open challenge in AI research. Although many attempts have been made, they have not scaled well [32, 27, 10, 31, 7]. However, two recent advancements have made significant headway: β -VAE and InfoGAN [5, 32].

3.4.1 InfoGAN

InfoGAN is an information theoretic extension of the Generative Adversarial Network, which has convincingly learnt generative factors of multiple data sets in an unsupervised manner, including MNIST, 3D faces and 3D chairs [5].

However, InfoGAN is sensitive to the choice of the prior distribution, and therefore requires a priori knowledge of the data set, as well as the number and type of generative factors [32]. Ideally, the number of generative factors would be inferred and not made explicit by the designer. On this basis, we can conclude that future work is needed.

3.4.2 β -VAE

β -VAE is the first method to overcome what InfoGAN could not: to learn the (unspecified amount of) generative factors of a data set in an unsupervised manner. First a formal derivation of β -VAE will be proposed, then a number of experimental results.

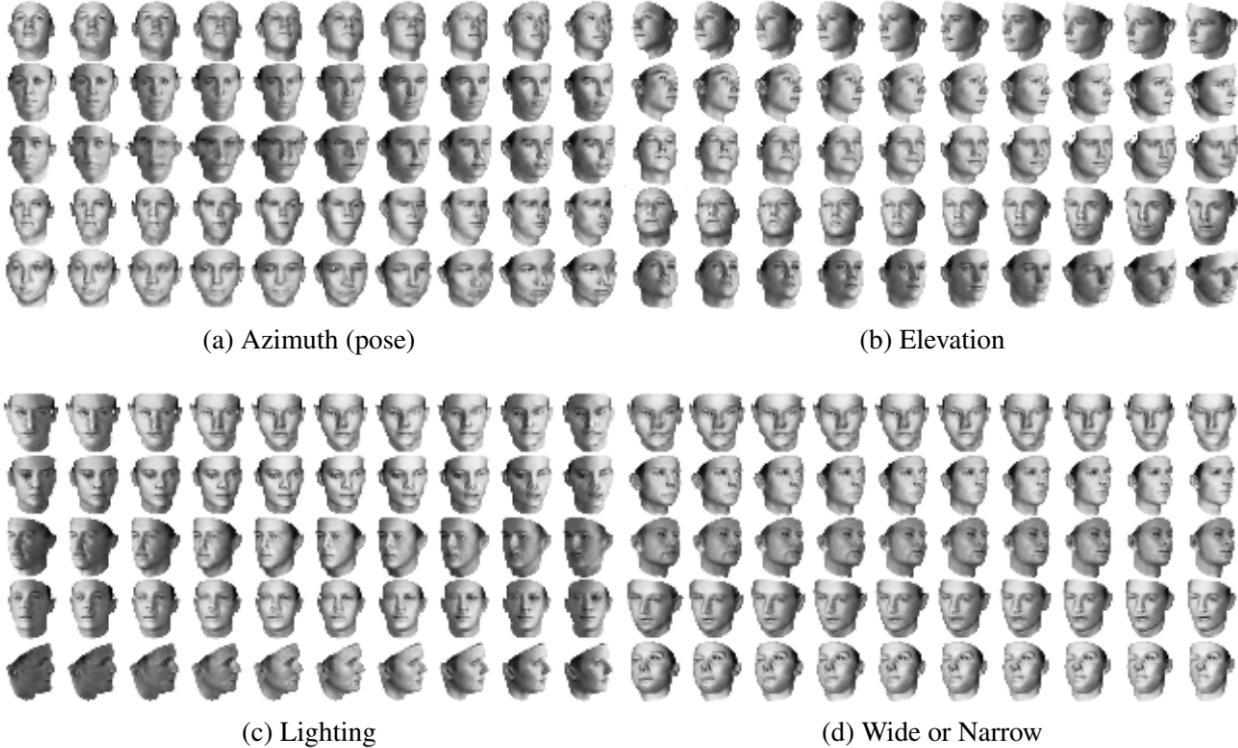


Figure 3.11: InfoGAN convincingly learns the underlying generative factors in the 3D face data set. Rows correspond to a data point and columns the value of the latent variable (varied from -1 to 1). Each section (a), (b), (c) and (d) consider a different latent variable. [5]

Derivation

Recall that for the variational autoencoder, we assumed the data point \mathbf{x} was generated by a process involving a latent variable \mathbf{z} and conditional $p(\mathbf{x}|\mathbf{z})$. For β -VAE we make a slightly different assumption: data point $\mathbf{x} \in \mathbb{R}^N$ was generated by conditionally independent factors $\mathbf{v} \in \mathbb{R}^K$ s.t. $p(\mathbf{v}|\mathbf{x}) = \prod_k p(v_k|\mathbf{x})$, conditionally dependent factors $\mathbf{w} \in \mathbb{R}^H$ and the conditional $p(\mathbf{x}|\mathbf{v}, \mathbf{w})$. In short, this is written as

$$p(\mathbf{x}|\mathbf{v}, \mathbf{w}) = Sim(\mathbf{v}, \mathbf{w}) \quad (3.32)$$

where Sim is short for a true world simulator. β -VAE seeks to represent these generative factors in the latent variable $\mathbf{z} \in \mathbb{R}^M$, such that

$$p(\mathbf{x}|\mathbf{z}) \approx p(\mathbf{x}|\mathbf{v}, \mathbf{w}) = Sim(\mathbf{v}, \mathbf{w}) \quad (3.33)$$

Note that, unlike InfoGAN, β -VAE does not specify the number independent generative factors K ; instead it is inferred. However, in order for $\mathbf{z} \in \mathbb{R}^M$ to learn the conditionally independent factors $\mathbf{v} \in \mathbb{R}^K$, we must assert that $M \geq K$.

As before, the posterior $p_{\theta}(\mathbf{z}|\mathbf{x})$ is intractable, which motivates the introduction of the probabilistic encoder $q_{\phi}(\mathbf{z}|\mathbf{x})$. A reasonable objective is to find the most likely parameters of the model over all latent variables that produced the observed data. This is summarised more precisely as the optimisation problem

$$\boldsymbol{\theta}^*, \boldsymbol{\phi}^* = \max_{\boldsymbol{\phi}, \boldsymbol{\theta}} \mathbf{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] \quad (3.34)$$

A constraint is added to control the capacity of information in the latent space and encourage the factors \mathbf{v} to be learnt in a disentangled manner. The probabilistic encoder $q_{\phi}(\mathbf{z}|\mathbf{x})$ is pressured to be close to an isotropic Gaussian, implicitly applying an independence pressure. By using the KL divergence, we form the optimisation problem

$$\begin{aligned} \boldsymbol{\theta}^*, \boldsymbol{\phi}^* &= \max_{\boldsymbol{\phi}, \boldsymbol{\theta}} \mathbf{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] \\ \text{s.t. } D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) &< \epsilon \end{aligned} \quad (3.35)$$

The Lagrangian of optimisation problem (3.35) is

$$\mathcal{F}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}) = \mathbf{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \beta(D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) - \epsilon) \quad (3.36)$$

Since $\beta, \epsilon > 0$,

$$\mathcal{F}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}) \geq \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}) = \mathbf{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \beta D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) \quad (3.37)$$

By applying pressure for $q_{\phi}(\mathbf{z}|\mathbf{x})$ to be close to the prior $p(\mathbf{z})$, the latent space learns the conditionally independent factors \mathbf{v} in a different subset of \mathbf{z} than the conditionally dependent factors \mathbf{w} . β controls the degree of this pressure. Also note that $\beta = 1$ yields the *ELBO* formulated earlier (3.22), and $\beta = 0$ yields the maximum likelihood case. It's hypothesised

that disentangled representations are learnt for $\beta > 1$.

As the prior needs to be isotropic, it's often chosen to be the standard multivariate Gaussian:

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (3.38)$$

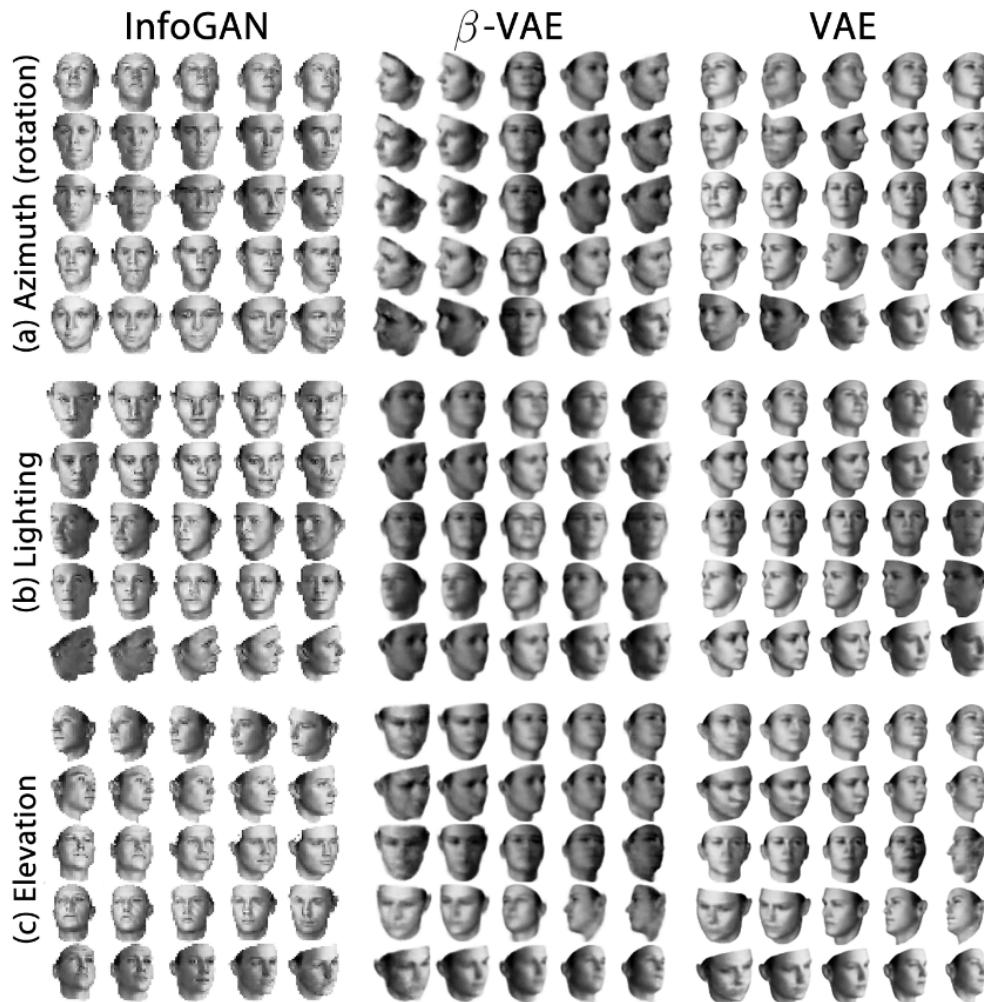


Figure 3.12: A comparison of InfoGAN, β -VAE ($\beta = 20$) and VAE on the 3D face data set. Different latent variables are varied for sections (a), (b) and (c). All models learnt lighting and elevation, but only InfoGAN and β -VAE managed to continuously vary the azimuth. [32]

3.5 Improving Sampling from Generative Autoencoders with Markov Chains

A generative autoencoder may be defined as an autoencoder that pressures its latent distribution $q_\phi(\mathbf{z}|\mathbf{x})$ to match a given prior $p(\mathbf{z})$ [8]. As discussed earlier, the KL divergence term in the *ELBO* applies this pressure:

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}) = -D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p_{\boldsymbol{\theta}}(\mathbf{z})) + \mathbf{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})] \quad (3.23)$$

In Section (3.2) it was assumed that the data set X was generated by

$$\mathbf{z}^{(i)} \sim p(\mathbf{z}) \quad \mathbf{x}^{(i)} \sim p_{\boldsymbol{\theta}^*}(\mathbf{x}|\mathbf{z}^{(i)}) \quad (3.39)$$

where $p(\mathbf{z})$ is a parameterless prior distribution and $p_{\boldsymbol{\theta}^*}(\mathbf{x}|\mathbf{z}^{(i)})$ is the true conditional. As the true set of parameters $\boldsymbol{\theta}^*$ is unknown to us, we may attempt to generate samples with learnt parameters $\boldsymbol{\theta}$:

$$\mathbf{z}^{(i)} \sim p(\mathbf{z}) \quad \mathbf{x}^{(i)} \sim p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z}^{(i)}) \quad (3.40)$$

Now suppose

$$\int q_\phi(\mathbf{z}|\mathbf{x})p(\mathbf{x})d\mathbf{x} = \hat{p}(\mathbf{z}) \quad (3.41)$$

where $p(\mathbf{x})$ is the distribution that generated data set X . The generative procedure (3.40) makes the erroneous assumption that $q_\phi(\mathbf{z}|\mathbf{x})$ perfectly matches the prior $p(\mathbf{z})$, that is, $p(\mathbf{z}) = \hat{p}(\mathbf{z})$ [8]. Since only pressure was applied to match the two, this is clearly not true in general. That is, in general,

$$p(\mathbf{z}) \neq \hat{p}(\mathbf{z}) \quad (3.42)$$

Therefore we have

$$\int p_{\theta}(\mathbf{x}|\mathbf{z})\hat{p}(\mathbf{z})d\mathbf{z} = p(\mathbf{x}) \neq \int p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z} \quad (3.43)$$

which suggests that the generative procedure (3.40) does not produce samples from $p(\mathbf{x})$ exactly.

To sample from $p(\mathbf{x})$, a new generative procedure is proposed [8]:

$$\mathbf{z}^{(i)} \sim \hat{p}(\mathbf{z}) \quad \mathbf{x}^{(i)} \sim p_{\theta}(\mathbf{x}|\mathbf{z}^{(i)}) \quad (3.44)$$

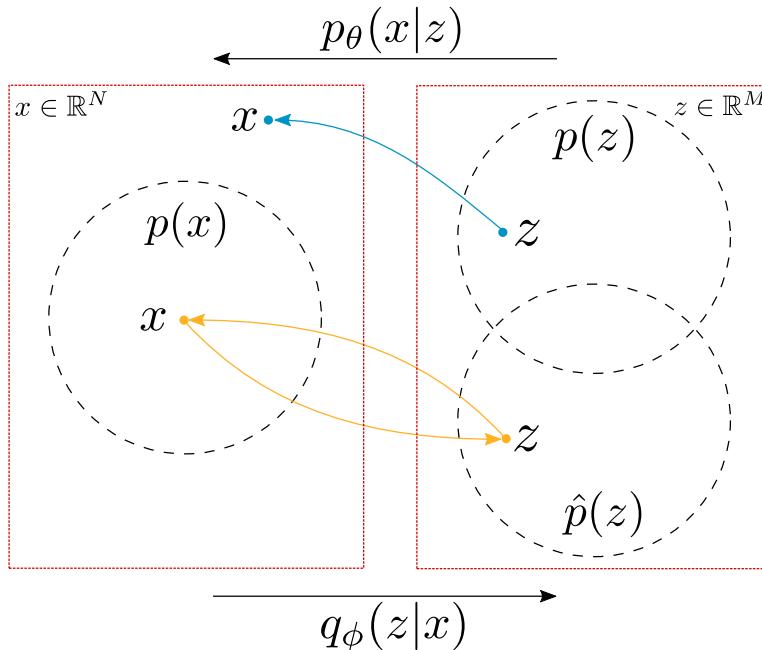


Figure 3.13: The probabilistic encoder $q_{\phi}(\mathbf{z}|\mathbf{x})$ maps a given data point \mathbf{x} to the unknown distribution $\hat{p}(\mathbf{z})$. The probabilistic decoder $p_{\theta}(\mathbf{x}|\mathbf{z})$ is trained to map samples from $\hat{p}(\mathbf{z})$ back to $p(\mathbf{x})$, since its inputs are drawn from the probabilistic encoder $q_{\phi}(\mathbf{z}|\mathbf{x})$. A sample from the prior $p(\mathbf{z})$ will not be mapped back by $p_{\theta}(\mathbf{x}|\mathbf{z})$ to $p(\mathbf{x})$ exactly if $p(\mathbf{z}) \neq \hat{p}(\mathbf{z})$. [8]

As $\hat{p}(\mathbf{z})$ is unknown, it is not obvious how to sample from generative autoencoders using procedure (3.44). However, it is possible to formulate a Markov chain Monte Carlo (MCMC) chain that starts with an arbitrary latent variable $\mathbf{z}_{t=0}$ and converges with $\mathbf{z}_{t \rightarrow \infty} \sim \hat{p}(\mathbf{z})$ [8]. This is done by successively encoding and decoding the same sample as follows:

$$\mathbf{x}_{t=k+1} \sim p_{\theta}(\mathbf{x}|\mathbf{z}_{t=k}) \quad \mathbf{z}_{t=k+1} \sim q_{\phi}(\mathbf{z}|\mathbf{x}_{t=k+1}) \quad (3.45)$$

where $\mathbf{z}_{t=0}$ is arbitrarily chosen.

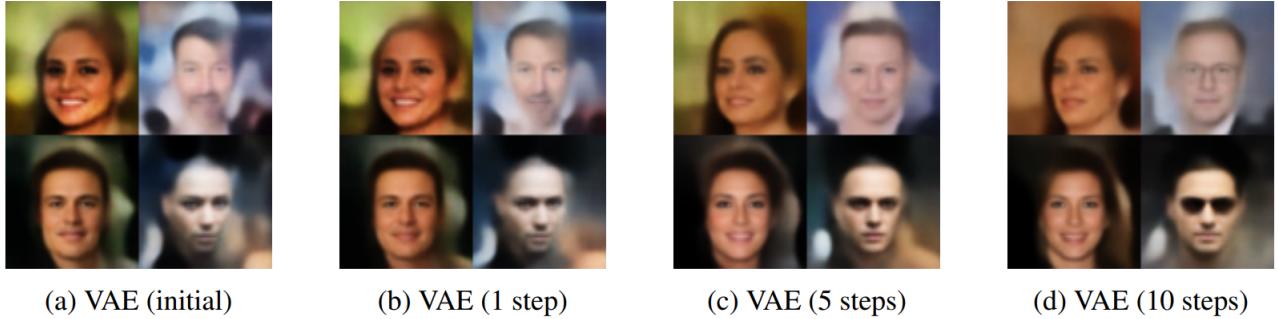


Figure 3.14: Samples from a variational autoencoder trained on the CelebA data set after $t = 0, 1, 5$ and 10 steps of the generative procedure (3.45). The MCMC chain was initialised with a sample from the prior $\mathbf{z}_{t=0} \sim p(\mathbf{z})$, which often improves the quality of the samples. [8]

3.6 Regularizing CNNs with Locally Constrained Decorrelations

[25]

3.7 Batch Normalisation

...

Chapter 4

Implementation

4.1 Dimensionality Reduction

Machine learning with high-dimensional data, such as images, is often computationally intensive. Atari 2600 frames are RGB images (3 channels) of size 210×160 , expressed in shorthand as $(3, 210, 160)$. (This data is much higher in dimensionality than MNIST, which has dimensions $(1, 28, 28)$). Considering the training sets considered in this project are of hundreds of thousands of images, data points of dimension $(3, 210, 160)$ are too computationally intensive with the best hardware available ($2 \times$ NVIDIA Tesla K80 GPU Accelerators). It is therefore necessary to reduce the dimensionality of our data set. Google DeepMind's *Human-level Control Through Deep Reinforcement Learning* [21] made use of Stella, as we have, and have convincingly struck a reasonable balance between resolution and dimensionality reduction. This section will be a short but necessary mentioning of the pre-processing pipeline used to generate the data sets in later chapters.

4.1.1 Pre-processing Pipeline

Ensuring Object Persistence

Atari 2600 games can only store a limited number of sprites per frame due to the limitations in hardware during its development [21]. This is an issue as some objects that appear in one frame fail to appear in the next. To solve this lack of object persistence, even and odd frames are combined by taking the maximum over each channel (RGB). By taking the maximum, we ensure that any object present in one frame is also present in the other.

Extracting Luminance and Cropping

The luminance Y is then extracted from the RGB image [29]:

$$Y = 0.2126 \times R + 0.7152 \times G + 0.0722 \times B \quad (4.1)$$

The resultant greyscale image of shape $(1, 210, 160)$ is cropped to $(1, 84, 84)$.

File Formats

Image formats were discovered to be more important than originally thought. Empirically, PNG or GIF formats were reasonable formats, but using JPEG often resulted in distortions near the perimeter of objects. An example of this effect is shown in Figure (4.2).

4.2 Qualitative Assessment Using GUIs

- To qualitatively assess the significance of a latent neuron in the final reconstruction, we change its value over a given range and inspect the reconstructions

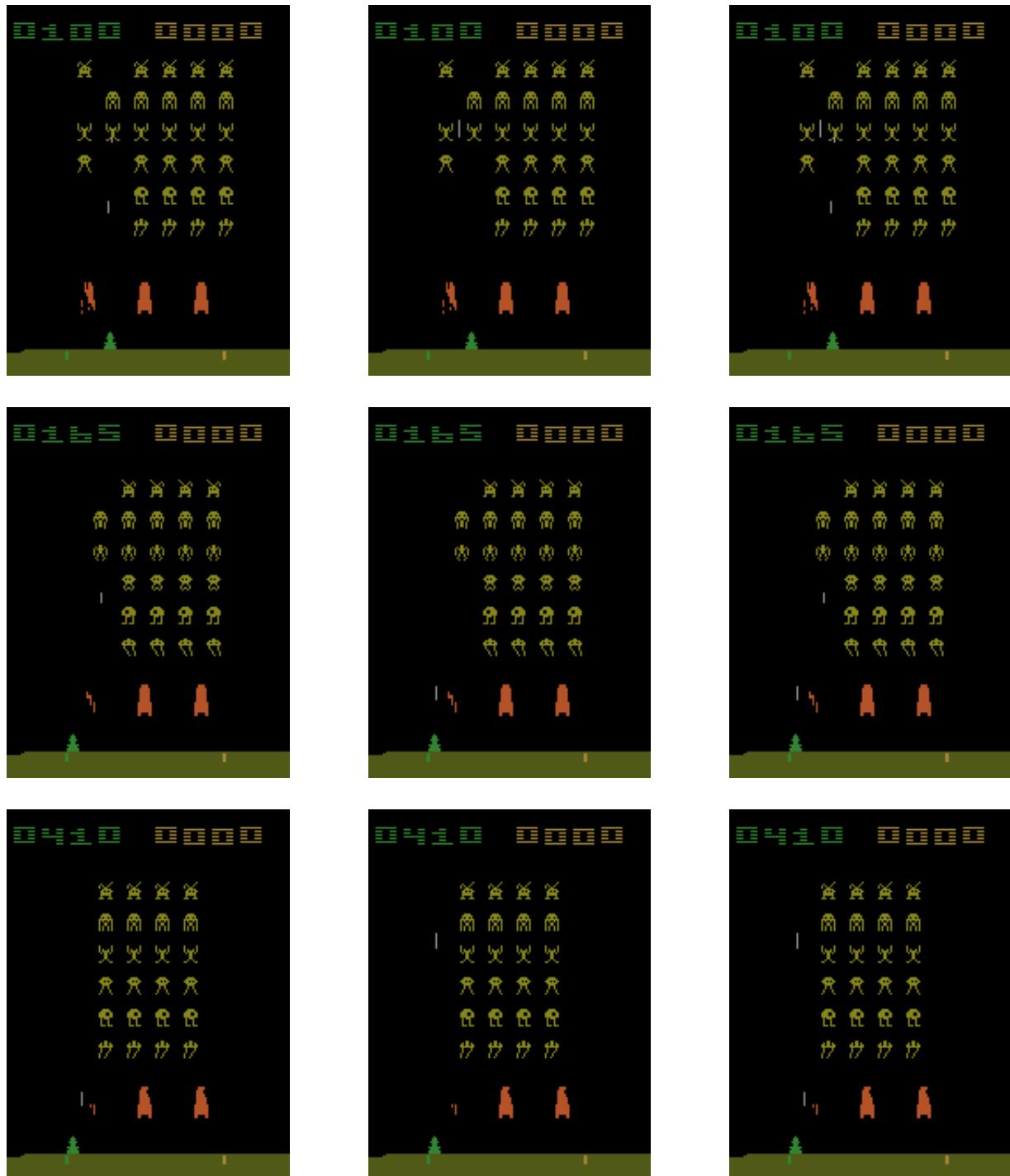


Figure 4.1: A collection of frames captured from Space Invaders emulated on Stella. **Left column:** an even frame. **Middle column:** the (odd) frame following. **Right column:** Combining the even and odd frames by taking the maximal value over each channel (RGB). Clearly the bullets visible in one frame fail to persist in the next. As mentioned, this is due to the limited number of sprites Atari 2600 can load in a single frame.

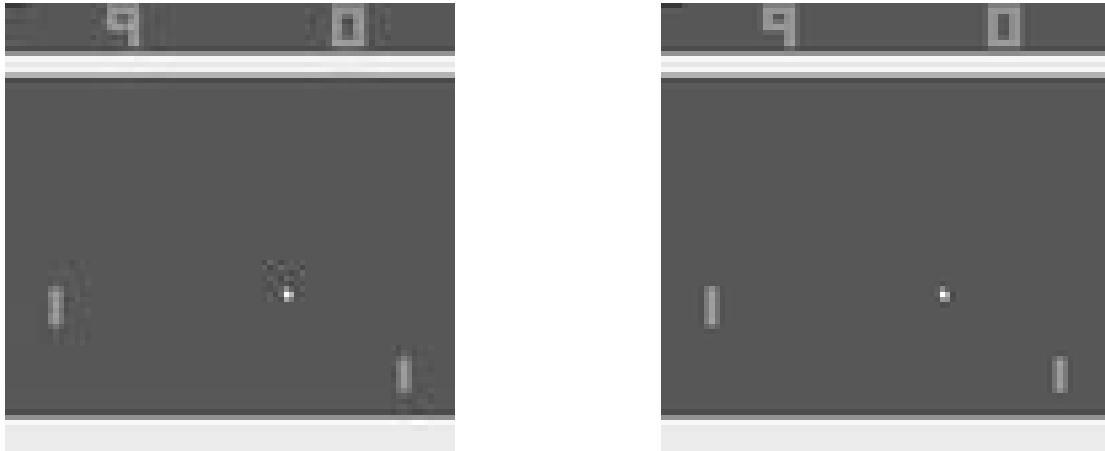


Figure 4.2: Pre-processed frames captured from Pong emulated on Stella. These frames were originally 84×84 , but are printed here as 168×168 to emphasise distortions. **Left:** The JPEG format distorts the ball, paddle and score sprites. **Right:** The PNG format displays the frame without such distortions.

- In order to speed up this process, Tkinter was used to develop a GUI with a sliding bar corresponding to the latent neuron's value
- This real-time reconstruction allows for much quicker feedback

4.3 Ensuring Numerical Stability in the Latent Space

- Use log variance instead of variance to spread values evenly, which helps with backpropagation, as it's easier to learn with an approximately equal step size between values rather than those that are either clumped or sparse
- Since log is monotonic, the value ordering is invariant, which allows back propagation to implicitly optimise the variance σ^2
- State implemented ELBO

4.4 Activation Functions in the Latent Space

- Conventionally activation functions are used in the hidden layers, such as sigmoid or, more commonly, ReLU.

- However, we must be careful not to mindlessly follow this trend.
- It makes no sense to use a ReLU activation function on the mean layer of the variational autoencoder. The mean $\mu \in \mathbf{R}$ is unconstrained, and hence a linear activation function is appropriate.
- Stuff about log variance - also use linear activation function

Chapter 5

Methods

- Introduce requirements
- Need to preserve spatial information for symbolic front-end
- Need disentangled latent space for symbolic front-end
- Explain degree of freedom here (disentangled filters, or disentangled neurons, or something else?)
- This is unexplored territory here - we decide what the reasonable architecture should be

5.1 Single Latent Filter

- As mentioned, the type and position of an object must be preserved in the latent space
- This may be achieved by using a single latent filter, with the object type corresponding to the value of the weighted sum of the neuron.

5.1.1 Architecture

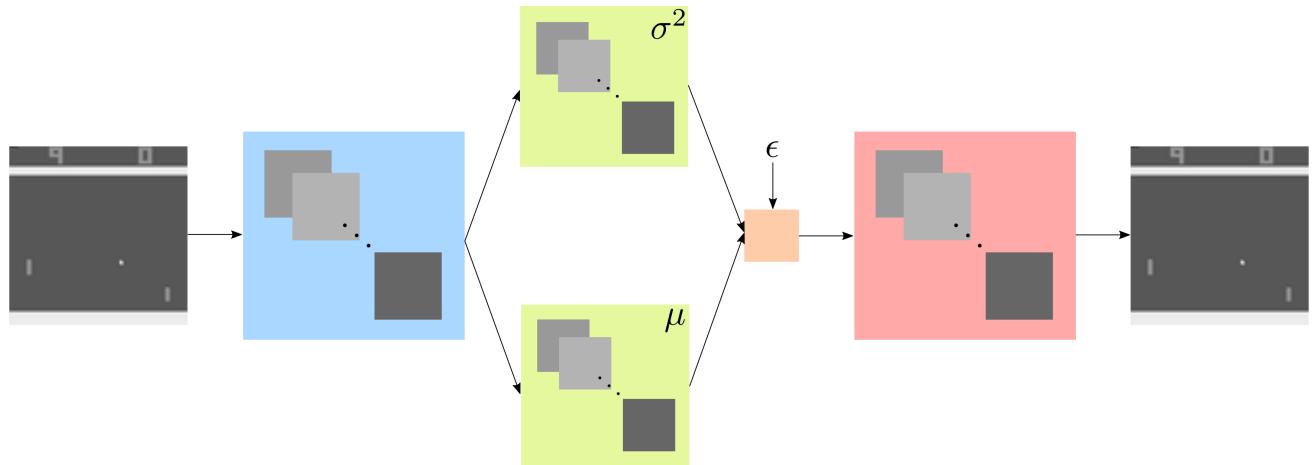


Figure 5.1: Caption.

5.1.2 Derivation

5.2 Disentangling Latent Neurons

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

5.2.1 Architecture

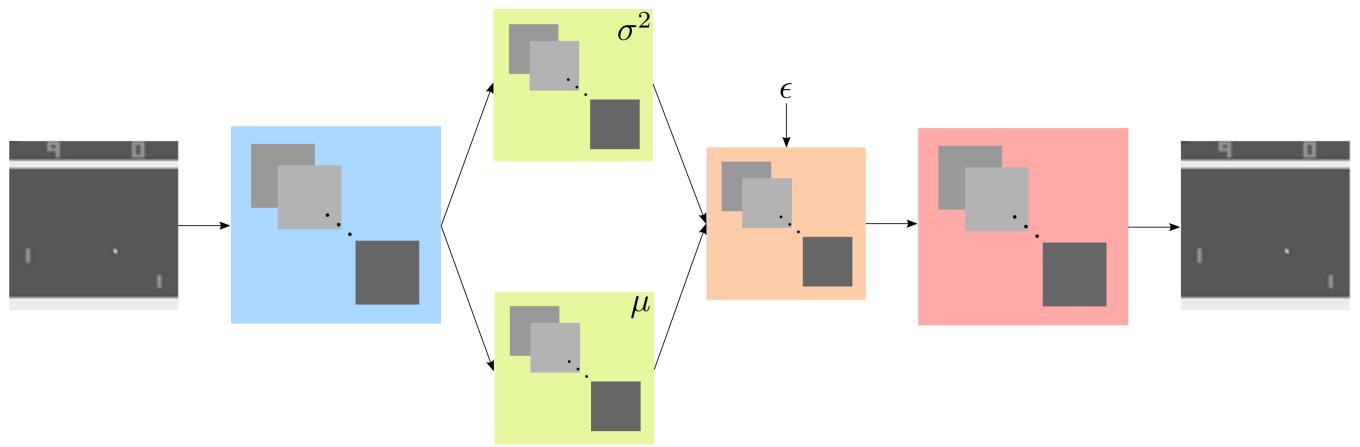


Figure 5.2: Caption.

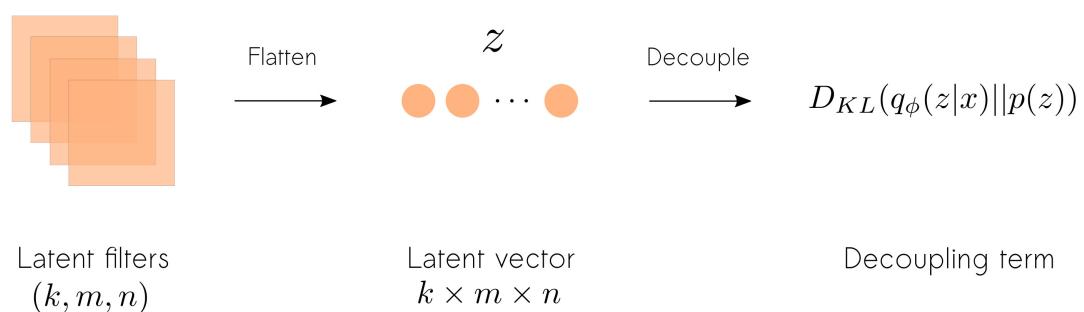


Figure 5.3: Caption.

5.2.2 Derivation

5.3 Disentangling Latent Filters Using Averages

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

5.3.1 Architecture

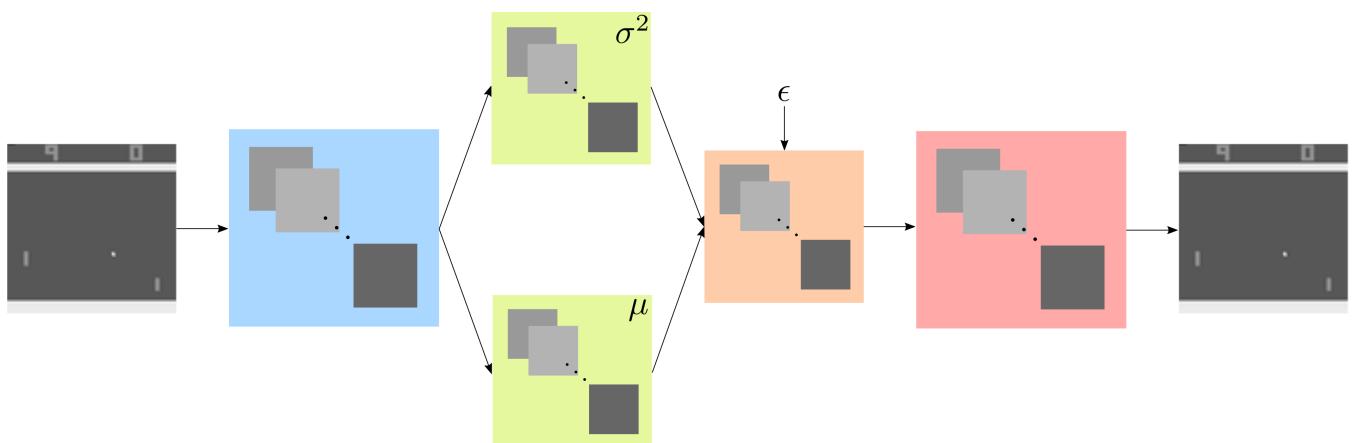


Figure 5.4: Caption.

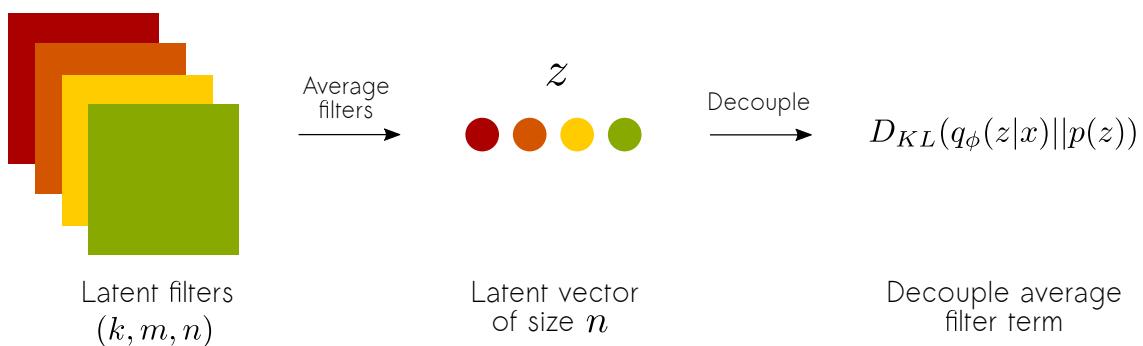


Figure 5.5: Caption.

5.3.2 Derivation

5.4 Decoupling Latent Filters Using Weighted-Averages

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

5.4.1 Architecture

5.4.2 Derivation

5.5 Separating Colour Spaces

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

5.5.1 Architecture

5.5.2 Derivation

5.6 Orthogonal Convolutions

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

5.6.1 Architecture

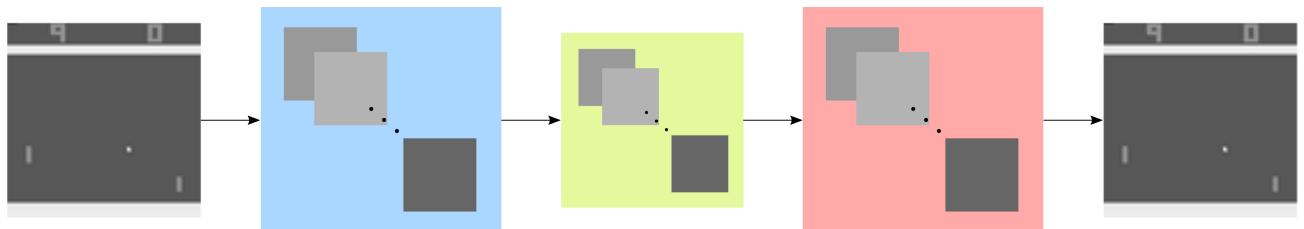


Figure 5.6: Caption.

5.6.2 Derivation

5.7 Winner Takes All

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis

natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

5.7.1 Architecture

5.7.2 Derivation

Chapter 6

Results

6.1 The No Free Lunch Relationship Between Reconstruction Loss and KL Divergence

6.2 Using Batch Normalisation With Convolutional Variational Latent Spaces

6.3 Using Batch Normalisation With Convolutional Variational Latent Spaces

Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetuer id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla,

malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

Chapter 7

Conclusion

7.1 Summary of Thesis Achievements

Summary.

7.2 Applications

Applications.

7.3 Future Work

Future Work.

Bibliography

- [1] M. G. Bellemare, Y. Naddaf, J. Veness, and M. Bowling. The arcade learning environment: An evaluation platform for general agents. In *IJCAI International Joint Conference on Artificial Intelligence*, volume 2015-Janua, pages 4148–4152, 2015.
- [2] D. M. Blei. Variational Inference, 2011.
- [3] D. M. Blei, A. Kucukelbir, and J. D. McAuliffe. Variational Inference: A Review for Statisticians. jan 2016.
- [4] K. Burnham and D. Anderson. *Model Selection and Multimodel Inference: A Practical Information-Theoretic Approach (2nd ed)*, volume 172. 2002.
- [5] X. Chen, Y. Duan, R. Houthooft, J. Schulman, I. Sutskever, and P. Abbeel. InfoGAN: Interpretable Representation Learning by Information Maximizing Generative Adversarial Nets. *arXiv:1606.03657 [cs.LG]*, pages 1–14, 2016.
- [6] F. Chollet. Keras: Deep Learning library for Theano and TensorFlow, 2015.
- [7] T. S. Cohen and M. Welling. Transformation Properties of Learned Visual Representations. dec 2014.
- [8] A. Creswell, K. Arulkumaran, and A. A. Bharath. Improving Sampling from Generative Autoencoders with Markov Chains. oct 2016.
- [9] B. Darrach. Meet Shakey, the First Electronic Person. *Life*, pages 58–68, 1970.
- [10] G. Desjardins, A. Courville, and Y. Bengio. Disentangling Factors of Variation via Generative Entanglement. oct 2012.

- [11] C. Doersch. Tutorial on Variational Autoencoders. *arXiv*, pages 1–23, 2016.
- [12] I. Dykeman. Conditional Variational Autoencoders, 2016.
- [13] M. Garnelo, K. Arulkumaran, and M. Shanahan. Towards Deep Symbolic Reinforcement Learning. *arXiv Preprint*, pages 1–13, 2016.
- [14] I. Higgins, L. Matthey, X. Glorot, A. Pal, B. Uria, C. Blundell, S. Mohamed, and A. Lerchner. Early Visual Concept Learning with Unsupervised Deep Learning. *arXiv*, 2016.
- [15] M. Hutter. *Universal Artificial Intelligence*, volume 1. 2005.
- [16] D. P. Kingma and M. Welling. Auto-Encoding Variational Bayes. *Iclr*, (MI):1–14, 2014.
- [17] W. Knight. AI Winter Isn’t Coming. *MIT Technology Review*, 2016.
- [18] A. Krizhevsky, I. Sutskever, and G. E. Hinton. ImageNet Classification with Deep Convolutional Neural Networks. *Advances In Neural Information Processing Systems*, pages 1–9, 2012.
- [19] F.-F. S. U. Li, J. S. U. Johnson, and S. S. U. Yeung. Convolutional Neural Networks for Visual Recognition, 2016.
- [20] C. Y. Liou, J. C. Huang, and W. C. Yang. Modeling word perception using the Elman network. In *Neurocomputing*, volume 71, pages 3150–3157, 2008.
- [21] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.
- [22] D. Ormerod. Lee Sedol plays the first move of game three against AlphaGo, 2016.
- [23] E. U. o. T. Reingold. PSY371F Higher Cognitive Processes, 2001.
- [24] S. R. Richter, V. Vineet, S. Roth, and V. Koltun. Playing for data: Ground truth from computer games. In *Lecture Notes in Computer Science (including subseries Lecture Notes*

- in Artificial Intelligence and Lecture Notes in Bioinformatics)*, volume 9906 LNCS, pages 102–118, 2016.
- [25] P. Rodríguez, J. Gonzàlez, G. Cucurull, J. M. Gonfaus, and X. Roca. Regularizing CNNs with Locally Constrained Decorrelations. nov 2016.
- [26] R. Rosen. IBM’s Deep Blue vs Gary Kasparov, 2012.
- [27] J. Schmidhuber. Learning Factorial Codes by Predictability Minimization. *Neural Computation*, 4(6):863–879, 1992.
- [28] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. van den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, S. Dieleman, D. Grewe, J. Nham, N. Kalchbrenner, I. Sutskever, T. Lillicrap, M. Leach, K. Kavukcuoglu, T. Graepel, and D. Hassabis. Mastering the game of Go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016.
- [29] M. Stokes, M. Anderson, S. Chandrasekar, and R. Motta. A Standard Default Color Space for the Internet - sRGB, 1996.
- [30] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, volume 07-12-June, pages 1–9, 2015.
- [31] Y. Tang and G. Hinton. Tensor Analyzers. *Proceedings of the 30th International Conference on Machine Learning (ICML-13)*, 28, 2013.
- [32] B. G. Thiagarajan, A. Member, and G. Z. Voyiadjis. beta-VAE: Learning Basic Visual Concepts with a Constrained Variational Framework. *ICLR’17 submission*, (July):1–13, 2016.
- [33] M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, volume 8689 LNCS, pages 818–833, 2014.

- [34] K. Zsolnai-Fehér. Computer Games Empower Deep Learning Research, 2016.